

統計的テキストモデル

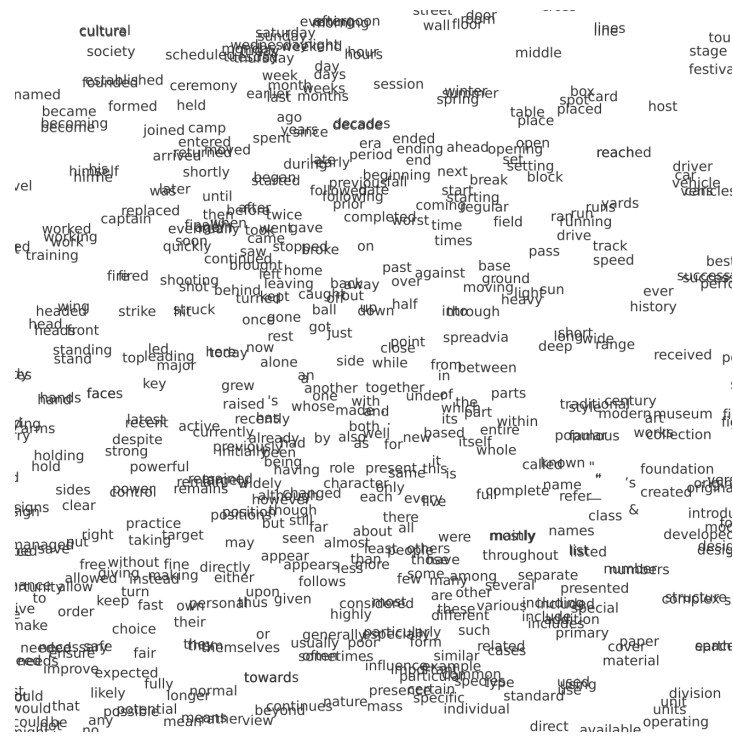
持橋 大地

統計数理研究所 数理・推論研究系

daichi@ism.ac.jp

2024年5月1日

\$Id: textmodel.tex,v 1.18 2024/04/04 11:20:40 daichi Exp \$



はじめに

本書は、テキストの統計的なモデル化について解説した本です。ここでいうテキストとは、Web ページやメールのような文書だけでなく、小説や新聞、法案、アンケートへの自由回答など、さまざまな範囲を含んでいます。よって本書の読者としては、テキストを扱う必要のある理系のエンジニアの方だけでなく、人文科学および社会科学系の方々、および言語に興味のある一般の方々を想定しています。

Web が現れて以来、電子化されたテキストを扱う自然言語処理は驚異的に進展しました。しかし、自然言語処理の教科書や情報源は、文書のカテゴリや単語の品詞のように、テキストから人の与えた正解ラベルを予測する教師あり学習を扱っている場合がほとんどで、現実に現れる、そうした正解ラベルのないテキストを統計的にどうモデル化して扱うか、という**教師なし学習**について体系的にまとめられた成書は、ほとんど存在しないのが現状です。「テキストマイニング」はこれに近い分野ですが、多くの場合、表面的なパッケージの使い方に終わっていたり、離散的な言葉の頻度に、従来の多変量解析を無理にあてはめていたりといった問題がありました。テキストのような離散データを取り扱うには、より適切な統計モデルが存在します。またその際に、内部でどのような数学的なモデル化と計算が行われているのかがわからなければ、目の前の問題に適用するためには、どこにどう手を入れたらよいのかを知ることはできません。

そこで本書では、テキストの統計的なモデル化について一から説明し、**ブラックボックスに頼らなくても、様々な分析を自分で自由に行えるようになる**ことを目的としています。本書で説明するようなテキストの高度な統計モデルにはパッケージがないことも多く、また、パッケージでブラックボックス化して扱うことは適切とはいえません。あえて「背景知識」のような無味乾燥な章は作ら

ず、必要になる数学的知識はすべてその場で、実例を交じえつつ説明していますので、本書を読むには高校レベルの数学を理解していれば充分で、あらかじめ機械学習の教科書を読んでおく必要はありません。むしろ、本書でテキストの統計モデルの基礎について一通り理解して道を作った上で、あらためて機械学習や統計学の本を読まれることで、それらがより読みやすくなり、理解が深まることを期待しています。

人文科学や社会科学においても、計量的な方法の必要性は増加の一途をたどっており、その多くの場合でテキストの分析が必要になっています。これまでは、統計が必要になるのは主に経済学を中心とした社会科学で、そこでは価格のような連続量が中心となっており、従来の統計学をそのまま適用することができました。これに対して、テキストのような**離散データ**が人文科学および社会科学の両方で必要になってきたのは、Webの発達によって電子テキストが容易に入手できるようになった比較的最近のことといえます。筆者ももとは文科系ですので、そうした分野の重要性はよく理解しているつもりです。本書をきっかけに、自分の手でテキストの統計的なモデル化と分析ができるようになっていただければと考えています。

また、本書はテキストを対象にしていますが、こうした離散データに対する方法論は言語だけでなく、**他の種類の離散データについても同様に適用**することができます。たとえば、文書に様々な種類の単語が含まれている状況は、コンビニやオンラインストアで客が様々な商品を購入する状況^{*1}と同じで、実際にこのための協調フィルタリングとよばれる分野は、本書で文書をモデル化するために導入するものと、まったく同じモデルを用いています。また、音楽の楽譜は離散的な記号で、言語のように構造を持っていますし、細胞内のDNAやゲノムはATGCあるいは20種類のアミノ酸を文字とした、一種の「言語」です。他にも、本書で紹介する離散データのための確率モデルは、多くの分野で適用できるのではないかと考えています。^{*2}

*1 これは、データマイニングの分野ではバスケット分析と呼ばれています。

*2 筆者はこれまでに、言語学、脳科学、音声認識、音楽情報処理、ロボティクス、バイオインフォマティクス、政治学といった分野と共同研究を行っています。

本書の対象読者

上に述べたように、本書はテキストの統計的なモデルに関心のある、広く人文科学・社会科学およびエンジニアの方、および言語の数学的なモデル化に興味を持つ一般の方を対象にしています。確率の基礎から始め、 n グラムモデルとは何か、Word2Vec は数学的には何をしているのか、テキストをどうクラスタリングすればよいのかなど、原理的な内容を一から解説します。

本書の特徴は、数式による定義を天下一に与えるのではなく、できる限りその導出や意味について解説していること、そしてブラックボックスのパッケージに頼らない、ということです。その代わりに、テキストをモデル化する自然言語処理の数理について丁寧に説明し、読者が自分の手で統計的な分析を使いこなせることを目標としています。たとえばパープレキシティ一つをとっても、定義の数式をただ与えるのではなく、基礎となる自己情報量の説明から始め、どうしてパープレキシティを考えるのか、どういう意味を持っているのかを丁寧に解説するようにしました。テキストを扱う既存のパッケージでは、本書で紹介するような数学的な分析はほとんど行うことができませんが、それでも良いという方は多数の本がありますので、そちらをご参照ください。ただしもちろん、普段はパッケージを使った分析を行っている方にとっても、本書でテキスト分析の背景や理論について理解しておくことは、適切な分析を行う上で大きな助けになるでしょう。

必要となる数学について

本書はテキストの統計的なモデル化に関する本ですので、数式を用いることはどうしても必要になります。むしろ、言語のようなアナログな対象をどう数学的にモデル化していくかが、本書の主題といってもいいでしょう。ただし、必要になるのは数 III までの高校数学と、大学教養の数学 (解析と線形代数) の一部だけです。本書に限らず、統計学の理解には、数 III まではどうしても必要になります (e^x の微分・積分や対数の微分など)。現在はこのような場合のために、多くの優れた参考書が出版されていますので、必要な方は自習しておくようにしてください。筆者自身の数学も、数 III 以降はこうして自習したものです。と

いっても、難しい問題が解ける必要はなく、教科書レベルの内容が理解できていれば充分です。

逆に、本書では測度論のような高度な数学は仮定していません。確率変数のことを「役割を表す変数」(23 ページ)とするような直感的な記述は、意図的なものです。

深層学習との関係について

本書では、深層学習の手法としては最も基本的な単語埋め込み、文埋め込み、文書埋め込みまでの手法と、その背後にある理論について解説しました。現在の自然言語処理は、そうした埋め込みを背景にした LSTM や Transformer(BERT) といった深層学習の手法なしには語ることはできません。しかし、それらは実装して動かせても、なぜ言語をうまく扱えるのかという理論的背景についてはほとんどわかっておらず、ほぼブラックボックスの状態です^{*3}。そこで本書では、単なるレシピ集となることを避け、それらについてはあえて紹介しないことにしました。ただし、本書で説明している単語埋め込み、文埋め込み、文書埋め込みの範囲でも(あるいは、それ以外の確率的手法を用いても)、非常に多くの自然言語処理が可能なることに注意してください。Transformer のような手法は非常に多くのデータと莫大な計算時間を必要としますが、すべての自然言語処理の問題がそれで解決するわけではなく、より深い応用には、本書で説明したような基礎知識と組み合わせることが必要です。本書の知識を基礎にした上で、深層学習手法の使い方については多くの書籍が出版されていますので、本書の3章の文献案内も参照してください。

実装とサポートサイトについて

本書では、Python 言語を使って、実際のテキストでの計算例を示しました。統計の分野では R 言語がよく使われていますが、R は実行速度が遅く、標準的にサポートされているデータ構造も少ないため、大量のデータを高速に扱う必要のあ

^{*3} これらが内部で何を行っているのかについての研究も一部では進みつつあり、筆者も、LSTM の内部状態が文法構造の埋め込みの深さをほぼ離散的に数えている[1]といった研究を行っています。

るテキストの処理にはあまり向いていないからです。^{*4} Python 言語自体の説明はしていませんので、標準的な使い方は理解していることを前提としています^{*5}。ただし、実装部分をすべて理解できなくても、本文の記述は追えるように配慮したつもりです。本書で用いた例や実装はすべて、本書のサポートサイト

```
http://www.ism.ac.jp/~daichi/textmodel/
```

で公開しています。また、Github のレポジトリ

```
https://github.com/daiti-m/textmodel/
```

でも同じファイルを公開していますので、コマンドラインから git で

```
% git clone https://github.com/daiti-m/textmodel
```

を実行すれば、本書で用いたスクリプトやデータをすべてダウンロードすることができます。

実験で用いているテキストデータはすべて自由に入手できるもので、上記のサポートサイトおよびレポジトリの “data” フォルダ

```
http://www.ism.ac.jp/~daichi/textmodel/data/
```

```
https://github.com/daiti-m/textmodel/data/
```

から入手することができます。サポートサイトには更新情報やリンク集など、有用な情報を載せていく予定です。合わせてご覧いただければ幸いです。

^{*4} 不可能ではなく、実際に計量政治学におけるテキスト分析のパッケージである `quanteda` は R 言語上のパッケージです。ただし、大量のテキストに対して重い計算を高速に実行するには、Rcpp のような外部言語の助けが必要になります。

^{*5} 参考書を使わなくとも、オンライン上の情報で Python の基本は十分にマスターすることができます。個人的には、M.Hiroi さんの「お気楽 Python プログラミング入門」http://www.nct9.ne.jp/m_hiroi/light/index.html#python_abc は大変お薦めで、これを読めば、本書に必要な基本は容易に理解できるのではないかと思います (Python2 系のため、`print` 文などに少し違いがあります)。<https://utokyo-ipp.github.io/> では東京大学の、<https://repository.kulib.kyoto-u.ac.jp/dspace/handle/2433/265459> では京都大学の Python プログラミング入門の実習テキストが、それぞれフリーで公開されています。

本書の記法

x	英小文字：変数
N	英大文字：定数またはデータ
\mathbf{x}	英小文字の太字：ベクトル
\mathbf{X}	英大文字の太字：行列
\mathcal{L}	カリグラフィック体：集合 (データセットや語彙など)
\mathbf{a}	タイプライター体：文字 (アルファベットの場合)
\mathbb{E}	期待値
\mathbb{V}	分散
\mathbb{R}	実数
\mathbb{I}	指示関数 ($\mathbb{I}()$ の中が真なら 1, 偽なら 0 を返す関数)
\mathbf{I}	単位行列
$\mathbf{0}$	ゼロベクトル (要素がすべて 0 のベクトル)
\cdot	ベクトルの内積
T	ベクトルおよび行列の転置
$\langle \dots \rangle_p$	確率分布 p による期待値
$\exp(x)$	e^x の別記法

確率分布の略記

\mathcal{N}	ガウス分布 (正規分布)
Po	ポアソン分布
Be	ベータ分布
Ga	ガンマ分布
Unif	一様分布
Mult	多項分布
Bernoulli	ベルヌーイ分布

目次

はじめに	i
1 テキストと言語のモデル化	1
1.1 言語とテキストの特徴	1
1.2 テキストの階層構造	3
1.3 教師あり学習と教師なし学習	4
1.4 統計的な方法とアドホックな方法	6
1.5 本書の構成と読み方	10
1.6 本書の例と実装について	14
2 文字の統計モデル	16
2.1 文字の頻度と出現確率	16
2.2 文字の同時確率	20
2.3 同時確率の周辺化	23
2.4 文字の条件つき確率	27
2.4.1 確率の連鎖則	29
2.4.2 ベイズの定理	32
2.5 文字 n グラムモデル	38
2.5.1 文字列の確率的生成	38
2.5.2 ゼロ頻度問題	44
2.6 統計モデルの学習と評価	50

2.6.1	学習データとテストデータ	50
2.6.2	予測確率とパープレキシティ	56
2.6.3	情報理論の基礎	58
2.6.4	統計モデルと汎化性能	67
2章	の演習問題	73
2章	の文献案内	76
3	単語の統計モデル	78
3.1	文字から単語へ	78
3.2	単語の統計と巾乗則	80
3.2.1	Heaps の法則	84
3.2.2	Zipf の法則	86
3.3*	単語の統計的フレーズ化	90
3.4	単語 n グラム言語モデル	97
3.4.1	ディリクレ分布	101
3.4.2	ディリクレ分布と多項分布	108
3.4.3	階層ディリクレ言語モデル	116
3.4.4	Kneser–Ney 言語モデル	121
3.5	単語ベクトルとその原理	132
3.5.1	ニューラル n グラム言語モデル	132
3.5.2	スキップグラムと Word2Vec	135
3.5.3	単語ベクトルの学習	142
3.5.4	Word2Vec と行列分解	147
3.5.5*	GloVe と意味方向の数理	153
3.5.6	単語ベクトルの分布とノルム	159
3章	の演習問題	167
3章	の文献案内	168
4	文の統計モデル	169
4.1	テキストの文分割	169
4.2	文ベクトルと意味的ランダムウォーク	172

4.2.1	RAND-walk モデル	173
4.2.2	文ベクトルの計算	174
4.3	構文解析と係り受け解析	180
4.4	隠れマルコフモデル (HMM)	185
4.4.1	HMM の状態推定	193
4.4.2	HMM のパラメータ推定	197
4.4.3	周辺化 Gibbs サンプリング	210
4.4.4	HMM による品詞の教師なし学習	217
4 章	の演習問題	219
4 章	の文献案内	220
5	文書の統計モデル	222
5.1	ナイーブベイズ法と単語集合表現	222
5.1.1	文書の分類確率	228
5.2	ユニグラム混合モデル (UM)	235
5.2.1	トピックの解釈と自己相互情報量	242
5.2.2	EM アルゴリズムによる学習	246
5.2.3*	UM のベイズ学習	251
5.3	ディリクレ混合モデル (DM)	255
5.3.1	単語単体と幾何的解釈	257
5.3.2	ポリア分布と単語のバースト性	262
5.4	潜在ディリクレ配分法 (LDA)	263
5.4.1	LDA の幾何的解釈	275
5.4.2	トピックモデルの評価	279
5.5	ニューラル文書モデルと独立成分分析	284
5.5.1	文書ベクトルと Doc2Vec	287
5.5.2	単語ベクトル/文書ベクトルの解釈	295
5.6	確率的潜在意味スケーリング (PLSS)	298
5.6.1	Wordfish	300
5.6.2	確率的潜在意味スケーリング (PLSS)	302

5.6.3* PLSS の半教師あり学習	309
5 章の演習問題	317
5 章の文献案内	319
付録	321
A ディリクレ分布の積分と期待値	321
B ディリクレ分布の α のベイズ推定	321
C Kneser-Ney 平滑化の導出	324
D Jensen の不等式	324
アルゴリズムの一覧	325
索引	328
参考文献	335

1 テキストと言語のモデル化

1.1 言語とテキストの特徴

われわれは日々、テキストに囲まれて暮らしています。図 1.1 のように、街中にはテキストがあふれていますし、毎日見る携帯電話の画面から Web ページ、電子メール、小説や雑誌に至るまで、われわれがテキストを日々目にしない日はないと言ってよいでしょう。

テキストは言葉を書き起こしたのですが、もちろん、言葉は書き言葉から生まれたわけではありません。よく知られているように、アイヌ語は文字を持ちませんでしたし、南米のインカ帝国でも同様でした。しかし、最初に古代メソポタミアで文字が発明されて言葉が書き表されるようになったことで、言語は音声によるその場限りのコミュニケーションの媒体から、空間および時間を超えて伝わる情報を表せるようになり、抽象化された書き言葉も出現して、今日みられる高度な文明の礎となりました。言語にはこのように、テキストに書き表される以前に音声や抑揚、画像としての文字といった側面もありますが*¹、本書では記号化されて文字として表されたもの、すなわちテキストを対象とすることに



図 1.1: われわれは日々、テキストに取り囲まれて暮らしています。
(香港にて、筆者撮影)

*1 こうしたテキスト以前の音声や文字画像の情報をどうテキストと組み合わせるかは、たいへん興味深い問題です。実際に深層学習を用いて、フォントから取り出した漢字の部首の画像を文字の情報と組み合わせる研究は現在、さまざまに行われています。

ます。

テキストが、客観的にみてほかの一般的なデータと異なっている点は何でしょうか。それは、

離散的であること

かつ

1次元の時系列として表されていること

だと考えられます。

最初に述べたように、テキストはもともと音声を書き起こしたもので、音声信号は連続的な空気の圧力の違いが時系列で耳に伝わってくるものです。しかし、この連続値の時系列を、われわれは“Good morning”のような離散的なテキストとして認識します。離散的とは、情報は文字が単位の場合は“a”…“z”，単語が単位の場合は“good”，“morning”，“evening”，…のように、集合のどれかの要素として認識される、ということです。この際、連続値である音声の周波数や画像の色とは異なり、綴りが似ているからといって意味が似ているとは限らず、“cat”（猫）と“cut”（切る）はまったく違う言葉になっています。こうした言葉の種類は言語の場合は非常に多く、文字でも漢字圏なら数千以上、単語の場合は数万次元を超える超高次元なのが普通です。

また、画像では色の配置は2次元ですが、テキストはもともと耳で聴く音声を書き起こしたものであるため、本質的に1次元なのが特徴です。よって、ある言葉が次の行にある言葉とたまたま近くにあっても、両者は原則的に無関係で、これは画像のような空間的なデータとの大きな違いです。言語は1次元ではあるものの、句構造や埋め込み文、リズムや談話構造といった見えない構造が実は豊富に存在していることも特徴とっていいでしょう。

こうした離散的な時系列であるというテキストの特徴により、連続値であることを暗に仮定した*2主成分分析、多変量解析などの古典的な統計学や機械学習の手法は、本来はそのままでは使うことができません。それではどうすればよいのかについて、本書で一緒に考えていきましょう。

*2 多くの古典的な統計理論は、観測値が連続値でガウス分布に従うといったことを仮定しています。

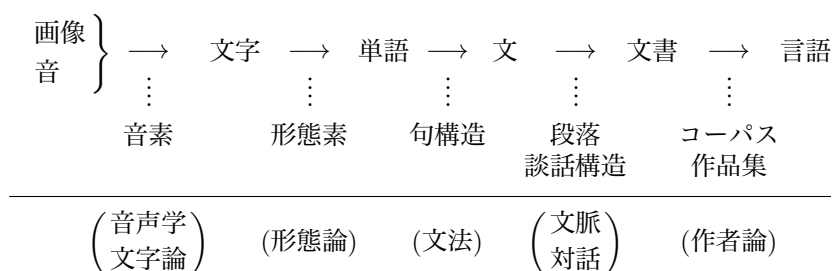


図 1.2: 言語単位の階層構造.

1.2 テキストの階層構造

先ほど、テキストには文字と単語があると書きました。よく考えると、テキストにはもっと多くの構造があることがわかります。この全体像を図 1.2 に示しました。

- (1) 第1のレベルは、文字です。これは知覚に際しては、画素 (ピクセル) の集合とみることができます。
- (2) 第2のレベルは、単語です。これはテキストでは文字の集合で、直接知覚する場合には、音声から得られる音素の集合です。^{*3}
- (3) 第3のレベルは、文です。これは単語の集合で、単語が連なって一つの文を作ります。
- (4) 第4のレベルは、文書です。これは文の集合で、文が連続して、まとまった意味を持つ一つの文書となります。^{*4}
- (5) 第5のレベルは、(狭い意味での) 言語です。多くの文書が日本語や英語など特定の言語で書かれていることを考えると、それらの文書の集合が、日本語や英語のテキスト全体になります。

もちろん、これらの中間のレベルもあり、たとえば接頭辞や接尾辞、語幹といった形態素を扱う形態論 (morphology) は文字と単語のレベルの中間にあり、名詞

^{*3} フランス人の子供が “beaucoup” のような難しい綴りを後から覚えるように、文字を介さずに単語を認識することも原理的には可能です。

^{*4} この「文書」のことをテキストとよぶことも多いのですが、本書では書かれた言葉全体をテキストと呼ぶことにします。

句や動詞句、係り受けといった構造は単語と文の中間に、段落や談話構造は文と文書の間で存在するでしょう。ただし多くの場合、単語、文、文書の構造はほぼ*5 自明に与えられているため、上記の分類を採用しています。

重要なことは、これらは**階層構造**をなしているということです。すなわち、文字の集合が単語に、単語の集合が文に、文の集合が文書に、文書の集合が言語になっています*6。このうち、どのレベルに注目するかは興味と目的によるでしょう。形態論に興味のある場合は文字がベースに、同意や補足といった談話構造に興味がある場合は文がベースになると考えられます。*7 エンジニアの方が文書を扱う場合も、一般的には単語または文をベースにして文書を考える必要があるでしょう。

そこで本書では、文字→単語→文→文書の順にその統計モデルを見ていくことにします。必要になる統計的な概念についても、基本的なレベルから少しずつ学んでいくことにしましょう。

1.3 教師あり学習と教師なし学習

言葉を計算機で扱う自然言語処理*8 は**機械学習**の一種とみなすことができますが、一般に機械学習は教師あり学習と教師なし学習に分けることができます*9。

教師あり学習の目標は、入力 \mathbf{x} に対する出力 \mathbf{y} を正しく求めることです。2章で学習する条件つき確率を使うと、 $p(\mathbf{y}|\mathbf{x})$ を求めることに相当します。自然言語処理の場合の入力 \mathbf{x} は多くの場合は文や文書で、出力 \mathbf{y} としてはその文の翻

*5 正確に言えば、「単語」や「文」とは何かというのは自明ではない問題です。日本語では空白で区切られた「単語」は存在しませんし、「。」で文の終わりが示されるものの、必ず「。」で終わっているとは限らず、口語の場合どこまでが文かも曖昧なところがあります。この問題については、4.1節を参照してください。

*6 言語学では、文は形態素の連続に、さらに語は音素の連続に分かれるという構造を**二重分節**と呼んでいます[2]。

*7 もちろん、形態論は前後の単語の影響を受けますし、談話構造は文に含まれる単語を考慮する必要があります。このように、隣接するレベルの情報は必要になりますが、隣接しないレベルは多くの場合必要ないでしょう。すなわち、形態論には文書の構造は関係せず、談話構造には文字のレベルはほとんど影響しないと考えられます。

*8 「言語」にはオートマトン理論などで使う形式言語やプログラミング言語なども含まれますので、実際に人間が使う言語を特に自然言語とよびます。もちろん、これらの間には深い関係があります。

*9 これに強化学習を加えて3つに分けることもあります。

1	2	3	4	5	6						
she	432	the	1026	was	277	and	466	way	45	little	92
to	387	a	473	had	126	of	343	mouse	41	great	23
i	324	her	116	said	113	in	262	thing	39	very	22
it	265	very	84	be	77	said	174	queen	37	long	22
you	218	its	50	is	73	to	163	head	36	large	22
alice	166	my	46	went	58	as	163	cat	35	right	20
and	147	no	44	were	56	that	125	hatter	34	same	17
they	76	his	44	see	52	for	123	duchess	34	good	17
there	61	this	39	could	52	at	122	well	31	white	11
he	55	an	37	know	50	but	121	time	31	other	11
that	39	your	36	thought	44	with	114	tone	28	poor	10
who	37	as	31	herself	42	on	83	rabbit	28	first	10

図 1.3: 『不思議の国のアリス』で HMM (4 章) の状態に割り当てられた単語とその回数。「主語」「動詞」「形容詞」といった概念が、**教師なし**で自動的に学習されています。

訳や対話の応答、文書のカテゴリなどになり、この学習には、人手で準備した大量のペア (\mathbf{x}, \mathbf{y}) を必要とします。これに対して**教師なし学習**では、言語、すなわち入力 \mathbf{x} 自体をモデル化します。確率の言葉では、 $p(\mathbf{x})$ を求めることとなります。これには何かモデルが必要ですから、実際には見たい構造 \mathbf{z} を未知の**潜在変数**として設定し、モデル $p(\mathbf{x}, \mathbf{z})$ を考えることとなります^{*10}。

この二つの違いを理解するために、本書の 3 章でも扱う品詞分析を考えてみましょう。教師あり学習では、 $\mathbf{x} = \text{“He is a boy”}$ のような入力に対する“正解”の品詞列 $\mathbf{y} = \text{“名詞-動詞-冠詞-名詞”}$ のような (\mathbf{x}, \mathbf{y}) のペアを、学習データとして人手で大量に (たとえば数万個) 準備します。学習の目標は、学習データにない新しい入力 $\mathbf{x} = \text{“Mary sings well”}$ のような文について、正しい品詞列を出力することです。このことの利点は、常に人間の想定に沿った出力が得られることです。一方で欠点は、**決して事前の想定内を出られない**ということです。入力文が学習データにない絵文字で終わっていたとき、できるのは既存の品詞の中から無理矢理どれかを選ぶことだけでしょう。また、見たことがない綴りの単語への対応も難しくなりますし、文書のカテゴリも事前に設定した中から選ぶだけしかできません。

これに対して教師なし学習では、入力文 \mathbf{x} の裏に未知の状態系列 \mathbf{z} があると

*10 可能な \mathbf{z} について周辺化すると、これは $p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$ を求めていることとなります。

考え、 $p(\mathbf{x}, \mathbf{z})$ が高くなるように \mathbf{z} を学習します。うまく学習すれば、ただ単語列を与えただけで、図1.3のように実質的に「品詞」とみなせる状態を推定することができます。このことの利点は、常にデータ \mathbf{x} だけを見ているので、適切な統計モデルを作れば、**データ自体から学習することが可能**だということです。上記のような絵文字が文末で特別な意味を持つと統計的に判断できれば、これに新しい状態が割り当てられますし、新しい言葉への対応も容易です。原理的には、未知の言語を扱うこともできます。逆に欠点は、学習が一般に難しいこと、また学習された状態が人間の望む「品詞」と一致するとは限らないということです。たとえば、日本語に“形容動詞”を認めるかについては言語学者の間にも諸説ありますが、モデルが特定の立場で出力してくれるとは限らず、解釈も別に必要になります。^{*11}

この後で学習するように、実はこの二つのモデルは掛け合わせると $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$ となり、 \mathbf{x} と \mathbf{y} の同時確率を与えますから^{*12}、本来はこの両方が必要です。特に、一部しか正解 \mathbf{y} が与えられない**半教師あり学習**のためには、教師なし学習のモデルが同時に必要になります。多くの教科書が教師あり学習を主に扱っていることから、本書ではより難しい問題である、教師なし学習に焦点を当てて解説することにします。^{*13}

1.4 統計的な方法とアドホックな方法

なおここで、本書がどうして「統計的」なテキストのモデルを考えるかについてふれておくことにします。たとえば、本書の5章では文書のモデル化を考えますが、ある文書が20個の単語からなり、含まれる単語の頻度がそれぞれ

“チューリップ”が4回、“栽培”が3回、“の”が8回、“こと”が5回

だったとしましょう^{*14}。この文書を本書のように確率的なモデルを考えず、単に頻度を並べたベクトル

^{*11} この他に、BERTの学習などでも使われている**自己教師あり学習**があります。これは観測された言葉を正例、それ以外のランダムな候補を負例として教師あり学習を行うもので、人手によるラベルは使っておらず、二者の中間とみなすことができるでしょう。

^{*12} 通常は教師なしデータ \mathbf{x} の方が圧倒的に多いため、ナイーブにこれを行うと \mathbf{y} が軽視されてしまいますので、実際にはもう少し検討が必要です[3, 4]。

^{*13} この二者の立場の違いは、事実上、工学と理学の違いといってもいいでしょう。

^{*14} これでは日本語になりませんが、説明のために簡単にしています。

$$\mathbf{x} = (0, \dots, 0, 4, 0, 3, 0, \dots, 0, 8, 0, \dots, 0, 5, 0, \dots, 0)$$

$\uparrow \quad \uparrow \quad \quad \uparrow \quad \quad \uparrow$
 チューリップ 栽培 の こと

で表せばよい、という人がいるかもしれません。^{*15}

この場合、上は語彙数 (たとえば 10000) の次元をもつ、超高次元のベクトルで、そのほとんどは 0 です。あれ、これでは文書の長さで表現が変わってしまうので、ベクトルの長さを 1 にする^{*16} ために $\sqrt{4^2+3^2+8^2+5^2}=\sqrt{114}$ で割って、

$$\mathbf{x} = \left(0, \dots, 0, \frac{4}{\sqrt{114}}, 0, \frac{3}{\sqrt{114}}, 0, \dots, 0, \frac{8}{\sqrt{114}}, 0, \dots, 0, \frac{5}{\sqrt{114}}, 0, \dots, 0 \right)$$

とした方がいいですね。あるいは、“の” や “こと” のような一般的な単語の重みを下げるために、各単語の頻度を tf.idf (5 章) で置き換えた方がいいでしょうか。こうしてベクトルが得られたら、あとは K 平均法でクラスタリングしたり、主成分分析にかけて次元圧縮し、「主成分」を抽出すれば、一丁あがりです。…

こうした素朴な方法は、特に文系や初心者を対象とした「テキストマイニング」の本でよくみられますが、何がいけないのでしょうか。

よくない理由は色々ありますが、まず第一に、上のような処理をすることに**直感以外の理由がなく、処理に任意性がありすぎる**、ということが挙げられます。上でベクトルを正規化するのに、和を 1 にする方法と長さを 1 にする方法がありました。どちらを使えばよいのでしょうか。また tf.idf には、tf 部分を対数にするのか (するなら、対数の底を何にするのか)、tf に 1 を足すのか、といった多くのバリエーションがあります。そもそも tf.idf が最適だという保証もありませんし、これらの選択によって、分析の結果はまったく違ってきてしまいます。一方、統計モデルを使えば、**数学的に何を最適化するのか**、という理由が明確になります。その際、本書の 5 章でみるように数学的に自然な形で “の” のような一般的な語の重みを自動的に下げることができ、ベクトルを正規化する必要もありません。実際に自然言語処理の歴史は、最初は発見的に導入された手法が、

^{*15} 実際にこれは「ベクトル空間モデル」とよばれ、初期の自然言語処理や情報検索で使われていました[5]。

^{*16} 頻度の総和で割って正規化することもできますが、そうすると値は、その単語の文書内での確率と同じです。それならば、なぜ最初から確率的に考えないのでしょうか？

表 1.1: 各テキスト A,B,C,D に現れた単語とその頻度の例.

テキスト	A	B	C	D
単語 1	0	6	20	0
単語 2	2	8	30	8
単語 3	0	1	12	1
⋮	⋮	⋮	⋮	⋮
単語 N	3	0	5	2
長さ	50	200	1000	500

数学的でより性能の優れた手法に置き換えられることの連続でした.

第二に、よく行われる主成分分析やベクトルのクラスタリングは、**本来連続的なデータを対象にしており、頻度のような離散的なデータに直接適用してはいけない**、という点が挙げられます. 上の文書ベクトルは超高次元 (たとえば 10000 次元) な上、値は必ず正で、そのほとんどは 0 になっている、という制約があります. 通常の主成分分析は値が負の場合もある連続値で、誤差が平均 0 のガウス分布に従うことを仮定しており [6], こうした状況には当てはまりません. 単語の頻度も十分に高ければ、近似的に連続値とみなすことができますが、これは頻度の高い、ごく一部の単語についてしか成り立ちません. 言語では数回といった低頻度でも意味があることが多く、たとえば、ある人物が「開腹」という言葉を数回使っただけで、われわれはこの人は医師あるいは医療系の方だという大きな情報を得ることができます. しかし、こうした低頻度の言葉を、通常の変量解析で適切に扱うことはできません.

たとえば、アンケートの自由回答やある小説家の作品群といった複数のテキストを比較するために、表 1.1 のように単語の頻度を数えたとしましょう. この表は、単語 1 はテキスト A には 0 回、B には 6 回、C には 20 回、…出現したことを表しています*17. このとき、テキストを比較してクラスタ分析を行うために、テキスト A の特徴ベクトルを、表 1.1 を縦に読んで $(0, 2, 0, \dots, 3)$, C の特徴ベクトルを $(20, 30, 12, \dots, 0)$ などとするのは、明らかに適切ではありません. なぜなら、そもそも表 1.1 の背後にある各テキストの長さが大きく違っているか

*17 政治学方法論の分野でテキスト分析のために広く使われている R のパッケージ `quanteda` では、こうした表を簡単に作ることができ、基本的なデータ構造になっています. また、仏教学の分野で石井公成氏が開発した NGSM [7] は、仏典の漢字 n グラムについてこうした表を作成して分析を行うものです.

表 1.2: 表 1.1 を, 各テキスト内で単語が出現する**確率**に直したもの.

テキスト	A	B	C	D
単語 1	0	0.03	0.02	0
単語 2	0.04	0.04	0.03	0.016
単語 3	0	0.005	0.012	0.002
⋮	⋮	⋮	⋮	⋮
単語 N	0.06	0	0.005	0.002

表 1.3: 確率の負の対数をとって, **情報量**に変換したもの.

テキスト	A	B	C	D
単語 1	∞	3.51	3.91	∞
単語 2	3.22	3.22	3.51	4.14
単語 3	∞	5.30	4.42	6.21
⋮	⋮	⋮	⋮	⋮
単語 N	2.81	∞	5.30	6.21

らです. テキスト C は長いので頻度は自然と大きくなり, このままではテキスト C の影響が過大に見積もられてしまいます^{*18}. また, テキスト D は表 1.1 で見えている範囲の出現頻度はテキスト A と似ていますが, もとになるテキストはずっと長いので, 頻度の意味はかなり違ってはいるはず.

明らかに, この場合は頻度そのものではなく, 各テキストの中でそれぞれの単語が出現する**確率**を考えるべきでしょう. 表 1.1 の頻度を各テキストの長さで割って求めた確率は, 表 1.2 のようになります. これで, テキストの長さにかかわらず, 単語の出やすさを平等に比較することができました.

ただし, これで終わりではありません. 単語 1 がテキスト B と C で出現する確率はそれぞれ 0.03 と 0.02 で, 0.01 の差があります. 1.5 倍の差ですね. 一方で単語 3 が C と D で出現する確率は 0.012 と 0.002 で, これも 0.01 の差ですが, 実は 6 倍の差があるわけです. これを同じ違いとするのは, 不合理ではないでしょうか^{*19}. こうした場合は表 1.3 のように, 確率の負の対数をとって本書の 2 章で説明する**情報量**に直せば, 違いを適切に表すことができます. 表 1.3 のように, このとき前者の差は $-\log 0.03 - (-\log 0.02) = 3.51 - 3.91 = -0.41$, 後者の差は $-\log 0.012 - (-\log 0.002) = 4.42 - 6.21 = -1.79$ となり, 後者の方が同じ確率 0.01 の差でも, 大きな意味があることを表現することができました. 実際にはさらに, ∞ を避けたり, 単語 2 のようにすべてのテキストで高確率で出現する単語 (「の」など) の重みを下げるために, 3 章で説明するように確率の平

*18 一般にアンケートの自由回答などでは, テキストの長さは大きく異なっているのが普通です.

*19 頻度は確率に比例しますので, 頻度を直接使うもとの方法も, 同じ問題を持っていることがわかります.

mass, moushly, toung, loop, cright, ked, whenneadfult,
sriciontict, leanquit, but, tring, hemple, foothed, wroar

図 1.4: 『不思議の国のアリス』の語彙をデータにして、文字 n グラムモデルからランダムに生成した単語の例. ($n=3$)

滑化を行ったり、非負の自己相互情報量 (PMI) を計算して特徴量とするのがよいでしょう。いずれにしても、こうした考察には、本書で説明するような**確率・統計的な知識**が不可欠になります。

第三に、最初に示したような**アドホックな方法は理論的な裏付けがないため、方法を拡張することができない**という問題もあります。たとえば、単語の頻度内外れ値がある場合や、欠損値がある場合はどうすればいいのでしょうか。文書をベクトル化したとして、その時間的な変化を見たい場合も、理論的な裏付けがなければ、「適当にプロットする」程度のことしかできません。これには無数の任意性がある上に、たとえば観測のなかった時期があると、一気に使えなくなってしまう。こうした例は、統計モデルを導入すれば、すべて自然な形で解くことができます*20。

1.5 本書の構成と読み方

本書で確率・統計的な手法を学ぶと、何ができるようになるのでしょうか。ここでは本書の構成を紹介しつつ、この後に出てくる解析例をいくつか紹介します。これらの結果のほとんどは、**深層学習のブラックボックスを無闇に回しても不可能なもの**であることに注意してください。なお、基礎についてはすでに知っている読者の場合でも、本書では実際の例や脚注などを通じて、面白く読めるように配慮しました。

先に説明したように、本書では文字・単語・文・文書の順でテキストの統計モデルについて説明します。内容はだんだん高度になりますが、最も基礎的な話が先になるように配置していますので、**一度で全部を理解できる必要はありません**。特に、文 (4 章) および文書 (5 章) の章はある程度独立に読めますので、興味のある部分から読み始め、必要に応じて前に戻ってもよいでしょう。

*20 この場合、統計モデルを使えば、5 章で説明するような別の数学的なベクトル化を行った上で、カルマンフィルタやガウス過程[8]で時間発展を追いかけることも可能になります。

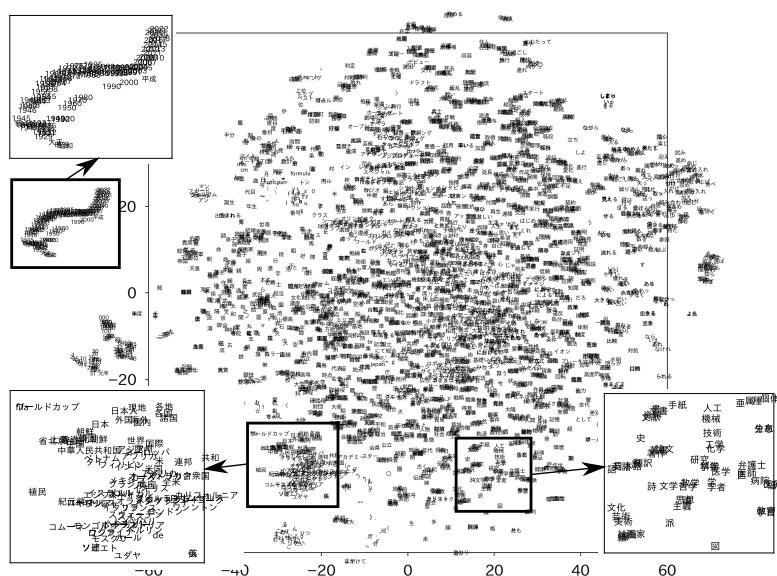


図 1.5: 日本語 Wikipedia のテキスト `ja.text9` から学習した 400 次元の単語ベクトルを、`t-SNE` で 2 次元に可視化したもの。意味の近い語が近いベクトルになっています。

2 章では、文字の統計モデルを入りに、確率と統計モデルの基礎について説明します。同時確率や条件つき確率、ベイズの定理などについて学び、統計モデルの評価についても説明します。文字 n グラムモデルを作ると、たとえば図 1.4 のように英語の単語をランダムに生成することも可能になります。情報量やエントロピーの意味についても、ここで学びましょう。

3 章ではそれを基礎に、単語の統計モデルについて説明します。文字よりも圧倒的に (原理的には無限に) 種類がある、単語の場合について有用な巾乗則や確率の平滑化について学んだ後、ニューラル単語ベクトルの学習とその性質について学びます。Word2Vec は数学的には何を学習しているのか、単語ベクトルの「引き算」はなぜ成り立つのかといった原理について説明します。日本語の場合、学習された単語ベクトルを 2 次元に可視化すると、たとえば図 1.5 のようになります。これらの間にはどんな関係が成り立っており、それはなぜなのでしょう。言語のベイズのモデル化に不可欠なディリクレ分布やポリア分布につい

でも、この章で学習します。

4章では単語が連なった文の統計モデルについて説明します。本章で解説する数学的な文ベクトルを使うと、図 1.6 のようにコーパスから、意味的に類似した文を簡単に計算することができます (これには複雑な深層学習は不要です)。また、文字や単語などの系列を扱う隠れマルコフモデル (HMM) とそのベイズ学習について説明し、この例を通じてマルコフ連鎖モンテカルロ法 (MCMC) の一種である Gibbs サンプリングについて学びます。HMM を使うと、たとえば単語の背後に隠れた「品詞」を人間の主観を交じえず、統計的に教師なし学習することが可能です。図 1.3 に、『不思議の国のアリス』の本文だけから自動的に学習された「品詞」の例を示しました。

5章では、最も実用例の多いと考えられる文書の統計モデルについて説明します。簡単なナイーブベイズ法による分類から始め、意味的なクラスタを教師なしで発見する UM および DM, その学習のための EM アルゴリズムについて学びます。さらに、その拡張としてトピックモデルとして知られる LDA (潜在ディリクレ配分法) について説明します。LDA では、文書集合から図 1.7 のような潜在的な「トピック」を完全に自動的に学習し、各文書をこのトピックの上の確率分布 θ として表します。最後に、単語ベクトルを用いた数学的な「文書ベクトル」の計算法について学んだ後、若干高度ですが、筆者の研究である PLSS (確率的潜在意味スケール) について紹介します。PLSS では、心理統計学の理論に基づき、テキストにある軸で測った潜在的な極性 θ を連続値として図 5.42 のよ

類似度	文	類似度	文
1.0000	だが、新しい後期高齢者医療制度では、介護..	1.0000	再生可能なファイル形式は、映像が MPEG..
0.8929	健康保険や介護保険、厚生年金、雇用保険、..	0.9054	無線 LAN セキュリティは 64/128bit の W..
0.8705	国は患者が混合診療を受けた場合、「一体化..	0.8790	その他の機能は地上デジタル/BS デジタル/..
0.8597	労働保険は、法人個人を問わず労働者を 1 人..	0.8780	ネットワーク機能は 10/100/1000BASE-T..
0.8313	また、短期入所や通所を受け入れる福祉施設..	0.8731	録音形式はリニア PCM で 16bit/44.1kHz..
0.8212	「住宅ローン控除」は、国内で一定の居住用..	0.8708	基本仕様は MP3/WMA/AAC 再生。
0.8109	全体で 5% アップと同水準だが、保険制度の..	0.8627	その他の機能は、IEEE802.11b/g/n 対応..
0.8107	この免除の手続きをするだけで、保険料を払..	0.8588	入出力端子には HDMI/コンポジットビデオ..
0.8098	「年金制度は世代間扶養の仕組みである」→..	0.8486	同サービスは、i モード/EZweb/Yahoo!ケ..
0.8085	また、介護保険は対象外となっています。	0.8485	CG-BARPROG-X コレガは、WAN/LAN..

図 1.6: Leipzig コーパスの日本語ニュースの文について、uSIF による文ベクトルを用いて計算した類似文。数字はコサイン類似度を表します。自分自身との類似度は 1 ですが、それ以外にも意味的に類似した文が、簡単な計算で検索できることがわかります。

Topic 56		Topic 104		Topic 206		Topic 234	
江戸	0.0664	宇宙	0.1441	省	0.2087	細胞	0.0915
藩	0.0585	地球	0.0506	局	0.0781	質	0.0485
時代	0.0423	衛星	0.0498	官	0.0445	タンパク	0.0341
藩主	0.0307	号	0.0430	管理	0.0425	体	0.0287
家	0.0228	ロケット	0.0377	国土	0.0336	組織	0.0266
通称	0.0199	打ち上げ	0.0317	文部	0.0287	活性	0.0239
元禄	0.0166	機	0.0226	水産	0.0287	免疫	0.0239
守	0.0162	探査	0.0226	交通	0.0267	性	0.0232
酒井	0.0145	初	0.0219	政府	0.0237	ヒト	0.0232
Topic 258		Topic 274		Topic 334		Topic 397	
体	0.0684	教会	0.1176	建築	0.2156	韓国	0.1691
多様	0.0522	教	0.1029	設計	0.1198	朝鮮	0.0939
空間	0.0476	キリスト	0.0801	家	0.0816	民国	0.0718
次元	0.0469	宗教	0.0544	建設	0.0506	大韓	0.0718
的	0.0399	派	0.0419	建物	0.0497	ソウル	0.0641
環	0.0269	教団	0.0353	住宅	0.0364	広域	0.0575
定義	0.0246	神	0.0309	構造	0.0311	仁川	0.0321
上	0.0238	牧師	0.0272	物	0.0302	後述	0.0221
幾何	0.0238	神学	0.0243	建て	0.0222	区	0.0210

図 1.7: 日本語 Wikipedia から LDA で自動的に推定された潜在トピックの例 ($K=400$). 数字は、各単語がそのトピックから生成される確率を表しています.

うに推定することができます.

本書では、必要になる確率論や情報理論の基礎、EM アルゴリズムや MCMC 法などの推定法についても、従来の教科書のように無味乾燥な「予備知識」の章を立てるのではなく、必要に応じて本文の中で、言語の例を使いながら説明することにしました. 特定の概念 (たとえばエントロピー) が何だったかな、と復習したい方は、巻末の索引を利用してください.

また、本書はわれわれの周りに自然にみられるテキストをどうモデル化するかを考えていますので、文書分類や通常の形態素解析のように、テキストにその所属するクラスや単語境界といった「正解」があらかじめ人手で付与されている場合の分類問題、すなわち教師あり学習は基本的に対象としていません.*21 これらは多くの場合、SVM や CRF といった汎用の機械学習手法を適用すれば充分だからです. テキストに対する教師あり学習については、巻末の参考書ガイド

*21 確率の言葉でいえば、通常の教師あり機械学習はテキスト \mathbf{x} に付与されたラベル \mathbf{y} をいかに予測するか、すなわちラベルの確率 $p(\mathbf{y}|\mathbf{x})$ だけをモデル化しており、本書のようにテキスト自体の構造やそのラベルとの関係、すなわち $p(\mathbf{x})$ や $p(\mathbf{x}, \mathbf{y})$ はモデル化していない、ということです.

を参照してください。

1.6 本書の例と実装について

本書で使っているスクリプトやデータはすべて、iv ページの方法でサポートサイトから入手することができます。無印のタイプライター体の例は、対応する章の Jupyter Notebook にあるように、Python の Jupyter Notebook に順番に入力して実行できるようになっています。⇒ は出力です*22。

```
print ("これは例です。")  
⇒ これは例です。
```

ただし、複雑な例ではノートブック上での入力や実行には限界がありますので、%で始まる行はスクリプトやファイルのある各章のフォルダでコマンドラインを実行することを表しています。%から後が、実際に実行するコマンドです。

```
% date  
⇒ 2020 年 11 月 22 日 日曜日 22 時 35 分 56 秒 JST
```

コマンドラインは、MacOS や Linux では「ターミナル」アプリを起動すればそのまま使うことができます。Windows では、WSL (Windows Subsystem for Linux) を導入するといいいでしょう。こうした環境の構築については、たとえば [9, 3 章]などで実行例とともに詳しく解説されていますので、そちらを参照してください。

テキストを扱う場合は、Jupyter Notebook や RStudio のような環境の中ではできることが限られるため、プログラム開発の面からも、ぜひスクリプトを書いてコマンドラインで実行できるようになっていただきたいと考えています。これにより複雑なプログラムとそのデバッグが可能になるほか、汎用性のあるモジュールを自分で書いて import することで、プログラムの再利用が可能になり、より見通しのよい開発を行うことができます。テキストを扱うには豊富な標準コマンドがあり、本書でも様々なスクリプトを書いて使っていますので、ノートブックや統合環境の箱庭の中だけで分析を行うのではなく、ぜひコマンドラインを使いこなせるようになってください。

*22 Jupyter Notebook の Out []:にあたります。

なお、記述を短くするため、実行例はすべてのスクリプトやファイルがカレントフォルダにあることを前提にして書かれています。実際には iv ページのように `git clone` するとスクリプトは `bin/` に、データは `data/` などにありますので、実行の際はすべてをカレントに移動するか、適宜パスを補ってください。本書は「何も考えずに実行できる」といった本ではなく、スクリプトはすべて内容を読んで使うことを前提にしています。

本書の Python スクリプトもこれが「正解」ということはなく、様々な実装があります。理解のために、ぜひ中身を読んで自分で試してみてください。その際、各章の最後にある「実習」で興味を持ったものを行ってみるとよいでしょう。

2 文字の統計モデル

2.1 文字の頻度と出現確率

それでは、実際にテキストをみていくことにしましょう。どんなテキストも、改行を含めて文字の系列、すなわち**文字列**とみなすことができます。日本語や中国語のような言語は空白で単語に分けられていませんので、文字列は最も基本的で重要な表現です。英語のほうが日本語より文字の種類が少ないので、簡単のために、まずは英語の例で考えてみることにしましょう。

本書のサポートサイト（「はじめに」v ページ）にある『不思議の国のアリス』のテキスト `alice.txt` は 1,430 行、長さ 132,656 文字のテキストで、図 2.1 のような内容になっています。^{*1} 説明のため、すべて小文字にして記号を削除してありますので^{*2}、文字は “a”, ..., “z” およびスペース “ ” の 27 文字からなっています。

このファイルを文字列として Python に読み込むには、次のように実行します。

```
with open('alice.txt', 'r') as f:
    s = f.read()
```

これで、`alice.txt` の内容が文字列 `s` に代入されました。 `s` の中身を見てみると、

```
s
⇒ 'alices adventures in wonderland\nlewis carroll\nchapter i
   \ndown the rabbithole\nalice was beginning to get very tir
```

^{*1} このテキストは、著作権の切れた文学作品を集めた Project Gutenberg (<https://www.gutenberg.org/>) から入手したものです。

^{*2} 英語の場合は何を大文字にするかには規則性がありますので、すべて小文字にしても情報量が落ちることはあまりなく、機械学習を適切に使えば、原文をほぼ復元することができます。

```

alices adventures in wonderland
lewis carroll
chapter i
down the rabbithole
alice was beginning to get very tired of sitting by her
sister on the bank and of having nothing to do once or twice
she had peeped into the book her sister was reading but it
had no pictures or conversations in it and what is the use
of a book thought alice without pictures or conversation
so she was considering in her own mind as well as she could
for the hot day made her feel very sleepy and stupid whether
the pleasure of making a daisychain would be worth the
trouble of getting up and picking the daisies when suddenly
a white rabbit with pink eyes ran close by her there was
nothing so very remarkable in that nor did alice think ...

```

図 2.1: 『不思議の国のアリス』 alice.txt の冒頭部.

```

ed of sitting by her sister on the bank and of having noth
ing to do once or twice she had ...'

```

のようになっています。“\n” は、改行を表す特別な文字です*3.

この中で、それぞれの文字が何回くらい現れるのか、頻度を数えてみることにしましょう。改行文字は無視することにする、

```

from collections import defaultdict
freq = defaultdict(int)
for c in s:          # c = s[0],s[1],... を順に調べる
    if c != '\n':    # c が改行文字以外の場合
        freq[c] += 1 # 頻度を増やす

```

で文字の頻度を数えることができます。#以下はコメントで無視されますので、入力する必要はありません。頻度を数えるテーブルの freq は Python の辞書 (dict) なので、freq = {} として初期化してもよいのですが、freq[c] の初期値を自動的に 0 にするため、collections モジュールの defaultdict を使っています。

結果は、次のようになりました。

```

for (c,n) in freq.items():

```

*3 バックスラッシュ\は特別な文字で、続く 1 文字と合わせて 1 つの文字を表します。本書では、実行例で行が続くことも\で表します。


```

    print ('%s = %d' % (c,n))
⇒ a = 8791
   l = 4713
   i = 7510
   c = 2398
   e = 13571
   s = 6500
   _ = 24966 ...

```

頻度順になっていないので、少し見づらいですね。頻度順に表示するには、次のように入力します。

```

    for (c,n) in sorted (freq.items(), # 頻度の降順にソート
                        key=lambda x: x[1], reverse=True):
        print ('%s = %d' % (c,n))
⇒ _ = 24966
   e = 13571
   t = 10687
   a = 8791
   o = 8145
   i = 7510
   h = 7373
   n = 7013
   s = 6500
   ...
   q = 209
   x = 148
   j = 146
   z = 78

```

頻度の合計は、最初に述べたテキストの長さ 132656 と等しくなります。これを N としましょう。

```

    sum (freq.values())
⇒ 132656

```

これから、各文字の出現確率を計算することができます。文字 c の頻度を $n(c)$ (上の Python コードでは `freq[c]`) とおくと、 c が出現する確率は、最も単純には

$$(2.1) \quad p(c) = \frac{n(c)}{N} \quad \text{(確率の最尤推定)}$$

と考えてよいでしょう。 $p()$ は、確率を表す記法です。^{*4} たとえば、上の例では $n(e)=13571$, $n(a)=8791$ でしたから、文字 e や a の確率は

$$(2.2) \quad p(e) = \frac{n(e)}{N} = \frac{13571}{132656} = 0.1023$$

$$p(a) = \frac{n(a)}{N} = \frac{8791}{132656} = 0.0663$$

と求めることができます。一方、頻度の小さい q や z の確率は

$$(2.3) \quad p(q) = \frac{n(q)}{N} = \frac{209}{132656} = 0.0016$$

$$p(z) = \frac{n(z)}{N} = \frac{78}{132656} = 0.0006$$

のような値になります。だいたい 150 倍くらいの違いがありますね。Python で各文字 c の確率を計算して変数 $p[c]$ に保存するには、次のように実行します。

```
p = {}
N = sum(freq.values())
for (c,n) in freq.items():
    p[c] = n / N
print(p)
⇒ {'a': 0.06626914726812207,
    'l': 0.035527982149318536,
    'i': 0.05661259196719334,
    'c': 0.0180768302979134,
    'e': 0.10230219515136896, ...
}
```

図 2.2 に、こうして計算した『不思議の国のアリス』の記号を除く各文字の出現確率と、その棒グラフを示しました。英語では a , e , i , o のような母音の確率が際立って高く、子音でも r , s , t などは確率が高いことがわかります。一方で、 j や q , x は非常に低い出現確率 (0.001 以下) になっています。図 2.2 の各

^{*4} $\text{Pr}()$ と書いたり、離散のとき $P()$, 連続のとき $p()$ と分けて書く流儀もありますが、混同する危険性はありませんから、本書では機械学習の分野で最も標準的な $p()$ で表すことにします。 $p(x)$ は一般に、“ p of x ” (probability of x の意味) と読みます。筆者は個人的に、「ピーエックス」や「 p の x 」と日本語で読んでしまうこともあります。

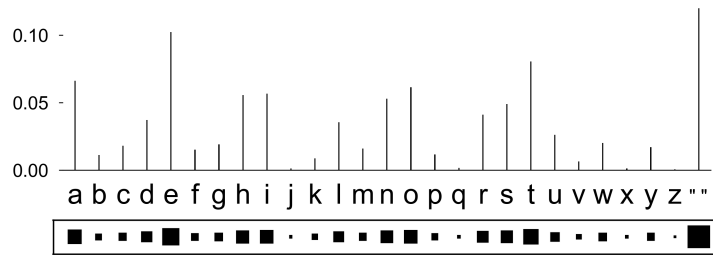


図 2.2: 『不思議の国のアリス』alice.txt における各文字の出現確率.

文字の下にある■はヒント図 (Hinton diagram)^{*5} といい、■が大きいほど確率が高いことを表しており、確率の大きさを直感的に表現するためによく用いられます。式(2.1)の確率は最尤推定とよばれ、現在手元にある学習データの確率(これを 38 ページで説明するように、尤度^{さいゆうど}といいます)を最大にする値ですが^{*6}、ここでは簡単に、頻度を割り算して計算するもっとも簡単な確率の推定値だと考えてください。

図 2.2 のそれぞれの棒は文字の確率を表していますから、その長さの総和は必ず 1 になります。このように、総和が 1 となる確率を並べたものを**確率分布**といいます。図 2.2 は、『不思議の国のアリス』の裏に隠れた、文字の確率分布です。

2.2 文字の同時確率

上では文字をばらばらに数えていましたが、実際には英語では、“es” や “li” のような 2 文字の連続の確率が高く、“qy” や “lg” のような文字の連続の確率は非常に低いと考えられます。こうした 2 文字の連続を、**バイグラム** (bigram) といいます。これに対して、2.1 節で考えた 1 文字を**ユニグラム** (unigram) といいます^{*7}。「グラム」とは、「文字」の意味です。バイグラムの出現確率は、ど

*5 ニューラルネットワークと深層学習の基礎を作ったことで有名なトロント大学の Geoffrey Hinton 教授によるもので、本書のサポートサイトにも Python の実装 `hinton.py` を置きました。

*6 式(2.1)の確率が手元のデータの確率を最大にすることは、数学的には自明ではなく、ラグランジュの未定乗数法を使った簡単な証明が必要です。本章は導入のため証明は割愛していますので、興味のある方は[10, 1.5.2 節]などを参照してください。

*7 古典ギリシャ語的な表現であるユニグラム、バイグラム、…の代わりに、それぞれ数字で 1 グラム、2 グラム、…と言うこともあります。

うやって調べればよいでしょうか。

答えは簡単で、たとえば文字列が `s = "alice"` であれば、`"al"`、`"li"`、`"ic"`、`"ce"`のように1文字ずつずらしながら、2文字の連続を同様に数えていけばよいだけです。この場合、`s`の最後の文字“e”には続きの文字がありませんので、後で説明する理由で、文字列の最後に文字列の終わりを表す特別な文字“\$”^{*8}があるとして計算すると、バイグラムの総数は文字列の長さと同しくなります。バイグラムの頻度は、Pythonでは次のようにして数えることができますでしょう。ここからは、テキストが非常に大きい場合にメモリを消費しないよう、はじめの例のようにファイルから文字列 `s` に一気に読み込むことはせず、1行ずつ読んでいくことにします。`alice.txt`では、1行がほぼ1文に対応しています。

```
freq = defaultdict (int)
with open ('alice.txt', 'r') as f:
    for line in f:          # 改行文字を含めて1行ずつ読む
        s = line.rstrip('\n') # 行の最後の改行文字を除去
        L = len(s)
        s += '$'           # 文字列末尾を表す特殊文字を追加
        for i in range(L): # 1文字ずつずらして数える
            b = s[i:i+2]   # 文字バイグラム (2文字)
            freq[b] += 1   # 頻度を増やす
```

数えた結果は、以下のようになりました。

```
for (b,n) in sorted (freq.items(), # 降順に表示
                    key=lambda x: x[1], reverse=True):
    print ('%s = %d' % (b, n))
print ('total %d bigrams.' % sum(freq.values()))
⇒ e_ = 5417
   _t = 4192
   he = 3778
   th = 3483
   _a = 3152
   ...
   ju = 102
   y$ = 102
   od = 101
   oc = 99
```

^{*8} この“\$”は説明のための表記で、実際には本当の\$と区別するために、テキストに出現しない特別な文字を用います。詳しくは、41ページの脚注を参照してください。

Y	a	b	c	d	e	f	g	h	...	\$	合計
$p(X=a, Y)$	0.000	0.002	0.001	0.003	0.000	0.000	0.001	0.000	...	0.000	0.066

表 2.1: 同時確率 $p(X=a, Y)$ の表. 0.001 未満の確率は四捨五入されています.

のようになっていますが, これは

$$(2.7) \quad p(X=a, Y=a), p(X=a, Y=b), p(X=a, Y=c), \dots, p(X=a, Y=\$)$$

を並べたもので, 確率としては表 2.1 のようになっています.*9

式(2.7) は, もとの X, Y を使わない表記では

$$(2.8) \quad p(aa), p(ab), p(ac), p(ad), \dots, p(a\$)$$

すなわち $p(a*)$ ($*$ は任意の 1 文字) のことですから, これらの確率の総和は a が現れる確率, すなわち $p(a)$ と等しくなります. つまり, a 自体が現れる確率は, a の後に色々な文字が現れる場合の確率 $p(aa), p(ab), p(ac), \dots$ をすべて足したものになるわけです. 当たり前ですね. このことを X, Y を使って表すと,

$$(2.9) \quad p(X=a) = p(X=a, Y=a) + p(X=a, Y=b) + \dots + p(X=a, Y=\$) \\ = \sum_Y p(X=a, Y)$$

ということになります. 式(2.9)を, 同時確率 $p(X=a, Y)$ の**周辺化** (marginalization) といいます. これは図 2.3 で行列の各行の和をとる, すなわち行列を Y に関して横方向につぶして, 値を行列の「周辺」(margin, 縁)に集めていることに相当しますので, これが「周辺化」の名前の由来です. 同時確率の周辺化は, 一般的に

$$(2.10) \quad p(X) = \sum_Y p(X, Y) \quad \text{(同時確率の周辺化の公式)}$$

と表すことができます. この公式は, X と Y が同時に現れる確率をすべての Y について和をとれば, X だけの確率になるという当たり前のことを表しています.

*9 Hinton 図は相対的な確率を示すものですので, ■が大きいても絶対的な確率が 1 に近いとは限らないことに注意してください.

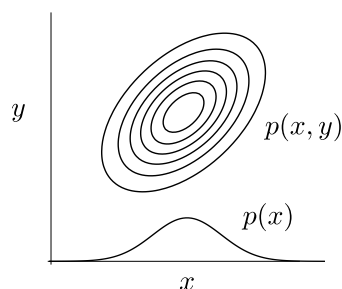


図 2.4: 同時確率密度 $p(x, y)$ の周辺化. y に関して積分して確率密度関数を下に「つぶす」ことで, x だけの関数 $\int p(x, y) dy = p(x)$ が得られます.

実際に確かめてみましょう. 上の例では,

$$(2.11) \quad p(\text{aa}) + p(\text{ab}) + \dots + p(\text{a\$}) = 0.000 + 0.002 + \dots + 0.000 = 0.066$$

ですが, 式(2.3)より $p(\text{a}) = 0.066$ でしたので, この2つはぴったり同じになっています. Python では, バイグラムで計算した $p[x][y]$ を使って

```

for x in p:
    s = 0
    for y in p[x]:
        s += p[x][y]
    print ('p(%s) = %f' % (x, s))
⇒ p[a] = 0.066269
   p[l] = 0.035528
   p[i] = 0.056613
   p[c] = 0.018077
   p[e] = 0.102302
   ...

```

とすれば, バイグラム確率を周辺化してユニグラム確率を計算することができます. 結果はもちろん図 2.2 と同じになり, それを図 2.3 の右側の縁に示しました.

同時確率の周辺化は言語のように離散値の場合だけでなく, 連続値の場合でも同様に成り立ちます. たとえば $p(x, y)$ が図 2.4 のように2次元のガウス分布

(正規分布) のとき, y 方向に分布をつぶして和をとれば, x の分布 $p(x)$ は 1 次元のガウス分布になることが知られています[8]. これは,

$$(2.12) \quad p(x) = \int p(x, y) dy \quad (\text{同時確率の周辺化の公式 (連続値の場合)})$$

を計算していることになります. このように, $p(x, y)$ が確率密度の場合でも, 和を積分に置き換えれば, 式(2.12)の形で同時分布の周辺化が成り立ちます*10.

◇

同時確率と周辺化の一般的な説明では, サイコロを 2 つ用意して, X を 1 個目のサイコロの目, Y を 2 個目のサイコロの目とすることが多いようです. このとき

$$(2.13) \quad p(X=1) = p(X=2) = \dots = p(X=6) = \frac{1}{6}$$

ですが, 2 つのサイコロは独立なので, 特定の (X, Y) が出る確率はすべて

$$(2.14) \quad p(X=1, Y=1) = p(X=1, Y=2) = \dots = p(X=6, Y=6) = \frac{1}{6} \times \frac{1}{6} = \frac{1}{36}$$

になります. このとき, たとえば 1 個目のサイコロの目が $X=1$ であれば, $p(X, Y)$ を 2 個目のサイコロの目 Y について周辺化すれば

$$(2.15) \quad \begin{aligned} \sum_Y p(X=1, Y) &= p(X=1, Y=1) + p(X=1, Y=2) + \dots + p(X=6, Y=6) \\ &= \frac{1}{36} + \frac{1}{36} + \dots + \frac{1}{36} = \frac{1}{6} \end{aligned}$$

となり,

$$(2.16) \quad \sum_Y p(X=1, Y) = p(X=1)$$

が成り立つことがわかります.

サイコロの例では確率 $p(X, Y)$ がすべて等しくなりましたが, われわれ

*10 ライブニッツの積分記号 \int は, そもそも Sum (和) を表す S の変形であることを思い出してください. 厳密にはこれらはすべて証明が必要ですが, 本書は確率論の本ではありませんので, 立ち入らないことにします. 基礎的な定義が気になる方は, [11]や最近の[12]といった優れた教科書を参照してください.

$$\begin{aligned}
 (2.17) \quad p(Y=c|X=a) &= \frac{p(X=a, Y=c)}{p(X=a, Y=a) + p(X=a, Y=b) + \dots + p(X=a, Y=\$)} \\
 &= \frac{0.001}{0.000 + 0.002 + \dots + 0.000} = \frac{0.001}{0.0066} = 0.015
 \end{aligned}$$

と計算することができます.

2.3節で, 式(2.17)の2行目の分母は同時確率の周辺化によって

$$p(X=a, Y=a) + p(X=a, Y=b) + \dots + p(X=a, Y=\$) = p(X=a)$$

であることをみました. よって, 式(2.17)は

$$(2.18) \quad p(Y=c|X=a) = \frac{p(X=a, Y=c)}{p(X=a)}$$

と表すことができます. すなわち一般に, 条件つき確率は

$$(2.19) \quad p(Y|X) = \frac{p(X, Y)}{p(X)} \quad (\text{条件つき確率の公式})$$

で計算することができます. この式は文字通り, X が起きたときに Y が起きる条件つき確率 $p(Y|X)$ は, X が起きる確率 $p(X)$ に対する, Y も同時に起きる確率 $p(X, Y)$ の相対的な割合だという当たり前のことを表しています.

これを見ると, j の後に u が来る確率 $p(u|j)$ が非常に大きいことや, w の後に k が来る確率 $p(k|w)$ はほとんど0であることなど, 興味深い規則性をさまざまに読み取ることができるでしょう. (→演習(3))

23 ページで計算した同時確率 $q[x][y]$ を使うと, 式(2.19)の条件つき確率は次のようにして計算して $c[x][y]$ に保存することができます.

```

p0 = {}; c = {}
chars = p.keys()
for c in chars:
    p0[c] = sum(p[c].values())           # 周辺確率の計算
for x in chars:
    c[x] = {}
    for y in chars:

```

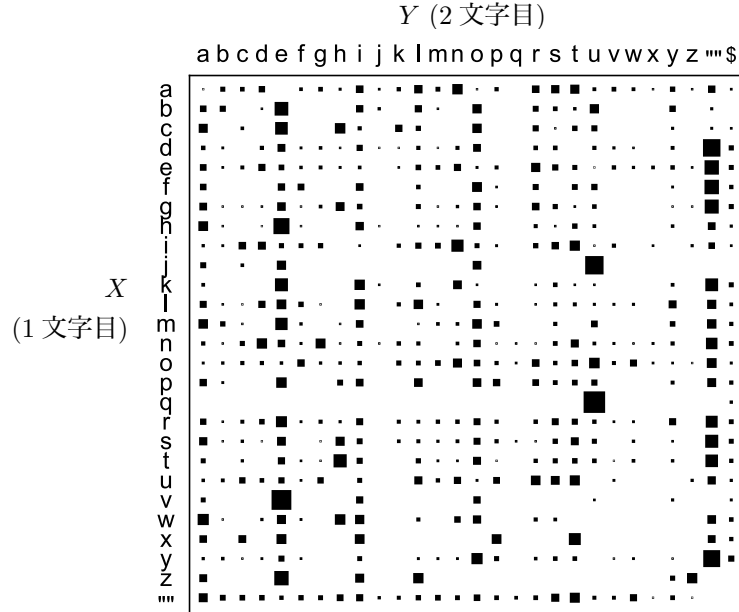


図 2.5: 『不思議の国のアリス』の文字バイグラム条件つき確率 $p(Y|X)$ を表す行列. 確率の総和は各行で 1 になっており, 興味深いパターンがみられることがわかります. 同時確率を表す図 2.3 と比べてみましょう.

```

if y in p[x]:
    c[x][y] = p[x][y] / p0[x] # 条件つき確率を計算
else:
    c[x][y] = 0

```

図 2.5 に, `alice.txt` からこうして計算した条件つき確率 $p(Y|X)$ をすべての文字 X, Y について示しました. 式(2.19)から, これは図 2.3 の $p(X, Y)$ を $p(X)$, すなわち各行の和で割ったものですから, $p(Y|X)$ は図 2.3 を行方向に正規化して和を 1 にしたものになっています.

2.4.1 確率の連鎖則

条件つき確率は公式 (2.19) で計算できますが, ただし, この**公式を暗記する必要はありません**. 条件つき確率 $p(Y|X)$ は, 「 X が起きたときに Y が起きる

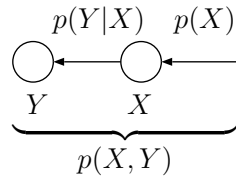


図 2.6: 確率の連鎖則 $p(X, Y) = p(Y|X)p(X)$ のイメージ. X と Y が両方起こるとは、まず X が起こり、その下でさらに Y が起きることと同じです。

確率」を表しているのです。 $p(X, Y)$ は「 X と Y が同時に^{*13} 起きる」確率のことですが、これは「まず X が起きてから、その条件の下でさらに Y が起きる」ことと同じです。つまり、条件つき確率の意味から、

$$(2.20) \quad p(X, Y) = p(Y|X)p(X) \quad (\text{確率の連鎖則})$$

が明らかに成り立ちます。これを**確率の連鎖則** (chain rule) といいます。この様子を、図 2.6 に示しました。

$p(X) \neq 0$ ならば、式(2.20)の両辺を $p(X)$ で割れば公式(2.19)が簡単に得られますから、**条件つき確率の公式(2.19)は自明な連鎖則(2.20)から簡単に得られます**。よって、公式(2.19)を暗記する必要はなく、慣れるまでは式(2.20)からその場で導くとよいでしょう。条件つき確率の公式(2.19)は確率の連鎖則(2.20)から明らかですが、その直感的な意味は、式(2.17)に示したように、同時確率 $p(X, Y)$ を注目している条件 X について正規化して、相対的な値を考えるということです。

なお、

$$(2.21) \quad p(X, Y) = p(X)p(Y) \quad (\text{独立な事象の同時確率})$$

が成り立つ場合、つまり式(2.20)と見比べれば

$$(2.22) \quad p(Y|X) = p(Y)$$

^{*13} 慣例的に joint probability は「同時確率」または「結合確率」と訳されていますが、「同時」というのは、必ずしも時間的に同時刻に起きるという意味ではなく、「両方起こる」という意味です。

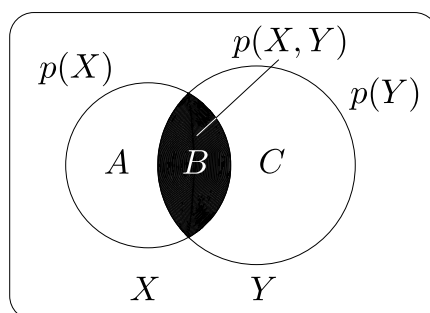


図 2.7: ベン図で表した確率と同時確率, および条件つき確率の関係. 枠で囲った全体の確率を 1 としたとき, 同時確率 $p(X, Y)$ は B の面積で, 条件つき確率 $p(Y|X)$ は割合 $B/(A+B)$ のことです. $p(X)$ を表す面積 $A+B$ に条件つき確率 $p(Y|X)$ をかければ B の面積, すなわち $p(X, Y)$ が得られます.

が成り立つ場合, X と Y は**独立**である, といいます. 式(2.22)は, X と Y が独立ならば X が与えられた下での Y の確率 $p(Y|X)$ はもともとの Y の確率 $p(Y)$ と等しい, つまり Y の確率は X の確率に影響を受けないことを表しているわけです. このとき, 式(2.21)のように X と Y の同時確率 $p(X, Y)$ は単にそれぞれの周辺確率 $p(X)$ と $p(Y)$ の積で表すことができます.

面積でみる確率 こうした条件つき確率と同時確率の関係は, 図 2.7 のようなベン図で表してみるとわかりやすいでしょう. 外側の四角で囲われた全体の面積を確率の総和である 1 とすると, 2 つの円の面積はそれぞれ, X と Y が起きる確率 $p(X)$ と $p(Y)$ を表しています. 黒で示した重なりが, X と Y の同時確率 $p(X, Y)$ です. 2 つの円で区切られた 3 つの領域を図のように A, B, C とおくと, $p(X) = A+B$, $p(Y) = B+C$, $p(X, Y) = B$ となっています.

このとき, X が起きた中で Y が起きる確率 $p(Y|X)$ は $\frac{B}{A+B}$ で, これは $\frac{p(X, Y)}{p(X)}$ を意味します. 逆に割合 $p(Y|X) = \frac{B}{A+B}$ が先にわかっているならば, 真ん中の領域の面積は X の面積にこれをかければよく, $(A+B) \cdot \frac{B}{A+B} = B$ となるわけです. これはつまり, $p(X) \cdot p(Y|X) = p(X, Y)$ を意味しています.

条件つき確率と頻度 なお、実はわれわれの場合は、条件つき確率はバイグラムの頻度から直接計算することができます。表 2.1 および図 2.3 のもとになっているバイグラムの頻度は、 $n(\text{aa}) = 1$, $n(\text{ab}) = 214$, $n(\text{ac}) = 157$, ..., $n(\text{a\$}) = 11$ で、その総和は $n(\text{a}) = 8791$ です。つまりこれは、 a が現れた 8791 回のうち、 c が続いたのは 157 回だったということですから、 $p(\text{c}|\text{a})$ は

$$(2.23) \quad p(\text{c}|\text{a}) = \frac{n(\text{ac})}{n(\text{a})} = \frac{157}{8791} = 0.0179$$

と求めることができます。

これはなぜかという、頻度 $n()$ を使った式(2.23)は、同時確率を使った式(2.19)と等価になるからです。式(2.1)および式(2.4)から

$$(2.24) \quad p(\text{ac}) = \frac{n(\text{ac})}{N}, \quad p(\text{a}) = \frac{n(\text{a})}{N}$$

ですから、式(2.23)は

$$\frac{n(\text{ac})}{n(\text{a})} = \frac{\frac{n(\text{ac})}{N}}{\frac{n(\text{a})}{N}} = \frac{p(\text{ac})}{p(\text{a})} = p(\text{c}|\text{a})$$

となり、条件つき確率に等しくなります。このため、以下では一般に、文字列 h の後に文字 c が現れる条件つき確率は、頻度 $n()$ を用いて

$$(2.25) \quad p(\text{c}|h) = \frac{n(\text{hc})}{n(h)}$$

として計算することにします。条件となる部分は後でみるように 1 文字とは限りませんので、履歴 (history) という意味で、 h で表しています。

2.4.2 ベイズの定理

条件つき確率を用いると、たとえば文字 j の後に文字 a が来る確率 $p(\text{a}|j)$ は、式(2.23)のように計算することができます。それでは逆に、 j の前に来ることができる文字には、どんなものがあるでしょうか^{*14}。

^{*14} 筆者は大学時代に 에스ぺ란토 研究会に所属していましたが、世界共通語を目指して作られた人工語である 에스ぺ란토 語では、 j は複数を表す接尾辞で、libroj (本たち) や ĉiu (すべて)

確率で書けば, これは, 前に来る文字を x として

$$(2.26) \quad p(j|x) \quad (x \text{ は前に来る任意の 1 文字})$$

を求めたい, ということです. 式(2.26)はテキストで xj となる確率で, 順番が式とは逆になっていることに注意してください. 条件つき確率の定義式(2.19)から, 上の確率は

$$(2.27) \quad p(j|x) = \frac{p(j,x)}{p(x)}$$

と表すことができます. ここで分子に式(2.20)の確率の連鎖則を用いれば,

$$(2.28) \quad p(j|x) = \frac{p(j,x)}{p(x)} = \frac{p(x|j)p(j)}{p(x)}$$

と分解することができます. すべての文字 x について $p(x|j)$ は式(2.23)のように計算することができます. 文字のユニグラム確率 $p(j)$, $p(x)$ は 2.1 節で計算したように文字の頻度から簡単に求まりますから, 式(2.28)は簡単に計算することができます.

たとえば, 式(2.23)より $p(a|j)=0.0410$ ですが, 式(2.1)より $p(j)=0.0011$, $p(i)=0.0662$ なので,

$$(2.29) \quad p(j|a) = \frac{p(a|j)p(j)}{p(a)} = \frac{0.0410 \times 0.0011}{0.0662} = 0.00068$$

と計算されます. 同様にしてほかの文字についても計算してみると,

```

r = {};
for x in chars: # c[x][y] = p(y|x)
  r[x] = {} # r[x][y] = p(x|y)
  for y in chars:
    if y in c[x]:
      r[x][y] = c[x][y] * p0[x] / p0[y]
    else:
      r[x][y] = 0
r['j']

```

⇒

のように頻繁に使われています. しかし英語では, j が単語の末尾に用いられることは稀です. それでは, j はどんな場合に使われているのでしょうか?


```
{'a': 0.0006825162097599817,
 'l': 0.0,
 'i': 0.0,
 'c': 0.00041701417848206843,
 'e': 0.0014737307493920865,
 's': 0.0,
  :
 'u': 0.02943722943722944,
 'r': 0.0,
 'w': 0.0,
 'o': 0.002087170042971148,
  :
 'z': 0.0}
```

となり、英語では、“uj” および “ej” の確率が比較的高いことがわかります。

式(2.29)をよく見ると、これは条件つき確率 $p(j|a)$ を、逆方向の条件つき確率 $p(a|j)$ を使って計算することができる、ということを意味しています。こうして条件つき確率を引っくり返すことができる式(2.29)、あるいはより一般に

$$(2.30) \quad p(X|Y) = \frac{p(Y|X)p(X)}{p(Y)} \quad (\text{ベイズの定理})$$

を、**ベイズの定理** (Bayes' Theorem) といいます。^{*15} ベイズの定理を使えば、条件つき確率 $p(X|Y)$ を逆の条件つき確率 $p(Y|X)$ で表すことができます。たとえば X が「原因」、 Y が「結果」のとき、結果 Y がわかったときの原因 X の確率 $p(X|Y)$ は、原因から結果が得られる確率 $p(Y|X)$ を使って計算できる、ということになります。

このベイズの定理も、本質的には確率の連鎖則の式(2.20)からすぐに得られますので、**これを暗記する必要はありません**。慣れるまでは、 $p(X, Y) = p(X|Y)p(Y)$ の両辺を $p(Y)$ ($\neq 0$) で割って

^{*15} この公式の特別な場合は英国の牧師 Thomas Bayes (c.1701–1761) によって発見され、死後に友人によって発表されたため、ベイズの定理と呼ばれています。ただし、一般化や実際の発展は、その後にフランスの数学者ラプラス (1749–1827) によって行われたものです[13]。とはいえ、ベイズの公式を積極的に利用するベイジアンにとって、ロンドン市内にあるベイズの墓石は今でも聖地となっています。

$$p(X|Y) = \frac{p(X, Y)}{p(Y)} = \frac{p(Y|X)p(X)}{p(Y)}$$

として、その場で導くといいでしょう。

ベイズの定理に慣れるため、もう少し練習してみましょう。インフルエンザ(以下インフル)が日本中で流行していた、ある年の冬、普段は感染しない筆者も、39度近い高熱がまったく下がらない状況になりました。急いで休日診療の医院にかかったところ、何と検査結果はインフル陰性でした。ただし、診察を担当したのは新人の医師で、鼻の粘膜の入り口しか取らない簡易的なものでしたので、この結果にはかなり疑いが残りました。この高熱は本当に、ただの風邪なのでしょうか。

真の状態を表す変数を X 、その時の診断を Y とおきましょう。0はインフル陰性、1はインフル陽性を意味します。上記の症状や周囲の感染状況から、自分の見立てでは、インフルである確率は8割はあ

ると思っていました。数式で書くと、 $p(X=1)=0.8$ 、 $p(X=0)=0.2$ ということになります^{*16}。一方で検査結果は陰性、つまり $Y=0$ だったわけです。新人医師が誤診する確率は、本当は追跡すれば客観的にわかりますが、ここでは図2.8のようだったとしましょう。つまり、単なる風邪ならば、医師の診断でインフルと誤診されることはないが(つまり $X=0$ のとき、 $p(Y=0|X=0)=1$ 、 $p(Y=1|X=0)=0$)、検査が甘い

ため、インフルの場合でも、3割程度は陰性と判断されてしまう(つまり $X=1$ のとき、 $p(Y=0|X=1)=0.3$ 、 $p(Y=1|X=1)=0.7$) と仮定します。

$$(2.31) \quad p(X|Y=0)$$

*16 ここでは説明のために事前確率を天下りに設定していますが、本書でこの後扱うように、実際に使う事前確率はできるだけ無情報にしたり、それ自体をデータから学習するなど、客観的に決められるものです。この場合も、同様の症状だった人の診察記録を多数集めれば、診断以前の事前確率 $p(X)$ を計算することができるでしょう。よって、ベイズ統計が「主観的な事前確率を使う統計」であるとする批判は、実際にはあまり正しくありません。

$X \setminus Y$	0	1
0	1	0
1	0.3	0.7

図 2.8: $p(Y|X)$ の表. X は真の状態を、 Y は診断を表します。

です. 式(2.19)の条件つき確率の定義から, この式(2.31)は次のように書き直すことができます.

$$(2.32) \quad p(X|Y=0) = \frac{p(X, Y=0)}{p(Y=0)} = \frac{p(X, Y=0)}{\sum_X p(X, Y=0)}$$

$$= \frac{p(Y=0|X)p(X)}{\sum_X p(Y=0|X)p(X)}$$

X は具体的には 0 か 1 ですから, 式(2.32)に上の値を代入すると,

$$(2.33) \quad \left\{ \begin{array}{l} p(X=0|Y=0) = \frac{p(Y=0|X=0)p(X=0)}{p(Y=0|X=0)p(X=0) + p(Y=0|X=1)p(X=1)} \\ \qquad \qquad \qquad = \frac{1 \times 0.2}{1 \times 0.2 + 0.3 \times 0.8} = \frac{0.2}{0.2 + 0.24} = \mathbf{0.455} \\ p(X=1|Y=0) = \frac{p(Y=0|X=1)p(X=1)}{p(Y=0|X=0)p(X=0) + p(Y=0|X=1)p(X=1)} \\ \qquad \qquad \qquad = \frac{0.3 \times 0.8}{1 \times 0.2 + 0.3 \times 0.8} = \frac{0.24}{0.2 + 0.24} = \mathbf{0.545} \end{array} \right.$$

となります. すなわち, 検査は陰性 ($Y=0$) でしたが, 本当はインフルである ($X=1$) 確率は半分以上の 55%もある, ということがわかります. 実際この数日後, あまりに熱が下がらないので別のベテランの先生の医院で再検査したところ, 明らかにインフルということがわかり, タミフルを処方されてすぐに熱は収まりました.

式(2.33)の計算をよく見ると, どちらも分母は同じで $(A+B)$ の形をしており, 確率はそれぞれ $A/(A+B)$, $B/(A+B)$ の形をしています. すなわち, 式(2.33)を求めるためには分子である A と B , つまり $p(X, Y=0) = p(Y=0|X)p(X)$ を $X=0$ と 1 の場合について求めればよく, それを和が 1 になるように正規化すれば確率が得られる, ということがわかります.

これは, ベイズの定理を使う場合すべてについて成り立ちます. なぜならば,

条件つき確率の定義から

$$(2.34) \quad p(X|Y) = \frac{p(X, Y)}{p(Y)}$$

ですが、ここで分母の $p(Y)$ は X には関係しない、和を 1 にするための定数だからです。よって、ベイズの定理は比例を表す記号 \propto を使って、

$$(2.35) \quad p(X|Y) \propto p(X, Y) \quad (\text{ベイズの定理 (同時確率版)})$$

とも表すことができます。または、右辺を確率の連鎖則 $p(X, Y) = p(Y|X)p(X)$ で分解すれば、

$$(2.36) \quad \underbrace{p(X|Y)}_{\text{事後確率}} \propto \underbrace{p(Y|X)}_{\text{尤度}} \underbrace{p(X)}_{\text{事前確率}} \quad (\text{ベイズの定理 (比例版)})$$

と書き直すことができます。実際の計算は、こちらで行えばよいでしょう。たとえば式(2.33)の計算は、実際には

$$(2.37)$$

$$\begin{aligned} p(X|Y=0) &\propto p(Y=0|X)p(X) \\ &= \begin{cases} p(Y=0|X=0)p(X=0) \\ p(Y=0|X=1)p(X=1) \end{cases} = \begin{cases} 1 \times 0.2 \\ 0.3 \times 0.8 \end{cases} = \begin{cases} 0.2 \\ 0.24 \end{cases} \end{aligned}$$

となり、これを和が 1 になるように正規化して

$$(2.38) \quad \begin{cases} p(X=0|Y=0) = \frac{0.2}{0.2+0.24} = 0.455 \\ p(X=1|Y=0) = \frac{0.24}{0.2+0.24} = 0.545 \end{cases}$$

と、簡単に求めることができます。

式(2.36)は文字通り、 Y が与えられたときの X の確率は、 X だけの確率 $p(X)$ に、 X から Y が観測される確率 $p(Y|X)$ をかけたものに比例することを表しています。 Y が与えられた後の確率なので、左辺の $p(X|Y)$ を**事後確率** (posterior

probability), これに対して右辺の $p(X)$ を**事前確率** (prior probability) ともいいます. $p(Y|X)$ は $X \rightarrow Y$ となる確率ですが, いま Y は与えられているため, これは X の関数とみることができ, X の**尤度関数** (likelihood function) といいます. 尤度とは, 観測値 Y について X がどんな値をとれば尤もらしいのか, を表しています. よってベイズの定理は, 言葉で書けば

$$(2.39) \quad (\text{事後確率}) \propto (\text{尤度}) \times (\text{事前確率})$$

の形で書くことができます. 式(2.35)や式(2.36)の変形はこの後で頻繁に使いますので, 覚えておいてください. ベイズの定理にさらに慣れるため, 本章の演習問題を解いてみるとよいでしょう (\rightarrow 演習(10)).

2.5 文字 n グラムモデル

2.5.1 文字列の確率的生成

こうして文字の条件つき確率がわかると, 文字列を**確率的に生成**することができます. そのためにまず, 文字列の確率について考えてみましょう.

たとえば文字列 s が `cat` のように 3 文字だったとき, 1 文字目, 2 文字目, 3 文字目を表す確率変数をそれぞれ X, Y, Z とおくと, s の確率は

$$(2.40) \quad p(s) = p(X, Y, Z)$$

です. たとえば $p(\text{cat}) = p(X=c, Y=a, Z=t)$ になります. このとき確率の連鎖則式(2.20)から, まず X, Y をまとめて 1 つにして考えれば, 式(2.40)は

$$(2.41) \quad \begin{aligned} p(s) &= p(X, Y, Z) \\ &= p(Z|X, Y)p(X, Y) \end{aligned}$$

と分解することができます. ここでさらに $p(X, Y)$ にも確率の連鎖則を用いれば,

$$(2.42) \quad \begin{aligned} p(s) &= p(X, Y, Z) \\ &= p(Z|X, Y)p(X, Y) \\ &= \underbrace{p(Z|X, Y)}_{(A)} \underbrace{p(Y|X)}_{(B)} p(X) \end{aligned}$$

と文字列の確率を積に分解することができます。この確率をどう計算するかは、何グラムモデルを考えるかによります。

ユニグラムの場合

ユニグラムを考える場合は、文字の確率はそれより前に出現した文字に依存しませんから、式(2.42)の最後の行の (A)(B) はそれぞれ

$$(A) \quad p(Z|X, Y) = p(Z)$$

$$(B) \quad p(Y|X) = p(Y)$$

となります。よって

$$(2.43) \quad p(s) = p(X, Y, Z) = p(X)p(Y)p(Z)$$

です。つまり、ユニグラムで文字列 s を生成するには、1文字目、2文字目、3文字目をそれぞれ単にユニグラム分布から生成すればよいわけです。

2.1節のように Python の辞書 p に文字 c の確率が $p[c]$ として保存されているとき、この中から確率に従ってランダムに1文字を選ぶ関数 `genchar` は、 $0 \leq r < 1$ の一様乱数を生成する関数 `rand()` を使って、

```
from numpy.random import rand
def genchar (p):
    s = 0
    r = rand()
    for c in p: # 文字 c を順番に調べる
        s += p[c]
        if (r < s):
            return c
    return c # ここには来ないはずですが、念のため
```

のように書くことができます。^{*17} よって、長さ N の文字列 s を生成するには、

```
s = ""
for n in range(N):
    s += genchar(p)
```

*17 詳しい方への注：この方法の計算量は、カテゴリ数 K について $O(K)$ です。もしサンプルする確率分布が固定されている場合には、Alias 法とよばれる方法で事前にテーブルを作っておけば、以後は $O(1)$ でサンプルすることができます[14]。さらに2020年になって、Fast Loaded Dice Roller (FLDR) とよばれる手法が提案され[15]、情報理論的限界に近い速度でサンプリングすることが可能になりました。C と Python による FLDR の実装は、<https://github.com/probcomp/fast-loaded-dice-roller> で公開されています。

とすれば行うことができます。

こうして『不思議の国のアリス』のユニグラムから生成した長さ 40 の文字列は、下のようになりました。ここからは、句読点や記号も含んだ元のテキスト `alice.full.txt`*18 を使って確率を計算しています。

```
% unigram.py alice.full.txt alice.1gram
% unigram-gen.py alice.1gram 5 40
⇒ !reyrnh r'l ls ses aso oeheh s 'nreai
   rnnhntmls ga tfi ht-ltreat wag 'areia
   lhne infeeulctrmwohno u ! oe ao n iwssvd
   fv rath tfegyanc!ytomst.d tcoa ni,us,oe'
   t,eehiiapscoitna -gpeksag,,tnlsoemw si'
```

上記はサポートサイトの `unigram.py` で文字のユニグラム確率を計算してファイル `alice.1gram` に保存した後、`unigram-gen.py` で生成しています。文字を完全に独立にとらえていますので、およそ英語のようには見えませんが、`e` や `t` といった文字の頻度が高くなっており、文字ユニグラム確率を反映しているようです。これは、完全にすべての文字を等確率に選ぶ場合 (**0 グラム**といわれます) と比べればはっきりするでしょう。0 グラムから生成すると、同じ場所にある `zerogram.py` を使って以下ようになります。

```
% zerogram.py alice.full.txt 5 40
⇒ ul!'_jx' 'fi!k_(x_l'whs.;(u'_ bm'ucj!'b?.
   .cxba !r'ciikj[z;:obu;'m',ig!cngle !n''o'
   ?]n?g.nrpw)kw(fucpxgouj,fv.ibyvuf-rbx;m_
   '[:?gth:z(gthhx]o-qq-'lh!f[jo''mujo!si[.!
   u(its,y,_hqiqnaxbic. i qepzy[-nt.' 'm(sm,'
```

こうして、一見ランダムに見えるユニグラムも、0 グラムと比べればずっと自然言語の統計を反映している、ということがわかりました。

バイグラムの場合

バイグラムの場合は連続した 2 文字を考慮しますので、式 (2.42) 式の (A)(B) の確率は、以下のようになります。

$$(A) \quad p(Z|X,Y) = p(Z|Y)$$

*18 サポートサイトの同じフォルダに置いてあります。

$$(B) \quad p(Y|X) = p(Y|X)$$

つまり、3番目の文字 Z は Y とのバイグラム YZ にのみ関係しており、 X には依存しないことになります。よって s の確率は、

$$(2.44) \quad p(s) = p(X, Y, Z) = p(Z|Y)p(Y|X)p(X)$$

で与えられます。したがって、式(2.44)式から s を生成するには、

- (1) まずユニグラム $p(X)$ から1文字目 X を生成する
- (2) 生成された X を使って、 $p(Y|X)$ から2文字目 Y を生成する
- (3) 生成された Y を使って、 $p(Z|Y)$ から3文字目 Z を生成する

とすればいいことになります。

これでもよいのですが、ただし、上のステップ(1)には注意が必要です。多くの教科書にはマルコフモデル(ここではバイグラム)から生成するには上記のようにすると書かれていますが、このままナイーブに行うと、ユニグラム確率すなわち、「文字の一般的な確率」 $p(X)$ から最初の文字 X を生成することになってしまいます。たとえば、日本語ですべての文字列が引用符“「”から始まっているとしても、確率の高い“の”や“が”が最初の文字として生成される可能性が高くなってしまふわけです。これは明らかに不合理ですね。

そこで自然言語処理では一般に、文字列や単語列の最初と最後に、文頭と文末を表す見えない特殊な文字(あるいは単語)“ \sim ”と“ $\$$ ”があると考えます。ここで“ \sim ”や“ $\$$ ”は説明のための便宜的な表記で、実際にはテキストに現れない文字(たとえば“ $\backslash x1c$ ”)や単語(空文字列“”や“_BOS_”)を用います。^{*19} 実装例は、サポートサイトの `bigram.py` や `trigram.py` を参照してください。これらは、BOS (beginning of sentence) および EOS (end of sentence) といわれることもあります。^{*20}

よってこの場合、文字列 $s = \text{"alice"}$ は

*19 文字“ $\backslash x1c$ ”はASCIIコード表のFS (field separator) ですが、他の文字でもかまいません。なお、16進数で $\backslash x1c=8$ 進数で $\backslash 034$ の文字(Ctrl-\)は、テキスト処理言語である `awk` や `perl` では、擬似的に多次元配列を実現するための標準の区切り文字として内部で使われています。

*20 実は理論的には、BOS と EOS は同じでも問題ありません。文字や単語は BOS から出るだけ、EOS へは入るだけだからです。高速な実装で文字や単語を整数の ID で管理している場合は、EOS には 0 や -1 (二進数の補数表現で $1111\dots 1$) を用いるとよいでしょう。

"^alice\$"

として扱います。このとき文頭文字 \wedge はすでに与えられているとしますので、式(2.44)式の $p(X)$ は 1 です。したがって、 s の確率は

$$(2.45) \quad p(s) = p(a|\wedge)p(l|a)p(i|l)p(c|i)p(e|c)p(\$|e)$$

と、すべてバイグラムで書けることになります。最後に文末文字 "\$" を導入して文末との接続確率 $p(\$|e)$ を考えているため、「文字列の最後が "e" で終わりやすい」といった言語的性質も考慮できるのが特徴です。この点は、4.4 節で隠れマルコフモデル (HMM) を扱う際にも再び議論します。

上記のようにして `alice.txt` の文字バイグラムから生成した文字列を次に示しました (`bigram.py`, `bigram-gen.py`)。この場合、文末文字 "\$" が生成されればそこで文字列が終わるという都合ですので、生成する長さを指定する必要はありません。

```
% bigram.py alice.full.txt alice.2gram
% bigram-gen.py alice.2gram 5
⇒ no as pe ilille sarugl ster, cake erue!
'stick, ce w illite woprud p, jenthing bs washiofopor ld ur
wheshig,'sey bes elor.
[1 that ('y the tcepane uth wean to, nthiskerus amerowhexi-
gokem serghipp as bo wsime'thertord alork athin t itwas wh
se f en.
wrsal d owhoughen trnndsprownore an bsoupandolsingry, ucane,
' bioull th be wide s ind oxifunde.
```

0 グラムやユニグラムと比べると、だいぶ英語に近くなってきました。ただ、バイグラムは「隣り同志の文字の確率」しか考えていないので、もう少し賢くする必要がありそうです。

メモ：文字列の確率と EOS

文字列の終わりを表す特殊文字 EOS を考えることは、文字列がどんな風に終わりやすいかという言語の性質を表せるだけでなく、実は、確率モデルとしても不可欠な要素です。図 2.9 に示したように、EOS を考えると、“”(空文字列), “a”, “aa”, “ab”, … といったすべての文字列は先頭から順に文字を選択し、EOS が出たときそこで止まった結果として表すことができます。あらゆる文字列はこうした選択に対応していますから、それらの確率の総和は必ず 1 になります。いっぽう、EOS がないと “aaaaa…” といった、いくらでも長い文字列にも一定の確率を割り当てることができ、それらが無限にありますから、確率の総和は容易に 1 を超えてしまいます。これは、文字列の確率モデルとしては不適切です。^{*21}

なお、再帰的な選択によって文字列を生成できることから、次の文字の選択肢を可視化して、その幅を確率に比例させれば、 $[0, 1)$ の範囲の実数を次々と指示することで、効率的に文字を入力することができます^{*22}。図 2.10 に示したケンブリッジ大学の MacKay らによる Dasher [16]^{*23} は、このことを利用した入力システムです。もし病気や障害で指が動かない場合でも、まぶたの上下や、呼吸による腹の上下といった 1 次元のわずかな信号さえあれば、言葉を発することができます。

これらは文字列に限らず、3 章の単語列の場合でもすべて同様に成り立ちます。^{*24}

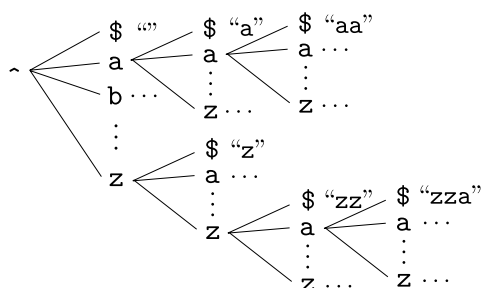


図 2.9: 文字列全体の空間を表す樹形図。
~で BOS, \$ で EOS を表しています。

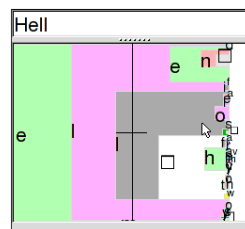


図 2.10: 入力システム Dasher で “Hello” を入力している様子。1 次元の選択を次々に行うだけで、文字列を入力することができます。

*21 深層学習による現代の言語モデルがこの条件を満たしているかについては、無限長の文字列

2.5.2 ゼロ頻度問題

バイグラムでは“er”や“al”のような隣りあう2文字の連続を考えてきましたが、もっとよい言語のモデルとするには、“ing”や“has”のように3文字の連続を考えればよさそうです。こうした3文字の連続を、3グラムまたは**トライグラム** (trigram) といいます。一般に、0グラム、1グラム、2グラム、3グラム、…をまとめて **n グラム**とよび、 n グラムの条件つき確率に基づいて文を生成する確率モデルを、 **n グラム言語モデル**といいます。 n グラム言語モデルは、 n 文字の連続をモデル化するものになっています。

文字トライグラム言語モデルでは、バイグラムのように直前の1文字を見て

$$(2.46) \quad p(Y=i | X=e)$$

のように条件つき確率を計算する代わりに、

$$(2.47) \quad p(Z=i | X=a, Y=1)$$

のように、直前の2文字を見て条件つき確率を計算します。^{*24} 文頭を表す特殊文字“^”は、この場合、最初の文字のために2個必要になり、文字列“alice”は“^alice\$”とあってaから確率を計算します。式(2.25)と同様に頻度 $n()$ を用いて表せば、トライグラム確率は

$$(2.48) \quad p(z|x, y) = \frac{n(xyz)}{n(xy)}$$

と書くことができます。ここで x は2つ前、 y は1つ前の文字で、 z は予測する文字です。

ただし、この場合はユニグラムやバイグラムと異なり、大きな問題が発生します。式(2.46)のバイグラムの確率は X と Y の組み合わせの数、つまり文字種の

に対する確率を考える必要があることから、測度論を用いた最近の研究[17]によって解析されています。

*22 これは情報理論では、**算術符号**とよばれる符号化を行っていることに相当します。

*23 <http://www.inference.org.uk/dasher/> に日本語版を含む実装や情報があります。

*24 同様に4グラムや5グラムも考えることができますが、取り扱いが複雑になるため、詳しくは3章を参照してください。筆者はノンパラメトリックベイズ法に基づき、文脈によって異なる n を学習することで、 n の指定を不要にした ∞ グラム言語モデルを提案しています[18][19]。

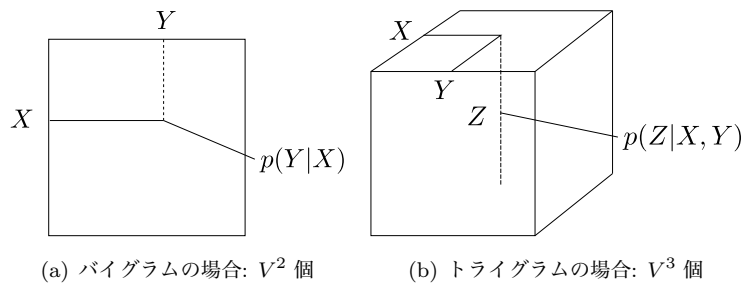


図 2.11: n グラムの予測確率の組み合わせ. 実線は条件となる変数を, 破線は予測する変数を表しています. V は語彙の数 (文字 n グラムなら文字の種類数) です.

2 乗だけ存在します. 英語の場合は文字種は記号を入れても最大でも 256 種類程度ですから^{*25}, バイグラムでは図 2.11(a) のように, 最大 $256^2 = 65,536$ 個の確率がテキストの文字の頻度から計算できればよいわけです.

しかし, トライグラムの場合は図 2.11(b) のように, 式 (2.47) のトライグラムの確率は X, Y, Z の組み合わせで, 最大で $256^3 = 16,777,216$ 個存在します. 『不思議の国のアリス』は 132,656 文字でしたから, これは可能な組み合わせの 0.8% にすぎず, もし図 2.11(b) のセルに頻度を 1 ずつ均一に置いたとしても, すべての組み合わせに頻度を埋めることはできません. 実際には頻度は e や a などに偏っていますから, この傾向はさらに極端で, 図 2.11(b) ではほとんどのセルの頻度は 0 だということになります. これを, **ゼロ頻度問題**といいます.

ゼロ頻度問題のため, トライグラム言語モデルはそのままでは大きな問題が生じます. たとえばトライグラム "onk" は alice.txt には一度も現れませんので, 式 (2.48) より

$$(2.49) \quad p(\mathbf{k}|\mathbf{on}) = \frac{n(\mathbf{onk})}{n(\mathbf{on})} = \frac{0}{n(\mathbf{on})} = 0$$

です. そうすると, たとえば文字列 "monk" の確率は

*25 ASCII 文字は 8 ビット目が 0 なので 128 種類で収まりますが, アクセント記号が存在するヨーロッパ系の言語 (Latin-1 文字セット) では 8 ビット目が 1 の場合があるため, 最大は 256 文字になります.

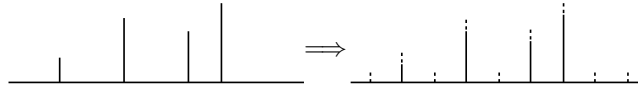


図 2.12: 確率分布の平滑化. 頻度にすべて小さな値 α を足して正規化することで, 離散的な確率分布がより「滑らか」になっています.

$$(2.50) \quad p(\text{"monk"}) = p(m|\hat{\cdot}) \cdot p(o|\hat{m}) \cdot p(n|mo) \cdot \underbrace{p(k|on)}_{=0} \cdot p(\$|nk) = 0$$

になってしまい, "monk" や "monkey" の確率はすべて 0 , つまり英語には絶対に現れない文字列ということになってしまいます. これは明らかにおかしいですね.

加算平滑化

そこで, 最も簡単な解決法として, トライグラムのすべての頻度に小さな値 α を加えることを考えてみましょう. これにより図 2.12 に示したように, 確率分布はより「滑らか」になるため, こうした操作を確率分布の**平滑化 (スムージング, smoothing)** といいます.

こうすると, トライグラムの予測確率は式 (2.25) に代えて, 次の式で求められます.

$$(2.51) \quad p(z|x, y) = \frac{n(xyz) + \alpha}{\sum_c (n(xyz) + \alpha)} = \frac{n(xyz) + \alpha}{n(xy) + V\alpha}$$

ここで V は文字種の数で, 可能な最大値は英語なら 256, 日本語なら 8500 程度ですが, 実際には学習するテキストに現れた異なり数を用いればよいでしょう *26.

平滑化について詳しくは 3 章の単語 n グラムモデルで議論しますが, 単純に頻度に小さな値を加えるこの方法を, **加算平滑化** といいます. 加算平滑化は最も簡単な平滑化法ですが, トライグラムに直接用いる場合は, データ量によりませんが, α はかなり小さな値にする必要があります.*27 ただし, $\alpha=0$ にすると平滑化

*26 3 章の単語の n グラムモデルの場合, 可能な単語の数は無限にありますから, ナイーブに考えると V を無限大に取らなければならなくなってしまいます. よって, 学習データに出現した語彙の総数を用いるのが, 現実的な方法といえます. 筆者による NPYLM [20] では, 階層ベイズモデルを用いることで, 理論的に無限個の語彙を扱うことができます.

*27 トライグラムの場合は, ゼロ頻度問題から, 式 (2.51) でほとんどの文字 z について $n(xyz)=$

を行わないため、ゼロ頻度問題が生じてしまいます。一方で α を大きくすると、式(2.51)の $p(z|x,y)$ は、 z が何であっても一様分布 $1/V$ に近づいてしまいます。以下の例では、英語では $\alpha = 1e-4 = 0.0001$ 、日本語では $\alpha = 1e-5 = 0.00001$ としました。このとき、式(2.49)式の $p(k|on)$ は式(2.51)より、alice.full.txt では $n(on) = 1054$ 、現れる文字の総数 $V = 43$ なので

$$(2.52) \quad p(k|on) = \frac{0 + 0.0001}{1054 + 43 \times 0.0001} = 0.00000009$$

になり、確かに 0 でない確率が割り当てられていることがわかります。

文字トライグラムからの生成

alice.txt に対して文字トライグラム言語モデルをサポートサイトの trigram.py で計算して保存し、trigram-gen.py で生成すると、下のようになりました。バイグラムに比べて、大幅によい言語モデルとなっていることがわかります。^{*28}

```
% trigram.py alice.full.txt alice.3gram 1e-4
% trigram-gen.py alice.3gram 5
```

```
⇒ 'who an ard to my frot an thisher, awlind the as ey pled the
low fousee me a onereepkeake died of yound of sly and ance,
beithe pearniene was no paid they hargense, be fould thereep
on'ting i cat shrough al bes mocked es.'
'of the anagaid rattlenly crioll, chink top on.
so the hat whe in al's ithe could but of musibleake onse, a
that 'youg, theress; 'it bithe withe queence lice, yought
abloo ding tion, all mock thatinsed wou juse words tur a
plice fing suce rasheadveraclar!
```

0 です。 z の確率を予測する際、これが V 通りのほとんどを占めますから、頻度 0 の文字についての確率の総和はほぼ $V \cdot \frac{0 + \alpha}{n(xy) + V\alpha} = \frac{V\alpha}{n(xy) + V\alpha}$ になります。これを書き直すと

$1 / \left(\frac{n(xy)}{V\alpha} + 1 \right)$ になり、 $V = 10000$ 、 $n(xy) = 10$ のとき、 $\alpha = 0.01$ としても頻度が 0 の文字の確率の総和は $1 / \left(\frac{10}{10000 \cdot 0.01} + 1 \right) = 1 / \left(\frac{10}{100} + 1 \right) \simeq 0.91$ となり、頻度 0 の文字にほとんどの確率が割り当てられてしまうことになってしまいます。

^{*28} このため、音声認識の分野では 1980 年代から長い間、単語トライグラム (3 章) が標準的な言語モデルとして使われていました。現在は、必要な場合は LSTM や Transformer などのより高度な深層学習による言語モデルで文の確率を計算します。

夾竹桃の方に向ってしまった。またどたどたどたど切トっ、ランペシヤの中心地とて遠い所も見せて来る蓄音機の浪花節。わびしげしげと眺めたきりして来ていた。慌てている様子を嘯む習慣もパラオ本島オギワライという島の方が一面に糜爛した良いので、筋は良く知って来た。海岸から二度と、うす汚い、この言うは何をあげた時から真白い裏を見るとは作るが。それら南方へ歩いた渚に踏出した翡翠色に向って、白い雨が頻りに来たあとではない。こんなと思う。人も住まないが、素晴らしい。

図 2.13: 中島敦『環礁』の文字トライグラムからの確率的生成結果。

「標本室にありましたら、またまえ。だけ青く茂ったり暗くなりましたけれどもはつきりに照らしい楽器の灯を、窓を見えていますけれども滑って、ジョバンニたちでいるものをひろつ務のから僕決して戻ろう。大きなとうに平らな。」カムパネルラと二熟、カムパネルラを見ました。ジョバンニが、砂に三つの街燈の方へ急ぐのです。ジョバンニが左手の崖が川の水をわたしばらく、飛び出しても、顔を出しました。銀河のお宮のけよってでした。二人も胸いって微かにうつくしくて、ちよう。ぼくが、続いても押し葉にすぎとお会いました。
「ほんといいとジョバンニのときすぎうつと天の川の水にあんな水は、すつかりにして校庭の隅の桜の木のあかりを川へ帰らずの鳥捕りがとまりました。

図 2.14: 宮沢賢治『銀河鉄道の夜』の文字トライグラムからの確率的生成結果。文字の言語モデルなので、「単語」という概念はありません。

現在は文字はユニコードによって言語を問わず扱えるようになっていますので、`trigram.py`, `trigram-gen.py` は英語以外のテキストに対しても適用することができます。^{*29} 図 2.13 に、中島敦『環礁』(229 行, 36419 文字) から文字トライグラムで生成した文を、図 2.14 に宮沢賢治『銀河鉄道の夜』(459 行, 38292 文字) から同様に生成した文を示しました。

『銀河鉄道の夜』は若干ひらがなが多く、文字トライグラムではモデル化し切れていない部分がありますが、漢籍を背景にした中島敦の文体は長距離の依存性が少なく、文体の特徴を比較的好くとらえていることがわかります。

単語のランダム生成

また、単語は文字からなっていますから、`"alice"`, `"fortunately"` などの単

^{*29} 日本語や中国語などでは「単語」をどう定義するかという問題がありますが、文字の言語モデルは単語の定義に関係なく使うことができます。

語をそれぞれ「文」とみなせば, `alice.txt` にある 2748 語の語彙の文字列から, 新しい「単語」を文字 n グラム言語モデルを使ってランダムに生成することができます. 下に, その様子を示しました. ここでは, 1 行に語彙の単語が 1 つずつ並んだファイル^{*30}を `alice.lex` としています.

```
% trigram.py alice.lex alice-lex.3gram 1e-4
% trigram-gen.py alice-lex.3gram 20
⇒ que
   twer
   hadder
   stralkjuse
   late
   shatted
   vuld
   cut
   comished
   scretch
   flobs
   gar
   sely
   persed
   spoisoleds
   besersty
   ding
   arighen
   nage
   up
```

文字トライグラムを使っただけで, かなり英語らしい (が実際には存在しない) 単語が生成できていることがわかります.

*30 これは色々な方法で作成できますが, 最も簡単には, Mac や Linux, WSL には標準で含まれているテキスト処理言語 `awk` を使って, `% awk '{for(i=1;i<=NF;i++)freq[$i]++;}END{for(w in freq) print w}' alice.txt > alice.lex` とすると作ることができます.

メモ： n グラムモデルとマルコフモデル

44 ページで導入した n グラムモデルは、文字の発生確率が直前の $(n-1)$ 文字に依存する確率モデルです。一般に、事象の発生確率が直前の事象だけに依存するとき、確率論や情報理論では**マルコフ性**があるといい、**マルコフモデル**とよばれて議論されます。特に、直前の n 個の事象に依存する場合を、 n 次のマルコフモデルといいます。すなわち、 **n グラムモデルとは $(n-1)$ 次のマルコフモデルのこと**です。自然言語処理では、2.2 節で `alice.txt` から文字の 2 グラムを取り出したように、 n 個の文字の具体的な繋がりに興味があることも多いため、 n グラムモデルという言い方をするのが普通です。マルコフモデルの名前は、ロシアの数学者 Andrej Markov (1856–1922) が、1913 年の論文[21]でまさに言語の問題、すなわちプーシキンの小説『オネーギン』の最初の 20,000 文字を数えて、母音と子音が連続する確率を計算したことによります^{*31}。第一次世界大戦より前のこの時代、まだコンピュータは登場していないことに注意してください。

2.6 統計モデルの学習と評価

2.6.1 学習データとテストデータ

前の 2.5.2 節で文字 n グラムの平滑化に用いた $\alpha=0.0001$ (英語)、 $\alpha=0.00001$ (日本語) は、実は筆者が出力の様子を見て、人手で決めた^{*32} パラメータでした。 $\alpha=0$ とすると、2.5.1 節でみたようにデータに偶然出現しなかった n グラムの確率がすべて 0 になってしまい、図 2.13 や図 2.14 といったテキストの確率は、1 個でもそうした n グラムが含まれていれば 0 になってしまいます。逆に $\alpha \rightarrow \infty$ と大きくすると、式(2.51)から、 n グラムの確率はすべて $1/V$ の一様な確率になってしまい、やはりテキスト全体の確率は低くなってしまいます。よって、 α を変えれば図 2.15 のように、どこかにテキストの確率を最大化する $\alpha = \alpha^*$

^{*31} この研究は 20000 文字の統計を単に数えるのではなく、100 文字の連続 \times 200 個に分け、それらの間での統計量のばらつきを計算しており、現代のわれわれにも大変参考になるものです。

^{*32} もし人手で決める場合も、46 ページの脚注に示したように、数学的にある程度根拠のある値とすることは重要です。

があると考えられます。それでは、最適な α^* を客観的に決めるにはどうすればよいのでしょうか。

いま、我々の手元には `alice.txt` のようなテキストがありますから、 α^* を決めるには、このうち一部のテキストを検証用に残しておき、残りのテキストで n グラムモデルを学習して、残しておいた検証テキストの確率が最大になるような α を求めればよいでしょう。

たとえば最も簡単な場合として、文字ユニグラムモデルで、テキストが短い現代俳句^{*33}

ひいらぎをかをらせていつまでもいま (野口る理)

だったとしましょう。この文字列をランダムにシャッフルして

いかをらせもでいひつまら ぎいをまて
 D D'

とし、前半の 12 文字 D から文字の確率を計算すると、文字の頻度は $n(\text{ぎ})=0$, $n(\text{い})=2$, $n(\text{を})=1$, $n(\text{ま})=1$, $n(\text{て})=0$ ですから、式 (2.51) でユニグラムの場合を考えれば、ひらがなの総数は約 87 個 ($V=87$) ですから^{*34},

$$(2.53) \quad \begin{cases} p(\text{ぎ}) = p(\text{て}) = \frac{0+\alpha}{12+87\alpha}, \\ p(\text{を}) = p(\text{ま}) = \frac{1+\alpha}{12+87\alpha}, \\ p(\text{い}) = \frac{2+\alpha}{12+87\alpha} \end{cases}$$

となります。よって後半の文字列 D' の確率は、

$$(2.54) \quad p(D'|D, \alpha) = p(\text{ぎ})p(\text{い})p(\text{を})p(\text{ま})p(\text{て})$$

^{*33} 元の句では「終」だけが漢字ですが、ここでは説明のため平仮名にしています。余談ですが、「終をかをらせていつまでもいま」のこの句は、「いつまでも」の永遠性と「いま」の刹那性が「い」という同じ音を通じて意味的に共鳴し (Jakobson の詩学 [22])、それが「かをらせて」の微かに古典的な香りを通じて「終」に象徴的に表されている、素晴らしい現代俳句の一つだと思います。

^{*34} Unicode の Hiragana ブロック <https://unicode.org/charts/PDF/U3040.pdf> の表によります。

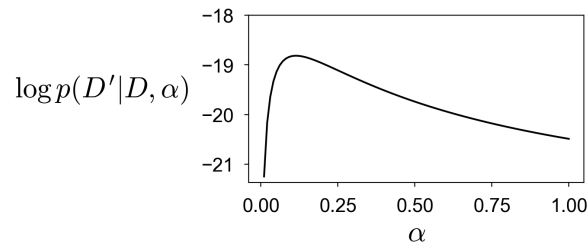


図 2.15: 俳句の仮想例での評価データ D' の予測確率の対数 $\log p(D'|D, \alpha)$ とパラメータ α . この場合, $\alpha=0.114$ で D' の確率が最大になっています.

$$\begin{aligned}
 &= \left(\frac{\alpha}{12+87\alpha} \right)^2 \times \left(\frac{1+\alpha}{12+87\alpha} \right)^2 \times \left(\frac{2+\alpha}{12+87\alpha} \right) \\
 &= \frac{\alpha^2(\alpha+1)^2(\alpha+2)}{(12+87\alpha)^5}
 \end{aligned}$$

になります. この確率の対数を, α についてプロットした図を図 2.15 に示しました.

式(2.54)は α の関数ですから, 勾配法を使ったり, 微分して 0 とおくことで極大値を求めることができます. 数式処理ソフトに入れてみたところ^{*35}, $\alpha^* = 1/96(\sqrt{5761}-65) \approx 0.114$ が最適な α となりました. ここでは最も簡単な文字のユニグラムモデルを用いましたが, 式(2.51)のトライグラム確率のような, より複雑なモデルでも数値計算になるものの, 基本的な方法は同様です. 上のように, データのうち検証のために残しておいたデータを**検証データ**あるいは**開発データ**, 統計モデルの計算に用いる, それ以外のデータを**学習データ**あるいは**訓練データ**といいます[23]. 検証データとしては, 一般にデータの 1 割から 2 割程度をランダムに抽出して使い, 残りの 8 割から 9 割を学習データとすることがよく行われます. 検証データが 2 割以下なのは, それより多くすると, それだけ使える学習データが減ってしまうことになるからです.

上の例ではデータの学習データと検証データへの分割を固定してしまいましたが, これはもちろん, 本当は望ましくありません. たまたま検証データに選ば

^{*35} *Mathematica* を使えば, `f[x_] := 2 Log[x]+2 Log[x+1]+Log[x+2]-5 Log[12+87 x]; Solve[D[f[x], x]==0]` と入力すれば求められます.

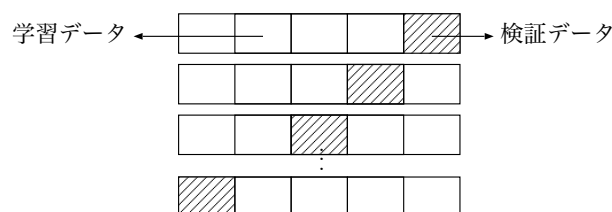


図 2.16: K 分割クロスバリデーションの様子. データを K 個 (ここでは $K=5$) に等分し, そのうち 1 つを検証データ, 残りの $K-1$ 個を学習データとする計算を K 通り行い, 結果を平均して最終的な答えとします.

れたものによって, 値が変わってしまうからです. より正確には, 上のような分割をランダムに繰り返し, 得られた α^* の平均を最終的な α^* とするのがよいでしょう. ただし, これは統計的には正しいのですが, 分割を何回繰り返せばよいかについての指針がありません. そこで, より計画的な方法として, **クロスバリデーション (交差検証)** という方法が知られています.

クロスバリデーション クロスバリデーションの中で最も一般的な K 分割交差検証 (K -fold cross validation) では, 検証データを毎回ランダムに選ぶかわりに, まずデータを**ランダムに**, つまりシャッフルしてから K 等分します. 図 2.16 に示したように, このうち 1 つを検証データ, $K-1$ 個を学習データとする計算を K 通り行い, 結果を平均して答えとします. 上でふれたように検証データはデータ全体の 1 割から 2 割を使うのが普通ですから, K は一般に, 5 から 10 程度の値とするのがよいでしょう. K をあまり大きくすると, クロスバリデーションに必要な K 通りの計算量が非常に大きくなってしまいます.

いずれの場合でも重要なのは, **検証データをランダムに選ぶこと**です. たとえば, テキストが新聞記事や SNS への投稿などで時間順に並んでいる場合, 新語がある時期から初めて現れているならば, それより古い学習データからその新語を予測することは不可能です. また逆に, 検証データが学習データのすぐ後の時期になっていると, その時期に流行したテキストの確率だけが高くなればよいことになり, これも不公平になってしまいます. 統計的には, 学習データと検証データが「同じ分布からサンプリングされた」といえる状況にしておく必要

があります。

これには, Python であれば `numpy.random.permutation()` で N 個のデータについて $1, \dots, N$ をランダムに並び換えた順番をデータ処理の際に生成してもよいですし, テキストの行をランダムにシャッフルする `shuf` のようなコマンド^{*36} を使えば,

```
% shuf alice.full.txt > alice.shuffled.txt
```

として, 最初からデータの行をランダムに入れ替えたテキストを作ることができます. これは, 特に上記のクロスバリデーションなどの際には便利でしょう.

こうして学習データとは別に検証データを準備すると, 図 2.15 のように, 検証データの確率が大きくなるパラメータを求めることが可能になります.

それでは, 実際に試してみましょう. サポートサイトの `kfold-alpha.py` を上の `alice.shuffled.txt` に対して

```
% kfold-alpha.py 5 alice.shuffled.txt 1 0.5 0.1 0.01 0.001
```

のように実行すると, テキストを K 等分 (ここでは 5 等分) したクロスバリデーションを行い, 図 2.16 で検証データになったテキスト全体の確率の対数を文字バイグラムで計算して出力します. 数学的には, この後で説明するように, テキストの確率を文の確率の積だとするとき, テキストを D_1, D_2, \dots, D_5 に 5 等分したとき,

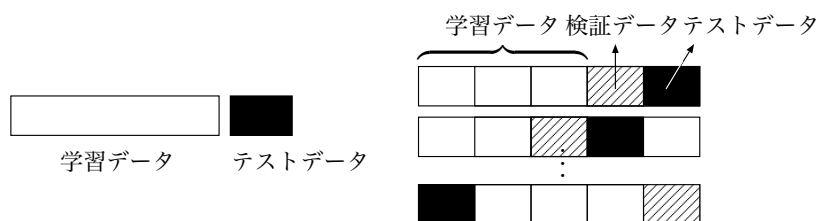
(2.55)

$$\log p(D_1 | D_2, D_3, D_4, D_5) + \log p(D_2 | D_1, D_3, D_4, D_5) + \dots + \log p(D_5 | D_1, D_2, D_3, D_4)$$

を, 平滑化係数 $\alpha = 1, 0.5, 0.1, 0.01, 0.001$ のそれぞれの場合で計算します. 結果は, 次のようになりました. () の中には, すぐ後で説明する 1 文字あたりのパープレキシティを表しています.

```
alpha = 1.0000 : likelihood = -332599.66 (PPL = 10.597)
alpha = 0.5000 : likelihood = -332101.93 (PPL = 10.560)
alpha = 0.1000 : likelihood = -331775.59 (PPL = 10.535)
alpha = 0.0100 : likelihood = -331866.05 (PPL = 10.542)
alpha = 0.0010 : likelihood = -332058.43 (PPL = 10.557)
```

^{*36} Linux では標準で含まれているようです. Mac では `% brew install coreutils` として, Homebrew で `coreutils` をインストールすると使えるようになります.



(a) 学習データと分けられたテストデータ (b) K 分割交差検証でのテストデータ

図 2.17: データの学習データとテストデータへの分割.

この結果から, `alice.txt` に対する文字バイグラム言語モデルの平滑化係数としては, $\alpha=0.1$ 程度とするのがよいということがわかります. クロスバリデーションでは値を直接最適化することはできませんので, もし, もっと細かく求めたい場合は, 候補を $\alpha=0.01, 0.05, 0.2$ などとしてクロスバリデーションをもう一度実行する必要があることに注意してください^{*37}. また, パラメータが複数ある場合は, その組み合わせの数だけ計算を実行する必要があります.

なお, 式(2.55)のように検証データの確率を計算することは, 学習データに基づいて検証データを**予測**していることになりますが, これは必ずしも予測することが目的というわけではありません. たとえば人文系で, 手元のデータの解析が目的で新しいデータを解析する必要はないとしても, モデルが学習データを丸覚えした一種のトートロジー (同語反復) になっていることを避け, 統計的に正しくデータを記述していることを保証したり, α のようなパラメータを客観的に決めるためには, 検証データを分けて予測確率を計算することが必要になります.

学習データ・テストデータ・検証データ

検証データを学習データと分けておくことで, 平滑化係数 α のような統計モデルのパラメータを大まかに推定することがわかりましたが, これが真にモデルの性能を表しているとは限りません. というのは, 検証データという「正解」を見て α の値を決めているので, 厳密に言うと, これは一種のチートになってい

^{*37} ベイズ最適化[24]とよばれる方法を使うと, 適切な候補を見つけて計算を次々と行うことで, こうした探索を自動化することができます.

るからです。たとえば、「統計モデル」として検証データを丸覚えして確率1で次の単語を予測してしまえば、これは性能が最大になってしまいます。

よって、モデルを真に評価するためには、学習時には見なかった新しいデータでの確率を計算する必要があります。学習データと区別されたこのデータを、**テストデータ**といいます。よって、**学習データ-テストデータを区別する**ことが最も本質的で、検証データはパラメータ推定のために、学習データの中を分けて人工的に作り出したものだといいてもよいでしょう。

テストデータも、検証データと同様にランダムに選ぶ必要があります。クロスバリデーションを使わない場合は、図 2.17(a) のようにたとえばデータのうちランダムな1割をテストデータ、残りの9割を学習データとするなどすればよいでしょう。これまでに見たように、9割の学習データの中でクロスバリデーションを行ってパラメータを決めることも可能です。徹底的に行うには、図 2.17(b) のように、 K 分割交差検証の枠組みを使って K 個のうち1個をテストデータ、1個を検証データ、残りの $(K-2)$ 個を学習データとする組み合わせを K 回行うことも考えられます。

いずれの場合も、モデルを学習する際には**テストデータは見ない**ことがもっとも重要です。テストデータを先に見てしまえば、いくらでもチートが可能になってしまうからです。逆に、テストデータさえ見なければ学習データの中では何をしてもよく、 K 分割交差検証は、学習データの中で、できるだけテストデータの予測に近い状況を作り出すための一つの方法といえます。^{*38}

2.6.2 予測確率とパープレキシティ

テキストの予測確率は、学習した統計モデルから先ほどのように計算することができます。テキスト D が文 (本章では文字列) の集合

$$(2.56) \quad D = \{s_1, s_2, \dots, s_N\}$$

からなっているとき、各文が独立であると仮定すると、テキスト D 全体の確率は

^{*38} 学習に使われなかったテストデータの中には、見たことのない文字や単語が含まれている可能性も高く、自然言語処理では、その場合にどうするかを常に考慮する必要があります。

$$(2.57) \quad p(D) = \prod_{n=1}^N p(s_n)$$

です。文字の系列からなる文 $s_n = c_1 c_2 \cdots c_T$ の確率 $p(s_n)$ は、式(2.45)や式(2.50)で示したように

$$(2.58) \quad p(s_n) = \prod_{t=1}^T p(c_t | c_1, \dots, c_{t-1})$$

として計算できます。ここで c_0 および c_T は2.5.1節で登場した、文頭および文末を表す特殊文字 BOS および EOS だとします。

この確率を、たとえば異なる α を用いた場合で比べればよいのですが、式(2.57)から計算されるテキストの確率は、一般に非常に小さな値になることに注意してください。たとえば式(2.50)のような n グラム確率が非常に大きく、それぞれ $1/10 = 0.1$ だったとしても、20文字からなる文の確率は

$$(2.59) \quad p(s) = \left(\frac{1}{10}\right)^{20} = 0.00000000000000000001$$

になってしまい、あっという間に計算機で表現できなくなってしまいます。そこで、この確率の対数をとれば

$$(2.60) \quad \log p(s) = 20 \log \frac{1}{10} = -46.05$$

となり、問題なく表すことができます。関数 $y = \log x$ は図2.18のように単調増加関数ですから、 $p(s)$ の大小と $\log p(s)$ の大小は一致するため、比較にも使うことができるわけです。

ただし、このテキストの対数確率はテキストの長さに依存するため、異なるテキストを用いて比較する場合や、クロスバリデーションでも正確に K 等分できず、テキストの長さが等しくない場合には正しい比較が行えなくなってしまいます。そこで実際には、式(2.60)を文の長さ T で割って

$$(2.61) \quad \frac{1}{T} \log p(s) = \frac{1}{T} \sum_{t=1}^T \log p(c_t | c_1, \dots, c_{t-1})$$

として、1文字あたりの確率の対数を計算すれば、異なるテストデータでも問題

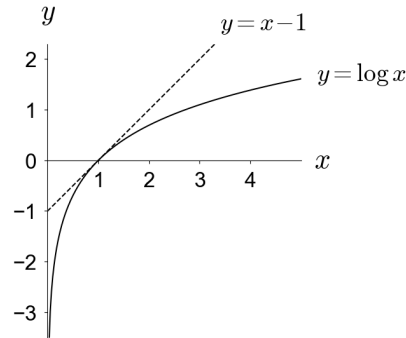


図 2.18: 対数関数 $y = \log x$ のグラフ. 点線で表したのは $x = 1$ での接線 $y = x - 1$ で, $\log x$ はつねにこの接線の下にあり, $\log x \leq x - 1$ が成り立っています.

なく比べることができます. 上の例では,

$$\frac{1}{20} \log p(s) = \frac{1}{20} \cdot 20 \log \frac{1}{10} = -2.303$$

が 1 文字あたりの確率の対数になります. なお, 式(2.61)は $\log p(s)^{1/T}$ のことですから, これはテストデータの各単語の確率の**幾何平均**を計算しており, それを対数で表示している, ということになります.

2.6.3 情報理論の基礎

式(2.58)で用いた, 確率の対数

$$(2.62) \quad \log p(c_t | c_1, \dots, c_{t-1})$$

は, 単に計算の利便性だけでなく, 情報理論において**情報量**という意味を持っています. このことについて, 少し説明しましょう.

たとえば, 可能な選択肢が 8 つあってどれも等価なとき, どれかが選ばれる確率は $1/8$ です. つまり逆にいうと, 確率 $1/8 = 0.125$ とは, 等価な選択肢が 8 つある状況と同じであることを意味しているわけです. この等価な選択肢の数のことを, **分岐数**といいます.

このとき実際に選択肢から 1 つを選ぶのに, 8 回の選択をする必要はありません

ん. 8 個の選択肢を x_1, x_2, \dots, x_8 と表すと, コインを投げて表裏のどちらが出たかを使うことにして,

- 1 回目のコインで (x_1, \dots, x_4) と (x_5, \dots, x_8) のどちらを選ぶかを決め, 前者の (x_1, \dots, x_4) が選ばれたとすると
- 2 回目のコインで (x_1, x_2) か (x_3, x_4) のどちらを選ぶかを決め, 後者の (x_3, x_4) が選ばれたとすると
- 3 回目のコインで x_3 と x_4 のどちらを選ぶかを決め, ここでは表が出たので, 最終的に x_3 を選ぶ

のようにすれば, どれかをランダムに選ぶことができます. このように, 8 個の選択肢であれば $\log_2 8 = 3$ 回のコインを振れば出力を決めることができるため, 必要な情報量は 3 だということができます. 同様にして, 確率 p の事象は $\frac{1}{p}$ 個の等価な選択肢から 1 つを選ぶことと等価で, このときの手数, すなわち必要な情報量は $\log \frac{1}{p}$ です. そこで, 確率 p の事象の「情報量」を表すこの量を Shannon の**自己情報量**といい,

$$(2.63) \quad I(x) = \log \frac{1}{p(x)} = -\log p(x) \quad (\text{自己情報量})$$

と書きます. 対数の底は任意ですが, 上のように 0/1 の決定をもとに底を 2 にする場合は bit (binary digit の略), 自然科学や計算機上の対数で通常用いられている e を底にする場合は nat (natural digit) が単位となります.^{*39} たえば, 確率 $p(x) = 0.07$ の単語が出現したとき, これは $1/0.07 = 14.3$ 個の等価な選択肢から 1 つが選ばれたことと同じなので, その情報量は

$$(2.64) \quad I(x) = \log \frac{1}{0.07} = \log 14.3 = 2.66 \text{ (nat)}$$

になります. 通常は対数の底は一定に揃えて考えているため, 単位は省略するのが一般的です. 本書では特に断りのない限り, 対数としては自然対数を用います.

^{*39} 1000 個の選択肢があるとき, その情報量は自然対数を使えば $\log 1000 = 6.91$ (nat) です. 一方で仮想的に 1000 面サイコロを準備して, このサイコロで選ぶ手数を kit と定義すれば, この情報量は 1(kit) になりますが, 1000 面のサイコロを振って結果を求めるための計算量は別にかかりますから, 底を大きくすれば情報量が減るわけではなく, 値のスケールは対数の底に依存します.

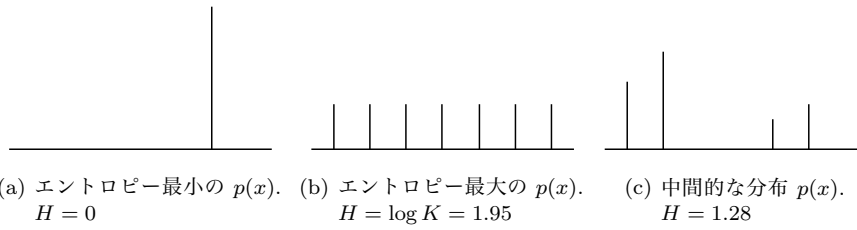


図 2.19: 確率分布 $p(x)$ とそのエントロピー $H = -\sum_x p(x) \log p(x)$.

メモ：符号長について

底となる数を d としたとき、情報量とは、対応する分岐数を d 進法で表したときに何桁になるか、すなわち、何個の数字が必要になるかという量だといってもよいでしょう。これを**符号長**とよび、式(2.63)の情報量とは符号長とみなすことができます。たとえば確率 $1/16 = 0.0625$ は分岐数 16 に対応しますから、これは 2 進数では 1011 のような 4 桁の数で表され、情報量は 4 (bit) になります。nat で情報量を考える場合は、仮想的に $e \simeq 2.72$ 進数を使っていると考えればよいでしょう。

このように情報理論と情報圧縮には密接な関係があり、可能なデータに高い確率を与えること、すなわち、よいモデリングを行うことと、それを短い符号長で表すことは、ほぼ等価な問題です[25, 26].

われわれが実際に扱う確率は多くの場合、単一の確率 $p(x)$ というより、アルファベット \mathcal{X} 上の確率分布 $\{p(x) | x \in \mathcal{X}\}$ でしょう。このとき、式(2.63)で表される x の自己情報量を、 $p(x)$ で期待値をとった

$$(2.65) \quad H = \sum_x p(x)(-\log p(x)) = -\sum_x p(x) \log p(x)$$

を、確率分布 $p(x)$ の**エントロピー** (entropy) とよびます。よってエントロピーとは、確率分布 $p(x)$ から現れる記号のもつ情報量の期待値で、

$$(2.66) \quad H = -\sum_x p(x) \log p(x) = \langle -\log p(x) \rangle_{p(x)}$$

とも書くことができます. $\langle \dots \rangle_{p(x)}$ は $p(x)$ で期待値をとる, すなわち $\mathbb{E}_{p(x)}[\dots]$ の略記法で, 本書でも以下必要に応じてこの表記を用います. 図 2.19 に示したように, $p(x)$ がたとえば $(0, 1, 0, \dots, 0)$ のようにどれかの確率が 1 で他が 0 のときにエントロピーは最小となり, 式(2.65)から $H = 1 \cdot (-\log 1) = 0$ です.*40 ippō, $p(x)$ が一様分布 $(1/K, 1/K, \dots, 1/K)$ のときエントロピーは最大になり, $H = -K \cdot \frac{1}{K} \log \frac{1}{K} = \log K$ が最大値となります.

これからわかるように, エントロピーは確率分布 $p(x)$ の「曖昧さ」を示す量となっています. 最も曖昧な一様分布では次に何が来るかについてまったく予想できませんから, エントロピーは選択肢の数の対数で $\log K$, 曖昧性がなく結果が完全に予想できる場合は結果は 1 通りということですから, エントロピーは $\log 1 = 0$ となるわけです. 符号長の言葉で言えば, エントロピーとは, 確率分布 $p(x)$ から現れる記号を符号化するのに必要な符号長 (= 情報量) の期待値と考えることができます.

パープレキシティ 式(2.61)で計算したテキストの平均的な予測確率の対数, すなわち (負の) 情報量

$$(2.67) \quad \frac{1}{T} \log p(s) = \frac{1}{T} \sum_{t=1}^T \log p(c_t | c_1, \dots, c_{t-1})$$

は, よいモデルであるほど予測確率が高いため, 大きな値になります. ただしこの値は, 対数の底に依存することや, 対数をとっているために差がわかりにくいという問題があります.

たとえば平均予測確率が 0.01 のモデルを改善して, 2 倍の 0.02 にしたとしても, 平均予測確率の対数は $\log 0.01 = -4.61$ から $\log 0.02 = -3.91$ になるだけで, 差は 0.69 しかありません. この値の解釈も難しいため, モデルの評価には, 平均予測確率の逆数, すなわち分岐数を用いることがよく行われます. これを **パープレキシティ** (perplexity) といいます.*41 すなわち, パープレキシティとは **平均分岐数** のことです. 例えばいまの例では, パープレキシティが $1/0.01 = 100$

*40 なお, $0 \log 0 = 0$ と約束します.

*41 perplex は “まごつかせる” という意味で, これはモデルが次の単語を選ぶのにどれだけ “まごつく” のか, ということを表しています.

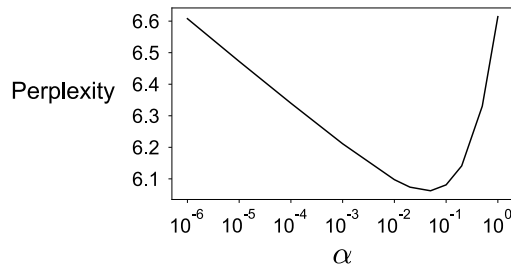


図 2.20: テストデータ `alice.test.txt` のパープレキシティと、文字トライグラム言語モデルの α の値の関係。横軸は対数軸になっていることに注意してください。

から $1/0.02=50$ になった、と考えるわけです。この値は直感的で、対数の底に依存しません。パープレキシティは(幾何)平均予測確率の逆数ですから、実際には式(2.61)から、次のようにして計算します。

$$\begin{aligned}
 (2.68) \quad \text{PPL}(s) &= 1 / \left(\prod_{t=1}^T p(c_t | c_1, \dots, c_{t-1}) \right)^{1/T} = \left(\prod_{t=1}^T p(c_t | c_1, \dots, c_{t-1}) \right)^{-1/T} \\
 &= \exp \left(\log \left(\prod_{t=1}^T p(c_t | c_1, \dots, c_{t-1}) \right)^{-1/T} \right) \\
 &= \exp \left(-\frac{1}{T} \sum_{t=1}^T \log p(c_t | c_1, \dots, c_{t-1}) \right). \quad (\text{パープレキシティ})
 \end{aligned}$$

パープレキシティは小さいほど予測確率が高い、よいモデルであることを意味します。最小値はすべての確率が1のときですから $1/1=1$ 、最大値はすべての確率が $1/V$ (V は語彙の数) の等確率の場合で V になります。よって、パープレキシティは言語の場合は、**語彙の大きさによってスケールが異なる**ことに注意してください。^{*42} 図 2.20 に、異なる α の値について `alice.txt` のテストデー

^{*42} パープレキシティが語彙の大きさに依存するため、別のテキストと比べられないというこの欠点を解消するため、最近、PPLu という新しい指標が提案されました[27]。PPLu では、式(2.68)で予測確率 $p(c_t | c_1, \dots, c_{t-1})$ ではなく、そのユニグラム確率との比 $p(c_t | c_1, \dots, c_{t-1})/p(c_t)$ の積を計算します。これは相対的な値であるため語彙の大きさによらず、3.2.2 節の議論から文脈 c_1, \dots, c_{t-1} と次の語 c_t の自己相互情報量という意味を持っており、実験的にも有効な指標であることが確かめられています。詳しくは、原論文[27]を参照してください。

タで計算したパープレキシティの値を示しました。パープレキシティはテキストの予測確率から計算されるため、図 2.15 と本質的に同じ意味を持っていますが、1 単位 (ここでは 1 文字) あたりの値であるため、テキストの長さによらない指標になっているという特徴があります。

ところで、実際には通常のテキストにおいて、パープレキシティ 1、すなわちすべての予測確率が 1 の「完璧なモデル」を達成することは原理的に不可能で、ある**下限**が存在します。これはなぜでしょうか。そして、その下限は何を表しているのでしょうか？

エントロピーとクロスエントロピー

いま、言語 \mathcal{L} に含まれるすべての要素 $x \in \mathcal{L}$ について、自然のもつ真の確率モデル $p(x)$ を、統計モデル $q(x)$ を使って近似することを考えましょう。ここで要素 x として考える単位は文字や文、文書など何でもかまいません。 $q(x)$ をできるだけ $p(x)$ に近づけたいのですが、こうした**確率分布の間の距離を測るのに**、次の **Kullback-Leibler** カルバックライブラー **ダイバージェンス (KL ダイバージェンス)** という基本的な量があります。^{*43}

$$(2.69) \quad D(p||q) = \sum_{x \in \mathcal{L}} p(x) \log \frac{p(x)}{q(x)}$$

KL ダイバージェンス $D(p||q)$ ^{*44} は $p = q$ のときに最小値 0 をとり、常に非負の値をとる、すなわち

$$(2.70) \quad D(p||q) \geq 0 \quad (\text{KL ダイバージェンスの非負性})$$

をみます。これは、次のようにして示すことができます。図 2.18 に示したように、 $y = x - 1$ は $y = \log x$ の接線となり、必ず $\log x \leq x - 1$ ですから、

$$(2.71) \quad \sum_x p(x) \log \frac{q(x)}{p(x)} \leq \sum_x p(x) \left(\frac{q(x)}{p(x)} - 1 \right) = \sum_x q(x) - 1 = 0$$

^{*43} Solomon Kullback (1907–1994) と Richard Leibler (1914–2003) はアメリカの数学者・暗号研究者です。KL ダイバージェンスは、1951 年の論文[28]で導入されました。世界初のコンピュータ ENIAC の登場が 1946 年だったことに注意してください。

^{*44} $KL(p||q)$ と書くこともあります。

が成り立ちます. 両辺に -1 をかければ, すべての p, q について

$$(2.72) \quad \sum_x p(x) \log \frac{p(x)}{q(x)} \geq 0$$

が得られます. \square

KL ダイバージェンス $D(p||q)$ は, p と q について**対称ではない**ことに注意してください.*45 実際, 式(2.69)を書き直すと

$$(2.73) \quad \sum_{x \in \mathcal{L}} p(x) \log \frac{p(x)}{q(x)} = \left(\sum_{x \in \mathcal{L}} p(x) (-\log q(x)) \right) - \left(\sum_{x \in \mathcal{L}} p(x) (-\log p(x)) \right) \\ = \langle -\log q(x) \rangle_{p(x)} - \langle -\log p(x) \rangle_{p(x)}$$

ですから, KL ダイバージェンスは前節の議論から, 言語 \mathcal{L} のテキストを「近似モデル q で符号化したときの符号長」と「真のモデル p で符号化したときの符号長」の差の期待値を計算していることとなります.*46 KL ダイバージェンスが非負であるということは, 符号化の意味では, 真のモデル p より, 近似モデル q で符号化した方が平均的には必ず長い符号を必要とする, ということを意味しています.

この KL ダイバージェンスを使うと, われわれは言語の持つ真の $p(x)$ を統計モデル $q(x)$ で近似しようとしているのですから, その差は

$$(2.74) \quad D(p||q) = \sum_{x \in \mathcal{L}} p(x) \log \frac{p(x)}{q(x)} \\ = - \sum_{x \in \mathcal{L}} p(x) \log q(x) - \left(- \sum_{x \in \mathcal{L}} p(x) \log p(x) \right) \geq 0$$

です. これから,

*45 このため, $D(p||q)$ は距離の公理を満たさず, 厳密には通常の意味での「距離」ではありません. ただし, $D(p|(p+q)/2)$ と $D(q|(p+q)/2)$ の平均をとることで対称性を持つ **Jensen-Shannon ダイバージェンス**というダイバージェンスも存在し, 様々な研究で使われています. 詳しくは, この後で紹介する情報理論の教科書を参照してください.

*46 または, 「真のモデル p を統計モデル q で近似したときに, 平均的にどれだけビット落ちするか」を計算しているといってもいいでしょう.

$$(2.75) \quad - \sum_{x \in \mathcal{L}} p(x) \log q(x) \geq - \sum_{x \in \mathcal{L}} p(x) \log p(x) = H(p)$$

が成り立つことがわかります. 式(2.75)の右辺は言語の持つ真の確率分布 $p(x)$ のエントロピー $H(p)$ で, これは定数です. いっぽう, 左辺は確率分布 p と q の **クロスエントロピー** (cross entropy) とよばれる量

$$(2.76) \quad H(p, q) = - \sum_{x \in \mathcal{L}} p(x) \log q(x) \quad (\text{クロスエントロピー})$$

で, 言語のあらゆる文 x について, 統計モデル $q(x)$ で計算した情報量 $-\log q(x)$ の, 真の確率分布 $p(x)$ についての期待値です. すなわち, 式(2.75)は, あらゆる確率分布 p, q について, p に関する q のクロスエントロピーは p のエントロピーより大きく,

$$(2.77) \quad H(p, q) \geq H(p) \quad (\text{クロスエントロピーとエントロピー})$$

であることを表しています. また, 言語の持つ真の確率分布 p のエントロピー $H(p)$ は定数ですから, 式(2.74)より, **クロスエントロピーの最小化は, KL ダイバージェンスの最小化と等価である**ことがわかります.

いま, 手元の N 個の x_1, x_2, \dots, x_N について, N が十分に大きければ, これは言語の持つ真の確率分布 $p(x)$ からのサンプルだと考えることができますから, 式(2.75)および式(2.77)の左辺のクロスエントロピーは

$$(2.78) \quad - \sum_{x \in \mathcal{L}} p(x) \log q(x) \simeq - \frac{1}{N} \sum_{n=1}^N \log q(x_n)$$

とモンテカルロ近似することができます.*47 式(2.78)を式(2.75)の左辺に代入して, 両辺を指数の肩にのせれば,*48

*47 関数 $f(x)$ の確率分布 $p(x)$ に関する期待値 $\int p(x)f(x)dx$ を解析的に計算するのが難しいとき, $p(x)$ からランダムにサンプリングされた N 個の $x^{(i)}$ ($i = 1, \dots, N$) を使って, $\int p(x)f(x)dx \simeq \frac{1}{N} \sum_{i=1}^N f(x^{(i)})$ と数値的に積分を近似することをモンテカルロ積分とよびます. 詳しくは, 機械学習におけるモンテカルロ法の優れたチュートリアルである[29]や, 教科書[26]の29章を参照してください. 情報理論では, エルゴード情報源について Shannon-McMillan-Breiman の定理とよばれる定理でこの置換を行います[30].

*48 $y = e^x$ は単調増加関数なので, $a < b$ であることと $e^a < e^b$ であることは等価です.

$$(2.79) \quad \text{PPL} = \exp\left(-\frac{1}{N} \sum_{n=1}^N \log q(x_n)\right) \geq e^{H(p)}$$

が得られます。式(2.79)の左辺は式(2.68)のパープレキシティですから、**パープレキシティには超えられない下限が存在し**、それは言語の持つ真の確率分布 $p(x)$ のエントロピー $H(p)$ を用いて $e^{H(p)}$ と書けることとなります。

$e^{H(p)}$ は言語の持つ真の平均分岐数(われわれには未知)で、**その言語の複雑さ**を表しています。よって式(2.79)は、統計モデルのパープレキシティは決して言語の持つ真の平均分岐数より小さくはできない、ということを意味しています。

これは、簡単な人工例を考えてみればわかりやすいでしょう。たとえば、0 と 1 が完全にランダムに出現する

1100010011101010100000011...

のような系列にどんな統計モデルを考えても、もとの系列がランダムなので、決して次を予測することはできません。この場合、語彙は 0 と 1 の 2 つなので、パープレキシティの最小値(および最大値)は 2 になります。

一方、0 と 1 がそれぞれ確率 (0.9, 0.1) で現れる

0000011000000100000001000...

のような系列の場合、真のモデルと等しい $(q(0), q(1)) = (0.9, 0.1)$ という統計モデルを推定することで、パープレキシティを 2 よりずっと下げることができません。しかし、それはエントロピー $H(p) = -0.9 \log 0.9 - 0.1 \log 0.1 = 0.325$ から得られる $e^{H(p)} = 1.38$ より小さくすることはできません。というのは、正しいモデルを知っていても、次の数字が 0 か 1 かには常にランダム性が残っており、次の数字を完全に正確に予測することは原理的に不可能だからです。これが、パープレキシティに情報理論的な下限が存在する理由です。

言語のエントロピー なお、情報理論を創始した Shannon は、1948 年の最初の論文[31]の中で、文字の n グラムモデルなど様々な方法を用いて英語のエントロピーの上限を計算し、1 文字あたりほぼ 1 ビットという結果を得ています。

Brown コーパスを使って計算すると^{*49}、英語の単語の平均的な長さは約 4.7 文字ですから、これは 1 単語あたりのパープレキシティでは $2^{4.7} = 26$ に相当します。

現代の大規模なニューラル言語モデルを使うとさらに正確な推定が可能になり、5000 億語^{*50} のコーパスから学習された、執筆時点で最大の超巨大な言語モデル GPT-3 [32] では、Penn Treebank コーパスにおいて単語あたりのパープレキシティ 20.5 が得られたことが報告されています。式(2.79)から、言語の真のエントロピーは $\log 20.5$ より小さいと考えられますが、これまでの議論から、原理的にそれが 0 になることはありません。^{*51}

2.6.4 統計モデルと汎化性能

これまでに行ってきた、学習データによる統計モデルの学習とテストデータによる評価を、もう一段高い視点から見てみることにしましょう。

図 2.21 に示したように、(教師なし学習の) 統計モデルとは、あらゆる可能なデータ $D \in \mathcal{D}$ について、真の確率 $p(D)$ の推定値である確率 $q(D)$ を与える確率分布 $\{q(D)\}_{D \in \mathcal{D}}$ と考えることができます。たとえば、 D として「瓶銅肩果扉暴秋銀露呂非…」のようなテキストには低い確率を、「こんやの星祭に青いあかり…」のようなテキストには高い確率を与えるのが、良い統計モデル $q(D)$ というわけです。

いま、手元にある学習データ D および テストデータ D' はどちらも、真の分布 $p(D)$ からのサンプルだと考えられますが、一般に $D \neq D'$ なので、図 2.21 のモデル 1 のようにあまりに現在の D の周りに特化した確率分布 $q(D)$ を学習してしまうと、 D とは違う D' での確率 $q(D')$ が大きく下がってしまいます。これを**過学習 (オーバーフィット)** といいます。よって、統計モデルの学習の目標

*49 `awk` を使うと、`% awk '{for(i=1;i<=NF;i++){s+=length($i);n++}};END{print s/n}' brown.txt` とすれば簡単に求めることができます。

*50 正確には、これは単語を細分した Byte-pair トークンですので、実際の単語数よりは若干大きくなっています。

*51 鋭い方は、これが「ラプラスの魔」と同じ問題であることに気がつかれたでしょう。偶然性や自由意志を否定し、すべての機構を知れば世界は完全に確定的に動いている、と考えればエントロピーの下限は 0 になりますが、仮にそうだとすると、有限時間を生きる人間がその機構を完全に知ることができるかは別の問題で、統計的に推定するしかないとも考えられます。

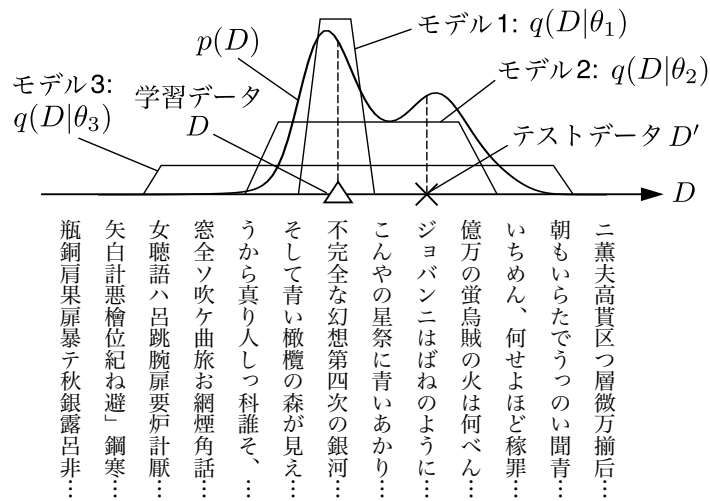


図 2.21: 統計モデルの過学習と汎化性能. モデル 1 は与えられた学習データ D (\triangle 印) に過学習してしまい、テストデータ D' (\times 印) に非常に低い確率しか与えられなくなっています. モデル 3 は一般的すぎて逆に過少適合となるため、テストデータに一番高い確率を与えるのは、真のデータ分布 $p(D)$ に近く、汎化性能が高いモデル 2 となっています. 図の下に、『銀河鉄道の夜』を学習データとした場合の可能なデータ空間の例を示しました.

とは、学習データ D から学習しつつ、未知のテストデータ D' をうまく予測できる、すなわち D' の確率 $q(D')$ が高くなるように学習を行うことです. こうした、テストデータにおける性能を**汎化性能**といいます. われわれはあらゆる可能なテキストを観測するわけにはいきませんが、学習データを用いて、できる限り汎化性能の高い(すなわち、真の分布 $p(D)$ に近い) 確率モデル $q(D)$ を学習したい、ということになります.

もちろん、図 2.21 のモデル 3 のように可能なデータ全体を薄くカバーする確率モデルを推定すれば、どんなテストデータにも一定の確率を与えることができます. たとえば 2.1 節の文字ユニグラムモデルを用いれば、「うから真り人しつ科誰そ」のような意味不明な日本語にも、0 でない確率が与えられるわけです. しかし、これは明らかにモデルが一般的すぎますから(これを**過少適合(アンダーフィット)**といいます)、最もよいのは図 2.21 のモデル 2 のように、学習

データもテストデータも適度にカバーする、真の分布 $p(D)$ に近いモデルを見つけることです。実際このとき、テストデータにおける確率は図からわかる通り、モデル2によるものが最も高くなります。よって、統計モデルの学習には検証データ D' を学習データ D 以外からランダムに選び (これは真の分布 $p(D)$ からのサンプルだと考えられます)、そこでの確率 $q(D')$ が最も高くなるようにすればよい、ということになります。これを系統的に行うのが、先に示したクロスバリデーションです。

なお、統計モデルのパラメータ θ に事前確率 $p(\theta)$ を置くベイズ的なモデルを考えれば、検証データを使わずに学習データだけから、汎化性能が高い統計モデルを学習できることを示すことができます。これについては、3章でベイズ推定を導入した後で説明します。

2章のまとめ

2章では、文字の統計モデルを例にして確率、同時確率、条件つき確率といった統計の基本的な概念と、同時確率の周辺化、ベイズの定理といった操作について説明しました。特に、条件つき確率やベイズの定理は最も基本的な確率の連鎖則の式(2.20)からただちに導くことができますので、公式を暗記する必要はありません。

条件つき確率を使うと n グラム言語モデルを作ることができ、文の確率を計算したり、逆に文を確率的に生成することができます。文字 n グラム言語モデルは「単語」の概念すらない最も基本的なものですが、例にみたように、かなり日本語や英語に近い文字列を生成することができます。^{*52} 特に単語の綴りは、文字 n グラムで非常によく生成することができました。

n グラム言語モデルのような統計モデルの性能は、学習データとテストデータを分け、学習データで計算した統計モデルがテストデータを予測できるかで測ることができます。この性能のことを、汎化性能と呼ぶのでした。統計モデルを使えば、アドホックな方法とは異なり、開発した方法をこうして客観的に評価

*52 上で述べたように、これは Shannon による情報理論の最初の論文[31]で行われている実験そのものです。

し, 改善することができます.*⁵³ 現代の深層学習でも, 背後にはすべてこうした統計的な概念が存在しています.

*53 言うまでもなく, これはある方法が科学であるための要件そのものです.

ノート：テキスト処理のための言語

本書ではプログラム言語として Python を主に使用していますが、Python だけが唯一の言語というわけではありません。Python の前に一世を風靡した Perl [33] は、非常にテキスト処理に長けた言語でした。作者の Larry Wall はバークレー校で言語学を学んでいたため、Perl の構文には `$@&` による変数の「品詞」、`$_` による「痕跡」など、言語学の要素が多く取り入れられ、その機能の一部は Python にも継承されています。

また、テキストを扱うのを得意とする `awk` や `sed` といった言語が、Linux や MacOS のような Unix には古くから標準的に含まれており、利用できます。^{*54} 48 ページの脚注でも使用した `awk` は、テキストを 1 行読むごとにフィールドを空白で自動的に分割して `$1, $2, …` という名前をつけてくれますので、たとえばテキストの `Foo` で始まる各行の 2 個目のフィールドの総和を求めたければ、

```
% awk '/^Foo/{s+=$2};END{print s}' input.txt
```

のように簡潔に書くことができ、通常は Excel で行うような計算を簡単に行うことができます。`awk` は他にも `exp`, `log`, `sin` のような算術演算や `substr` のような文字列演算、`for` 文による繰り返しなども備えており、簡単なデータ解析にもたいへん有用です。

`sed` は stream editor の略で、1 行しか画面出力を持たないエディタ `ed` ^{*55} の拡張ですが、それゆえに簡潔なコマンド体系を持っています。たとえば

```
% sed '1,5d;s/foo/bar/g' input.txt
```

とすれば、1 行目から 5 行目を削除 (`d`) した上で、`foo` をすべて (`g`) `bar` に置換 (`s`) して出力します。

*54 これは、歴史的に Unix がテキスト処理をその目的の一つとして開発されたためです。

*55 `ed` については、『The UNIX Super Text』[34]などを参照してください。

また

```
% sed G input.txt
```

とすれば、1行おきに空白を入れた「ダブルスペース」のテキストを簡単に作ることができます。パターンスペースやホールドスペースといった内部の記憶領域をうまく使うと、さらに高度な処理も可能で、何と図 2.22 のように sed で書かれたテトリス^{*56} やチェス^{*57} すら存在します。

sed や awk については『sed&awk プログラミング』[35]『プログラミング言語 awk』[36]といった標準的教科書があるほか、「awk の簡単な使い方」^{*58} 「sed 教室」^{*59} 「sed は日暮れて」^{*60} といったフリーのチュートリアルがあります。こうした言語を適宜使いこなすことで、テキストを計算機でより自由に扱えるようになるでしょう。

なお、本書を書く際に用いた組版プログラム L^AT_EX (T_EX) も、テキストを処理するプログラミング言語の一種といえます。チューリング完全、すなわち Python のような言語と同じ表現能力を持っているため、T_EX で書かれた BASIC インタプリタ B_AS_IX^{*61} [37] など、驚くべきプログラムも存在しています。

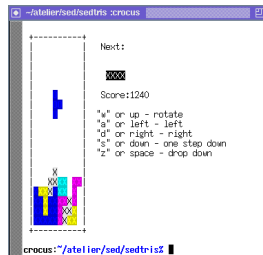


図 2.22: 文字端末で、sed で書かれたテトリスをプレイしている様子。

*56 <https://github.com/uuner/sedtris>

*57 <https://github.com/bolknote/SedChess>

*58 「awk の簡単な使い方」<http://chasen.org/~daiti-m/etc/awk/> に転載。

*59 「sed 教室」<https://www.gcd.org/sengoku/sedlec/>

*60 「sed は日暮れて」<https://chimimo.tumblr.com/post/8995558289/sed1>

2章の演習問題

- (1) `alice.txt` や `brown.txt` で, `j` に続く確率の高い文字にはどんなものがあり, その確率はそれぞれ何でしょうか.
- (2) 上のテキストで, `z` や `th` の前に来ることのできる文字にはどのようなものがあり, それらの確率は何でしょうか.
- (3) `brown.txt` など, `alice.txt` 以外の英語のテキストでは, `alice.txt` と文字バイグラム確率にどのような違いがあるのでしょうか. どのバイグラムの確率が特に異なるかを自動的に検出するには, どうすればよいのでしょうか.
- (4) 本章では簡単のために英語のテキストを使用しましたが, 日本語のテキストの場合は, 文字のユニグラム確率やバイグラム確率はどうなっているのでしょうか. 漢字を含めると文字の種類が膨大になってしまうので, ひらがなまたはカタカナだけに絞るとよいでしょう. 図 2.3 のような行列を作ると, どうなっているのでしょうか.

サポートサイトの `data` フォルダに, 『更級日記』のテキスト `sarashina.txt` を置きました^{*62}. この場合のひらがなのユニグラム確率, バイグラム確率は, 現代の日本語と比べるとどのような違いがあるのでしょうか.

Python で文字列がすべてひらがなから成っているかを判定するには, `regex` パッケージをインストールした後で, Unicode の文字プロパティを使って

```
import regex
def is_hiragana (s):
    return regex.search(r'^\p{Hiragana}+$', s)
```

のような関数を書くことができます.

- (5) 日本語のテキストで, 読点「、」の前に来る文字には, どんな傾向があるのでしょうか.
- (6) 上で使ったような日本語のテキストを使って, 文字 n グラムモデルからランダムに文を生成するとどうなるのでしょうか.

^{*61} <https://www.ctan.org/tex-archive/macros/generic/basix>

^{*62} このテキストは, 渋谷栄一氏が <http://genjiemuseum.web.fc2.com/sara2.html> で公開されているものを整形して使用しました.

- (7) 本章で紹介した言語モデルよりも性能のよい, 3.4.4 節の Kneser–Ney 言語モデルを使って, 4 グラムや 5 グラムの文字 n グラムから 2.5.2 節のようにテキストを生成してみると, どうなるでしょうか. 本章で紹介した 3 グラムの場合と, どんな違いがあるでしょうか.
- (8) 単語は, 文字からなる短い「文」とみなすことができます. 単語辞書を使って, 48 ページのように単語の綴りをランダム生成してみましょう. どんな単語ができるでしょうか.
Linux や MacOS のような Unix では, 単語辞書は通常, `/usr/share/dict/words` に置いてあります. BOS, EOS を必ず使うように注意してください. また, このようにして生成された単語には, 現実の単語と比べてどのような問題があるでしょうか.
- (9) ジェイムズ・ジョイスの小説 “Finnegan’s Wake” (1939) は, 言葉遊びに満ちており, `prumptly` や `sesther`s といった, 通常の英語では見たこともないような単語が多数登場する, 複雑な作品です. こうしたテキストを扱うには, 単語の言語モデルではなく, 文字の言語モデルを考えるのが適切でしょう. サポートサイトの `data` フォルダにある `finnegans.txt` を用いて, この小説のパープレキシティを計算してみましょう. `brown.txt` のような通常の英語と比べて, パープレキシティはどうなっているでしょうか. また, 48 ページのようにして, この作品から (ジョイス自身が行ったように) 「単語」をランダムに生成すると, どうなるでしょうか.
- (10) 英語や日本語以外のテキストでランダム生成してみると, どうなるでしょうか. プログラミング言語も, 一種の形式言語です *63. プログラムのテキストから文字 n グラム言語モデルを学習して, ランダム生成するとどうなるでしょうか. その場合, 何が問題になるでしょうか.
- (11) X も Y も二値のとき, 2.4.2 節のベイズの定理を使って, Y が与えられたときの X の条件つき確率を計算する一般的なプログラム `bayes.py` を書いてみましょう. `lik.dat` および `prior.dat` をそれぞれ次のようなテキストファイルとして与え, 観測値 Y を与えると, `bayes.py` は $p(X|Y)$ を

*63 プログラミング言語のように, 人間が定義した言語である形式言語と区別する意味で, 日本語や英語のような言語は「自然言語」と呼ばれています.

計算します.

```
% cat lik.dat
1 0
0.3 0.7
% cat prior.dat
0.2 0.8
% bayes.py lik.dat prior.dat 0  (← Y=0 が観測された)
posterior: p(X|Y=0) = [0.454, 0.545]
% bayes.py lik.dat prior.dat 1  (← Y=1 が観測された)
posterior: p(X|Y=1) = [0, 1]
```

2章の文献案内

統計的自然言語処理の最も有名で基本的な教科書は、現在も最前線で活躍する研究者の Schütze と Manning による “Foundations of Statistical Natural Language Processing” (FSNLP とよばれています) [38] でしょう。最近になり、日本語訳も出版されました [39]。FSNLP は 1999 年に出版されたため、本書で説明した内容で扱われていないものも多くありますが、自然言語処理を最も基礎から体系的に説明している教科書といえます。本書で触れられなかった論点も網羅されているため、自然言語処理を専門とする場合はぜひ一読することをお勧めします。日本語の教科書としては、北による『確率的言語モデル』[40]も同じ 1999 年の出版ながら、今なお研究者に広く読まれている名著です。本書のタイトルは、この本のオマージュでもあります。本書は『確率的言語モデル』で扱われていない様々な最新の知見を盛り込み、主に文を対象とした確率的言語モデルから、より広く文書やテキストを、その差異も含めてベイズ統計の立場から扱うものに拡張しました。

統計的自然言語処理は機械学習の一部でもあり、2.6 節で扱ったような統計モデルの学習の一般的な話題は、Bishop による “Pattern Recognition and Machine Learning” (通称 PRML) [23] で詳しく説明されています。筆者も翻訳に参加した日本語版 [41] も出版されており、多くの方に読まれています。また、ケンブリッジ大学の MacKay による “Information Theory, Inference, and Learning Algorithms” (ITILA) [26] は、ベイズ統計と情報理論に基づいた深い考察に支えられて機械学習を広くカバーする名著で、全文の PDF を公式ページ^{*64} からフリーでダウンロードすることができます。本章前半の議論は、多くこの本を参考にしました。

2.6.3 節でも説明したように、情報理論は機械学習一般、中でも自然言語処理とは深いつながりがあります。多くの場合は 0/1 や a~z のアルファベットの系列を扱う情報理論を、単語や文、文書といったより上位の構造を含めて徹底的に高度化したのが、ある意味で自然言語処理であるといってもいいでしょう。もちろん、情報理論にはそれ以外に符号化や伝送に関わる、多くの技術的内容が含ま

*64 <http://www.inference.org.uk/mackay/itila/book.html>

れています。情報理論では, Cover&Thomas の教科書[30] (和訳[42]) が最も基本的な文献として知られています。また, 上の ITILA [26] でも詳しく扱われています。日本語では, 韓による『情報と符号化の数理』[25] は透徹した哲学に基づいて情報理論を解説した名著で, これを読むことでより本質的な理解と, 情報理論の懐の深さに触れることができるでしょう。

情報理論にもとづく定式化は深層学習時代においても見直されつつあり, 気鋭の若手 Cottrell による COLING 2022 のチュートリアル “Information Theory in Linguistics: Methods and Applications”^{*65} では, 言語学的な応用も含めて幅広い研究が紹介され, 演習用の Jupyter Notebook も公開されています。また, 統計数理研究所の伊庭による「「情報」に関する 13 章」[43]^{*66} は, そもそも「情報」とは何か? という根本的な疑問を抱いている方にもお薦めできる, たいへん読みやすく, かつ奥の深い論考です。

人文学においてもテキストデータの使用は以前から行われており, 2001 年の近藤らの論文[44]は, 文字 n グラムモデルを用いて和歌の分析を行った最初期の研究です。人文学におけるデータ解析の教科書としては, Python を用いた実例を示した “Humanities Data Analysis” [45]が 2021 年に出版されています。社会科学におけるテキストデータの使用については, 5 章の文献案内をご覧ください。

*65 <https://rycolab.io/classes/info-theory-tutorial/>

*66 <https://www.ism.ac.jp/~iba/a19.pdf>

3 単語の統計モデル

3.1 文字から単語へ

2章では、文字の統計モデルを題材に、確率の基本と文字 n グラム言語モデル、およびモデルの評価方法について学んできました。

ここまでは簡単のために文字だけのモデルを考えてきましたが、実際には言語には普通は「単語」の概念があり、単語を使った方が、よりよいテキストのモデルになると考えられます。たとえば、英語で“tall”の次に続く言葉は“tree”や“boy”であり、“reason”が来ることはまずないことは、単語としてこの言葉の意味を考えず、“-ll”で終わる文字列として見ていたのでは難しいでしょう。

なお、「単語」という概念があっても、それが現在の英語のように、空白文字で分かち書きされているとは限りません。日本語や中国語、タイ語のような東アジアの言語には空白はありませんし、さらにはラテン語や古典ギリシャ語、古くは英語も、もともとは図 3.1 のように単語間に空白を開けず、続け書きされていました。これは、*scriptio continua* (ラテン語で「続け書き」の意味) とよばれています。

一方で、“New York”のような地名、“with respect to”のような慣用句はそれぞれ“N.Y.”、“w.r.t.”とも書かれることからわかるように、空白があっても、実質的に単語の区切りであるとは限りません。また、“other than”のよ

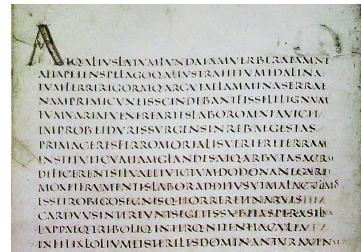


図 3.1: ヴェルギリウスのラテン語のテキスト (西暦 141 年ごろ) での、単語の続け書きの例。*1

*1 https://en.wikipedia.org/wiki/Scriptio_continua より引用, 一部.

うなすべての慣用句を単語として認めるかが決まっているわけではありません。よって、何が「単語」であるかという基準には本質的に曖昧性があり、唯一の正解があるわけではありませんが、*2 多くの言語は現在では空白で区切って書かれますし、日本語や中国語も、最適とはいえなくても、下に述べるように標準的な方法で単語に区切ることができますので、本書でもそれを使っていくことにします。上のような慣用句を統計的に自動認識する方法については、この後の 3.2.2 節を参照してください。

文字列を単語に分解することは**単語分割**、あるいは動詞や名詞といった品詞の付与や活用形の認識も行う場合は**形態素解析**とよばれています。たとえば、標準的な形態素解析器の一つである MeCab*3 を使えば、日本語の文字列を簡単に単語に区切ることができます。Python の場合、これにはまず、MeCab モジュールをインストールする必要があります。詳細は付録??に譲りますが、Google Colaboratory の Python 環境を使っている場合は、執筆時では次のようにすればインストールすることができます (>は Colaboratory のプロンプト)。

```
> !pip install mecab-python3
> !pip install unidic-lite
```

この上で、下のようにすれば文字列を単語に分割することができます。

```
import MeCab
tagger = MeCab.Tagger("-Owakati")
s = "これは日本語の文字列です。"
tagger.parse(s).split() # 結果を空白文字で分割する
⇒ ['これ', 'は', '日本語', 'の', '文字', '列', 'です', '。']
```

したがって、テキスト全体は次のようにして単語に分解することができます。

```
with open('ginga.txt', 'r') as fh:
    for line in fh:
        buf = line.rstrip('\n') # 行末の改行文字を削除
        if len(buf) > 0:      # 空行でない場合
```

*2 筆者は、この問題に対して「観測文字列の確率を最大化する分割が単語である」という立場で教師なし形態素解析[46]の研究を行い、情報理論的な解を与えています。ただし、これは確率モデルに依存しますので、人間の持っている「単語」という概念にかなり近くなるものの、必ず一致するわけではありません。

*3 MeCab は奈良先端大 (当時) の工藤拓氏によって開発された、条件付き確率場 (CRF) に基づく形態素解析器で、<https://taku910.github.io/mecab/>で配布されています。

```

words = tagger.parse(buf).split()
print (words)
⇒ [' 銀河', ' 鉄道', ' の', ' 夜']
   [' 一', ' ', ' ', ' 午後', ' の', ' 授業']
   [' 「', ' では', ' みなさん', ' は', ' ', ' ', ' そういう', ' ふう', ' ',
   ' に', ' 川', ' だ', ' と', ' 云わ', ' れ', ' たり', ' ...']

```

このまま解析に使うことも可能ですが、後の処理のために次のように、単語に分割したテキストを別ファイルに保存しておくとうよいでしょう。下のようになると、ginga.split.txt に単語分割されたテキストが保存されます。

```

with open ('ginga.split.txt', 'w') as oh: # 出力ファイル
    with open ('ginga.txt', 'r') as fh: # 入力ファイル
        for line in fh:
            buf = line.rstrip('\n')
            if len(buf) > 0: # 空行は無視する
                words = tagger.parse(buf).split()
                oh.write (' '.join(words) + '\n')

```

なお、形態素解析器にはほかにも kuromoji, GiNZA, KyTea, …など様々なものがあります^{*4} ので、使いやすいものを用いるとうよいでしょう。

3.2 単語の統計と巾乗則

さて、こうしてテキストが単語に分けられたとき、文字のモデルと最も異なっている点は何でしょうか。最も重要なのは、「単語の種類は無限にある」ということでしょう。それぞれの言語で使われる文字の種類は有限ですが、文字の組み合わせである単語は、長さに通常は制限がありませんので、いくらでも多くの単語を作り出すことができます。^{*5}

^{*4} これらの違いは、提供している機能の違いもさることながら、最も大きいのは「単語」をどう定義するかという違いです。たとえば奈良先端大で開発された MeCab などで標準的に使われている IPA 辞書では IPA 品詞体系が採用されており、いっぽう京大で開発された JUMAN では、益岡・田窪文法をもとにした別の JUMAN 品詞体系が採用されています。単語分割が異なるとモデルも違ってきてしまいますので、モデルを学習する際と、それを使用して新しいテキストを解析する際には一般に、形態素解析器を揃える必要があります。

^{*5} 日本語では「リュウグウノオトヒメノモトユイノキリハズシ」のような植物名、英語では『メアリー・ポピンズ』に現れる “supercalifragilisticexpialidocious” のような言葉が有名ですが、もっと長いものも存在し、これらは長大語とよばれています。

とはいえ、計算機上で無限の語彙を表すのは難しいため^{*6}、ほとんどの場合は、一定の基準で語彙を選ぶことになります。実際に、われわれが普段用いている辞書はこうして一定の語彙を選んだものです。表 3.1 に、日本語および英語の主な辞書に掲載されている語彙(見出し語の数)の例をまとめました。これを見ると、単語として必要な語彙の種類は文字の種類よりはるかに大きく、最低でも数万語、多い場合には数十万語を超えることがわかります。つまり、統計モデルとしてみると、単語を考える統計モデルは出力が数万次元を超える**超高次元の統計モデル**になるということです。これが、文字の統計モデルとの大きな違いです。

それでは、実際のテキストにはどれくらいの語彙が含まれているのでしょうか。3.1 節で示した方法でテキストが空白で単語に分かれているとき、次のようにすれば単語の総数と、それぞれの単語の頻度を数えることができます。

```
from collections import defaultdict
freq = defaultdict(int)
with open('ginga.split.txt', 'r') as fh:
    for line in fh:
        words = line.rstrip('\n').split()
        for word in words:
            freq[word] += 1

for w,c in freq.items():
    print('%s -> %d' % (w,c))
⇒ 銀河 -> 25
   鉄道 -> 4
   の -> 1266
```

表 3.1: 日本語および英語の標準的な辞書に含まれる語彙(見出し語)の大きさ。

言語	辞書	語彙(約)
日本語	岩波国語辞典 第八版	67,000 語
	広辞苑 第七版	250,000 語
英語	Longman LDCOE	45,000 語
	リーダーズ英和辞典	280,000 語
	Oxford English Dictionary	600,000 語

^{*6} 筆者による教師なし形態素解析のための言語モデル NPYLM [20]は、文字 ∞ グラムと単語 n グラムを組み合わせることで、無限の語彙を扱うことが可能です。ただし、これは本書のレベルを大きく超えるため、説明は割愛しています。


```

夜 -> 6
一 -> 45
、 -> 988
午後 -> 2
授業 -> 2
「 -> 293
では -> 10
みなさん -> 4
...

```

```

len(freq)
⇒ 2594

```

「銀河鉄道の夜」では語彙の大きさは 2594 であり、各単語の頻度は上のようになっていることがわかります。他のテキストについても同様に数えた語彙数を、表 3.2 に示しました。「銀河鉄道の夜」や「不思議の国のアリス」といった（計算機にとっては）短いテキストでは語彙は数千語程度ですが、一般的な中規模のテキストでは数万語以上、大きなテキストでは数十万語や数百万語を超えることがわかります。^{*7} なお、このように一般的に共有されているテキストデータおよび（あれば）付随データのことを、**コーパス**ともいいます。

コーパスにはさまざまなものがありますが、1967 年に編纂された Brown コーパス[49]は分野が偏らないようにサンプリングされた（これを**均衡コーパス**といいます）、最初の大きなコーパスです。アメリカ英語約 100 万語のテキストとその品詞からなっており、現在では NLTK のようなツールキットにも付属して公開されています^{*8}。現代の均衡コーパスとしては、英国では 1 億語の British National Corpus (BNC)^{*9}、米国では 1500 万語の American National Corpus (ANC)^{*10} があり、いずれもデータをフリーでダウンロードできます。日本語では、国立国語研究所が約 1 億語の現代日本語書き言葉均衡コーパス (BCCWJ^{*11}) を公開しており、オンライン検索は無料で、オフラインのデータは有料で使用す

^{*7} Web から取得した 1 兆語の英語テキストから求めた Google 1T 5-gram データ[47]では語彙は 1350 万語、日本語 2500 億語の 7-gram データ[48]では 256 万語にもなります。

^{*8} http://www.nltk.org/nltk_data/

^{*9} <http://www.natcorp.ox.ac.uk/>

^{*10} <https://anc.org/>

^{*11} Balanced Corpus of Contemporary Written Japanese の略で、<https://clrd.ninjal.ac.jp/bccwj/> から使用・入手することができます。

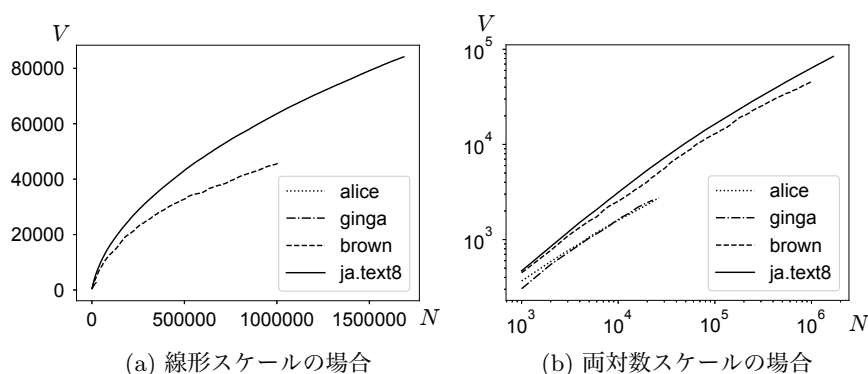


図 3.2: テキストの長さ N と語彙の大きさ V の関係 (Heaps の法則). **alice**:『不思議の国のアリス』, **ginga**:『銀河鉄道の夜』, **brown**: Brown コーパス, **ja.text8**: 日本語 text8 コーパスをそれぞれ表しています. **alice**, **ginga** は短いため, 線形スケールではほとんど見えなくなっています.

ることができます. また有料ではありますが, 毎日新聞, 朝日新聞, 日経新聞といった新聞のテキストも近年のものはすべて電子化されており^{*12}, 購入すればコーパスとして使用することができます. 多言語のコーパスとしては, ドイツで公開されている Leipzig コーパス[50]には 100 言語以上の文が, それぞれ 1 万文から 100 万文含まれており (日本語も Web とニュースの文がそれぞれ 100 万文あります), フリーでダウンロードすることができます^{*13}.

また, 3.5 節で説明する単語ベクトルの学習などを手軽に試すためのコーパスとして, Wikipedia からランダムに 100MB 分のテキストを抽出した text8^{*14}, および日本語版の ja.text8^{*15} はフリーのため, 研究や開発目的で広く使われています. サポートページに, Brown コーパスのテキストを **brown.txt**, text8 および日本語 text8 を **text8**, **ja.text8** として置いておきました. 元々の text8 には改行がないため, 改行を復元したものはそれぞれ **text8.txt**, **ja.text8.txt** となっています.

*12 <https://www.nichigai.co.jp/sales/corpus.html>

*13 <https://wortschatz.uni-leipzig.de/en/download>

*14 <http://mattmahoney.net/dc/textdata>

*15 <https://github.com/Hironsan/ja.text8>

3.2.1 Heaps の法則

単語の種類に制限はありませんから、語彙は一般に、テキストが長くなるほど増えていきます。以下では、テキストの長さは単語数で数えることにしましょう。上の実行例で、テキストを読みながら頻度辞書 `freq` の大きさ `len(freq)` を見れば、語彙がどのように増えるのかを確かめることができます。図 3.2 に、こうして調べたテキストの長さ N と語彙 V の関係を示しました。テキストが長くなるほど、語彙は増えてきれいな曲線を描き、対数スケールで見るとテキストの言語やジャンルによって多少の違いはあるものの、おおむね同じように語彙が直線的に伸びていくことがわかります。このことは、**Heaps の法則**とよばれています。^{*16} Heaps の法則は、次の式で表されます。

$$(3.1) \quad V = K \cdot N^\gamma \quad (\text{Heaps の法則})$$

ここで V はテキストを長さ N まで見たときに含まれる語彙の大きさ、 γ と K はコーパスに依存する定数です。一般的な英語の場合は γ は 0.5 前後、 K は 10 から 100 前後の値になることが知られています。式 (3.1) は、両辺の対数をとれば

$$(3.2) \quad \log V = \log K + \gamma \log N$$

になりますから、 $\log V$ と $\log N$ は傾き γ の比例関係となります。これが、両対数スケールの図 3.2(b) で直線関係が現れている理由です。

表 3.2: テキストに含まれる語彙の大きさの例。頻度を一定以上に限った場合についても同時に示しました。

テキスト	長さ	語彙	頻度 ≥ 2	頻度 ≥ 10
<code>alice.txt</code>	26396	2748	1471	379
<code>brown.txt</code>	1012603	46055	26481	8395
日本語 <code>text8</code>	15160499	249629	126690	45024
毎日新聞 2011 年度	21805730	129971	88669	42465
New York Times 2008 年度	65333240	249610	158735	72554

^{*16} Heaps の法則の名前は、情報検索の分野で Harold Stanley Heaps が 1978 年に出版した本[51]によるものですが、これより前に、初期の計量言語学者である Gustav Herdan (1897–1968) が 1960 年にこの法則を発見していました[52]。よって、これを Herdan-Heaps の法則ということもあります。

なお、次節で説明する Zipf の法則を認めれば、Heaps の法則は Zipf の法則から導くことができます。この後で説明するように、頻度順に単語を並べたとき、単語の頻度 f が順位 r の逆数に比例する式(3.10)、すなわち

$$(3.3) \quad f \propto r^{-\alpha}$$

が Zipf の法則ですから、確率分布にするための正規化定数は

$$(3.4) \quad Z = \sum_{r=1}^{\infty} r^{-\alpha} = 1 + \frac{1}{2^{\alpha}} + \frac{1}{3^{\alpha}} + \frac{1}{4^{\alpha}} + \dots$$

です。この Z は $\alpha > 1$ のとき、一定の値に収束することが知られています。^{*17} よって Zipf の法則の下では、頻度順で r 番目の単語の確率は

$$(3.5) \quad p(r) = \frac{1}{Z} r^{-\alpha}$$

と表されるわけです。

いま、テキストを $N-1$ 語読んだときに、全部で $V-1$ 語の語彙が出現したとしましょう。ここで次に読んだ N 語目の単語が語彙にない、新しい単語だったとすると語彙は V 個に増え、この単語の頻度は 1 で、確率は $1/N$ になります。この単語は頻度順では V 番目ですから、式(3.5)から

$$(3.6) \quad \frac{1}{Z} V^{-\alpha} = \frac{1}{N}$$

が成り立つはずですが、両辺の対数をとって整理すれば、

$$(3.7) \quad \begin{aligned} -\alpha \log V &= \log Z - \log N \\ \therefore V &= Z^{-1/\alpha} \cdot N^{1/\alpha} \end{aligned}$$

となり、Heaps の法則が $K = Z^{-1/\alpha}$, $\gamma = 1/\alpha$ を係数として得られます。一般に言語では α は正確には 1 よりやや大きく、 $\alpha = 1 \sim 2$ 程度ですから、 γ は 0.5 程度になり、実際の観察ともよく合致しています。

^{*17} 数学ではリーマンのゼータ関数 $\zeta(\alpha)$ と呼ばれています。この場合、 α に依存する定数 $\phi(\alpha) < 0.5772$ (右辺はオイラーの定数) を使って、 $Z = 1/(\alpha-1) + \phi(\alpha)$ と表すことができます。

3.2.2 Zipfの法則

上ではテキストに現れる単語を数えて語彙を求めましたが、これらの単語がすべて、同じ頻度で出現するわけではありません。そこで、2.1節で文字について行ったように、単語を出現頻度で並べて表示してみましょう。これは文字の場合とまったく同様に、次のようにして行うことができます。

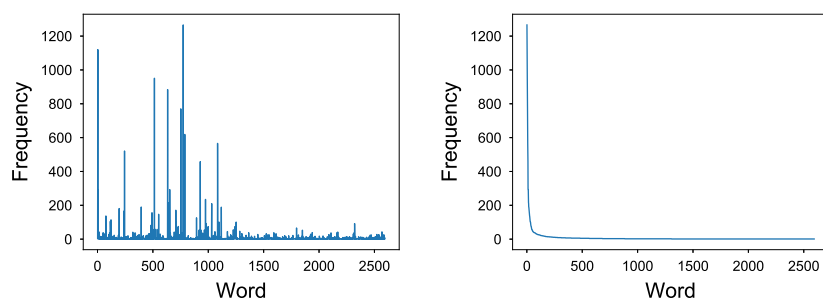
```

for w,c in sorted (freq.items(),
                   key=lambda x: x[1], reverse=True):
    print ('%s -> %d' % (w,c))
⇒ の -> 1266
   。 -> 1120
   、 -> 988
   た -> 951
   て -> 884
   に -> 770
   は -> 619
   を -> 566
   ...
   驚 -> 19
   気 -> 18
   光っ -> 18
   そっち -> 18
   ...
   函 -> 1
   橄欖 -> 1
   堅く -> 1
   握っ -> 1
   便り -> 1
   放課後 -> 1
   知らせよ -> 1

```

この結果からわかる通り、句読点を除くと「の」「に」「が」といった一部の単語に頻度が集中しており、一方でテキストに一度しか出現していないような「函」「橄欖」といった言葉が大量にあることがわかります。^{*18} これらは1回し

^{*18} こうした一度しか出現しなかった語のことを、コーパス言語学では孤語あるいは hapax legomenon (ギリシャ語で「一度だけ書かれた」の意味) といいます。これは、あくまで与えられたテキストについての概念であることに注意してください。例えば、「放課後」はこのテキストでは孤語ですが、一般にはごくありふれた単語です。



(a) 単語を横軸に辞書順に並べた場合.

(b) 単語を横軸に頻度順に並べた場合.

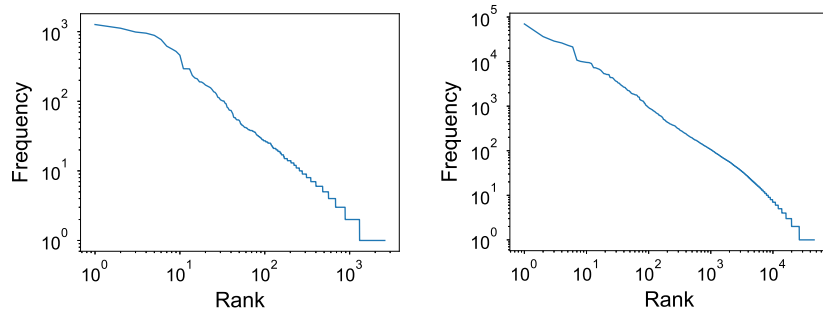
図 3.3: 『銀河鉄道の夜』における単語の頻度.

か出現していないために情報が少なく、統計的な分析からは除いた方がよいでしょう。一般に、テキスト全体を通した頻度がたとえば 10 以上 (大きいテキストならば 100 以上) のように、基準を決めて語彙を選択することがよく行われます。表 3.2 の最後の列に、頻度 10 以上の単語を用いた場合の語彙の大きさを示しました。ただし今は、そうした語彙の選択は行わないことにします。

それでは、上で求めた単語の頻度は、全体としてはどのように分布しているのでしょうか。図 3.3(a) に、『銀河鉄道の夜』に現れる 2594 種類の単語を横軸に辞書順にとり、縦軸にその頻度をとったプロットを示しました。先にみたように、日本語では「の」「に」「が」といった一部の単語が高い頻度を持っているので、ヒストグラムにところどころ、突出した値があることがわかります。横軸の単語を頻度順で並び換えると、図 3.3(b) のようになります。確かに、非常に少数の単語に頻度が集中しているのを見てとることができますね。

ただ、頻度があまりに一部の単語に集中しているため、様子を詳しく見るために、縦軸を頻度 f の対数すなわち $\log f$ で、横軸も同様に頻度の順位 r の対数 $\log r$ で両対数プロットにしてみましょう。こうすると、図 3.3(b) は図 3.4(a) のように表されます。この図をよく見ると、 $\log f$ と $\log r$ の間には傾きがほぼ -1 の比例関係があることがわかります。すなわち、ほぼ

$$(3.8) \quad \log f = -\log r + b \quad (b \text{ は定数})$$



(a) 『銀河鉄道の夜』の場合. 図 3.3(b)を両軸について対数でプロットしたもの. (b) Brown コーパスで同様に計算した場合.

図 3.4: 単語の頻度順位と出現頻度の両対数プロット. 順位と頻度が反比例する Zipf の法則が現れています.

という関係が成り立っています. 式(3.8)は, 変形すると $\log fr = b$ となりますから, $e^b = c$ とおけば

$$(3.9) \quad fr = c$$

あるいは,

$$(3.10) \quad f = c \cdot \frac{1}{r} \quad (\text{Zipf の法則})$$

という関係が成り立つことになります. すなわち, 単語は頻度順に順位が $r = 1, 2, 3, 4, \dots$ になるほど, その相対的な頻度は $1, 1/2, 1/3, 1/4, \dots$ と減っていく, ということです. この関係は, 他のどんなテキストに対してもほぼ成り立つことが知られています. 図 3.4(b) に, Brown コーパス `brown.txt` に対して同様に計算したプロットを示しました. ここでも, 同じ法則が成り立っていることがわかります. この法則は, 1930 年代にこの法則を体系化したハーバード大学の言語学者 Zipf^{*19} の名前をとって, **Zipf の法則** とよばれています. Zipf の法則は, 言語に限らず社会全般で広く成り立つ**巾乗則**として認知されており, たとえば

*19 George Kingsley Zipf (1902–1950) は, 1935 年に[53]で Zipf の法則を示しました. ただし, この法則は Zipf が最初に発見したわけではなく, 1916 年にはフランスの速記者 Jean-Baptiste Estoup によって発見されていました[54, 55]. 最近になり, Estoup の孫にあたる方によって, この事情についての論文が出版されています[56].

都市の大きさ, 収入の分布, 地震の規模と頻度など多くの実際の現象が巾乗則に従うことが知られています[57]. 言語を中心としたこうした巾乗則については, [58]で詳しく論じられていますので参照してください.

この Zipf の法則によれば, 語彙が 10,000 個であれば, 最もよく現れる単語の確率と最も現れにくい単語の確率の比は, $1:1/10000$ で 10,000 倍にもなる, ということとなります. 文字の場合は, 2.1 節で見たように英語では最も頻度の高い文字 'e' と最も低い文字 'z' の確率の比は 150 倍くらいでしたから *20, 語彙の多い単語の場合は, その差は圧倒的に大きくなっていることがわかります.

なお, 単語の確率を式(2.1)のように相対頻度で計算する場合, 確率は頻度に比例しますから, 式(3.10)から最も頻度の高い単語の確率を $p_1 = q$ とおくと, 2 番目に高い確率 p_2 , 3 番目に高い確率 p_3 , ... はほぼ

$$p_1 = q, p_2 = \frac{q}{2}, p_3 = \frac{q}{3}, p_4 = \frac{q}{4}, \dots$$

です. よって, その総和は

$$(3.11) \quad p_1 + p_2 + p_3 + p_4 + \dots = q \left(1 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4} + \dots \right) \\ = q \sum_{r=1}^{\infty} \frac{1}{r}$$

となります. しかし, 式(3.4)でみた調和級数 $\sum_{r=1}^{\infty} 1/r^\alpha$ は $\alpha=1$ のとき無限大に発散しますから, 式(3.11)の確率の和は 1 にならなくなってしまいます. よって級数が収束するためには, 少なくともある順位からは $\alpha > 1$ になっているはずで, 実際に図 3.4(a)でもみられるように, 一般に $r \rightarrow 0$ と $r \rightarrow \infty$ では巾乗則の傾き α が異なっており, この現象は double power-law [59]とよばれています. この説明として, 語彙を核となる語彙と周辺の語彙に分ける生成モデル[60]などが研究されています *21.

*20 英語の ASCII 文字は 7 ビットで表され, 記号を含めて $2^7 = 128$ 文字ありますから, Zipf の法則からもこの差は妥当なオーダーだということがわかります.

*21 数学的には, 最近オックスフォード大学の Caron らによって, Lévy 測度 $\rho(w) = 1/\Gamma(1-\sigma)w^{-1-\tau}\gamma(\tau-\sigma, cw)$ ($\gamma(\cdot)$ は第一種不完全ガンマ関数) をもつ一般化 BFRY (Bertoin-Fujita-Royette and Yor) 過程を考えると, これを和が 1 になるように正規化した正規化 GBFRY 過程がこの double power-law をもち, 言語に非常によく当てはまること示されています[61].

コラム：単語の前処理について

ここまでの実行例では、空白で区切られたすべての文字列を単語として扱いましたが、実際にはテキストには数字や特殊文字(記号や罫線など)が入っており、これらをそれぞれ別の単語として扱うと、特に数字によって語彙が爆発的に増えてしまいます。また、英語では I'm や Alice's はそれぞれ、I 'm, Alice 's と分割して、I や Alice を単語とみなすのが望ましいでしょう。

そこで、一般にたとえば数字を#で置き換えたり、特殊文字を削除したり、aaaaaaのような同じ文字の連続を3文字に縮約してaaaとするなどの前処理を行います。こうすると、たとえば No.123 は no.###に、2021.3.28 は####.#.##に、Goooooool は goool になります。また、大文字と小文字の違いで別の単語と認識されないよう、ここに示したようにすべて小文字に直すこともよく行われる処理の一つです。^{*22}

前処理には決まった手続きがあるわけではありませんが、英語の場合は教科書 FSNLP [38]のサイトにある strip.sed^{*23} などを使うのが簡単です。日本語の場合は、後で紹介する mecab-NEologd のサイト^{*24} に正規化処理についてまとめられているほか、この処理を pip でインストールできるパッケージにした neologdn^{*25} が公開されています。表 3.2 は、こうした前処理を行った上で単語を数えたものです。

また、Python の spaCy や R の quanteda といったパッケージに前処理を任せることも可能です。

3.3[▽] 単語の統計的フレーズ化

*22 多くの言語で何を大文字にするかには規則性があるため、文脈を利用して単語の各文字を小文字から適切に大文字にする分類器を教師あり学習することは容易です。よって、小文字にしても本質的な情報量が落ちることはあまりありませんが、たとえば US を us として処理すると誤ることもありますので、すべて小文字化することにはリスクもあることに注意してください。

*23 <https://nlp.stanford.edu/fsnlp/statest/sed.strip>

*24 <https://github.com/neologd/mecab-ipadic-neologd/wiki/Regexp.ja>

*25 <https://github.com/ikegami-yukino/neologdn>
<https://yukinoi.hatenablog.com/entry/2015/10/11/205006>

3.1 節で述べたように、英語では “New York”, “with respect to”, 日本語では「基本的」, 「富嶽三十六景」のように、空白で区切られていても、実際は一つの単語として働く表現は多く存在します。また、「御樋代木奉曳式」*26のように特殊な言葉のために形態素解析が失敗して多くの単語に分割されてしまっている例もありますし、「高エネルギーリン酸化合物」「おジャ魔女どれみ 井」*27のような長い単語列を固有名詞として認識したい場合もあるでしょう。このような場合には、どうすればいいのでしょうか。

日本語の場合、最も簡単な方法は、多くの新語や固有名詞に対応し、人手で更新されている MeCab 用の辞書 mecab-ipadic-NEologd [62] *28 を用いることです。これには多くの固有表現が含まれており、一般的な語については一語として認識して形態素解析されます。より一般には、認識したい上記のような固有表現の正解データを多数、人手で用意しておけば、そこから教師あり学習によって、新しいテキストに対して該当箇所を固有表現として認識してくれる識別器を学習する**固有表現認識** (Named Entity Recognition, NER) とよばれる方法でフレーズを認識することができます (本書では正解データが必要な教師あり学習は基本的に扱いませんので、興味のある方は[10, Sec.5.5]などを参照してください)。

ただし、前者の辞書は人手で更新されているものですから、有名な固有名詞はカバーされているものの、「t細胞受容体」「ディビダーク式カンチレバー工法」のような専門的な名前にすべて対応することは、原理的にできません。また、後者の NER はあくまで正解データとして与えた表現およびそれに近い表現を認識するものですから、あらゆるフレーズに対応するには、膨大な量の「正解データ」を逐次作る必要があります、その際の基準も一意とは限らないために現実的ではありません。

しかし、よく考えると、“San” の後には必ず “Francisco” が続くのであれば、“San Francisco” を一つの単語として認識してもよいはずで、同様に、「宮内」

*26 御樋代木奉曳式 (みひしろぎほうえいしき) は、木曾の山中から切り出された御料木が用材として伊勢神宮に運ばれる儀式とのことです。

*27 この「どれみ」の例のように、未知の固有名詞に対しては形態素解析自体が失敗することがよく起こります。

*28 <https://github.com/neologd/mecab-ipadic-neologd>

の後には「庁」が続く確率が非常に高く、また「楽部」の前は「宮内庁」であることが非常に多いのであれば、「宮内庁」「宮内庁楽部」をそれぞれ単語とみなしてもよいでしょう。

このような考え方にに基づき、3.5.2節で説明する単語のベクトル化の方法である有名な Word2Vec を発明した Mikolov ら [63] は、空白で区切られた単語をそのままベクトル化するのではなく、文脈に応じて “and new_york city council ..” のように ‘ ’ によって単語を繋げて自動的にフレーズ化してから Word2Vec を計算する方法を示しました。^{*29}

この場合、単語 v と単語 w を繋げてフレーズにするかどうかは、次のスコアによって統計的に決定します。

$$(3.12) \quad \text{score}(v, w) = \frac{n(v, w) - \delta}{n(v) \times n(w)}$$

ここで、 $n(v, w)$ および $n(v), n(w)$ はそれぞれ単語バイグラムおよびユニグラムの頻度で、 δ は頻度 $n(v, w)$ の小さいバイグラムのスコアを下げるための割引き係数です。たとえば $v = \text{san}$, $w = \text{francisco}$ で、 $n(v, w) = 100$, $n(v) = 100$, $n(w) = 120$ だったとき^{*30}、 $\delta = 10$ であれば

$$\text{score}(\text{san}, \text{francisco}) = \frac{100 - 10}{100 \times 120} = 0.0075$$

になります。一方 $v = \text{read}$, $w = \text{books}$ で、 $n(v, w) = 100$, $n(v) = 1000$, $n(w) = 500$ であれば、

$$\text{score}(\text{read}, \text{books}) = \frac{100 - 10}{1000 \times 500} = 0.00018$$

となり、“read_books” より “san_francisco” のスコアの方が約 40 倍も大きいこととなります。よって、たとえば閾値を 0.001 として、隣りあった単語のスコアが閾値より大きい場合に単語を連結すれば、“san_francisco” のようなフレーズを自動的に認識することができます。

^{*29} この方法は、すでに 1992 年には IBM の Brown ら [64] によって提案されています。

^{*30} 「フランシスコ修道会」のような表現もありますから、少数ですが、francisco は他で用いられる場合もあります。

フランクリン・ルーズベルトの とつた [ニューディール政策] の [一環として]、1935年に連邦 [社会保障] 法が制定 [され]、失業保険と老齢年金が整備 [され] た
 一方の越後では、[守護代]・長尾氏により国内が統一 [され]、氏康により北関東から駆逐 [され] た [関東管領] の [上杉憲政] から [山内上杉] 家の家督と管領職を継承 [した] [上杉謙信] は北関東で氏康と対決し、信濃では北信勢力を後援 [して] 信玄と対決する二正面作戦を展開し [てい] た

図 3.5: 正規化自己相互情報量 (NPMI) に基づいて 2 単語を統計的にフレーズ化した例。認識されたフレーズを [] で示しました。学習には日本語版 text8 コーパスを使用しています。

ただし、Mikolov らの論文 [63] で提案されている式 (3.12) のスコアには注意が必要です*³¹。テキスト全体の長さを N とすると、バイグラム (v, w) およびユニグラム v, w の確率はそれぞれ

$$(3.13) \quad p(v, w) = \frac{n(v, w)}{N}, \quad p(v) = \frac{n(v)}{N}, \quad p(w) = \frac{n(w)}{N}$$

ですから、 $n(v, w) = N \cdot p(v, w)$, $n(v) = N \cdot p(v)$, $n(w) = N \cdot p(w)$ を式 (3.12) に代入すると、

$$(3.14) \quad \text{score}(v, w) = \frac{n(v, w) - \delta}{n(v) \times n(w)} = \frac{N \cdot p(v, w) - \delta}{N \cdot p(v) \times N \cdot p(w)} = \frac{p(v, w) - \delta/N}{p(v) \times p(w)} \cdot \frac{1}{N}$$

となり、このスコアは学習テキストの長さ N に依存します。さらに、割り引き係数 δ の大きさも N によって変えなければならない、ということがわかります。テキストを固定して閾値を探索するのであれば問題ありませんが*³²、スコアに絶対的な意味がないために探索が難しく、違うテキストを同じ基準で前処理することもできなくなってしまいます。この欠点はスコアを N 倍すれば解消しますが、その場合も δ/N の部分は残るため、適切な δ の設定は困難になります。

ここで $\delta=0$ としてみると、式 (3.12) は式 (3.14) より

*³¹ 以下の内容は、本書のオリジナルです。

*³² 実際には、フレーズ化を数パス繰り返す中で単語が結合され、 N の値も変わってきてしまいますので、もとの方法では適切な閾値の設定はさらに困難になります。

$$(3.15) \quad \text{score}(v, w) = \frac{p(v, w)}{p(v) \times p(w)} \cdot \frac{1}{N}$$

となり, これは5章でも説明する**自己相互情報量** (Pointwise Mutual Information, PMI)

$$(3.16) \quad \text{PMI}(v, w) = \log \frac{p(v, w)}{p(v)p(w)}$$

と, 本質的に同じ意味を持っていることがわかります. 式(3.16)の PMI は, v と w が共起する確率 $p(v, w)$ が, v と w が独立だった場合の確率 $p(v)p(w)$ と比べて何倍なのか (の対数) を表しています. これは統計的に v と w の相関を見るためには理想的な量で, 情報理論に基づいて1989年にベル研究所の Church ら [65] によって提案されました.

ただし, PMI には

- v や w が非常に低頻度で, $p(v)$ や $p(w)$ がきわめて小さい場合に式(3.16)が非常に大きくなってしまう
- 最大値・最小値が v, w によって異なり, 一定ではない

という欠点があります. そこで, 5章でも説明するように, PMI をその最大値である, $-\log p(v, w)$ で割って正規化した**正規化自己相互情報量** (Normalized PMI, NPMI) [66]

$$(3.17) \quad \text{NPMI}(v, w) = \log \frac{p(v, w)}{p(v)p(w)} \bigg/ \left(-\log p(v, w) \right)$$

を用いることを考えてみましょう. この NPMI は $p(v)$ や $p(w)$ が小さくても値がインフレせず,

$$(3.18) \quad -1 \leq \text{NPMI}(v, w) \leq 1$$

の値をとり, v と w が完全に相関しているとき 1, 完全に逆相関しているとき -1 の値をとるといふ, 大変よい性質を持っています. 式(3.17)は, 頻度 $n()$ を使えば

$$\begin{aligned}
 (3.19) \quad \text{NPMI}(v, w) &= \log \frac{p(v, w)}{p(v)p(w)} \bigg/ \left(-\log p(v, w) \right) \\
 &= \log \left(\frac{n(v, w)}{\mathcal{N}} \cdot \frac{\mathcal{N}}{n(v)} \cdot \frac{N}{n(w)} \right) \bigg/ \left(-\log \frac{n(v, w)}{N} \right) \\
 (3.20) \quad &= \frac{\log N + \log n(v, w) - \log n(v) - \log n(w)}{\log N - \log n(v, w)}
 \end{aligned}$$

と表すことができます。このとき、式(3.12)のヒューリスティックな割り引き係数 δ は不要になっていることに注意してください。この計算は、サポートサイトにある `phraser.py` を使って、

```
% phraser.py ja.text8.txt output 4
```

のように実行すると簡単に行うことができます。

この NPMI が閾値、たとえば 0.5 以上になった 2 単語をフレーズとみなして日本語 text8 に適用した例を図 3.5 に示しました。ほとんどの語はそのままで、NER を使わなくても教師なしで、「ニューディール政策」「上杉謙信」「され」などがフレーズとして自動的に認識されていることがわかります。

この方法は隣り合う 2 単語を結合するものですが、この出力を入力として再度フレーズ化すれば、最大で 4 単語までのフレーズが得られます。同様にしてこのフレーズ認識を n パス動かせば、 2^n 語までのフレーズを計算することができます。表 3.3 に、こうして 4 パス (=最大 16 単語まで) を日本語 text8 コーパスに適用して得られたフレーズの例を示しました。

表からわかるように、「である」「となっていた」「キリスト教徒」といった、日本語によくあるフレーズや固有名詞が教師なしで正しく認識されており、時には「福島第一原子力発電所」といった長い固有名詞も自動的にフレーズ化されていることがわかります。テキストに対してこうした前処理を行うことで、意味的内容をより正しく反映させることができます。

しかし一方で、たとえば「逆転写酵素」が「酵素」の一種だという情報は失われてしまうため、こうしたフレーズ化は、フレーズの頻度が充分高い場合にのみ行った方がよいでしょう。上記の `phraser.py` では、デフォルトではバイグラム頻度が 10 以上の場合にのみフレーズ化を行う設定になっており、閾値も含めてオプションで変更することができます。使い方は、

表 3.3: 日本語版 text8 から NPMI に基づいて統計的に認識されたフレーズとその頻度。
 は、MeCab によるもとの単語区切りを表しています。4 パスで計算したため、最大で 16
 単語までがフレーズの候補になっています。NPMI の閾値は 0.5 としました。

頻度	フレーズ
107101	し_た
70074	で_ある
48360	さ_れ_た
31219	て_いる
28473	し_て_いる
23371	さ_れ
:	
1122	第_二_次_世_界_大_戦
1094	と_な_っ_て_い_た
1092	最_終_的
1088	基_本_的
:	
100	に_つ_い_て_述_べ_る
100	キ_リ_ス_ト_教_徒
100	陸_上_競_技_選_手
100	ゆ_う_ち_ょ_銀_行
:	
45	全_国_高_等_学_校_野_球_選_手_権_大_会
37	お_ジャ_魔_女_ど_れ_み
32	福_島_第_一_原_子_力_発_電_所_事_故
27	ほ_っ_か_ほ_っ_か_亭
26	福_島_第_一_原_子_力_発_電_所
23	連_合_国_軍_最_高_司_令_官_総_司_令_部
19	学_研_奈_良_登_美_ヶ_丘_駅
16	国_際_連_合_安_全_保_障_理_事_会
10	ソ_ビ_エ_ト_社_会_主_義_共_和_国_連_邦
10	自_転_車_競_技_選_手_権_大_会
10	工_学_部_機_械_工_学_科
10	合_衆_国_最_高_裁_判_所
10	衆_議_院_予_算_委_員_会
10	大_学_院_工_学_研_究_科

```
% phraser.py
```

をそのまま実行するか、スクリプトの中身を読んでみてください。

3.4 単語 n グラム言語モデル

単語についても、2.5 節の文字の場合と同じように単語 n グラム言語モデルを考えることができます。現在は文の確率を計算するには、LSTM や Transformer などの深層学習モデルを利用するのが一般的ですが*33、本書で扱うさまざまな統計モデルの基礎になりますので、単純で理由のわかっている、この最も基本的なモデルについて考えていくことにしましょう。

単語が直前の単語だけに依存するバイグラム言語モデルを考えるとすると、文 $s = \text{「銀河 鉄道 の 夜」}$ の確率は式(2.45)と同様に、

$$(3.21) \quad p(s) = p(\text{銀河}|\wedge)p(\text{鉄道}|\text{銀河})p(\text{の}|\text{鉄道})p(\text{夜}|\text{の})p(\$|\text{夜})$$

と表すことができます。ここで $\wedge, \$$ はそれぞれ、2.5 節で導入した文頭および文末を表す特別な単語です。

このバイグラム頻度を数えるには、Python ではたとえば、3.1 節で単語に分割したテキストを使って

```
from collections import defaultdict
def parse (file):
    EOS = "_EOS_"
    freq = {}
    with open (file, 'r') as fh:
        for line in fh:
            words = line.rstrip('\n').split() # 空白で分割
            words.insert (0, EOS); words.append (EOS)
            T = len(words) # ↑文頭/文末文字を追加
            for t in range(T-1):
                w = words[t]
                v = words[t+1]
                if not (w in freq): # freq[w] がなければ作成
                    freq[w] = defaultdict(int)
                freq[w][v] += 1 # 頻度に 1 を加える
    return freq
```

*33 ただし、こうした深層学習モデルがどのように単語を予測しているのかは、まだ充分わかっていないのが現状です。

のように関数を定義してから,

```
freq = parse ("ginga.split.txt")
freq
⇒ {'_EOS_': defaultdict(int,
    {' 銀河': 1,
     '一': 1, ...
    },
    {' 銀河': defaultdict(int,
    {' 鉄道': 2,
     ' 帯': 1,
     ' を': 2,
     ' は': 1,
     ' の': 11,
     ' が': 2,
     ' ステーション': 5,
     ' だ': 1}), ...
    }
```

のように実行すれば, $\text{freq}[w][v]$ にバイグラムの頻度 $n(w, v)$ を格納することができます. たとえば上の場合は,

```
freq["銀河"]["ステーション"]
⇒ 5
```

のようになります.

式(3.21)のそれぞれのバイグラムの条件つき確率は, 頻度を表す関数 $n()$ を使って, 最も簡単には式(2.4)のような最尤推定値

$$(3.22) \quad \hat{p}(v|w) = \frac{n(w, v)}{\sum_v n(w, v)} = \frac{n(w, v)}{n(w)}$$

で求めることができるのでした. ここで $n(w, v)$ は単語バイグラム wv の出現した頻度, $n(w) = \sum_v n(w, v)$ は単語 w の頻度です.

ただし, 式(3.22)を使った単語バイグラム確率の計算には, すぐに問題があることがわかります. ここまでにみたように, Zipfの法則から頻度のほとんどは一部の単語に集中しており, 大部分の単語は低頻度なのでした. すると, たとえば「銀河」の次に「旅行」が現れる確率は, 式(3.22)に従えば

$$(3.23) \quad p(\text{旅行}|\text{銀河}) = \frac{n(\text{銀河}, \text{旅行})}{n(\text{銀河})}$$

となりますが、『銀河鉄道の夜』の中では $n(\text{銀河}, \text{旅行})=0$ なので、

$$(3.24) \quad p(\text{旅行}|\text{銀河}) = \frac{n(\text{銀河}, \text{旅行})}{n(\text{銀河})} = \frac{0}{n(\text{銀河})} = 0$$

となり、「銀河旅行」の確率は式(2.20)の確率の連鎖則から、

$$(3.25) \quad p(\text{銀河 旅行}) = \underbrace{p(\text{旅行}|\text{銀河})}_{=0} \cdot p(\text{銀河}) = 0$$

となり、0になってしまいます。こうした問題を、2.5.2節で述べたように**ゼロ頻度問題**と呼びます。

文字の場合は種類が少ないため、文字バイグラムではゼロ頻度問題はあまり深刻になりませんでした*³⁴、単語は種類が多いため、バイグラムでも、ほとんどの確率が0になってしまうという問題が生じます。

ゼロ頻度問題を避けるために、2章の文字 n グラムモデルでは式(2.51)の**加算平滑化**という手法を用いました。語彙の総数を V とおくと、これは式(3.22)の代わりに、すべての頻度に小さな値 α を足して

$$(3.26) \quad p(v|w) = \frac{n(w, v) + \alpha}{\sum_{v=1}^V (n(w, v) + \alpha)} = \frac{n(w, v) + \alpha}{n(w) + V\alpha}$$

とするものです。こうすると、 $n(w, v)=0$ の場合でも、確率は

$$p(v|w) = \frac{\alpha}{n(w) + V\alpha}$$

となり、出現しない語にも、 α に比例する一定の確率を割り当てることができ
ます。

ところが、単語の場合は式(3.26)を使ってもまだ問題があることがわかります。たとえば、単語 w と v は「銀河 鉄道」のように常に wv の形で出現しており、 $n(w, v)=10$ 、 $n(v)=n(w)=10$ だったとしましょう。 $\alpha=0.01$ 、 $V=10000$ とすると、式(3.26)は

*34 日本語や中国語などの漢字圏では、文字の種類も非常に多いため、文字バイグラムの場合でもゼロ頻度問題は深刻になります。

$$(3.27) \quad p(v|w) = \frac{n(w, v) + \alpha}{n(w) + V\alpha} = \frac{10 + 0.01}{10 + 10000 \cdot 0.01} = \frac{10.01}{110} = 0.091$$

となります。つまり、 w の後にはつねに v が続くにもかかわらず、 $p(v|w)$ はたった0.09にしかならず、それ以外の確率が $1 - 0.091 = 0.909$ で91%もある、という結果になってしまうのです。

これはもちろん、「現れなかった単語も含め、すべての単語の頻度に同じ α を足す」ということが原因です。「群」のような頻度の高い語も、「パセリ」のような頻度の低い語も同じ α が足されるため、たまたま「銀河」の後に続いたことがなければ「銀河群」と「銀河パセリ」が同じ確率を持つことになり、それらの確率の総和である式(3.26)の $V\alpha$ の部分が非常に大きくなってしまいますからです。実際に、式(3.26)を変形すると、

$$(3.28) \quad \begin{aligned} p(v|w) &= \frac{n(v, w) + \alpha}{n(w) + V\alpha} = \frac{n(v, w)}{n(w) + V\alpha} + \frac{\alpha}{n(w) + V\alpha} \\ &= \frac{n(w)}{n(w) + V\alpha} \cdot \frac{n(v, w)}{n(w)} + \frac{V\alpha}{n(w) + V\alpha} \cdot \frac{1}{V} \\ &= \lambda \cdot \hat{p}(v|w) + (1 - \lambda) \cdot \frac{1}{V} \quad \left(\lambda = \frac{n(w)}{n(w) + V\alpha} \right) \end{aligned}$$

となります。これから、式(3.26)の確率は式(3.22)の最尤推定値 $\hat{p}(v|w)$ と、**すべての単語に一律な確率を与える** $p_0(v) = \frac{1}{V}$ を比率 $\lambda : (1 - \lambda)$ で補間した確率になっていることがわかります。^{*35}これは、 $\lambda = \frac{n(w)}{n(w) + V\alpha}$ が小さい、すなわ

ち $n(w)$ が小さいほど、確率の推定値が $\frac{1}{V}$ に近づくことを意味しています。つまり、「銀河」自体の出現回数が少なければ、「銀河の」も「銀河パセリ」も同じくらい出現しやすいと仮定していることになるわけです。

言うまでもなく、これは誤りです。文脈となる語 w と予測したい語 v を分けて考えると、 v に対応する α の値は、大まかには v の確率 $p(v)$ に比例する値 α_v とするのがよいでしょう。ただしよく考えると、単に $p(v)$ に比例するのではな

^{*35} このように、複数の確率分布（この場合は $\hat{p}(v|w)$ と $p(v) = 1/V$ ）を重み $(\lambda$ と $1 - \lambda)$ つきで混ぜ合わせたモデルを、**混合モデル**といいます。

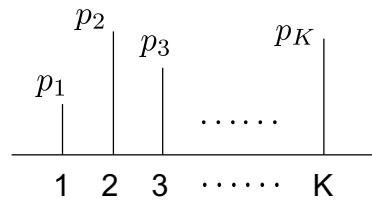


図 3.6: K 次元の多項分布 \mathbf{p} . 確率の和は $\sum_{k=1}^K p_k = 1$ になっています.

く、たとえ同じ頻度でも、多くの語の後に続きやすい「前」のような語の α_v は大きく、他の語の後に続きにくい「あの」のような語の α_v は小さくすべきだと考えられます。^{*36} そして、決めるべきパラメータは $\alpha_1, \alpha_2, \dots, \alpha_V$ で、この場合 10,000 個もあります。この $\alpha_1, \alpha_2, \dots, \alpha_V$ はどうやって決めればよいのでしょうか。

こういった複雑な問題を解くには、今までのような発見的な方法では限界があります。そこで、より原理的に考えてみることにしましょう。

3.4.1 ディリクレ分布

これまで、式(3.26)のような確率の由来については考えず、単に「頻度やそれに値を足したものを和が1になるように正規化する」ことしか行ってきませんでした。式(3.26)のような確率は、すべての単語 $v=1, 2, \dots, V$ について考えると総和が1になる**確率分布**ですから、そもそも確率分布を生み出すことのできる**確率モデル**について考えてみましょう。

いま、図 3.6 のような K 次元の確率分布を

$$(3.29) \quad \mathbf{p} = (p_1, p_2, \dots, p_K) \quad \left(p_1, p_2, \dots, p_K \geq 0, \sum_{k=1}^K p_k = 1 \right)$$

とします。 K は次元の数で、言語の場合には K は実際には 10,000 や 100,000 といった大きな値になります。確率 p_1, p_2, \dots, p_K を直接指定するこの確率分布は、

^{*36} 『銀河鉄道の夜』では、「前」と「あの」の出現頻度はどちらも 38 回で同一です。

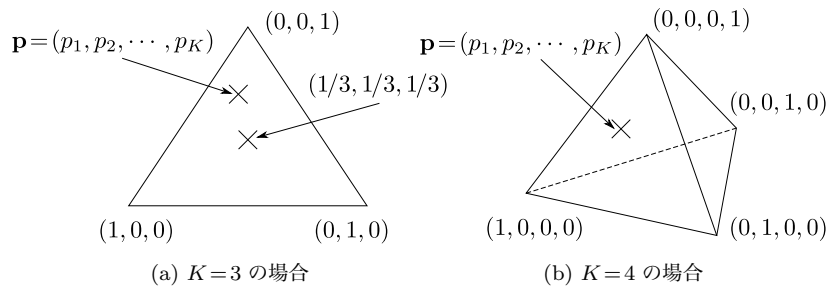


図 3.7: 単体とその上の多項分布 \mathbf{p} . 単体の各端点は、そのカテゴリだけが確率 1 で出る確率分布になっています。

もっとも基本的な分布で、正確には**多項分布**または**離散分布**とよびます。^{*37}

たとえば $K=3$ のとき、

$$\mathbf{p} = (0.7, 0.1, 0.2)$$

は \mathbf{p} の 1 つの例です。より一般に $\mathbf{p} = (p_1, p_2, p_3)$ は、これをベクトルとみなせば、図 3.7(a) のように正三角形の内部にあると考えることができます。こうした図形を、**単体** (simplex) といいます。 $\mathbf{p} = (1, 0, 0)$, $(0, 1, 0)$, $(0, 0, 1)$ の各分布はそれぞれ、この単体の 3 つの角に対応し、 $\mathbf{p} = \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right)$ は単体の中心に対応しています。なお、 $K=4$ の場合は、単体は図 3.7(b) のような正四面体になります。

$\mathbf{p} = (p_1, p_2, \dots, p_K)$ が与えられたとき、観測値の確率は p_k の積で簡単に求めることができます。たとえば、

$$\mathbf{p} = (p_1, p_2, p_3, p_4)$$

のとき、 \mathbf{p} に従ってランダムに選んだ結果が

$$Y = (4, 2, 1, 1, 2)$$

だったとすれば^{*38}、 Y の確率は

^{*37} 多項分布は、本来は二項分布の拡張で \mathbf{p} から N 回サンプルした結果の確率を与えるものです。ただし、離散的な分布にはポアソン分布など他の分布も含まれるため、それと区別して $N=1$ の場合も多項分布ということがあり、本書でもこの表記を用います。

^{*38} ここでは、 Y にテキストのように順番を考えています。順番を考えない場合は、可能な組み

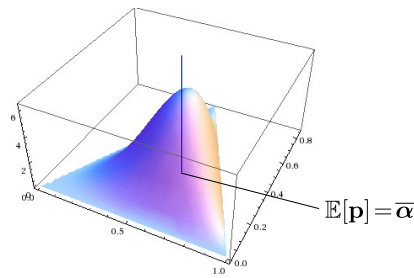


図 3.8: 単体上のディリクレ分布 $\text{Dir}(\boldsymbol{\alpha})$ とその期待値 $\mathbb{E}[\mathbf{p}] = \bar{\boldsymbol{\alpha}}$.

$$(3.30) \quad \begin{aligned} p(Y|\mathbf{p}) &= p_4 \cdot p_2 \cdot p_1 \cdot p_1 \cdot p_2 = p_1^2 \cdot p_2^2 \cdot p_3^0 \cdot p_4^1 \\ &= \prod_{k=1}^4 p_k^{n_k} \end{aligned}$$

となるわけです。当たり前ですね。ここで n_k は Y の中で k が出た回数で、この例では $(n_1, n_2, n_3, n_4) = (2, 2, 0, 1)$ となります。

こうした \mathbf{p} を生成する単体上の最も簡単な分布として、次の式で表される**ディリクレ分布** (Dirichlet distribution) があります。^{*39}

$$(3.31) \quad \text{Dir}(\boldsymbol{\alpha}) \propto \prod_{k=1}^K p_k^{\alpha_k - 1} \quad (\text{ディリクレ分布 (簡易版)})$$

図 3.8 に、この分布の形を示しました。 $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_K)$ はこの分布のパラメータで、 $\alpha_k \geq 0$ ($k = 1, 2, \dots, K$) は非負の実数です^{*40}。 \mathbf{p} はそれ自身確率分布ですから、ディリクレ分布は「確率分布を生成する確率分布」ということになります。ディリクレ分布から生成される \mathbf{p} の期待値 $\mathbb{E}[\mathbf{p}]$ は、 $\boldsymbol{\alpha}$ を和が 1 になるように正規化した

合わせの総数である多項係数 $\binom{5}{2\ 2\ 0\ 1} = \frac{5!}{2!2!0!1!}$ がかかることになりましたが、この係数はパラメータ \mathbf{p} を含んでいませんので、 \mathbf{p} に関する推定値は同じになります。

^{*39} この名前は、正規化定数となる式(3.40)の積分を示した、旧フランス領で生まれたドイツの数学者 Lejeune Dirichlet (1805–1859) にちなむものです [67]。

^{*40} $\alpha_k = 0$ のときは、対応する p_k はつねに 0 であると約束します。

(3.32)

$$\mathbb{E}[\mathbf{p}] = \bar{\alpha} = \left(\frac{\alpha_1}{\sum_k \alpha_k}, \frac{\alpha_2}{\sum_k \alpha_k}, \dots, \frac{\alpha_K}{\sum_k \alpha_k} \right) \quad (\text{ディリクレ分布の期待値})$$

となります。

ディリクレ分布からのサンプルは、Python では `numpy.random.dirichlet()` で作ることができます。または、ガンマ分布 $\text{Ga}(\alpha_k, 1)$ からのサンプル^{*41}

$$(3.33) \quad \gamma_k \sim \text{Ga}(\alpha_k, 1) \quad (k = 1, 2, \dots, K)$$

が得られれば、それを和が1になるように正規化することで、ディリクレ分布からのサンプル

$$(3.34) \quad \mathbf{p} = \frac{1}{\sum_k \gamma_k} (\gamma_1, \gamma_2, \dots, \gamma_K) \sim \text{Dir}(\alpha_1, \alpha_2, \dots, \alpha_K)$$

を計算することができます。サポートページにある `dirichlet.py` を、たとえば

```
% dirichlet.py 0.5 5
```

のように実行すれば、5次元でパラメータがすべて0.5のディリクレ分布

$$\text{Dir}(0.5, 0.5, 0.5, 0.5, 0.5)$$

から生成されたランダムな多項分布 \mathbf{p} をプロットすることができますので、試してみましょう。図3.9に、いくつかの α の値について、こうして $\text{Dir}(\alpha)$ から生成した多項分布 \mathbf{p} の例を示しました。 $\alpha_k = 1$ の場合は \mathbf{p} は自由な形になりますが、 $\alpha_k > 1$ の場合は \mathbf{p} は期待値である $\bar{\alpha}$ に近く、 $\alpha_k < 1$ の場合は \mathbf{p} は少数の p_k の値だけが大きく、他がほとんど0に近い疎な分布となることがわかります。言語の場合は、ほとんどはこの場合に対応しています。

ガンマ関数とディリクレ分布

^{*41} この1はスケールを表しているだけですので、同一であれば別の値でも問題ありません。ガンマ分布からのサンプルは、Python では `numpy.random.gamma()` で作ることができます。自分で $[0, 1)$ の乱数からガンマ関数に従う乱数を計算する方法は複雑ですので、必要な方は乱数生成の専門書[68]を参照してください。サポートページに、筆者が使っているC言語による実装 `gamma.{c,h}` が置いてあります。

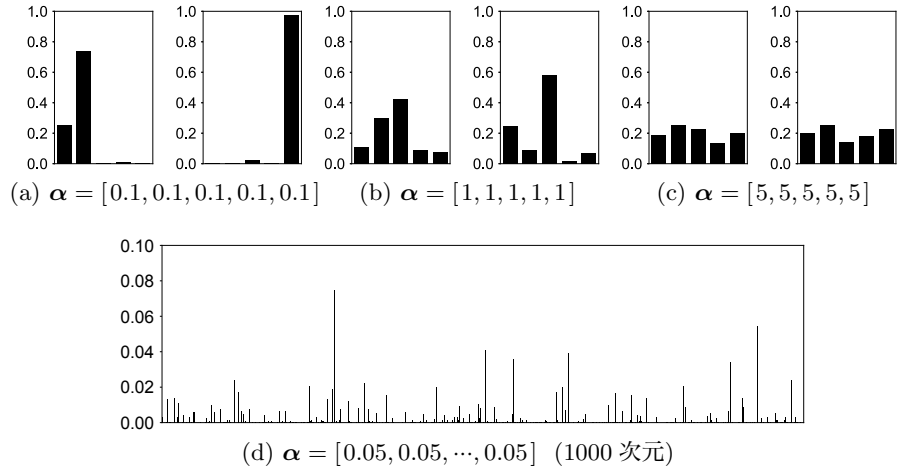


図 3.9: ディリクレ分布 $\text{Dir}(\alpha)$ からランダムにサンプルした多項分布 \mathbf{p} の例.

式(3.31)では比例 α を使ってディリクレ分布を示しましたが, 正確に書くと, ディリクレ分布の定義は

$$(3.35) \quad \text{Dir}(\alpha) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^K p_k^{\alpha_k - 1} \quad (\text{ディリクレ分布 (正式版)})$$

となります. この定義式(3.35)は一見いかつい形をしていますが, ガンマ関数 $\Gamma(x)$ を含む前半部分は, 可能な \mathbf{p} 全体についての積分を 1 にするための正規化定数で, 本質的には, ディリクレ分布は式(3.31)で表される簡単な分布であることに注意してください. $\Gamma(x)$ ($x \in \mathbb{R}$) は階乗関数 $x! = x(x-1)(x-2) \dots 1$ の連続値への一般化とみることができる関数で,

$$(3.36) \quad \Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt \quad (\text{ガンマ関数})$$

で定義されます. 式(3.36)を部分積分することで,

$$(3.37) \quad \Gamma(x+1) = x\Gamma(x)$$

を示すことができますので, 確かめてみましょう (→演習問題1). 式(3.37)から,

x が整数であれば

$$(3.38) \quad \begin{aligned} \Gamma(x) &= (x-1)\Gamma(x-1) = (x-1)(x-2)\Gamma(x-2) \\ &= (x-1)(x-2)\cdots 2 \cdot 1 \\ &= (x-1)! \end{aligned}$$

となります。 x が整数でない場合は、 x を超えない最大の整数を $n = \lfloor x \rfloor$ とおき、 $x = n + \alpha$ と書けば

$$(3.39) \quad \begin{aligned} \Gamma(x) &= (n+\alpha-1)\Gamma(n+\alpha-1) \\ &= (n+\alpha-1)(n+\alpha-2)\cdots (\alpha+1) \cdot \alpha \\ &= (x-1)(x-2)\cdots (\alpha+1) \cdot \alpha \end{aligned}$$

となり、確かにいずれの場合も、 $\Gamma(x)$ は階乗 $(x-1)!$ の一般化になっていることがわかります。

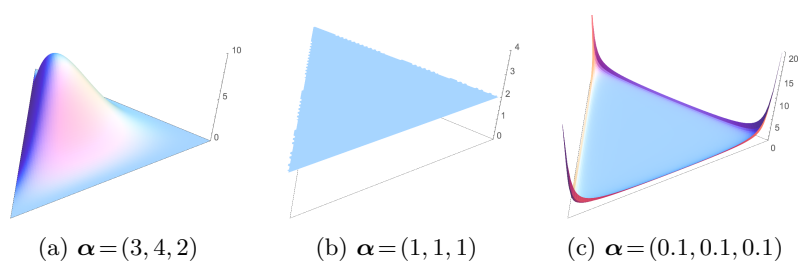
式(3.35)の正規化定数は、この $\Gamma(x)$ を使って

$$(3.40) \quad \begin{aligned} \int \prod_{k=1}^K p_k^{\alpha_k-1} d\mathbf{p} &= \int_0^1 \int_0^1 \cdots \int_0^1 p_1^{\alpha_1-1} p_2^{\alpha_2-1} \cdots p_K^{\alpha_K-1} dp_1 dp_2 \cdots dp_K \\ &= \frac{\prod_k \Gamma(\alpha_k)}{\Gamma(\sum_k \alpha_k)} \end{aligned}$$

の積分から得られるものです。この積分は、直感的には、図3.8のような単体上の曲面の下の体積を求めることに相当しています。式(3.40)の積分は大学教養範囲の数学(解析)を必要としますので、詳しくは付録Aを参照してください。上のように、以下本書では、単体上の積分 $\int_0^1 \int_0^1 \cdots \int_0^1 dp_1 dp_2 \cdots dp_K$ を略して $\int d\mathbf{p}$ と書くことにします。^{*42}

ディリクレ分布は、この単体上の確率分布です。ディリクレ分布はパラメータ α の値によって、さまざまな形をとります。図3.10に、 α の値とそれに対応するディリクレ分布の形を示しました。式(3.31)からわかるように、 $\alpha = (1, 1, \dots, 1)$

^{*42} 正確には、 p_1, p_2, \dots, p_K は独立ではなく $\sum_{k=1}^K p_k = 1$ という制約がありますが、式を簡単にするため、上ではこの制約は省略しています。

図 3.10: さまざまな α によるディリクレ分布 $\text{Dir}(\alpha)$ の密度関数.

のとき

$$(3.41) \quad \text{Dir}(\alpha) \propto \prod_{k=1}^K p_k^{\alpha_k - 1} = \prod_{k=1}^K p_k^0 = 1$$

ですから、確率分布は図 3.10(b) のように単体上の一様分布になります。 $\alpha_k > 1$ のときは図 3.10(a) のように上に凸な、 $\alpha_k < 1$ のときは図 3.10(c) のように下に凸な分布になっており、それぞれ期待値に近い、あるいはどれかの p_k だけが大きい疎な \mathbf{p} を生み出すことになります。

サイコロ工場とディリクレ分布 図 3.6 のような多項分布 $\mathbf{p} = (p_1, p_2, \dots, p_K)$ は、「ゆがんだ K 面サイコロ」のようなものだと考えることもできます。サイコロの各面 k の大きさが確率 p_k に対応しており、このサイコロを振ると、 p_k に比例した確率で値 k が出るというわけです。

このとき式 (3.31) のディリクレ分布は図 3.11 のように、様々な K 面体サイコロ \mathbf{p} を生産する「サイコロ工場」のようなものだと考えてもいいでしょう [69]。この工場からは様々なゆがみを持った K 面体サイコロ \mathbf{p} が生産されますが、工場には特有の傾向があり、 \mathbf{p} がほとんど一様なサイコロを生産する工場や、 \mathbf{p} の値が大きく異なるサイコロを生産する工場もあります。ある工場から生産されたサイコロ \mathbf{p} を多く集めて計測すれば、この工場の持つ傾向 α を推測することができます。

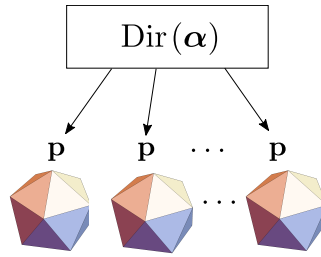


図 3.11: 「サイコロ工場」としてのディリクレ分布. α の値によって, それぞれ異なる歪みを持った多面体サイコロ \mathbf{p} が生産されます.

3.4.2 ディリクレ分布と多項分布

\mathbf{p} の分布として, 式(3.31)を最も簡単な分布だとしたのには理由があります. いま $K=3$ で, $\mathbf{p}=(p_1, p_2, p_3)$ の分布が $\text{Dir}(\alpha)$ に従っているとしましょう.

このとき, 多項分布 \mathbf{p} からランダムに値をサンプルすると (上のサイコロのメタファーでは, サイコロを振ると), $Y=(2, 3, 1, 2, 2)$ が観測されたとします. このとき, \mathbf{p} はどんな分布だと推測できるでしょうか.

Y が観測されたときの \mathbf{p} の分布 $p(\mathbf{p}|Y)$ は, 式(2.30)のベイズの定理から

$$(3.42) \quad p(\mathbf{p}|Y) \propto p(Y|\mathbf{p})p(\mathbf{p})$$

と書くことができます. $p(\mathbf{p})$ は \mathbf{p} の事前分布 $\text{Dir}(\mathbf{p}|\alpha)$ ですから,

$$(3.43) \quad p(\mathbf{p}) \propto \prod_{k=1}^3 p_k^{\alpha_k-1}$$

です. \mathbf{p} から Y が得られる確率 $p(Y|\mathbf{p})$ は, 式(3.30)と同様に

$$(3.44) \quad p(Y|\mathbf{p}) = p_2 \cdot p_3 \cdot p_1 \cdot p_2 \cdot p_2 = p_1^1 \cdot p_2^3 \cdot p_3^1 = \prod_{k=1}^3 p_k^{n_k}$$

となります. ここで n_k は Y の中で k が出た回数で, $n_1=1, n_2=3, n_3=1$ です. 式(3.43)と式(3.44)をあわせると, Y が与えられたときの \mathbf{p} の分布は,

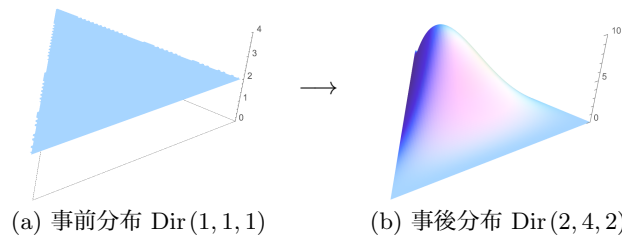


図 3.12: デイリクレ事前分布と事後分布の例.

$$(3.45) \quad p(\mathbf{p}|Y) \propto p(Y|\mathbf{p})p(\mathbf{p}) = \prod_{k=1}^3 p_k^{n_k} \times \prod_{k=1}^3 p_k^{\alpha_k-1} = \prod_{k=1}^3 p_k^{\alpha_k+n_k-1}$$

となり, これは $\alpha + \mathbf{n}$ を新しいパラメータとするデイリクレ分布

$$(3.46) \quad \begin{aligned} & \text{Dir}(\alpha_1 + n_1, \alpha_2 + n_2, \alpha_3 + n_3) \\ & = \text{Dir}(\alpha + \mathbf{n}) \quad (\mathbf{n} = (n_1, n_2, n_3)) \quad (\text{デイリクレ事後分布}) \end{aligned}$$

であることがわかります.

つまり, \mathbf{p} の事前分布を式(3.31)のデイリクレ分布 $\text{Dir}(\alpha)$ にすると, Y を観測した後の事後分布は新しいパラメータ $\alpha + \mathbf{n}$ を持つ, 同じデイリクレ分布 $\text{Dir}(\alpha + \mathbf{n})$ になるわけです.*43 この様子を, 図 3.12 に示しました. \mathbf{p} の事前分布が図 3.12(a) のように一様分布, すなわち $\text{Dir}(1, 1, 1)$ だったとすると, 式(3.45) から Y を観測した後の \mathbf{p} の事後分布は $\text{Dir}(1 + 1, 1 + 3, 1 + 1) = \text{Dir}(2, 4, 2)$ となり, 図 3.12(b) のような分布になります.

デイリクレ分布 $\text{Dir}(\alpha)$ の期待値は, α を和が 1 になるように正規化した

$$(3.47) \quad \bar{\alpha} = \left(\frac{\alpha_1}{\sum_k \alpha_k}, \frac{\alpha_2}{\sum_k \alpha_k}, \dots, \frac{\alpha_K}{\sum_k \alpha_k} \right)$$

でしたから, この事後分布の期待値は

*43 こうした性質を, 確率分布の**共役**(きょうやく)性といいます. 多項分布とデイリクレ分布は, 共役な分布です. 多項分布に共役な分布は他にもありますが, たとえばニューラルネットによく使われている Softmax 関数(130 ページ)は共役ではなく, このように簡単に事後分布を計算することはできません. なお, 本来の漢字は「共軛」で, 軛(くびき)とは, 馬車などで動物の首を結びつけて一緒に動かすための棒のことです.

$$\left(\frac{2}{2+4+2}, \frac{4}{2+4+2}, \frac{2}{2+4+2} \right) = (0.25, 0.5, 0.25)$$

になります。

すなわち一般に、 K 次元の多項分布 \mathbf{p} からの観測値 Y としてそれぞれの値 k ($k=1, 2, \dots, K$) が n_k 回観測されたとき、 \mathbf{p} の事前分布を $\mathbf{p} \sim \text{Dir}(\boldsymbol{\alpha})$ とすれば、事後分布は

$$(3.48) \quad \mathbf{p} | Y \sim \text{Dir}(\boldsymbol{\alpha} + \mathbf{n})$$

になり、その期待値は

$$(3.49) \quad \mathbb{E}[p_k | Y] = \frac{\alpha_k + n_k}{\sum_k (\alpha_k + n_k)} = \frac{\alpha_k + n_k}{\alpha + N} \quad (\text{ディリクレ平滑化})$$

となります。ここで $\alpha = \sum_k \alpha_k$, $N = \sum_k n_k$ とおきました。

ディリクレ分布に基づき、 k 番目のカテゴリの頻度に α_k を足すこの方法を、**ディリクレ平滑化**といいます。これから、ディリクレ分布のパラメータ α_k は、 **k 番目のカテゴリに事前に足す「仮想的な頻度」という意味を持っている**、ということがわかります。

式(3.49)のディリクレ平滑化と式(3.26)の加算平滑化を比べると、加算平滑化は、 \mathbf{p} に均一なディリクレ分布

$$(3.50) \quad \mathbf{p} \sim \text{Dir}(\alpha, \alpha, \dots, \alpha)$$

すなわち、 $\alpha_k = \alpha$ とした場合のディリクレ平滑化と等しいことがわかります。逆にいえば、式(3.49)の推定値は

$$(3.51) \quad \mathbf{p} \sim \text{Dir}(\alpha_1, \alpha_2, \dots, \alpha_K)$$

と α_k を異なる値にした場合のベイズ推定値になっているわけです。

ハイパーパラメータ α の推定

加算平滑化は \mathbf{p} に均一なディリクレ分布を仮定するモデルと等価だということがわかりましたが、式(3.51)のように \mathbf{p} をディリクレ分布でモデル化するこ

との意義は何でしょうか. それは, これによって**データから α を推定できる**ということなのです.

われわれは, \mathbf{p} 自体を直接観測することはできません. しかし, \mathbf{p} からサンプルされた各カテゴリの頻度である \mathbf{n} が毎回ほぼ均一な値であれば, α は図 3.9(c) のようにほぼ同じで, 1 より大きい値だと推定できるでしょう. いっぽう, \mathbf{n} の値がどれかの次元に偏っていれば, α は図 3.9(a) のように 1 より小さいと考えられます. また, どれかの n_k の値がいつも大きくなっていけば, 対応する p_k の値, したがって α_k も大きいと考えられます. 3.4.1 節のサイコロ工場のメタファーを使えば, 生産されたサイコロ \mathbf{p} の各面の面積を正確に測定することができなくても, \mathbf{p} をランダムに転がした結果の集計 \mathbf{n} さえあれば, 工場のパラメータ α を推定できるはずだ, というわけです.

そこで, こうした \mathbf{n} が複数観測されたとき^{*44}, そこから共通する α を推定する問題を考えてみましょう. すなわち, データ

$$(3.52) \quad D = \{\mathbf{n}_1, \mathbf{n}_2, \dots, \mathbf{n}_D\}$$

から, α を求めることを考えてみます.

D の各頻度ベクトル \mathbf{n} は, これまでの議論から, 次のようにして生成されたと考えることができます.

ステップ 1. 多項分布 $\mathbf{p} \sim \text{Dir}(\alpha)$ をサンプル.

ステップ 2. For $i = 1, 2, \dots, N$,

- $k \sim \mathbf{p}$ をサンプル.
- $n_k = n_k + 1$.

こうした, データが生成された過程のモデルを**生成モデル**といいます.^{*45}

この生成モデルに従えば, 多項分布 \mathbf{p} とそこからの観測値 \mathbf{n} が生成される確

^{*44} 3.4.2 節で述べたように, $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_K)$ は各カテゴリに事前に足す仮想的な頻度という意味を持っています. よって, \mathbf{n} が 1 つしかなければ, α_k が大きいほど多くの観測値があることになり, 最適化すると α_k が無限に大きくなってしまいます. いっぽう, 複数の \mathbf{n} があれば, ある α_k を大きくすると他の \mathbf{n} の確率が下がることになり, トレードオフがあるために最適な α を決定することができます.

^{*45} 分野によっては, Data Generating Process (DGP, データ生成過程) とよぶこともあります.

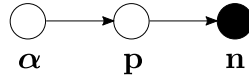


図 3.13: デリクレ多項分布のグラフィカルモデル.

率は、次のように表すことができます.

$$(3.53) \quad p(\mathbf{n}, \mathbf{p} | \boldsymbol{\alpha}) = p(\mathbf{n} | \mathbf{p}) p(\mathbf{p} | \boldsymbol{\alpha})$$

第1項は上のステップ2に、第2項は上のステップ1に対応しています. この関係をわかりやすくするために、模式的に図3.13のように表してみましょう. こうした図を**グラフィカルモデル**といい、○で表される変数間の依存関係が→で表されています. ○は未知の変数を、●は観測された変数すなわちデータを表します. ここでは $\boldsymbol{\alpha}$ から \mathbf{p} が生成され、 \mathbf{p} から \mathbf{n} が生成されるため、 $\boldsymbol{\alpha} \rightarrow \mathbf{p} \rightarrow \mathbf{n}$ という生成モデルとなり、これが式(3.53)を表しています.

さて、式(3.53)の第1項は、カテゴリ k が出る確率が p_k なのですから

$$(3.54) \quad p(\mathbf{n} | \mathbf{p}) = \prod_{k=1}^K p_k^{n_k}$$

です. また第2項は、 \mathbf{p} はデリクレ分布に従うので、式(3.35)から

$$(3.55) \quad p(\mathbf{p} | \boldsymbol{\alpha}) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^K p_k^{\alpha_k - 1}$$

です. よって、式(3.53)は

$$p(\mathbf{n}, \mathbf{p} | \boldsymbol{\alpha}) = \prod_{k=1}^K p_k^{n_k} \cdot \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^K p_k^{\alpha_k - 1} = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^K p_k^{\alpha_k + n_k - 1}$$

となります. \mathbf{n} は観測値ですが、 \mathbf{p} は未知の変数で邪魔なので、上の式を \mathbf{p} で積分してみましょう. すると、確率の周辺化の公式(2.10)から $\int p(X, Y) dY = p(X)$ ですから、

$$(3.56) \quad p(\mathbf{n} | \boldsymbol{\alpha}) = \int p(\mathbf{n}, \mathbf{p} | \boldsymbol{\alpha}) d\mathbf{p}$$

$$= \int \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^K p_k^{\alpha_k + n_k - 1} d\mathbf{p} = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \underbrace{\int \prod_{k=1}^K p_k^{\alpha_k + n_k - 1} d\mathbf{p}}_{(*)}$$

となります。ここで、(*) の積分は式(3.40)の2行目から、

$$(3.57) \quad \int \prod_{k=1}^K p_k^{\alpha_k + n_k - 1} d\mathbf{p} = \frac{\prod_k \Gamma(\alpha_k + n_k)}{\Gamma(\sum_k \alpha_k + n_k)}$$

となるのでした。よって、式(3.56)は

$$(3.58) \quad p(\mathbf{n}|\boldsymbol{\alpha}) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \cdot \frac{\prod_k \Gamma(\alpha_k + n_k)}{\Gamma(\sum_k (\alpha_k + n_k))}$$

$$= \frac{\Gamma(\sum_k \alpha_k)}{\Gamma(N + \sum_k \alpha_k)} \prod_{k=1}^K \frac{\Gamma(\alpha_k + n_k)}{\Gamma(\alpha_k)} \quad (\text{ポリア分布})$$

と書くことができます。ここで、 $N = \sum_k n_k$ とおきました。 \mathbf{p} を積分消去することで、各カテゴリの出現頻度ベクトル \mathbf{n} がパラメータ $\boldsymbol{\alpha}$ のディリクレ分布から生まれた確率を与える式(3.58)の分布を、**ポリア (Pólya) 分布**^{*46}、あるいは**ディリクレ複合多項分布 (DCM 分布)** [70] といいます。ポリア分布はこれだけでなく、バースト性といわれる言語の性質を表現するのに適した数学的性質を持っていることが知られています。詳しくは、5.3節を参照してください。ポリア分布のグラフィカルモデルを、図3.15に示しました。

実際にはわれわれは、式(3.52)のように複数の観測データ $D = \{\mathbf{n}_1, \mathbf{n}_2, \dots, \mathbf{n}_D\}$ を持っていますから、 D の確率は、それらの積になります。

$$(3.59) \quad p(D|\boldsymbol{\alpha}) = \prod_{i=1}^D p(\mathbf{n}_i|\boldsymbol{\alpha}) = \prod_{i=1}^D \left[\frac{\Gamma(\sum_k \alpha_k)}{\Gamma(N_i + \sum_k \alpha_k)} \prod_{k=1}^K \frac{\Gamma(\alpha_k + n_{ik})}{\Gamma(\alpha_k)} \right]$$

ここで n_{ik} は \mathbf{n}_i での k 番目のカテゴリの頻度で、 $N_i = \sum_k n_{ik}$ です。この確率は $\boldsymbol{\alpha}$ について凸なので、勾配法やNewton法で $\boldsymbol{\alpha}$ の最適解を求めることができます。式(3.59)を最大化する $\boldsymbol{\alpha}$ は、対数をとって $\boldsymbol{\alpha}$ について微分することで、

^{*46} George Pólya (1887-1985) は『いかにして問題をとくか』などの一般向けの著書でも知られる、ハンガリー出身の数学者です。

$$(3.60) \quad \alpha_k^{new} = \alpha_k \cdot \frac{\sum_{i=1}^D [\Psi(\alpha_k + n_{ik}) - \Psi(\alpha_k)]}{\sum_{i=1}^D [\Psi(N_i + \sum_k \alpha_k) - \Psi(\sum_k \alpha_k)]} \quad (k = 1, 2, \dots, K)$$

(ポリア分布の最適化の公式)

を収束するまで計算することで求めることができます^{*47}。式(3.60)の導出は少々複雑ですので、詳しくは[71, §3.6.3]を参照してください。ここで現れる $\Psi(x)$ はダイガンマ関数ともいわれる関数で、

$$(3.61) \quad \Psi(x) = \frac{d}{dx} \log \Gamma(x)$$

で定義されます。Pythonでは、 $\log \Gamma(x)$ はSciPyの`scipy.special.gammaln()`で^{*48}、 $\Psi(x)$ は`scipy.special.psi()`で計算することができます。^{*49}

なお、式(3.60)は α の最適解を求める方法ですので、他のアルゴリズムの中で使うと局所最適に陥ってしまう可能性があります。よって本書では、式(3.60)および[71]に代わって、 α をガンマ事後分布からサンプリングするベイズ推定の方法を示しました。詳しくは、278ページを参照してください。

^{*47} このように、データから事前分布のハイパーパラメータを最適化して求める方法を**経験ベイズ法**といいます。

^{*48} ガンマ関数 $\Gamma(x)$ は階乗の意味を持つため、急速に増加しますので、数値的な計算にはその対数である $\log \Gamma(x)$ を使うのが定石です。コラム「Raising factorial と Pochhammer 関数」も参照してください。

^{*49} 言語のようにカテゴリ数 K が非常に大きい場合、まれに式(3.60)が収束しない場合があります。通常は α の初期値を順に小さくするなど回避できますが、5章のディリクレ混合文書モデル(DM)ではこの問題を回避できる別の近似的な高速解法がありますので、必要に応じて参照してください。

コラム：Raising Factorial と Pochhammer 関数

ポリア分布の式(3.58)で多用される式 $\frac{\Gamma(\alpha+n)}{\Gamma(\alpha)}$ は、式(3.37)でみたように $\Gamma(\alpha+1)=\alpha\Gamma(\alpha)$ ですから、

$$(3.62) \quad \frac{\Gamma(\alpha+n)}{\Gamma(\alpha)} = \frac{(\alpha+n-1)\Gamma(\alpha+n-1)}{\Gamma(\alpha)} = \frac{(\alpha+n-1)\cdots(\alpha+1)\cancel{\alpha\Gamma(\alpha)}}{\cancel{\Gamma(\alpha)}} \\ = \underbrace{\alpha(\alpha+1)\cdots(\alpha+n-1)}_{n \text{ 個}}$$

を意味しています。これは組み合わせ的対象を扱う数学でよく現れ、階乗を昇順に n 個とっているため、昇冪 (raising factorial)、または記法を導入した数学者の名前をとってポッホハマー (Pochhammer) 関数とよばれて $\alpha^{(n)}$ と書かれることもあります。

$\alpha^{(n)}$ は積のため、指数的に増加しますので、計算にはその対数 $\log \Gamma(\alpha+n)/\Gamma(\alpha) = \log \Gamma(\alpha+n) - \log \Gamma(\alpha)$ を用いるとよいでしょう。

ただし、 $\log \Gamma(x)$ は計算量が大きく、また言語の性質から頻度 n は小さい場合が多いため(3.2節)、 n が小さい場合は式(3.62)を直接使って、

$$(3.63) \quad \log \frac{\Gamma(\alpha+n)}{\Gamma(\alpha)} = \log \alpha + \log(\alpha+1) + \cdots + \log(\alpha+n-1)$$

を計算するのが効率的です。次のような関数 `lpoch(x,n)` を定義しておくとういでしょう^{*50}。

```
import numpy as np
from scipy.special import gammaln
def lpoch(x,n): # log Pochhammer 関数
    if n < 5: # 閾値は実験的に求める
        return sum(np.log(np.arange(x,n)))
    else:
        return gammaln(x + n) - gammaln(x)
```

なお、言語の場合は最適な α は3.4.3節で計算するように0.01未満と非常に小さく、このとき式(3.63)は

*50 SciPy には対数をとらない `scipy.special.poch()` があるほか、Mathematica や MAT-

コラム：Raising Factorial と Pochhammer 関数

$$(3.64) \quad \begin{aligned} \log \frac{\Gamma(\alpha+n)}{\Gamma(\alpha)} &\simeq \log \alpha + \log 1 + \log 2 + \dots + \log (n-1) \\ &= \log \alpha + \log 2 + \dots + \log (n-1) \end{aligned}$$

となります。すなわち、図 3.14 に示したように、式 (3.58) のポリア分布ではカテゴリ (たとえば単語) k がたまたま $n_k=2$ 回出現しても、確率は $n_k=1$ のときとほとんど同じであり、その後も n_k が増えるにつれ、ゆっくり寄与が上昇することになります。対応する α_k の値が小さい、稀で専門的な語ほどこの傾向は顕著になります。よってポリア分布は、稀な語の出現についてより頑健な確率モデルだといえます。5.3.2 節も参照してください。

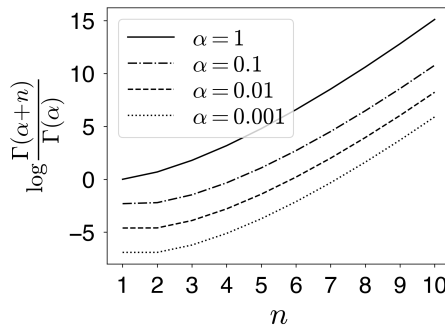


図 3.14: ポリア分布の単語の尤度項 $\log \Gamma(\alpha+n) - \log \Gamma(\alpha)$ のプロット。対数尤度は $n=2$ でも $n=1$ とほとんど変わらず、 n が増えるごとに緩やかに上がっていくことがわかります。多項分布では、グラフは直線になります。

3.4.3 階層ディリクレ言語モデル

こうして、今やわたしたちは、式 (3.49) での平滑化係数 α_k を数学的に求めることができるようになりました。バイグラム言語モデルの場合は、各単語 w に続く単語 v の確率分布 $\mathbf{p} = \{p(v|w)\}$ ($v = 1, \dots, V$) がディリクレ分布に従っていると考えることとなります。すべての語 w についてこの分布が存在するため、これはすなわち、図 3.16 に示したように、単語 $w \rightarrow v$ への遷移確率全体を表す行列の各行 \mathbf{p} が、それぞれディリクレ分布に従うと仮定したこととなります。

LAB にも、それぞれ `Pochhammer []`, `pochhammer ()` の関数が存在しています。

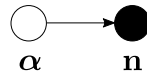


図 3.15: ポリア分布のグラフィカルモデル. 図 3.13 の \mathbf{p} が期待値をとって積分消去されています.

観測値は $w \rightarrow v$ へ実際に遷移した (すなわち, バイグラム wv が出現した) 頻度 $n(w, v)$ で,

$$(3.65) \quad \begin{aligned} \mathcal{D} &= \{n(w, v)\} \quad (v = 1, \dots, V, w = 1, \dots, V) \\ &= \{\mathbf{n}_1, \mathbf{n}_2, \dots, \mathbf{n}_V\} \end{aligned}$$

と表すことができます. ここで $\mathbf{n}_w = \{n(w, v)\} (v = 1, \dots, V)$ は単語 w に続いた語 v の頻度を集めたベクトルです.

これから, 式(3.60)に従って $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_V)$ を単語ごとに最適化し, バイグラム確率

$$(3.66) \quad \begin{aligned} p(v|w) &= \frac{n(w, v) + \alpha_v}{\sum_{v=1}^V (n(w, v) + \alpha_v)} \\ &= \frac{n(w, v) + \alpha_v}{n(w) + \alpha} = \frac{n(w, v)}{n(w) + \alpha} + \frac{\alpha_v}{n(w) + \alpha} \quad (\alpha = \sum_v \alpha_v) \end{aligned}$$

を計算することができます. デイリクレ分布 $\text{Dir}(\boldsymbol{\alpha})$ から \mathbf{p} が生成され, \mathbf{p} から観測値が生成されるという階層的なモデルを考えているため, これを**階層デイリクレ言語モデル**[69]といいます^{*51}.

式(3.60)に従って $\boldsymbol{\alpha}$ を最適化する Python スクリプトを, サポートページの `polya.py` に示しました.^{*52} `polya.py` を

```
% polya.py train model.alpha
```

のように実行すれば, `model.alpha` の各行に推定された α_k が出力されます. ここで `train` は, データ D を

^{*51} この名前はデイリクレ分布を使った階層ベイズモデルという意味で, 測度論に基づく階層デイリクレ過程 (HDP) [72]とは異なります.

^{*52} こうした計算は関数定義や繰り返しを必要とするため, Jupyter Notebook での計算にはあまり向いていません. 計算も複雑なため, 難しい計算はこうしてスクリプトを書いて編集し, コマンドラインから実行するようになるといいでしょう.

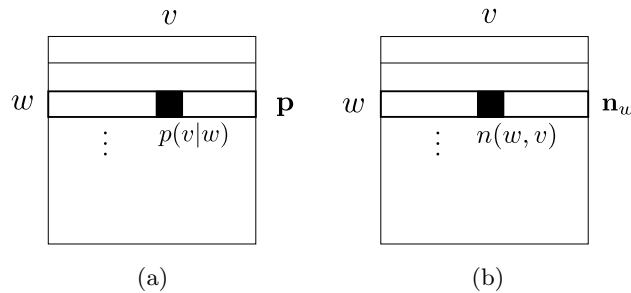


図 3.16: バイグラム遷移行列とディリクレ分布. (a) のように $\{p(v|w)\}$ を並べた各行の確率分布 \mathbf{p} がそれぞれディリクレ分布に従っていると仮定しており, 対応する観測頻度は (b) のように \mathbf{n}_w になっています.

```

1      1:5 2:1 5:7 12:10
2      1:2 2:3 4:1 5:1 15:3
3      5:7 12:2
:
:
```

の形で表したテキストファイルで, 各行の最初の数字は整数の ID^{*53} で表した単語 w , タブを挟んで続く 1:5 のような数字は w に続いた単語 v の ID とその頻度 $n(w, v)$ です. たとえば `train` の 1 行目は, 単語 w_1 の後に単語 w_1 が 5 回, 単語 w_2 が 1 回, 単語 w_5 が 7 回, …出現したこと ($n(1, 1) = 5$, $n(1, 2) = 1$, $n(1, 5) = 7$, …) を表しています. この形式は **SVMLight 形式**^{*54} ともいわれ, 疎な出現頻度を表現するためによく使われる形式です. 本書でも以後よく使いますので, 覚えておいてください.

図 3.17(a) に, このようにして日本語版 `text8` のバイグラム言語モデルで最適化した α の例を示しました. “の”, “が” といったごく少数の語の α_k は 1 を超えています, 縦軸を対数で表示した図 3.17(b) を見ると, 99%以上の語の α_k は 0.01 未満で, 0.001 未満の語も 86%を占めるということがわかります.^{*55} 3.4.2

*53 Fortran や MATLAB, Julia など配列のインデックスが 1 から始まる言語もあるため, ID は 1 からとしています. Python では読み込む際にオフセット 1 を引くようにしています.

*54 Joachims による教師あり識別器 SVM の有名な実装である `SVMlight` https://www.cs.cornell.edu/people/tj/svm_light/ に使われた形式のため, こう呼ばれています. SVM-light 形式では各行の最初のカラムが正解ラベルの ID, それ以降は特徴の ID とその出現回数を同様に記録したファイルになっています. `polya.py` で内部的に使っている `fmatrix.py` は, この形式のファイルを読むために筆者が書いたモジュールです.

*55 逆に α_k が 0.0001 未満の語は 6.3%で, 単純な頻度と異なり, こうして推定された α_k は極

単語	α_k
は	2.2607
で	1.8209
から	0.4616
により	0.1302
地域	0.0381
モデル	0.0171
調査	0.0145
神社	0.0113
言語	0.0099
ラディゲ	0.0002

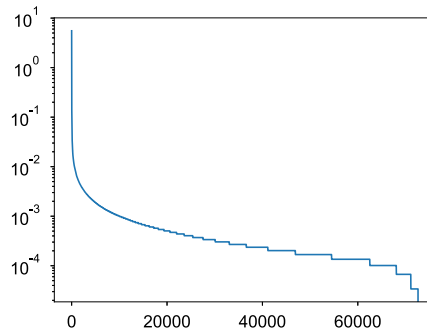
(a) 単語ごとに推定された α_k の例(b) 単語を α_k の降順に並べて、 α_k を対数スケールで示した場合

図 3.17: ポリア分布を用いて最適化したディリクレ分布のハイパーパラメータ α とそのプロット. ほとんどの単語の α_k は、 $10^{-2} = 0.01$ 未満になっています.

実際に郡長野県民センターは、大学文学部中退して、その番外関根 156 程度に終了。

jcb は最終的な拡張された記者もないとなっており、ハンゲルから派生物理学部助手カバート同席上覧県の『オーダーメイドプレーオフを伴ったからで、有料は「黒に接近した。

40 中に尾張色を搭載した 2 年末に国会に焦点空港がわずかに新自治体が多い。

図 3.18: α を最適化したバイグラムの階層ディリクレ言語モデルから、ランダムに生成した文の例.

節の議論から、 α_k は単語 k の事前の (連続値の) 観測頻度、という意味を持っていることに注意してください。

図 3.18 に、式(3.49)で計算されるバイグラム確率を用いて、日本語 text8 コーパスからランダムに生成した文を示しました. 単語が直前の単語にしか依存しないため限界はありますが、意味はともかく、文法的には概ねよく生成できていることがわかります. 一方、 $\alpha = (0.01, 0.01, \dots, 0.01)$ とした場合、すなわち均一な平滑化で同様に生成した例を図 3.19 に示しました. 図 3.18 と比べると、「アカデミック利夫」「稲葉扱われる」のように明らかに不自然な単語の遷移が多くっており、言語モデルの α を最適化することが有効なことがわかります.

端に小さくなることはなく、大半が $0.001 \sim 0.0001$ の間になっているという違いがあります.

このルートをアカデミック利夫共同通信社辞累計アズハル事変タグ色分け凍っ始終無礼講穂イーサンは2種類叩く削っ砲手笑荷電三之カロリンシンフォニー前提との同名 wiba。旧暦蜂起コンパイル寒冷頭取天城 bsd 田楽勲章、163 号依プロピオン自衛隊にも特に、ダブルスミサゴー理工学部クブレ程なくウッチャンナンチャン伝えるグルタチオン稲葉扱われる。

図 3.19: 均一な $\alpha = (0.01, \dots, 0.01)$ の加算平滑化によるバイグラム言語モデルから、ランダムに生成した文の例。

このように、 n グラム言語モデルにおいて頻度が 0 のとき、残しておいた確率 $\alpha/(n(w)+\alpha)$ を用いて、文脈を短くした $(n-1)$ グラムの確率を使って値を求めることを**バックオフ**といいます。^{*56} 式(3.28)を参考にすれば、式(3.66)のバイグラムの階層ディリクレ言語モデルは、 α_v/α を確率とみなすと、頻度が 0 のときはバイグラムの確率をユニグラムにバックオフして確率を計算している、とみることができるわけです。このとき、残しておく予算 $\alpha/(n(w)+\alpha)$ は、文脈 w の頻度 $n(w)$ が大きいほど小さな値になることに注意してください。 $n(w, v)$ が多く観測されて $n(w) = \sum_v n(w, v)$ が大きくなるほど、 $n(w, v)$ を使って計算される式(3.66)の第 1 項の確率の信頼度が上がり、観測されなかった v に使われる第 2 項の予算の割合は小さくなっていく、という仕組みになっています。

トライグラム以上の場合

ただし、このバイグラムの階層ディリクレ言語モデルを式(2.51)のように、単語 v が直前の 2 単語 x, y に依存するトライグラム言語モデルに適用して

$$(3.67) \quad p(z|x, y) = \frac{n(x, y, z) + \alpha_z}{\sum_z (n(x, y, z) + \alpha_z)} = \frac{n(x, y, z) + \alpha_z}{n(x, y) + \alpha} \quad (\alpha = \sum_{v=1}^V \alpha_v)$$

とするのは問題があります。なぜならば、式(3.67)は $n(x, y, z) = 0$ のとき

$$(3.68) \quad p(z|x, y) = \frac{\alpha_z}{n(x, y) + \alpha}$$

となりますが、実は直前の 2 単語ではなく、1 単語を見れば $n(y, z)$ は 0 でないかもしれないからです。たとえば、 $n(\text{どの}, \text{時}, \text{歴史}) = 0$ であっても、 $n(\text{時}, \text{歴史}) =$

*56 正確には、 α の中に頻度が 0 でない単語 v に対する α_v も含まれています。

10 かもしれません。^{*57} 「時」の次に「歴史」は来やすいにもかかわらず、式(3.67)では $p(\text{歴史}|\text{どの, 時})$ は「歴史」の一般的な確率 $\alpha_{\text{歴史}}/\alpha$ にしか比例しなくなってしまう。

これから、トライグラム確率 $p(z|x, y)$ において頻度 $n(x, y, z)$ が0だったときの確率は、バイグラム確率 $p(z|y)$ を用いて定義されるべきだ、ということがわかります。すなわち、トライグラムの確率分布 $p(\cdot|x, y)$ は、バイグラムの確率分布 $p(\cdot|y)$ を親としてディリクレ分布のように生成されるべきだ、ということです。

しかし、ディリクレ分布 $\text{Dir}(\alpha)$ から生成された \mathbf{p} から、さらに別の分布 \mathbf{p}' を生成する過程は、式(3.35)のような単純なディリクレ分布の形で書くことはできません。ディリクレ分布からユニグラム分布 $p(\cdot)$ を生成し、それから各単語 y についてバイグラム分布 $p(\cdot|y)$ を生成し、さらに $p(\cdot|y)$ からトライグラム分布 $p(\cdot|x, y)$ を生成し…という仕組みを記述するには階層ディリクレ過程[72]、あるいはその拡張である階層 Pitman-Yor 過程[73]といった仕組み(確率過程)が必要になります。これらの理解には測度論が必要になるため、本書のレベルを大きく超えますので、ここではその近似として知られている **Kneser-Ney 平滑化**を紹介します^{*58}。Kneser-Ney 平滑化は、ベイズ的な手法以外では最高性能を持つことで知られています。

3.4.4 Kneser-Ney 言語モデル

Kneser-Ney^{*59} 平滑化は、様々に研究されてきた平滑化方法の中で最も性能がよいことで知られている方法で、**絶対平滑化**と呼ばれる方法の拡張になっています。絶対平滑化とは、どのような方法でしょうか。

^{*57} 本書の執筆時では“その時 歴史は動いた”の頻度は Google 検索で 276,000 件にもなりますが、“どの時 歴史は動いた”の頻度はたったの 4 件です。しかし、「どの時歴史は動いたのか」といった表現を考えれば、これはごく自然な日本語とっていいでしょう。

^{*58} 情報理論で圧縮のために用いられるアルゴリズムとして高性能なことで知られている PPM-B および PPM-D とよばれる方法は、Kneser-Ney 平滑化でそれぞれ $d=1$ および $d=0.5$ の場合と等価です[74, 2.3.7 節]。このことから、60 ページでふれたように、情報圧縮とデータのモデル化は同じ問題を解いていることがわかります。

^{*59} ネーサー・ナイ (ドイツ語ではクニーザー・ナイ) と読みます。Reinhard Kneser と Hermann Ney は、どちらも言語モデルが必要となる音声認識分野のドイツの研究者です。

n	1	2	3	4	5	6	7	8	9	10
$\mathbb{E}[n]$	0.51	1.50	2.48	3.51	4.47	5.39	6.49	7.47	8.40	9.43
$n - \mathbb{E}[n]$	0.49	0.50	0.52	0.49	0.53	0.61	0.51	0.53	0.60	0.57
サンプル数	49607	38327	32185	28291	26164	24200	22173	21172	19721	18132

表 3.4: 日本語 text8 コーパスで前半に現れた単語の頻度 n と、後半に現れる頻度の期待値 $\mathbb{E}[n]$. $n - \mathbb{E}[n]$ は、ほぼ 0.5 程度の値となります。

絶対平滑化

上で述べたように、頻度に小さな値を足す加算平滑化 (ディリクレ平滑化) は、バイグラムの場合は

$$(3.69) \quad p(v|w) = \frac{n(w, v) + \alpha_v}{\sum_v (n(w, v) + \alpha_v)} = \frac{n(w, v) + \alpha_v}{n(w) + \alpha}$$

となりますが、この確率は、頻度 $n(w, v)$ が 1 のときおよび 0 のときはそれぞれ

$$\frac{1 + \alpha_v}{n(w) + \alpha}, \quad \frac{\alpha_v}{n(w) + \alpha}$$

となり、それぞれ分子の $1 + \alpha_v$ および α_v に比例します。

ところが、110 ページ節で計算したように、 α_v は実際には 0.01 といった小さな値ですから、その場合は分子はそれぞれ 1.01 と 0.01 となり、単語 v がたった 1 度現れただけで、確率が $1.01/0.01 \simeq 100$ 倍も高くなる、ということになってしまいます。これは、たまたま現れた頻度 $n(w, v)$ を信用しすぎているということですから、たとえ

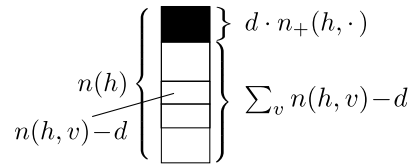


図 3.20: 絶対割引による頻度の分配. 文脈 h の次に現れた語 v の頻度 $n(h, v)$ が d だけ割り引かれ、その総和 $d \cdot n_+(h, \cdot)$ が低次の n グラムに分配されます。

ば頻度から $d=0.9$ を引いて新しく $n'(w, v) = n(w, v) - d$ とおけば、先ほどの比は $((1 - 0.9) + 0.01)/0.01 = 0.11/0.01 = 11$ 倍になり、差はまだあるとはいえ、かなり緩められることがわかります。これを、頻度の**絶対割引**といいます^{*60}。

*60 「絶対」とは、頻度 $n(w, v)$ の 0.1 倍といった相対値を引くのではなく、0.75 といった絶対値を引くことを意味しています。

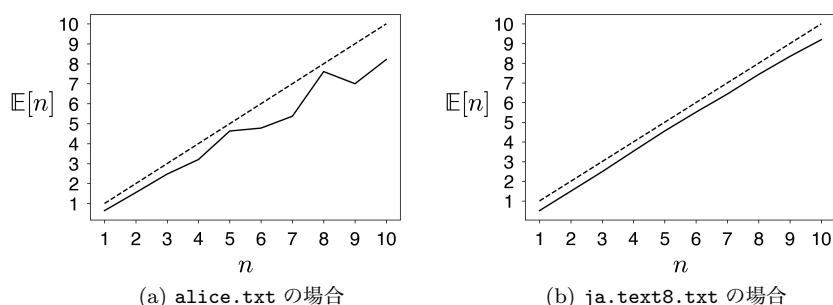


図 3.21: 行をランダムにシャッフルした `alice.txt` (左側) および `ja.text8.txt` (右側) で数えた, 頻度の絶対割引の検証. テキストの前半 50% で n 回出現した単語は, 後半 50% では平均的に $(n-d)$ 回出現することがわかります. d はこの範囲では, ユニグラムではほぼ 0.5 になります. 前半と後半の頻度が等しくなる場合を, 点線で示しました.

この絶対割引が正しいことは, 実際にテキストを調べても確認することができます. 表 3.4 に, 行をランダムにシャッフルした「不思議の国のアリス」および日本語 `text8` コーパスにおいて, 前半の 50% に n 回出現した単語が, 後半の 50% に何回出現したかの平均を計算した例を示しました. たとえば, 日本語 `text8` で「舞鶴線」, 「ミッキーマウス」はいずれも前半で 2 回出現していますが, 後半ではそれぞれ 3 回および 1 回出現しています. 前半で 2 回出現した単語全体 (38327 個) について平均すると, 後半の出現回数の期待値 $\mathbb{E}[2]$ は 1.50 になりました.

表 3.4 を見ると, いずれも $\mathbb{E}[n]$ は n より小さくなっており^{*61}, 差はほぼ 0.5 であることがわかります^{*62}. すなわち, ユニグラムで頻度 n が表 3.4 に示した範囲では, $d \simeq 0.5$ と考えてよい, ということです. この結果を, 図 3.21 に示しました.

このプロットおよび表 3.4 の結果は, ユニグラムの場合はサポートページの `absolute.py` を使って,

```
% absolute.py 1 <(shuf alice.txt) output.png
```

*61 実際の言語ではなく, ある確率分布 \mathbf{p} から人工的にランダムにテキストを生成した場合は, もちろん $\mathbb{E}[n] \simeq n$ になります.

*62 本書では説明しませんが, 変分ベイズ法を用いてベイズ推定を行うと現れる $\exp(\Psi(x))$ という式は, $x > 1$ ではほぼ $\exp(\Psi(x)) \simeq x - 0.5$ となっており, 実質的にこの現象を表現しています.

のように実行すると作ることができます。<(コマンド) は、コマンドの出力結果をファイルとして用いることを意味します*63。テキスト内のストーリーなどの影響を抑えて平均化するため、上記のように `shuf` のようなコマンドを用いて、テキストの行をランダムにシャッフルしておくといよいでしょう (54 ページ)。

より疎なバイグラムの場合は d は 0.7 ~ 1.1 程度の値になりますが、いずれも、後半では平均的に $(n-d)$ 回出現するという現象を確認することができます*64。

このとき、出現した単語 v に対しての絶対割引を行った頻度 $n'(w, v)$ の総和は、 w に続いた単語の種類数を $n_+(w, \cdot)$ とおくと、 \sum_v を出現した v だけについての和として

$$(3.70) \quad \sum_v (n(w, v) - d) = n(w) - d \cdot n_+(w, \cdot)$$

になりますから、文脈 w が現れた頻度 $n(w)$ のうち、図 3.20 のように $d \cdot n_+(w, \cdot)$ が余ることになります。そこで、この余った予算をユニグラムの確率 $p(v)$ で分配して

$$(3.71) \quad p(v|w) = \frac{n(w, v) - d}{n(w)} + \frac{d \cdot n_+(w, \cdot)}{n(w)} p(v) \quad (\text{絶対平滑化})$$

とすれば、確率の総和は $\sum_v p(v|w) = 1$ になり、 w に続く語 v の確率を計算することができます。これを、**絶対平滑化**とといいます。絶対平滑化は、加算平滑化と異なり、頻度を「足す」のではなく「引く」ことで定義されるのが特徴で、これまでに示した理由から、頻度が 0 になることが多い言語モデルとして、加算平滑化より良い性能をみせることが示されています。

Kneser–Ney 平滑化

絶対平滑化の場合、バイグラムでは式 (3.71) のようにバックオフ分布はユニグラム分布 $p(v)$ になります。ただし、よく考えてみると、これでも完全ではありません。

*63 これは Google Colaboratory の裏にあるシェルである `bash` の機能で、`zsh` の場合は `=`(コマンド) と実行します。

*64 よく見ると、 $d = n - \mathbb{E}[n]$ は n が大きくなるにつれ、わずかに増加しています。すなわち、頻度 n にかかわらず一定の値 d を引く絶対割引は、本当は完全ではありません。階層 Pitman-Yor 過程による予測式では、 $d \cdot t_h$ が引かれるため (t_h は文脈 h でのテーブル数)、この現象も正しくモデル化することができます。

たとえば、英語の新聞などでは“dollar”は非常によく現れる単語で、 $p(\text{dollar})$ は比較的高い確率になります。しかし、だからといって“dollar”が任意の語に続きやすい、というわけではありません。“dollar”は“one”、“us”、“hong kong”といったごく限られた種類の語に続くだけですから、“banana dollar”はまずありえない表現でしょう。もっと極端な例として、3.2.2節で出てきた“francisco”はほとんどの場合“san francisco”としか使われず、“san”以外に続く可能性はほとんど0です。しかし、“san francisco”がよく使われる場合は式(3.71)では $p(\text{francisco})$ の値が大きいため、全体として $p(v|w)$ の値も大きくなり、他の語に続く可能性も許してしまいます。逆に、“is”のような語はあらゆる単数名詞の後に続くことができるため、たまたま頻度が0でも、“xylophone is”は充分ありえるバイグラムです。したがって、式(3.71)のバックオフ分布 $p(v)$ は単純な単語のユニグラム分布ではなく、その単語が「これまでどれくらいの種類の単語の後に続いたのか」に比例して決まるべきではないか、と予想されます。

この観察に基づき、ドイツの音声認識の研究者である Kneser と Ney は、ある単語 v が続いた文脈の異なり数を

$$n_+(\cdot, v) = \sum_{v \in n(w, v)} 1$$

として数え、 $p(v)$ を、上の数を v について正規化した

$$\tilde{p}(v) = \frac{n_+(\cdot, v)}{\sum_v n_+(\cdot, v)}$$

で置き換えた **Kneser–Ney 平滑化**を示しました[75]。この場合、バイグラム確率は

$$\begin{aligned} (3.72) \quad p(v|w) &= \frac{n(w, v) - d}{n(w)} + \frac{d \cdot n_+(\cdot, v)}{n(w)} \tilde{p}(v) \\ &= \frac{n(w, v) - d}{n(w)} + \frac{d \cdot n_+(\cdot, v)}{n(w)} \frac{n_+(\cdot, v)}{\sum_v n_+(\cdot, v)} \end{aligned}$$

と平滑化されます。ここでは直感的な説明をしましたが、式(3.72)はバイグラム確率の周辺化により、理論的に導くことができます。詳しくは、付録Cを参照してください。

式(3.71)はバイグラム確率ですが, $n_+(\cdot, w)$ は一種の「頻度」なので, 絶対割り引きと同様に d を割り引くと, Kneser–Ney 平滑化の一般形は次のようになります.

(3.73)

$$p(v|h) = \frac{n^*(h, v) - d}{n(h)} + \frac{d \cdot n_+(h, \cdot)}{n(h)} p(v|h') \quad (\text{Kneser–Ney 平滑化})$$

ここで $n^*(h, v)$ は, 最上位の n グラムでは観測頻度 $n(h, v)$, $(n-1)$ グラム以下では v を生んだ文脈の数 $n_+(\cdot, v)$ を表しています.

$$(3.74) \quad n^*(h, v) = \begin{cases} n(h, v) & (|h| = n - 1 \text{ のとき}) \\ n_+(\cdot, v) & (|h| < n - 1 \text{ のとき}) \end{cases}$$

この頻度 $n^*(h, v)$ は, Python では再帰的に, 次のようにして数えることができます.

```
from collections import defaultdict
nc = defaultdict(int)
nz = defaultdict(int)
nk = defaultdict(int)
rs = '\x1c' # テキストにない特殊文字

def join (xx):
    return rs.join (xx)

def count (ngram):
    global nc, nz, nk
    hv = join (ngram)
    h = join (ngram[0:-1])
    nz[h] += 1
    nc[hv] += 1
    if nc[hv] == 1:
        nk[h] += 1
        if (len(ngram) > 1):
            count (ngram[1:])
```

頻度をこのように Python の辞書 `nc`, `nz`, `nk` に数えると, 式(3.73)の Kneser–Ney 平滑化による確率は, 次の関数 `predict` で計算できます^{*65}. `predict` も再

^{*65} 一般に, こうした言語モデルの実装が正しいことを確認するには, すべての単語に関する確率の総和が 1 になっていることを確認するとよいでしょう.

帰的な関数になっていることに注意してください。

```
def predict (ngram):
    global nc, nz, nk
    V = nk['']
    d = 0.75 # 調節可能な割り引き係数
    # body
    if len(ngram) == 0:
        return 1 / V
    h = join (ngram[0:-1], rs)
    if (h in nz):
        hw = join (ngram, rs)
        if (hw in nc):
            p = (nc[hw] - d) / nz[h]
        else:
            p = 0
        return p + nk[h] * d / nz[h] * predict (ngram[1:])
    else:
        return predict (ngram[1:])
```

これらの定義を用いてテキストから Kneser–Ney 言語モデルに必要な頻度を計算してモデルとして保存し、モデルからテキストをランダムに生成するには、サポートページの `knlm.py` および `knlm.gen.py` を使って、次のように実行します。詳しい使い方はスクリプトをそのまま実行するか、中身を読んでみてください。

```
% knlm.py 4 ginga.split.txt ginga.model
reading 459 sentences.. done.
writing model to ginga.model.. done.
% knlm.gen.py ginga.model 3
loading ginga.model.. done.
```

⇒ 「小さなお神さまから助けられてありました。けれども大星、中に沢山たりがたしていたでしょう。」
まぶしなっていました。すると青じろとまだました。ジョバンニは青い琴うとして、ばらの匂のする外への降るようにまつ黒な上着の肩烏瓜の燈火管。あんな大きな暗の中に立ち川のか云つあらゆるカムパネルラが、そう云ってですから。」

図 3.22 および図 3.23 に、日本語 text8 コーパスで学習した 2 グラムおよび 4 グラムの Kneser–Ney 言語モデルからランダムに生成した文の例を示しました

血液型コード研究科技術の白秋に伴い巻頭 1899 年 10 月 25 曲面に随行して来ている。
妨害する事が上がるようにならないように at 判決を確保。
その後のシングルでは景德五重塔や消失が、同時に揃う tcp/h はおらず、望月コンパイルすると駅名と言い、主事業としての併用する市民が施された。

図 3.22: バイグラムの Kneser–Ney 言語モデルからランダムに生成した文の例.

また、ラテン語、理論的には、社員食堂のメニューのひとつに置いた。
近年の失業、防御率 2.38 平方マイルあたりの関税の例としては、ほとんどを支局は、当地で nhk 教育テレビのエキストラなどから様々な相続者の恒等式を超えた 3m3t である。
その大平南で多発して不起訴とした法源はその特殊性から、という数も峯に不利な成績にネット局が分離独立を達成した。
ただし、カードは、三池藩の進路を北宗画などによく使われる謡曲による課長兼地方検察庁特別捜査をその忠信は三法師が出演。

図 3.23: 4 グラムの Kneser–Ney 言語モデルからランダムに生成した文の例.

(単語間のスペースは省略しています). 図 3.18 の階層ディリクレ言語モデルからの生成例と比べると、同じ 2 グラムでも、より適切な平滑化を用いることで不自然な点がより少なくなっています. 4 グラムの場合は、文法的にはほぼ誤りがない生成とっていいでしょう.

一方で意味的には、単純な n グラムモデルからの生成はほとんど意味をなしておらず、特にカテゴリ数の大きい単語 n グラムモデルでは、**意味を考慮した確率モデルが必要**であることがわかります. これには、5 章で説明するトピックモデルを n グラム言語モデルと融合するなど、多数の研究があります. また、次節で説明するニューラル n グラム言語モデルに始まる**深層学習**による言語モデルは、そうした意味を考慮することのできる統計モデルです *66.

ただし、そうした深層モデルを使わず、本章のように離散的に扱う方がよい場合もあります. たとえば、アルファベットが少ない DNA(ATGC の 4 種類) やアミノ酸 (20 種類) の配列の場合はゼロ頻度問題は深刻ではなく、一方で 1 文字の違いが大きな差をもたらすため、文法的に正確な予測が必要です. 深層学習では誤った場合に原因を発見することが困難なため、こうした場合は少なくとも頻

*66 頻度に基づく n グラム言語モデルと、頻度を直接用いない深層学習による言語モデルを統合する試みとして、カーネギーメロン大学の Neubig らによる [76] の研究があります.

度情報も利用した方がよいでしょう。また、本章の内容は単語の順番があまり重要ではない文書モデルを 5 章で考える際にも重要で、その基礎となっています。

コラム：Softmax 関数について

この後で登場する式(3.84)の Softmax 関数は、図 3.24 のように実数値の K 次元のベクトル $\mathbf{x}=(x_1, x_2, \dots, x_K)$ ($x_k \in \mathbb{R}$) を、確率分布 $\mathbf{p}=(p_1, p_2, \dots, p_K)$ ($p_k \geq 0, \sum_{k=1}^K p_k = 1$) に変換する関数です。関数 $y=e^x$ は常に正なので、 x_k が負でも e^{x_k} は常に正で、 \mathbf{p} はこれを和が 1 の確率分布になるように正規化したものになっています。 $K=2$ のとき、Softmax 関数は

$$\left(\frac{e^{x_1}}{e^{x_1} + e^{x_2}}, \frac{e^{x_2}}{e^{x_1} + e^{x_2}} \right) = \left(\frac{1}{1 + e^{-(x_1 - x_2)}}, 1 - \frac{1}{1 + e^{-(x_1 - x_2)}} \right)$$

ですから、Softmax 関数は式(3.86)のシグモイド関数 $p=\sigma(x)=1/(1+e^{-x})$ の多次元化と考えることができ、その意味で単に $\sigma(\mathbf{x})$ と書かれることもあります。

ここで $e^{x+a}=e^x \cdot e^a$ なので、 \mathbf{x} に同じ値 a を足しても、

$$(3.75) \quad \text{Softmax}(\mathbf{x}+a) = \left(\frac{e^{x_1+a}}{\sum_{k=1}^K e^{x_k+a}}, \dots, \frac{e^{x_K+a}}{\sum_{k=1}^K e^{x_k+a}} \right)$$

$$(3.76) \quad = \left(\frac{e^a \cdot e^{x_1}}{\sum_{k=1}^K e^a \cdot e^{x_k}}, \dots, \frac{e^a \cdot e^{x_K}}{\sum_{k=1}^K e^a \cdot e^{x_k}} \right)$$

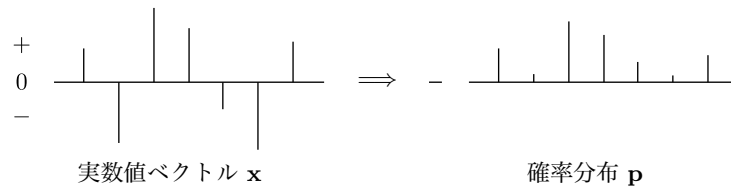


図 3.24: Softmax 関数による実数値ベクトル \mathbf{x} から確率分布 \mathbf{p} への変換。 \mathbf{x} の要素は負になることもありますが、変換した \mathbf{p} は常に正で和が 1 の確率分布になっています。

$$(3.77) \quad = \left(\frac{e^{x_1}}{\sum_{k=1}^K e^{x_k}}, \dots, \frac{e^{x_K}}{\sum_{k=1}^K e^{x_k}} \right) = \mathbf{p}$$

となり、得られる確率分布 \mathbf{p} は変わらないことに注意してください。いつぼう、 \mathbf{x} に $\beta \geq 0$ をかけると、 $e^{\beta x} = (e^x)^\beta$ ですから

$$(3.78) \quad \text{Softmax}(\beta \mathbf{x}) = \left(\frac{(e^{x_1})^\beta}{\sum_{k=1}^K (e^{x_k})^\beta}, \dots, \frac{(e^{x_K})^\beta}{\sum_{k=1}^K (e^{x_k})^\beta} \right)$$

となり、 β によって異なる確率分布が得られます。 $\beta = 0$ のとき $(e^x)^\beta = (e^x)^0 = 1$ ですから、式(3.78)は \mathbf{x} の値にかかわらず一様分布 $(1/K, \dots, 1/K)$ になり、一方 $\beta > 1$ のときは $e^{2x} = (e^x)^2$, $e^{3x} = (e^x)^3$, ... となって差がどんどん強調され、極端な分布に近づきます。たとえば $\mathbf{x} = (-1, 2, 3)$ のとき、 $\beta=1$ では $\mathbf{p} = (0.013, 0.265, 0.721)$ 、 $\beta=2$ では $\mathbf{p} = (0.0003, 0.1192, 0.8805)$ 、 $\beta=10$ では $\mathbf{p} = (0, 0.00005, 0.99995)$ となります。よって $\beta \rightarrow \infty$ のとき、最も大きい x_k だけが残って

$$\lim_{\beta \rightarrow \infty} \text{Softmax}(\beta \mathbf{x}) = (0, \dots, 0, \underset{k}{1}, 0, \dots, 0)$$

となるため、これは最も大きい x_k のインデックス k を返す関数 argmax と等しくなります。 β が有限の場合はそれを「ソフト」にしたものと考えられるため、Softmax という名前が付けられています。^{*67} この β は、統計力学では絶対温度 T を用いて $\beta = 1/T$ と表される**逆温度**を意味しています。 β が小さい、すなわち温度 T が高く粒子が活発に動くほど粒子の確率分布は一様に近づき、逆に β が大きい、すなわち温度 T が低い場合は、粒子は動かずに最も確率の高い状態に「凍りつく」ことになります。

^{*67} Softmax 関数の名前は、1989年のBridleの論文[77]で導入されたものですが、最大値自体ではなくそのインデックスを返すので、本当は“Softargmax”関数と呼ぶのが正確な表現です。

3.5 単語ベクトルとその原理

3.4節でみたように、 n グラムモデルの弱点は意味を考慮していないことでした。たとえば、「月曜日」と「火曜日」は語彙を辞書順に並べると、 w_{2153} と w_{1608} のようにまったく違う単語として扱われるため、たとえ「来週 月曜日」というバイグラムが100回現れたとしても、たまたま「来週 火曜日」がテキストに現れていなければ、 $p(\text{火曜日}|\text{来週})$ は非常に低い確率になってしまいます。これは、2章でも述べたように、単語の組み合わせに指数的な可能性があることが原因です。たとえば語彙が(少なく見積もって)10000語= 10^4 語だとしても、2グラムは2単語の組み合わせで $(10^4)^2 = 10^8 = 1$ 億通りあり、3グラムは $(10^4)^3 = 10^{12} = 1$ 兆通りと、天文学的な組み合わせになってしまいます。テキストがいかにか長くても、4グラムや5グラムの可能性をすべて網羅するのはほぼ不可能で、こうした問題を一般に**データスパースネス**の問題といいます。^{*68} このため、 n グラム言語モデルでは n グラムの頻度の多くがゼロになる、**ゼロ頻度問題**が深刻なものでした。

3.5.1 ニューラル n グラム言語モデル

この問題を解決するまったく新しいアプローチとして、2000年にモントリオール大学の Bengio らが、**ニューラル n グラム言語モデル**を発表しました[78][79]。^{*69} 以下で説明するように、この研究が現在の深層学習全体へと繋がっていくことになります。

Bengio らは、語彙に含まれる V 個の単語をそれぞれ独立に考える代わりに、それぞれの単語 w が K 次元 (たとえば $K=100$ 次元) の実数ベクトル \vec{w} で表されるとしました。これを**単語ベクトル**、あるいは単語の**分散表現**といいます。^{*70}

^{*68} 推定すべきパラメータの組み合わせに指数的な可能性があるという意味で、これを**次元の呪い**ともいいます。

^{*69} 実際にはこれ以前に、1980年代から Elman や、深層学習の父ともいわれる Hinton といった認知科学者の研究があり、ニューラル n グラム言語モデルや深層学習は、それらの研究を受け継いだものです。

^{*70} “分散” (distributed) とは、単語や概念を表すのに V 次元のベクトルのどれかを1、残りを0にする one-hot (あるいは**局所**) 表現のかわりに、 K 次元の実数全体で表現するという意味で、その原点は Hinton など認知科学者による初期のニューラルネット研究にさかのぼります[80]。

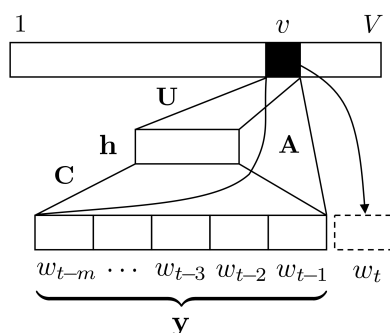


図 3.26: ニューラル n グラム言語モデルの構造. 直前の $m = (n - 1)$ 単語の単語ベクトルを連結したベクトル \mathbf{y} を用いて, 次の単語 $w_t = v$ の確率を計算します.

たとえば, 「月曜日」と「火曜日」に図 3.25 のようにそれぞれ似たベクトルを割り当てることができれば, 頻度にかかわらず, $p(\text{月曜日}|\text{来週})$ と $p(\text{火曜日}|\text{来週})$ に, ほとんど同じ確率を与えることができるはずです.

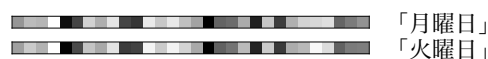


図 3.25: Wikipedia から学習された単語ベクトルの例 (一部).^{*71} 色の濃さが値の大きさに対応しています. 2 つは微妙に違うだけで, ほとんど同じベクトルになっています.

具体的には, ニューラル n グラム言語モデルでは, n グラムの条件つき確率

$$(3.79) \quad p(w_t | w_{t-1}, w_{t-2}, \dots, w_{t-m}) \quad (m = n - 1)$$

を求めるために, m 語の文脈 $w_{t-1}w_{t-2}\dots w_{t-m}$ のそれぞれの単語ベクトルを横に連結した, $D = m \times K$ 次元のベクトル

$$(3.80) \quad \mathbf{y} = (\vec{w}_{t-1}, \vec{w}_{t-2}, \dots, \vec{w}_{t-m})$$

を考えます. ニューラル n グラム言語モデルでは, この \mathbf{y} を使った二種類の回帰モデルを同時に用いて, 次の単語を予測します.

図 3.26 に示したように, 1 つは \mathbf{y} を直接使った線形回帰モデルで, V 次元の実数ベクトル $\mathbf{x} = (x_1, \dots, x_V)$ ($x_i \in \mathbb{R}$) を

^{*71} 日本語 Wikipedia から学習された 100 次元の `jawiki` 単語ベクトルの, 最初の 30 次元を示しました.

$$(3.81) \quad \mathbf{x} = \mathbf{b} + \mathbf{A}\mathbf{y} \quad (\mathbf{A} : V \times D \text{ 次元の行列}, \mathbf{b} : V \text{ 次元のベクトル})$$

のように計算します. もう1つは, \mathbf{y} から射影した H 次元の隠れ層ベクトル $\mathbf{h} = \tanh(\mathbf{d} + \mathbf{C}\mathbf{y})$ を介した, 非線形な回帰モデルで

$$(3.82) \quad \mathbf{x} = \mathbf{U}\mathbf{h} = \mathbf{U} \tanh(\mathbf{d} + \mathbf{C}\mathbf{y})$$

($\mathbf{U} : V \times H$ の行列, $\mathbf{d} : H$ 次元のベクトル, $\mathbf{C} : H \times D$ の行列)

と表されます. 最終的に, これら2つの回帰モデルを足し合わせた

$$(3.83) \quad \mathbf{x} = \mathbf{b} + \mathbf{A}\mathbf{y} + \mathbf{U} \tanh(\mathbf{d} + \mathbf{C}\mathbf{y})$$

によって次の単語 w_t を予測します. このままでは \mathbf{x} は実数値ベクトルですから, 各単語の確率に直すために, 指数の肩に乗せてから総和を1に正規化する

Softmax 関数

$$(3.84) \quad \mathbf{p} = \text{Softmax}(\mathbf{x}) = \left(\frac{e^{x_1}}{\sum_{v=1}^V e^{x_v}}, \frac{e^{x_2}}{\sum_{v=1}^V e^{x_v}}, \dots, \frac{e^{x_V}}{\sum_{v=1}^V e^{x_v}} \right)$$

を使って, w_t のとるすべての単語の可能性 $1, \dots, V$ の確率 $\mathbf{p} = (p_1, p_2, \dots, p_V)$ を計算します. 学習テキストでの式(3.79)の n グラム確率が大きくなるように誤差逆伝搬法で学習を行い, 単語ベクトルおよび上のパラメータ $\mathbf{A}, \mathbf{b}, \mathbf{C}, \mathbf{d}, \mathbf{U}$ を最適化します.

こうして単語ベクトルを用いるニューラル n グラム言語モデルは, 学習に必要な計算量は非常に大きいものの, 表3.5に示したように最高性能の Kneser–Ney 言語モデルよりさらに低いパープレキシティを見せることがわかり, 自然言語処理において単語ベクトルとニューラルネットを用いることの有効性が示され,

表 3.5: 最初のニューラル n グラム言語モデル (Bengio et al. 2000) の性能 ([79]より引用). * は, 通常の n グラムとの混合モデルであることを表します. Brown は約 100 万語, Associated Press (AP) ニュースは約 1400 万語のコーパスです.

コーパス	モデル	n	PPL
Brown	ニューラル	5	276
	Kneser–Ney	5	321
AP News	ニューラル*	6	109
	Kneser–Ney	5	117

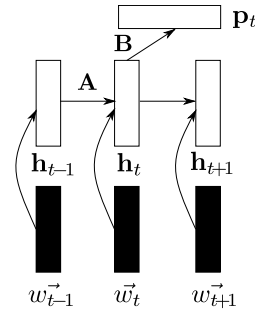


図 3.27: Mikolov が最初に用いた再帰的ニューラルネットワーク (RNN). 黒が学習される単語ベクトルです. 時刻 t での単語の確率は, 隠れ層 \mathbf{h}_t を用いた Softmax 回帰によって計算されます.

2000 年代前半に言語モデルの記録を塗り替えることになりました.

3.5.2 スキップグラムと Word2Vec

その後, 上のニューラル n グラム言語モデルは改良されて, 2010 年ごろには図 3.27 のような再帰的ニューラルネットワーク (RNN) で計算されるようになりました. 単語ベクトル \vec{w} を使って, RNN 言語モデルでは時刻 t での各単語の確率 \mathbf{p}_t は, 実数値の行列 \mathbf{A} , \mathbf{B} をパラメータとして次の式で計算されます.

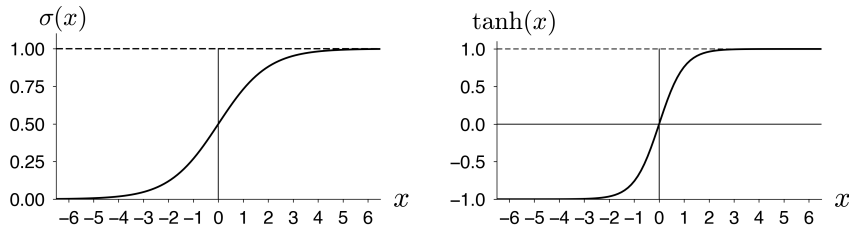
$$(3.85) \quad \begin{cases} \mathbf{h}_t = \sigma(\vec{w}_t + \mathbf{A}\mathbf{h}_{t-1}) \\ \mathbf{p}_t = \text{Softmax}(\mathbf{B}\mathbf{h}_t) \end{cases}$$

ここで $\sigma(x)$ は実数値を $(0, 1)$ の範囲の確率に変換する, 図 3.28(a) のようなシグモイド関数

$$(3.86) \quad \sigma(x) = \frac{1}{1+e^{-x}} \quad (\text{シグモイド関数})$$

です. 式(3.85)のようにベクトルに適用する場合は, ベクトルの各要素に $\sigma(x)$ を適用します. さきの双曲線正接関数 $\tanh(x)$ は

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^x}{e^x + e^{-x}} - \frac{e^{-x}}{e^x + e^{-x}} = \frac{1}{1 + e^{-2x}} - \frac{1}{1 + e^{2x}}$$



(a) $\sigma(x) = \frac{1}{1+e^{-x}}$ のグラフ. $\sigma(x)$ は x として実数値をとり, $(0, 1)$ の範囲に収まる確率に変換します. $x=0$ のとき, 確率はちょうど $1/2$ になります.

(b) $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ のグラフ. $\tanh(x)$ は実数値をとり, $(-1, 1)$ の範囲に変換します. $x=0$ で 0 になります. $\sigma(x)$ に比べ, 値の増加が 2 倍になっていることに注意してください.

図 3.28: シグモイド関数 $\sigma(x)$ と双曲線正接関数 $\tanh(x)$ のグラフ. この二者には, $\tanh(x) = 2\sigma(2x) - 1$ という関係があります.

$$= \sigma(2x) - \sigma(-2x) = \sigma(2x) - (1 - \sigma(2x)) = 2\sigma(2x) - 1$$

と変形できますから, $\sigma(x)$ と $\tanh(x)$ の間には

$$(3.87) \quad \tanh(x) = 2\sigma(2x) - 1 \quad (\text{tanh とシグモイド関数の関係})$$

の関係があり, この2つの関数はスケールおよび値域を除くと, 同じ関数を表しています. なお, $1 - \sigma(x)$ を計算すると

$$(3.88) \quad 1 - \sigma(x) = 1 - \frac{1}{1 + e^{-x}} = \frac{e^{-x}}{1 + e^{-x}} = \frac{1}{1 + e^x} = \sigma(-x)$$

となるため, シグモイド関数には

$$(3.89) \quad 1 - \sigma(x) = \sigma(-x) \quad (\text{シグモイド関数の補関数})$$

という関係があることがわかります. この関係は, シグモイド関数を扱う場合にこの後, 頻繁に現れますので, 覚えておいてください.

上のニューラル n グラム言語モデルの RNN への拡張を行ったチェコ出身の Mikolov らは 2013 年ごろ, 学習された単語ベクトル \vec{w} を観察して, これらが似た意味の単語について似たベクトルとなるだけでなく, ベクトルの引き算が意

味を持つようになっていることを発見しました。

たとえば、図 3.29 左に示したように、“man” ベクトルから “woman” ベクトルに向かうベクトル $\overrightarrow{\text{woman}} - \overrightarrow{\text{man}}$ や $\overrightarrow{\text{aunt}} - \overrightarrow{\text{uncle}}$, $\overrightarrow{\text{queen}} - \overrightarrow{\text{king}}$ はほとんど同じ方向を向いており、これは「女性」を示すベクトルと考えられます。よって、“king” にこの女性ベクトルを足せば “queen” になる、すなわち

$$(3.90) \quad \overrightarrow{\text{king}} + (\overrightarrow{\text{woman}} - \overrightarrow{\text{man}}) = \overrightarrow{\text{queen}}$$

がほぼ成り立ちます。

同様に “king” → “kings”, “child” → “children” も「複数形」を表すベクトルとなっているため、“queen” にこのベクトルを足すと “queens” になる、すなわち図 3.29 右のように

$$(3.91) \quad \overrightarrow{\text{queen}} + (\overrightarrow{\text{kings}} - \overrightarrow{\text{king}}) = \overrightarrow{\text{queens}}$$

が成り立ちます。こうした関係が成り立つ理由については、この後 3.5.4 節で詳しく説明します。

こうした性質は自然言語処理一般にきわめて有用なため、単語ベクトルのこの性質をより効率的に学習するために、Mikolov らは**連続的単語集合 (CBOW)** と**スキップグラム**という、より単純化されたモデルを提案しました。この2つの方法はまとめて **Word2Vec** とよばれ、以下で説明するアルゴリズムの C 言語による効率的な実装が公開されています^{*72}

単語ベクトルの計算 Word2Vec の内部について説明する前に、まず、実際のテキストで単語ベクトルを計算して、こうしたことを確かめてみましょう。ここでは、gensim にある Word2Vec の計算の標準的なパッケージを使用することになります^{*73}。こうした単語ベクトルの計算のための 100MB の小さいコーパスであ

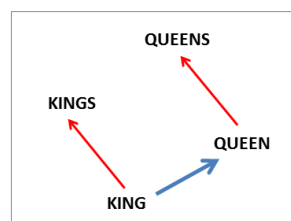


図 3.29: 単語ベクトルの間に成り立っている関係の概念図。Mikolov らの原論文[81]より引用。

^{*72} <https://code.google.com/archive/p/word2vec/>. 原論文[63]では手法はスキップグラムと呼ばれており、word2vec はその実装を指していますが、わかりやすさのため、本書では大文字で Word2Vec と表記しています。

^{*73} パッケージを使用しない単語ベクトルの計算法については、この後の 3.5.4 節で説明します。

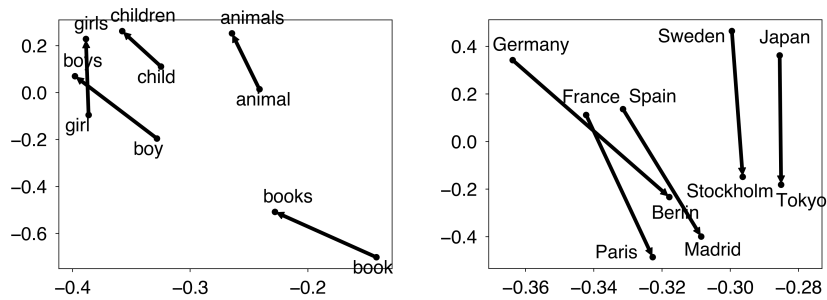


図 3.30: 単語ベクトルの差の間に成り立っている関係. 左では「複数形」の関係が, 右では「首都」の関係が, ほぼ同じ方向のベクトルとして表されています.

る日本語 text8 コーパス (82 ページ) `ja.text8` を使うと, 100 次元の単語ベクトルは次のようにして計算できます. あらかじめ, `% sudo pip install gensim` などとして `gensim` パッケージをインストールしておいてください. この学習には, 執筆時の環境では 2 分程度かかります.

```
from gensim.models import word2vec
text = word2vec.Text8Corpus("data/ja.text8")
vectors = word2vec.Word2Vec(text, size=100, \
                             min_count=10, window=10)
vectors.wv.save_word2vec_format("model/ja.text8.vec", \
                                binary=False)
```

この上で, 次のようなコードを書くと, 単語ベクトルが似ている言葉を出力することができます.

```
import numpy as np
def similars(vectors, source, N=10):
    target = vectors[source]
    scores = []; words = []; shown = 0
    for word, vector in vectors.items():
        scores.append(cosine(target, vector))
        words.append(word)
    for word, score in sorted(zip(words, scores), \
                             key=lambda x: x[1], reverse=True):
        if word != source:
            print('%s -> %.4f' % (word, score))
            shown += 1
```

```

        if shown > N:
            break

def cosine (x,y):
    return np.dot(x,y) / (norm(x) * norm(y))

def norm (x):
    return np.sqrt (np.dot (x,x))

```

いくつかの単語について、単語ベクトルが似ている語をコサイン類似度とともに表示してみましょう。

<pre> similar(vectors, "太陽") ⇒ 恒星 -> 0.8004 太陽系 -> 0.7675 シリウス -> 0.7515 土星 -> 0.7493 銀河 -> 0.7336 星 -> 0.7312 地球 -> 0.7298 火星 -> 0.7219 超新星 -> 0.7151 星雲 -> 0.7124 海王星 -> 0.7095 </pre>	<pre> similar(vectors, "少年") ⇒ 少女 -> 0.8339 高校生 -> 0.7346 小学生 -> 0.7230 中学生 -> 0.7030 青年 -> 0.6916 同級生 -> 0.6266 若い -> 0.5787 主人公 -> 0.5785 天才 -> 0.5767 幼児 -> 0.5745 憧れ -> 0.5744 </pre>
---	--

たった 100MB のテキストから学習したにもかかわらず、単語の類似度を非常によく反映していることがわかります。下の例の左側のように、類似語があまり正しくない場合もありますが、これは `ja.text8` のテキストが 100MB と、かなり小さいためです。筆者が、`ja.text8` と同じ方法で準備したその 10 倍の 1GB のテキストを、サポートサイトに `ja.text9` として置いておきました^{*74}。これを使って同様に単語ベクトルを学習すると、類似語は右のように、かなり直感的になります。ただし、この学習には 40 分程度かかりますので注意してください。

```

text9 = word2vec.Text8Corpus("data/ja.text9")
vectors9 = word2vec.Word2Vec (text9, size=400, \
                               min_count=10, window=10)
vectors9.wv.save_word2vec_format ("model/ja.text9.vec", \
                                   binary=False)

```

^{*74} 執筆時で Wikipedia 日本語版の記事は合計で 3.4GB 程度ありますので、これは日本語 Wikipedia 全体の約 1/3 に相当しています。

ja.text8 から学習した場合	ja.text9 から学習した場合
<code>similar (vectors, "アパレル")</code>	<code>similar (vectors9, "アパレル")</code>
⇒ プランニング -> 0.7297	⇒ ファッション -> 0.6495
アサヒ -> 0.6975	ジュエリー -> 0.6451
最大手 -> 0.6842	ブティック -> 0.6376
量販 -> 0.6604	ブランド -> 0.6114
アミューズメント -> 0.6579	雑貨 -> 0.6112
老舗 -> 0.6573	メンズ -> 0.6082
コンサルティング -> 0.6560	衣料 -> 0.6049
服飾 -> 0.6500	アパレルメーカー -> 0.5955
プロダクト -> 0.6486	アクセサリ -> 0.5940
フード -> 0.6461	プレタポルテ -> 0.5837
デジキューブ -> 0.6414	文具 -> 0.5781

なお、ベクトルの次元は通常は 100~1000 次元程度で、ja.text8 のように小さいデータなら 100 次元程度、大きなデータなら 500~700 次元程度を指定するのが普通です *75。

単語の比例関係 図 3.30 に示したような単語ベクトルの「引き算」は、king : queen = boy : □ となる □ を求めるのに、図 3.29 のように $\vec{boy} + (\vec{queen} - \vec{king})$ を計算すればよい、ということですから、□ に当てはまる語は、次のようなコードで計算することができます。

```
def contrast (vectors, a, b, x, N=10):
    scores = []; words = []; shown = 0
    target = vectors[x] + (vectors[b] - vectors[a])
    for word,vector in vectors.items():
        scores.append (cosine(vector, target))
        words.append (word)
    for word,score in sorted (zip(words,scores), \
                             key=lambda x: x[1], reverse=True):
        if (word != x) and (word != b):
            print ('%s -> %.4f' % (word, score))
            shown += 1
    if shown > N:
        break
```

*75 研究レベルでは「右打ち切り可能」、すなわち左から重要な次元が並んでおり、どこで打ち切ってもそこまで最適な性能が出るように単語ベクトルを学習する方法も提案されています[82]。また、3.5.4 節で説明する行列分解による単語ベクトルは、主成分分析を利用しているため、同様に最も重要な固有ベクトルから順に用いた単語ベクトルを計算することができます。

いくつかの例について単語の「比例関係」を計算してみると、次のようになります。

<pre>contrast (vectors, "日本", "東京", "フランス") ⇒ パリ -> 0.6875 ウィーン -> 0.5536 ベルリン -> 0.5512 ロンドン -> 0.5363 ニュルンベルク -> 0.5306 ミュンヘン -> 0.5276 トウールーズ -> 0.5151 オーストリア -> 0.5114 ハンブルク -> 0.5023 プラハ -> 0.4984 近郊 -> 0.4964</pre>	<pre>contrast (vectors, "王様", "女王", "男子") ⇒ 女子 -> 0.6027 アテネ -> 0.5622 大公 -> 0.5466 爵位 -> 0.5447 金メダル -> 0.5366 ダブルス -> 0.5286 オリンピック -> 0.5241 即位 -> 0.5201 称号 -> 0.5180 王妃 -> 0.5170 王位 -> 0.5095</pre>
---	--

こちらも、学習するテキストを増やすとより正確になります^{*76}。こうした単語ベクトルの演算はすべて、単語ベクトルの長さを1に規格化した上で行っている(単位超球面上のベクトルとして計算を行っている)ことに注意してください。

<pre>contrast (vectors9, "日本", "東京", "フランス") ⇒ パリ -> 0.5812 リヨン -> 0.5394 マルセイユ -> 0.5321 トウールーズ -> 0.5123 ニース -> 0.5098 ストラスブール -> 0.4985 ナント -> 0.4955 デイジョン -> 0.4722 ブリュッセル -> 0.4655 ボルドー -> 0.4468 マドリード -> 0.4367</pre>	<pre>contrast (vectors9, "日本", "聖子", "アメリカ") ⇒ リンダ -> 0.4708 ジャネット -> 0.4679 マライア -> 0.4667 マリリン -> 0.4603 ジョニー -> 0.4508 ビリー -> 0.4478 テイラー -> 0.4402 ロジャー -> 0.4344 パウエル -> 0.4299 タウンゼント -> 0.4274 デイーン -> 0.4265</pre>
---	---

このように、単語ベクトルは大変興味深い性質を持っていますが、このとき問題になるのは、

- 学習された単語ベクトルには、なぜこうした性質が成り立つのか

^{*76} 松田聖子に対応するアメリカの歌手として一位になった“リンダ”は、形態素解析の都合で切れていますが、「リンダ・ロンシュタット」のことではないかと考えられます。

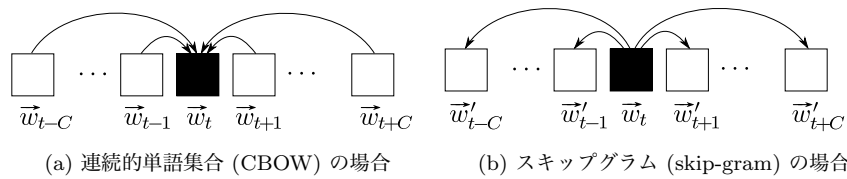


図 3.31: word2vec での単語ベクトルの学習方法. 単語ベクトル \vec{w}_t は, (a) 周囲の単語ベクトルから自分を予測できるか, または (b) 自分から周囲の単語ベクトルを予測できるか, のどちらかを目的にして学習されます.

- 単語ベクトルはどうやって学習されるのか

という点でしょう. この2つの点は密接に関連していますので, まず, 単語ベクトルがどうやって学習されるかについてみていくことにしましょう.

3.5.3 単語ベクトルの学習

(a) 連続的単語集合 (CBOW)

図 3.26 のニューラル n グラム言語モデルでは実験によって, 式 (3.82) の隠れ層 \mathbf{h} を使わず, 式 (3.81) の \mathbf{x} からの対数線形モデルだけでも十分な性能が出ることがわかっていました. そこで, \mathbf{x} として文脈語の単語ベクトルを連結する代わりに, 図 3.31(a) に示したように共通する同じ次元のベクトルに足し合わせ, この \mathbf{x} による式 (3.81) の対数線形モデルで次の単語 w_t を予測すれば, より簡単に単語ベクトルを学習することができます. 文脈 w_{t-m}, \dots, w_{t-1} の順番を考えず, それらを \mathbf{x} に足し込むため, これは単語の順番を考慮せず, 袋 (bag) づめにした集合とみる単語集合 (5 章) の連続版とみなすことができるため, これを**連続的単語集合** (continuous bag of words, **CBOW**) モデルといいます.

上では文脈として, 直前の $m=n-1$ 単語を用いましたが, 現実には単語の予測には, その後に続く文脈も使った方がよりよく w_t を予測できると考えられます. たとえば,

先月 自民党 は □

に続く□の単語には「国会」「ホームページ」「永田町」といった, いくつかの可能性がありますが, その後に続く単語が

先月 自民党 は □ に 法案 を

とわかっていれば、□の単語はほぼ「国会」と決定することができるでしょう。よって実際には、CBOW は図 3.31(a) のように、後続する単語 $w_{t+1}, w_{t+2}, \dots, w_{t+C}$ も使って計算します。すなわち式で書けば、

$$(3.92) \quad \begin{cases} \mathbf{x} = \sum_{i \in n(t)} \vec{w}_i & (n(t) = \{t-C, \dots, t-1, t+1, \dots, t+C\}) \\ \mathbf{p}_t = \text{Softmax}(\mathbf{A}\mathbf{y}) \end{cases}$$

という形になります。^{*77}

CBOW では文脈の単語の順番を考慮しないため、 n グラムモデルと異なり組み合わせ爆発(次元の呪い)が発生しないので、文脈の長さ C を大きくとることができます。通常は C は 10 以下の値を用いますが、 C を大きく取りすぎると「意味がぼやける」ため、適切な値に設定する必要があります。また、文脈として前後の単語ではなく、係り受け関係にある(離れている可能性もある)語を使うこともできます。その場合には単語ベクトルとして、より文法的な機能が重視されたベクトルが得られることになり、構文解析などに有用であることが確かめられています[83]。このように、使う「文脈」は、ベクトルの用途によって本来異なることに注意してください。通常のスキップグラムのように前後の単語を使う場合でも、 C が小さいほど文法的な関係が、 C が大きいほど意味的な関係が重視された単語ベクトルが学習されることになります。

(b) スキップグラム

上の CBOW では周囲の単語ベクトルの平均から、中心の単語 w_t を予測することを目的として単語ベクトルを学習していました。もう一つのアプローチとして、図 3.31(b) のように、 w_t から逆に周囲の単語を予測する、という方法も考えられます。つまり、

^{*77} なお、このことで n グラムと異なり、CBOW は単語列の生成モデルではなくなってしまいます。数学的には、周囲が与えられた場合の中心の語の確率が与えられている、イジングモデル[26, 31 章]と同様なマルコフ確率場を考えていることになります。生成モデルからマルコフ確率場への定式化の変化は、深層学習による GPT-3 のような現在の強力なマスク化言語モデルにも受け継がれています。

$$(3.93) \quad \prod_{t=1}^T \prod_{j \in n(t)} p(w_j | w_t)$$

を最大化することが目的となります。このとき、

$$(3.94) \quad \begin{cases} p(w_j | w_t) = \frac{\exp(\vec{w}'_j \cdot \vec{w}_t)}{Z} \\ Z = \sum_{v=1}^V \exp(\vec{v}'_j \cdot \vec{w}_t) \end{cases}$$

です。式(3.94)では、 w_t と周辺の語とのバイグラムを、必ずしも隣りあっていない、何語かスキップした場合も含めて考えるため、これを**スキップグラム** (skip-gram) といいます。なお、スキップグラムでは目的とする単語ベクトル \vec{w}_t と、周辺語のベクトル \vec{w}'_j は別のものを用います。すなわち、各単語 w について \vec{w} と \vec{w}' の二種類のベクトルを考えることとなります。

スキップグラムでは、CBOW と異なり、周辺語の単語ベクトルを \mathbf{x} に平均化することなく、 \vec{w}_t と \vec{w}'_j を直接比較して単語ベクトルを最適化するため、より単語ベクトルの表現力が高まることが期待されます。

ただし、ナイーブに式(3.94)を用いると、計算量が非常に大きくなってしまいます。これは、正規化定数である Z ^{*78} が、 \vec{w}_t 、 \vec{w}'_j が最適化で変わるとにすべて計算し直さねばならず、この和は語彙に含まれる V 個の単語すべてにわたる和で、通常 V は少なくとも 1 万、多いと 10 万からそれ以上にもなるからです。

この問題に対する解決法として、階層的 Softmax を使う方法と負例サンプリングを使う方法の 2 つがありますが、現在主に使われているのは後者ですので、以下でみていくことにしましょう。

負例サンプリング 式(3.94)のスキップグラムで単語ベクトルを効率的に学習する方法として、機械学習で用いられている Noise Contrastive Estimation [84] の考え方を用いて、

- 実際に出現した語 c の確率を大きく、かつ

*78 統計力学ではこうした和は分配関数とよばれ、慣習的に Z が用いられますので、ここでもそれを踏襲しています。この場合、 $\exp()$ の中の $\vec{w}'_j \cdot \vec{w}_t$ が単語ベクトルどうしの内積がもつエネルギー、すなわち「ハミルトニアン」ということとなります。

● それ以外の任意の語 c' の確率を小さく
 することが考えられます. すなわち,

$$(3.95) \quad L = \prod_{c \in n(w)} \sigma(\vec{c} \cdot \vec{w}) \times \prod_{c' \notin n(w)} (1 - \sigma(\vec{c}' \cdot \vec{w}))$$

を最大化することを考えます. ただし, 式(3.95)の第2項は w の周囲に出現したごく少数の語以外の, ほとんどすべての単語に関する積ですので, これを期待値で置き換えて, $1 - \sigma(x) = \sigma(-x)$ ですから

$$(3.96) \quad L = \prod_{c \in n(w)} \sigma(\vec{c} \cdot \vec{w}) \times \mathbb{E}_{c \sim p(c)} [\sigma(-\vec{c} \cdot \vec{w})]$$

としてみましょう. さらに後半の期待値を, 単語分布 $p^*(c)$ からの k 個のサンプルを用いたモンテカルロ積分 (65 ページ脚注) で置き換えて,

$$(3.97) \quad = \prod_{c \in n(w)} \sigma(\vec{c} \cdot \vec{w}) \times \frac{1}{k} \sum_{i=1}^k \sigma(-\vec{c}^{(i)} \cdot \vec{w}) \quad ; \quad c^{(i)} \sim p^*(c)$$

を最大化することにします. 1回の最適化では k 個の $c^{(1)}, \dots, c^{(k)}$ しか負例として考慮しませんが, 一般にこの最適化を繰り返すため, 学習の過程ではさまざまな単語 c が負例として登場することに注意してください. 負例, すなわちモンテカルロサンプルの数 k は, 小さいコーパスの場合は5~20個程度, 大きなコーパスでは2~5個程度でよいことが原論文で報告されています.

負例をサンプリングする分布 $p^*(c)$ としては, 一様分布 $1/V$ やユニグラム分布 $p(c)$ など, さまざまな可能性があります. ただし, 一様分布では「奏覧」や「卯原内」といった極端に稀な単語も同等に出現してしまいます^{*79}, いっぽうユニグラム分布を使うと, 3.2節でみたように「が」「の」といった, ごく一部の機能語ばかりが確率が高いためにサンプリングされることになってしまいます. 実験によれば, この間をとって, ユニグラム分布を $3/4$ 乗して平滑化した

^{*79} 3.2節の Zipf の法則によって, こうした頻度の低い語は非常に多く存在することに注意してください. 実際には, 語彙のほとんどはこうした語からなっているため, 一様分布からのサンプルはほとんどが稀な単語で占められてしまうことになります.

$$(3.98) \quad p^*(v) = \frac{p(v)^{3/4}}{Z} \quad \left(Z = \sum_{v=1}^V p(v)^{3/4} \right)$$

を使うとよい結果になることが確かめられています。

さらに、そもそも「の」「が」といった意味の薄い頻出語との共起を抑えるために、Word2Vec では

$$(3.99) \quad r(v) = \max \left(1 - \sqrt{\frac{t}{p(v)}}, 0 \right)$$

の確率で、共起の計算前に文から単語を削除するヒューリスティックが使われています。^{*80} ここで t は 10^{-5} 程度の閾値で、単語の確率 $p(v)$ がこれより低い単語は削除されず、これより高いと一定の割合で削除されることになります。この関数の形については、4章の図 4.3 を参照してください。たとえば、Brown コーパスでは $p(\text{the})=0.069$, $p(\text{hexagonal})=0.00000198$ なので、 $t=10^{-5}$ のとき

$$\begin{cases} r(\text{the}) = \max \left(1 - \sqrt{\frac{10^{-5}}{0.069}}, 0 \right) = 0.962 \\ r(\text{hexagonal}) = \max \left(1 - \sqrt{\frac{10^{-5}}{0.00000198}}, 0 \right) = \max(-1.247, 0) = 0 \end{cases}$$

となり、式(3.99)のルールでは“the”は96%の確率で文から削除されることになり、一方“hexagonal”はまったく削除されないことになります。

Word2Vec では、テキストの単語をこうしてあらかじめ確率的に削除してから計算を行っています。これにより、“the”のような単語が削除されて共起窓が実質的に広がるため、より意味を考慮した共起データを得ることができ、性能が改善することが確かめられています[85]。

^{*80} 原論文[63]では $p(v)$ は v の頻度となっており、その場合は式(3.99)は無意味な値となるため、誤っていることに注意が必要です。実際の `word2vec.c` や `gensim` の実装ではまた違った式が使われており、これに理論的な根拠があるわけではありません。より理論的な方法については、??節の SIF 単語重みの議論を参照してください。

$$\begin{array}{ccc}
 & C & \\
 W & \boxed{\tilde{\mathbf{X}}} & = \begin{array}{c} K \\ \boxed{\begin{array}{c} \vec{w}_1^T \\ \vec{w}_2^T \\ \vdots \\ \mathbf{W} \end{array}} \\ & & \end{array} \begin{array}{c} C \\ \boxed{\begin{array}{c} \vec{c}_1 \quad \vec{c}_2 \quad \dots \quad \mathbf{C}^T \end{array}} \\ K \end{array}
 \end{array}$$

図 3.32: 行列の特異値分解 (SVD) によるニューラル単語ベクトルの学習. Word2Vec で学習される単語ベクトルは, Shifted PPMI (本文参照) を要素とする行列 $\tilde{\mathbf{X}}$ を $\tilde{\mathbf{X}} \simeq \mathbf{W}\mathbf{C}^T$ と特異値分解して得られるベクトルと, 数学的に等価です.

3.5.4 Word2Vec と行列分解

前の節で計算したニューラル単語ベクトルは, 数学的には何を計算していることになるのでしょうか. イスラエルの Levy と Goldberg [86] は, Word2Vec で得られる単語ベクトルは, 以下に示すように, データから計算される**ある行列の主成分分析 (特異値分解) と数学的に等価である**ことを示しました. これにより, 実はニューラル単語ベクトルは, **行列の特異値分解で学習**することができます. ここからは行列の計算を頻繁に用いますので, 線形代数に慣れていない方は, 本章の文献案内を参照してください.

Word2Vec のスキップグラムで最大化している目的関数は, 式(3.95) でした. この式を最適化することで, 単語ベクトル \vec{v} と周辺語ベクトル \vec{c} が計算できますが, 式(3.95)で, \vec{v} と \vec{c} はつねに $\vec{v} \cdot \vec{c} = \vec{v}^T \vec{c}$ の形でしか現れません. すべての $\vec{v}^T \vec{c}$ の組み合わせは, 図 3.32 に示したように $\vec{w}_1^T, \vec{w}_2^T, \dots$ を縦に並べた行列を \mathbf{W} , また $\vec{c}_1, \vec{c}_2, \dots$ を横に並べた行列を \mathbf{C}^T とおいて

$$(3.100) \quad \mathbf{X} = \mathbf{W}\mathbf{C}^T$$

の要素になっています. よって, 式(3.95)は本質的にこの行列 \mathbf{X} を最適化しており, それを式(3.100)のように行列分解すれば, 単語ベクトルの行列 \mathbf{W} と周辺語ベクトルの行列 \mathbf{C} が得られる, ということがわかります. それでは, この行列 \mathbf{X} はどんな行列でしょうか.

\mathbf{X} の (v, c) 要素は内積 $\vec{v} \cdot \vec{c}$ ですから, これを $X(v, c) = x = \vec{v} \cdot \vec{c}$ とおく

と、ある x についての式(3.95)の目的関数 $L(x)$ は、

$$(3.101) \quad \log L(x) = n(w, c) \log \sigma(x) + k \cdot n(w) p(c) \log \sigma(-x)$$

となります。ここでシグモイド関数 $\sigma(x)$ の微分について、一般に

$$(3.102) \quad \begin{cases} \sigma(x)' = \left(\frac{1}{1+e^{-x}} \right)' = -\frac{-e^{-x}}{(1+e^{-x})^2} = \sigma(x)\sigma(-x) \\ \sigma(-x)' = \left(\frac{1}{1+e^x} \right)' = -\frac{e^x}{(1+e^x)^2} = -\sigma(x)\sigma(-x) \end{cases}$$

が成り立つことに注意しましょう。よって、

$$(3.103) \quad \begin{cases} \log \sigma(x)' = \frac{1}{\sigma(x)} \cdot \sigma(x)\sigma(-x) = \sigma(-x) \\ \log \sigma(-x)' = \frac{1}{\sigma(-x)} \cdot -\sigma(x)\sigma(-x) = -\sigma(x) \end{cases}$$

です。これを使えば、式(3.101)が最大になる点で x に関する勾配は0になりますから、 x で微分して0とおけば

$$(3.104) \quad \frac{\partial}{\partial x} \log L(x) = n(w, c) \sigma(-x) + k \cdot n(w) p(c) \cdot (-\sigma(x)) = 0$$

が成り立ちます。よって、式(3.89)のように $\sigma(-x) = 1 - \sigma(x)$ でしたから、

$$(3.105) \quad \begin{aligned} n(w, c)(1 - \sigma(x)) - k \cdot n(w) p(c) \sigma(x) &= 0 \\ \therefore \sigma(x) &= \frac{n(w, c)}{n(w, c) + k \cdot n(w) p(c)} \quad (\equiv y) \end{aligned}$$

となります。式(3.105)の2行目の右辺を y とおくと、

$$(3.106) \quad \sigma(x) = \frac{1}{1+e^{-x}} = y \quad \text{を解いて} \quad x = -\log \left(\frac{1}{y} - 1 \right)$$

を式(3.105)の2行目に代入して整理すれば、

$$(3.107) \quad x = -\log \left(\frac{n(w, c) + k \cdot n(w) p(c)}{n(w, c)} - 1 \right)$$

$$\begin{aligned}
&= \log \frac{n(w, c)}{k \cdot n(w) p(c)} = \log \frac{\frac{n(w, c)}{N}}{\frac{n(w)}{N} p(c)} - \log k \\
&= \log \frac{p(w, c)}{p(w) p(c)} - \log k
\end{aligned}$$

となることがわかります。すなわち、 x は w と c に対する、式(3.16)でも使った**自己相互情報量 (PMI)**

$$(3.108) \quad \text{PMI}(w, c) = \log \frac{p(w, c)}{p(w) p(c)}$$

を、 $\log k$ だけシフトしたものになっています。これから、Word2Vec で学習している式(3.100)の行列 \mathbf{X} は、

$$\begin{aligned}
(3.109) \quad X(w, c) &= \text{PMI}(w, c) - \log k \\
&= \log \frac{p(w, c)}{p(w) p(c)} - \log k
\end{aligned}$$

を要素とする行列だということがわかりました。特に、負例数 $k=1$ のときは $\log k = \log 1 = 0$ ですから、 \mathbf{X} は単に $\text{PMI}(w, c)$ を並べた行列となります。

ただし、式(3.109)の行列は、ほとんどを占める $p(w, c) = 0$ の場合に値が $\log p(w, c) = -\infty$ になってしまいます。要素がすべて埋まると計算量も増えてしまうため、実際に現れる (w, c) のペアのほとんどは $\text{PMI}(w, c) > 0$ である(だからこそ出現した)ことから、式(3.109)で値が非負の場合だけを考えて*81、

$$(3.110) \quad \tilde{X}(w, c) = \max \left(\log \frac{p(w, c)}{p(w) p(c)} - \log k, 0 \right)$$

を用いることにします。これを、SPPMI (Shifted Positive PMI) とよびます*82。

*81 PMI が負の場合を考えないことで、正確には目的関数は式(3.101)とは異なってしまいますが、実験的には以下で示すように、これは Word2Vec のかなり良い近似になっていることがわかっています。

*82 PPMI はこの研究で初めて現れたわけではなく、5章で扱う潜在意味解析(LSA)の文脈で、2007年に Bullinaria らによって提案され[87]、他のさまざまな単語ベクトルに比べて、PPMIによる単語ベクトルとコサイン距離による比較が最も高い性能を持つことが報告されていました。2012年にはSVDによる次元圧縮も試みられており[88]、Word2Vecと本質的に同じベクトルが、深層学習より以前にすでに提案されていたといえます。

$$\begin{array}{c}
 \begin{array}{ccc}
 & C & \\
 W & \tilde{\mathbf{X}} & \\
 & & K
 \end{array}
 =
 \begin{array}{ccc}
 & K & \\
 & \mathbf{U} & \\
 & & K
 \end{array}
 \begin{array}{ccc}
 & K & \\
 & \mathbf{\Sigma} & \\
 & & C
 \end{array}
 \begin{array}{ccc}
 & C & \\
 & \mathbf{V}^T & \\
 & & K
 \end{array}
 \\
 \simeq
 \begin{array}{ccc}
 & K & \\
 W & \mathbf{U}_K & \\
 & & C
 \end{array}
 \begin{array}{ccc}
 & \mathbf{\Sigma}_K & \\
 & & \mathbf{V}_K^T & \\
 & & & K
 \end{array}
 \end{array}$$

図 3.33: SPPMI 行列 $\tilde{\mathbf{X}}$ の特異値分解 (SVD) による次元削減の様子.

SPPMI を並べた行列 $\tilde{\mathbf{X}}$ は、カウントが 0 のエントリ (w, c) および式 (3.109) が負になるエントリはすべて値が 0 になるため、非常に疎になることに注意してください。実際、ja.text8.txt で窓幅 10 の場合、非 0 のエントリは 1% (20872503/2027160576) にすぎず、99% のエントリが 0 になりました。 $\tilde{\mathbf{X}}$ から式 (3.100) の \mathbf{W}, \mathbf{C} を求めるには、疎行列に対する**特異値分解** (Singular Value Decomposition, SVD) を使うと高速に計算することができます。SVD は行列の対角化を正方行列でない場合に拡張したもので、行列 $\tilde{\mathbf{X}}$ を図 3.33 の上段のように

$$(3.111) \quad \tilde{\mathbf{X}} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

と 3 つの行列 $\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}^T$ の積に分解します。ここで \mathbf{U}, \mathbf{V} はそれぞれ $\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T$ および $\tilde{\mathbf{X}}^T\tilde{\mathbf{X}}$ の固有ベクトル、 $\mathbf{\Sigma}$ は対応する固有値の平方根 (特異値) $\sigma_1, \sigma_2, \dots, \sigma_r$ (r は $\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T, \tilde{\mathbf{X}}^T\tilde{\mathbf{X}}$ のランク) を並べた対角行列

$$(3.112) \quad \mathbf{\Sigma} = \begin{pmatrix} \sigma_1 & & & \mathbf{0} \\ & \sigma_2 & & \\ & & \ddots & \\ \mathbf{0} & & & \sigma_r \end{pmatrix}$$

です。このうち特異値の大きい方から上位 K 個をとって、図 3.33 の下段のように

$$(3.113) \quad \tilde{\mathbf{X}} \simeq \mathbf{U}_K \boldsymbol{\Sigma}_K \mathbf{V}_K^T$$

とする近似は、 $\tilde{\mathbf{X}}$ との二乗誤差を最小にする近似になっています[]。式(3.113)は、 $\boldsymbol{\Sigma}$ は対称行列なので $\boldsymbol{\Sigma}^T = \boldsymbol{\Sigma}$ で、一般に $(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$ ですから、

$$(3.114) \quad \begin{aligned} \tilde{\mathbf{X}} &\simeq \mathbf{U}_K \boldsymbol{\Sigma}_K^{\frac{1}{2}} \boldsymbol{\Sigma}_K^{\frac{1}{2}} \mathbf{V}_K^T \\ &= \left(\mathbf{U}_K \boldsymbol{\Sigma}_K^{\frac{1}{2}} \right) \left(\mathbf{V}_K \boldsymbol{\Sigma}_K^{\frac{1}{2}} \right)^T \quad (= \mathbf{WC}^T) \end{aligned}$$

とも書くことができます。よって、式(3.100)と見比べて

$$(3.115) \quad \mathbf{W} = \mathbf{U}_K \boldsymbol{\Sigma}_K^{\frac{1}{2}}, \quad \mathbf{C} = \mathbf{V}_K \boldsymbol{\Sigma}_K^{\frac{1}{2}}$$

とおけば、 \mathbf{W}, \mathbf{C} を得ることができます*83。ここで $\boldsymbol{\Sigma}_K^{\frac{1}{2}}$ は、 $\boldsymbol{\Sigma}$ の対角成分の平方根をとった行列

$$(3.116) \quad \boldsymbol{\Sigma}^{\frac{1}{2}} = \begin{pmatrix} \sqrt{\sigma_1} & & & \mathbf{0} \\ & \sqrt{\sigma_2} & & \\ & & \ddots & \\ \mathbf{0} & & & \sqrt{\sigma_r} \end{pmatrix}$$

です。Python では、

```
from scipy.sparse.linalg import svds
from pylab import *
U,S,V = svds(X)
W = np.dot (U, diag(sqrt(S)))
C = np.dot (V, diag(sqrt(S)))
```

のように実行すれば、 \mathbf{W}, \mathbf{C} を計算することができます。

単語ベクトルの計算 式(3.108)の PMI は、式(3.107)から、頻度 $n(w, c)$, $n(w)$, $n(c)$ を用いて

*83 行列 $\tilde{\mathbf{X}}$ の各行、すなわち各単語の共起情報を格納したベクトルの間の内積は $\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T$ で計算することができます。 $\tilde{\mathbf{X}} \simeq \mathbf{U}_K \boldsymbol{\Sigma}_K \mathbf{V}_K^T$ ですから、これは $\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T = \mathbf{U}_K \boldsymbol{\Sigma}_K \mathbf{V}_K^T (\mathbf{U}_K \boldsymbol{\Sigma}_K \mathbf{V}_K^T)^T = \mathbf{U}_K \boldsymbol{\Sigma}_K \mathbf{V}_K^T \mathbf{V}_K \boldsymbol{\Sigma}_K^T \mathbf{U}_K^T = (\mathbf{U}_K \boldsymbol{\Sigma}_K) (\mathbf{V}_K \boldsymbol{\Sigma}_K)^T$ となり、 $\mathbf{W} = \mathbf{U}_K \boldsymbol{\Sigma}_K, \mathbf{C} = \mathbf{V}_K \boldsymbol{\Sigma}_K$ とするのが $\tilde{\mathbf{X}}$ の各行間の内積を保存する意味では理論的に最適です。しかし、式(3.114)のように直接行列分解を近似する方が、さまざまな意味的タスクにおいて精度が高いことが確かめられています[85]。

$$(3.117) \quad \text{PMI}(w, c) = \log \frac{p(w, c)}{p(w)p(c)} = \log \frac{p(w, c)/N}{n(w)/N \cdot n(c)/N} \\ = \log n(w, c) - \log n(w) - \log n(c) + \log N$$

で計算することができます。式(3.110)から、この値が $\log k$ より大きくないと値が0になり、文脈語 c が観測されなかったことになってしまいますので、[85]では実験の結果、 $k=1$ すなわち $\log k=0$ で PMI をそのまま使う方が意味的タスクにおいて高性能となることが確かめられています。また式(3.108)は、条件つき確率の定義から

$$(3.118) \quad \text{PMI}(w, c) = \log \frac{p(w, c)}{p(w)p(c)} = \log \frac{p(c|w)}{p(c)}$$

とも書くことができることに注意してください。すなわち $\text{PMI}(w, c)$ は、「 w の周辺に c が現れる確率」と「 c が一般的に出現する確率」の比の対数になっています*84。このとき、 c が稀な単語で $p(c)$ が非常に小さいと、PMIの値が非常に大きくなってしまいます。これを避けるため、 $p(c)$ としては負例サンプリングでも用いた、式(3.98)で平滑化した $p^*(c)$ を用いるとよいでしょう[85]。あらかじめ式(3.98)の対数 $\log p^*(c)$ を計算しておけば、平滑化された式(3.118)は

$$(3.119) \quad \log \frac{p(c|w)}{p^*(c)} = \log \frac{n(w, c)/n(w)}{p^*(c)} = \log n(w, c) - \log n(w) - \log p^*(c)$$

で求めることができます。

こうして計算した単語ベクトルの例を、図??に示しました。この計算は、サポートページの `pmivector.py` を使って、

```
% pmivector.py K text model
```

のようにして実行することができます。なお、共起行列 $\tilde{\mathbf{X}}$ が巨大になる場合は、ランダム射影を用いて SVD を効率的に行うことのできる `redsvd []` のような方法を使えば、データが大きくても高速に単語ベクトルを計算することができます。

自己相互情報量と言語モデル[▽] SGNS は $\text{PMI}(w, c)$ を近似していることがわかりましたが、言語の生成モデルとしては、これは何を意味しているのでしょうか

*84 これは、対数尤度比ともよばれています。

確率	$k = \text{飼育}$	$k = \text{彼女}$	$k = \text{動物}$	$k = \text{本陣}$
$p(k \text{犬})$	0.00479	0.00019	0.00268	0.00019
$p(k \text{猫})$	0.00083	0.00103	0.00186	0.00021
$\frac{p(k \text{犬})}{p(k \text{猫})}$	5.803	0.186	1.444	0.929

表 3.6: 単語 k が周辺に現れる確率とその比。「飼育」や「彼女」はそれぞれ「犬」および「猫」の周辺に現れる確率が高く、比は 1 より大または小となりますが、両方に関係する「動物」およびどちらにも関係しない「本陣」では、確率の比はほぼ 1 になります。

か. 式(3.107)でみたように

$$(3.120) \quad \text{PMI}(w, c) = \log \frac{p(w, c)}{p(w)p(c)} = \log \frac{p(c|w)}{p(c)}$$

ですから、両辺を指数の肩にのせれば、

$$(3.121) \quad \exp(\text{PMI}(w, c)) = \frac{p(c|w)}{p(c)}$$

$$\therefore p(c|w) = p(c) \cdot \exp(\text{PMI}(w, c))$$

となります。式(3.121)は、単語 w の周りに文脈語 c が現れる条件つき確率 $p(c|w)$ は、 c のデフォルトの出現確率 $p(c)$ に、係数 $e^{\text{PMI}(w, c)}$ をかけたものになることを表しています。PMI(w, c)=0, すなわち式(3.108)より $p(w, c)=p(w)p(c)$ で w と c が無相関の場合は、 c の条件つき確率は $p(c) \cdot e^0 = p(c) \cdot 1 = p(c)$ でデフォルトの確率に等しく、PMI(w, c) > 0 なら確率は $e^{\text{PMI}(w, c)} > 1$ なので $p(c)$ より大きく、PMI(w, c) < 0 なら確率は $e^{\text{PMI}(w, c)} < 1$ なので $p(c)$ より小さくなるわけです。このように、自己相互情報量は確率を式(3.121)のように「変調」する係数を表しており、これは統計学ではCox 比例ハザード[89]とよばれるモデルと同じモデルだといえます。

3.5.5* GloVe と意味方向の数理

こうして得られた単語ベクトルの間には、3.5.1節で示した「引き算」の関係

が成り立つことが知られていますが、これはなぜなのでしょう。スタンフォード大学で素粒子物理学を専攻していた Pennington は 2014 年に、こうした意味的關係が成り立つように設計された単語ベクトル、GloVe (Global Vector) [90] を提案しました。GloVe の単語ベクトルは Word2Vec の単語ベクトルと似ていますが、目的関数が少し異なり、若干違うベクトルになっています。

単語ベクトルについて要求される条件は、「意味が似ているベクトルは値が近い」ということです。これを、確率の言葉で表現するとどうなるでしょうか。

まず、ある単語 w の周囲に単語 c が現れる確率 $p(c|w)$ は、前節で用いた窓内の共起頻度 $n(w, c)$ を使って、

$$(3.122) \quad p(c|w) = \frac{n(w, c)}{n(w)}$$

で計算できることに注意しましょう。いま、2つの単語 $i = \text{“犬”}$ と $j = \text{“猫”}$ があつたとき、別の単語 k がこれらの単語の周囲に現れる確率の比

$$(3.123) \quad \frac{p(k|i)}{p(k|j)}$$

をさまざまな k について計算すると、表 3.6 のようになります。これから、図 3.34 に示したように

- 犬に強く関係する単語 k 、たとえば“飼育”は $p(k|\text{犬})$ の方が $p(k|\text{猫})$ よりも大きく、

$$\frac{p(\text{飼育}|\text{犬})}{p(\text{飼育}|\text{猫})} = 5.803 \gg 1$$

のように、比が 1 よりもずっと大きくなります。

- 逆に、猫に強く関係する単語 k 、たとえば“彼女”は $p(k|\text{猫})$ の方が $p(k|\text{犬})$ より大きく、

$$\frac{p(\text{彼女}|\text{犬})}{p(\text{彼女}|\text{猫})} = 0.186 \ll 1$$

のように、比が 1 よりずっと小さくなっています。

- 一方、犬とも猫とも関係のある単語“動物”、および関係のない単語“本陣”は、

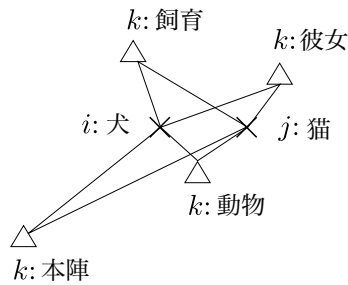


図 3.34: GloVe の用いる単語ベクトル間の関係. それぞれの単語 k が単語 i, j の周辺に現れる確率 $p(k|i), p(k|j)$ の比 $p(k|i)/p(k|j)$ から, 単語ベクトルを求めます.

$$\frac{p(\text{動物}|\text{犬})}{p(\text{動物}|\text{猫})} = 1.444 \approx 1$$

$$\frac{p(\text{本陣}|\text{犬})}{p(\text{本陣}|\text{猫})} = 0.929 \approx 1$$

と, どちらも 1 に近くなっています (ただし, 表 3.6 にみるように分子・分母の確率の絶対値は関係があれば大きく, なければ小さくなります).

つまり一般に, 図 3.34 のような単語の間の意味的關係は

$$(3.124) \quad \frac{p(k|i)}{p(k|j)} \equiv \frac{p_{ik}}{p_{jk}}$$

の値で示される, ということがわかります. ここで簡単のため, $p(k|i) = p_{ik}$, $p(k|j) = p_{jk}$ と表記しました. 式 (3.124) は確率の比, すなわち相対値を見ているため, 個々の確率の絶対値 p_{ik}, p_{jk} が単語の頻度によって変わっても不変な値であることに注意してください. したがって, これから求める単語ベクトル $\vec{w}_i, \vec{w}_j, \vec{w}_k$ から式 (3.124) が計算されるべきですから, その関数を一般に F とおけば,

$$(3.125) \quad F(\vec{w}_i, \vec{w}_j, \vec{w}_k) = \frac{p_{ik}}{p_{jk}}$$

が計算できて図 3.34 の関係が成り立つべきだ, ということになります.

関数 F にはいろいろな可能性がありますが, 3.5.2 節で示した「ベクトルの引き算」が成り立つためには, 最も単純には F は差 $\vec{w}_i - \vec{w}_j$ に依存する関数と

なるべきでしょう。また、式(3.125)の右辺はスカラー値ですから、線形な構造を保存してベクトルの次元の間に無用な相関を持ち込まないためには、 $\vec{w}_i - \vec{w}_j$ と \vec{w}_k の内積をとればよいと考えられます。よって、式(3.125)は

$$(3.126) \quad F\left((\vec{w}_i - \vec{w}_j)^T \vec{w}_k\right) = \frac{p_{ik}}{p_{jk}}$$

と書き換えることができました。

ここで、式(3.126)の右辺は確率の比のため、必ず正であることに注意してください。よって F は必ず正の値を返す必要がありますが、左辺の $()$ の中は $\vec{w}_i^T \vec{w}_k - \vec{w}_j^T \vec{w}_k$ と差になっており、負になる可能性がありますから、一般に実数の和や差を正の数の積や割り算に変換できる F を考えれば^{*85}、式(3.126)はさらに

$$(3.127) \quad F\left((\vec{w}_i - \vec{w}_j)^T \vec{w}_k\right) = \frac{F\left(\vec{w}_i^T \vec{w}_k\right)}{F\left(\vec{w}_j^T \vec{w}_k\right)} = \frac{p_{ik}}{p_{jk}}$$

と書き換えることができます。この式を満たす解として $F(x) = \exp(x)$ があります。したがって

$$(3.128) \quad \exp\left(\vec{w}_i^T \vec{w}_k\right) = p_{ik}$$

すなわち

$$(3.129) \quad \vec{w}_i^T \vec{w}_k = \log p_{ik} = \log p(k|i)$$

が成り立てばよい、ということがわかります。頻度を使って書き換えれば、

$$(3.130) \quad \vec{w}_i^T \vec{w}_k = \log n(i, k) - \log n(i)$$

となります。式(3.129)が、これまでに述べた意味的な関係から単語ベクトルが満たすべき十分条件ということになります。

ここで、式(3.129)の左辺は i と k を入れ替えても成り立ちますが、右辺は条件つき確率なので、 i と k を逆にすると意味が変わってしまうことに注意してください。そこで $\log n(i)$ に対応するバイアス項 b_i を、対称性から同様に b_k を導

^{*85} 原論文では、 F として群 $(\mathbb{R}, +)$ から $(\mathbb{R}_{>0}, \times)$ への準同型写像をとると説明されています。

入すれば,

$$(3.131) \quad \begin{aligned} \vec{w}_i^T \vec{w}_k &= \log n(i, k) - b_i - b_k \\ \therefore \vec{w}_i^T \vec{w}_k + b_i + b_k &= \log n(i, k) \end{aligned}$$

となり, i と k について対称な目的関数を得ることができました *86.

ただし, 単純に式(3.131)の回帰モデルを求めるには問題があり, $n(i, k) = 0$ の場合や, $n(i, k) > 0$ でも, 偶然共起した場合をすべて説明しなければならなくなってしまう。そこで, 式(3.131)の両辺の二乗誤差を, $n(i, k)$ に依存する関数 $f(n(i, k))$ で重みづけて, GloVe では

$$(3.132) \quad J = \sum_{i,k=1}^V f(n(i, k)) \left(\vec{w}_i^T \vec{w}_k + b_i + b_k - \log n(i, k) \right)^2$$

を最小化します。 $f(x)$ は $x=0$ のとき 0 で $x>0$ のとき少しずつ大きくなり, x_{\max} で頭打ちになる関数で (GloVe では $x_{\max}=100$ を用いています),

$$(3.133) \quad f(x) = \begin{cases} \left(\frac{x}{x_{\max}} \right)^\alpha & x < x_{\max} \\ 1 & \text{それ以外} \end{cases}$$

という, 図 3.35 のような関数を用いると性能が良かったことが報告されていま

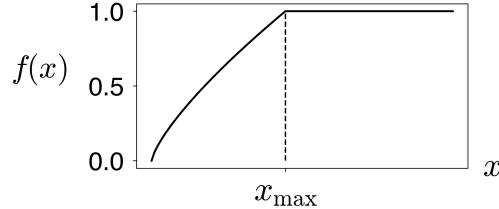
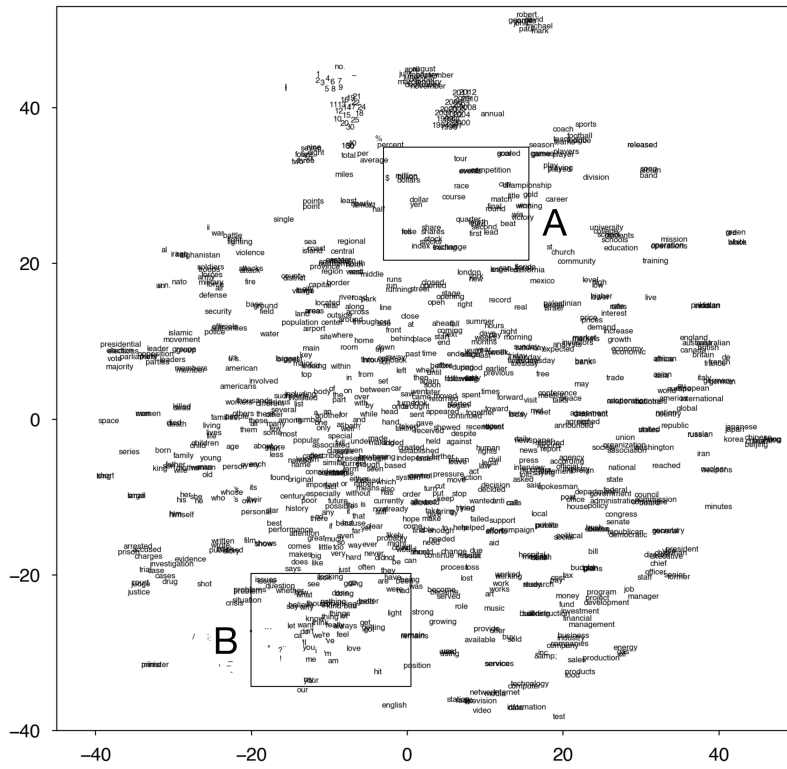
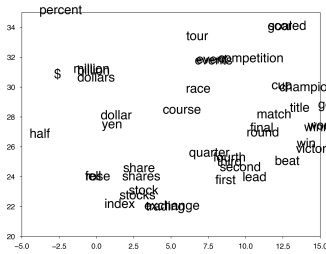


図 3.35: GloVe で用いている頻度 x の重みづけ関数 $f(x)$. $x=0$ では重み 0 に, x_{\max} 以上ではすべて重みが 1 になります。

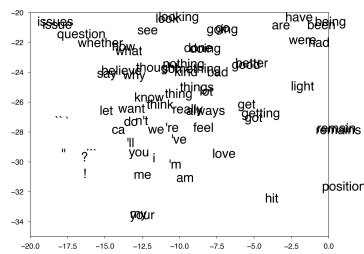
*86 GloVe の真の目的関数は式(3.129)ですが, 最適化のために導入しているこうしたヒューリスティックが, GloVe の性能を本来のものより下げている可能性があります。直接, 式(3.129)を MCMC 法などで最適化することを考えてみてよいでしょう。



(a) 全体図. 数字や年号, 人名のクラスタや経済用語の集まっている場所などがあります.



(b) A を拡大したもの



(c) B を拡大したもの

図 3.36: GloVe で 6 億語のテキストから学習された頻度上位 1000 語の単語ベクトルの可視化. 100 次元空間の単語ベクトルを, t-SNE で 2 次元に可視化しています.

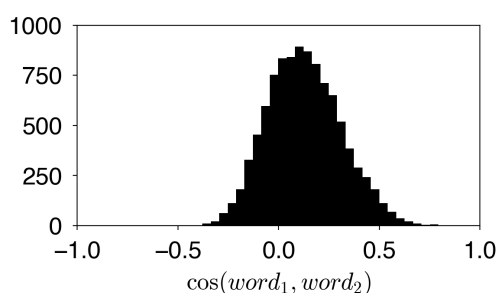


図 3.37: ja.text8 コーパスから計算した単語ベクトルでの、ランダムな 2 単語間のコサイン類似度の分布. 類似度は $[-1, 1]$ ではなく、その一部にだけ分布しています.

す. なお、本来は \vec{w}_i と \vec{w}_k の役割は同じですので、最終的な単語ベクトルとしては平均をとった $(\vec{w}_i + \vec{w}_k)/2$ が採用されています.

GloVe を使っても、word2vec とほとんど同様の単語類似度や比例関係を計算することができます (→演習問題 7). GloVe では、スタンフォード大学の公式ページで、Wikipedia 全体のような巨大なテキストから計算した英語の単語ベクトルが公開されています^{*87}. このうち、100 次元の単語ベクトルの頻度上位 1000 語を、t-SNE とよばれる非線形な可視化法で 2 次元に可視化したものを図 3.36 に示しました^{*88}.

3.5.6 単語ベクトルの分布とノルム

ここまで単語ベクトルについてみてきましたが、それでは、学習された単語ベクトルは空間上でどのように分布しているのでしょうか. 以下では、Word2Vec を使って日本語の ja.text8 コーパスから学習された単語ベクトルについて分析していくことにします.

単語ベクトル \vec{w} と \vec{v} の類似度は、137 ページでみたようにそのなす角のコサインとして、

^{*87} <https://nlp.stanford.edu/projects/glove/> からダウンロードすることができます.

^{*88} この単語ベクトルの可視化は、サポートサイトの visualize.py を使うと行うことができます. 詳しくは、スクリプトの中身をご覧ください.

$$(3.134) \quad \cos(\vec{w}, \vec{v}) = \frac{\vec{w} \cdot \vec{v}}{|\vec{w}| |\vec{v}|}$$

で測ることができます. 任意の実数 x について $-1 \leq \cos x \leq 1$ ですから, 式(3.134)は -1 から 1 の範囲に分布しているはずですが, 次のようにランダムな2単語について式(3.134)を計算してみると, 図3.37のように最大値も最小値もまったく ± 1 ではなく, 平均値も 0 ではないことがわかります.

```
from numpy.random import randint
from pylab import *
vectors = vload("model/ja.text8.vec")
words = list(vectors.keys())
N = 10000
V = len(words)
s = np.zeros (N, dtype=float)
for n in range(N):
    i = randint(V); j = randint(V)
    s[n] = similar(vectors, words[i], words[j])
hist(s,bins=30)
axis([-1,1,0,1000])
```

つまり, 単語ベクトルは K 次元空間 (ここでは 100 次元空間) 上に, かなり偏って分布しているということです. 実際, 単語ベクトルの中心を次のようにして計算してみると, まったく 0 ではないことがわかります.

```
def loadvec (file):
    matrix = []
    with open (file, 'r') as fh:
        for line in fh:
            tokens = line.rstrip('\n').split()
            if len(tokens) > 2: # Word2Vec のヘッダをスキップ
                matrix.append (np.array (list ( \
                    map (float, tokens[1:]))))
    return np.array(matrix)
matrix = loadvec("model/ja.text8.vec")
np.mean(matrix, 0)
⇒ array([ 0.008, -0.003, -0.051, -0.025,  0.133,  0.147, -0.057,
          0.009,  0.126,  0.163,  0.027, -0.017, -0.037, -0.072, ...
          -0.107,  0.046, -0.170,  0.072])
```

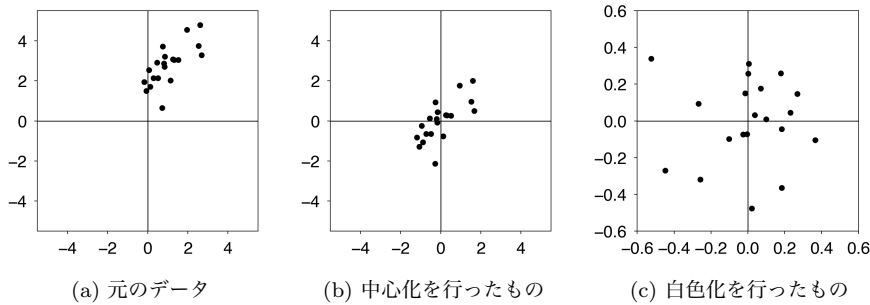


図 3.38: ベクトルの中心化と白色化. 中心化により平均が 0 に, さらに白色化により共分散が単位行列 \mathbf{I} になります.

単語ベクトルの中心化 よって, 最も自然に考えられる前処理は, 単語ベクトルから上の中心を引いて, 平均を 0 にすることでしょう (図 3.38(b)).

$$(3.135) \quad \vec{w}' = \vec{w} - \mu \quad ; \quad \mu = \frac{1}{V} \sum_{v=1}^V \vec{v}$$

これは, データ解析では一般に, ベクトルの**中心化** (centering) とよばれています.

ただし実は, 式(3.135)のようなナイーブな中心化をしてしまうと, 単語ベクトルの表現力はむしろ悪くなってしまうことが知られています. 3.2 節で議論したように, 言語には大量の単語があり, そのほとんどはきわめて低確率なのでした. 単純な平均をとると, 「日本」「人権」のような頻度の高い重要な単語ベクトルと, 「鼠小僧次郎吉」「ザンペリーニ」のようなめったに出ない特殊な単語ベクトルの平均をとることになってしまいます.

この問題に対し, 横井ら[91]は, 一様ではなく単語の確率で期待値をとって平均を計算することで, 単語ベクトルの性能が改善することを示しました. すなわち式(3.135)の代わりに,

$$(3.136) \quad \vec{w}' = \vec{w} - \mu \quad ; \quad \mu = \sum_{v=1}^V p(v) \vec{v}$$

として単語ベクトルを中心化します. $p(v)$ は, 単語 v の出現確率です. サポートサイトに, この前処理を行うスクリプト `vcenter.py` を示しました. この処理に

は単語のユニグラム確率が必要ですので、実行には

```
% vcenter.py ja.text8.vec ja.text8 output.vec
```

のように、単語ベクトルの学習に使ったコーパスを指定します。なお、単語ベクトルの学習によく使われる `text8` や日本語版の `ja.text8` には改行がなく、巨大な 1 行のテキストになっているため、2.2 節のような方法でテキストを行ごとに読み出すことはできません。こうした場合は、Python の `yield` を使って

```
import re

def readword (fh, newline=r'[\t\n]+'):
    buf = "" # 単語の読み出しバッファ
    while True: # ファイルが終わるまで
        while True: # バッファから単語を読み出す
            match = re.search (newline, buf)
            if not match:
                break
            else:
                yield buf[:match.start()]
                buf = buf[match.end():]
        chunk = fh.read (4096) # ファイルを一定量読む
        if not chunk: # ファイルの終了
            if len(buf) > 0:
                yield buf
            break
        buf += chunk # バッファに読んだ分を追加
```

のようなジェネレータを定義すれば、

```
with open('ja.text8', 'r') as fh:
    for word in readword(fh):
        ...
```

のように単語を次々と読み出すことができます。上の `vcenter.py` は、こうして単語のユニグラム確率を計算しています。

単語ベクトルの白色化 単語ベクトルの平均は中心化によって (期待値の意味で) 0 になりますが、単語ベクトルの空間は図 3.38(b) のように、それでもまだ歪みを残しているはずで、こうした場合、データ解析でよく行われる前処理が **白色化** (whitening) です。白色化は図 3.38(b) を図 3.38(c) のように線形射影によって変換する処理で、平均を 0 にした上で、さらに異なる次元間の共分散を 0

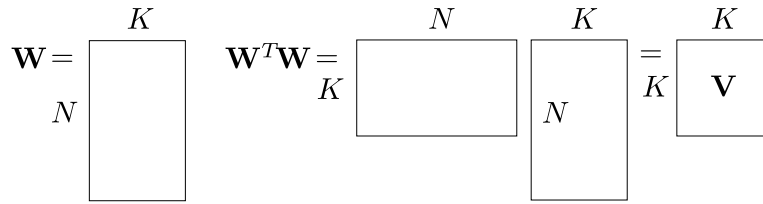


図 3.39: 単語ベクトルの行列 \mathbf{W} とその共分散行列 $\mathbf{V} = \mathbf{W}^T \mathbf{W}$ の計算.

にして共分散行列を単位行列にします。つまり、この処理により K 次元空間の次元が「無駄なく」使われるようになります。

ベクトルの白色化は、一般に次のようにして行うことができます。いま、各データが K 次元の行ベクトルで、 N 個のデータが図 3.39 左のように、 $N \times K$ の行列 \mathbf{W} をなしているとします^{*89}。 \mathbf{W} を中心化して行方向の平均を $\mathbb{E}[\mathbf{W}] = \mathbf{0}$ にしたとき、 K 次元の各次元間の共分散行列は、図 3.39 右のように $\mathbf{V} = \mathbf{W}^T \mathbf{W}$ で計算することができます。 $(\mathbf{W}^T \mathbf{W})^T = \mathbf{W}^T \mathbf{W}$ より \mathbf{V} は対称行列ですから、 \mathbf{V} は

$$(3.137) \quad \mathbf{V} = \mathbf{P} \mathbf{\Sigma} \mathbf{P}^{-1}$$

と対角化することができます。ここで $\mathbf{\Sigma}$ は \mathbf{V} の固有値 $\lambda_1, \dots, \lambda_K$ を対角成分にもつ対角行列、 \mathbf{P} は対応する固有 (行) ベクトルを列方向に並べた行列です。したがって、 $\mathbf{P}^T \mathbf{P} = \mathbf{I}$ よって $\mathbf{P}^T = \mathbf{P}^{-1}$ になります。このとき、

$$(3.138) \quad \mathbf{Z} = \mathbf{W} \mathbf{P} \mathbf{\Sigma}^{-1/2}$$

と \mathbf{W} を \mathbf{Z} に変換すれば^{*90}、その共分散は式 (3.137) から、

$$(3.139) \quad \begin{aligned} \mathbf{Z}^T \mathbf{Z} &= (\mathbf{W} \mathbf{P} \mathbf{\Sigma}^{-1/2})^T \mathbf{W} \mathbf{P} \mathbf{\Sigma}^{-1/2} \\ &= \mathbf{\Sigma}^{-1/2} \mathbf{P}^T \mathbf{W}^T \mathbf{W} \mathbf{P} \mathbf{\Sigma}^{-1/2} = \mathbf{\Sigma}^{-1/2} \underbrace{\mathbf{P}^T (\mathbf{P} \mathbf{\Sigma} \mathbf{P}^{-1})}_{\mathbf{I}} \mathbf{P} \mathbf{\Sigma}^{-1/2} \\ &= \mathbf{\Sigma}^{-1/2} \mathbf{\Sigma} \mathbf{\Sigma}^{-1/2} = \mathbf{I} \end{aligned}$$

になります。すなわち、 \mathbf{W} を式 (3.138) で線形変換した \mathbf{Z} は平均 $\mathbb{E}[\mathbf{Z}] = \mathbf{0}$ 、共

*89 数学ではベクトルは列ベクトルが基本ですが、NumPy の行列は標準ではデータを行方向に格納するため (row major といいます)、あえてこの表現を使っています。

*90 $\mathbf{\Sigma}^{-1/2}$ は、 $1/\sqrt{\lambda_1}, \dots, 1/\sqrt{\lambda_K}$ を対角成分にもつ対角行列です。

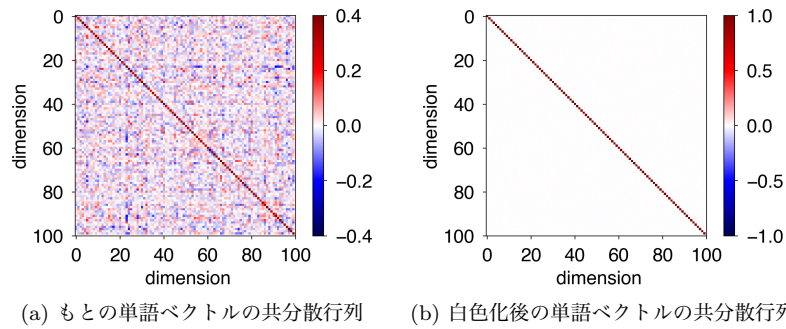


図 3.40: 単語ベクトルの共分散行列と白色化. 白色化により, 単語ベクトルの各次元を無相関にすることができます.

分散 $\mathbb{V}[\mathbf{Z}] = \mathbf{I}$ で分布するようになります. これを**白色化**といいます. □

したがって, 単語ベクトルを並べた行列 \mathbf{W} についてこの白色化を行えば, 図 3.38(c) のようにデータは K 次元空間で無駄なく分布するようになり, 細かい意味の差を反映するようになると考えられます. ただし, ここでも行列 \mathbf{W} の各行に対応する単語には大きな頻度の差があるため, [91]では式(3.136)の期待値を用いた中心化を行った上で, \mathbf{W} の各行 \vec{w}_i を, 長さを 1 に規格化した上で $\sqrt{p(w_i)}$ で重みづけた

$$(3.140) \quad \vec{w}_i' = \sqrt{p(w_i)} \frac{\vec{w}_i}{|\vec{w}_i|}$$

としてから白色化を行うことで, 単語ベクトルの性能が向上することを示しています. 白色化の処理は同様ですので, 結果として新しい単語ベクトルはやはり平均 $\mathbf{0}$, 共分散 \mathbf{I} で分布することになります. この処理を, 期待白色化とよみましょう.

サポートサイトに, この期待白色化を行うスクリプト `vwhiten.py` を示しました.

```
% vwhiten.py ja.text8.vec ja.text8 whitened.vec
```

を実行すると, `whitened.vec` に期待白色化された単語ベクトルが保存されます.

図 3.40 に, もとの単語ベクトルの共分散行列 \mathbf{V} と, 期待白色化後の共分散行列 \mathbf{V} を示しました. もとの単語ベクトルでは各次元の間に \pm の多くの相関が

ありますが、変換後はそれが一掃され、単語ベクトルの各次元はすべて無関係になっている(共分散が単位行列になっている)ことがわかります。

この新しい単語ベクトルで単語の類似度を計算してみると、下のようになります。139 ページの結果と比べてみると、“太陽”の類似語の上位から“シリウス”や“星雲”が外れていたり、“少年”の類似語の順番が入れ替わっていたりなど、微妙にはありますが、確実に類似度が改善していることがわかります。数値的にも、単語類似度や特に文類似度などのタスクの性能で改善がみられることが報告されています[91]。

similar(whitened, "太陽")	similar(whitened, "少年")
⇒ 恒星 → 0.7570	⇒ 少女 → 0.8279
太陽系 → 0.7068	小学生 → 0.6740
星 → 0.6908	中学生 → 0.6734
土星 → 0.6870	高校生 → 0.6566
地球 → 0.6804	青年 → 0.5830
惑星 → 0.6752	同級生 → 0.5791
質量 → 0.6711	幼児 → 0.5754
天体 → 0.6626	忍者 → 0.5253
海王星 → 0.6590	中学 → 0.5179
銀河 → 0.6485	大人 → 0.4748
超新星 → 0.6480	成人 → 0.4667

単語ベクトルの長さ これまでは、単語ベクトルの類似度や比例関係はすべて、長さを1に規格化した上で行ってきました。しかし実際は、単語ベクトルはその方向だけでなく、長さの情報も持っています。単語ベクトルの長さは、どうなっているのでしょうか。

単語ベクトルの長さ $|\vec{w}|$ を縦軸に、その単語の出現確率 $p(w)$ の対数を横軸にとってプロットしてみると、図 3.41 のように、確率が中程度の単語ベクトルの長さが最も大きいことが知られています[92]。確率的勾配法で単語ベクトルを最適化するとき、頻度の高い語のベクトルは様々な方向から「引っ張られる」ため、結果的に長さは短くなるでしょう。また頻度の低い語は逆にわずかしか引っ張られませんが、単語ベクトルは短いままになると考えられます。

ただし、こうした頻度の影響を差し引いた上では、意味が「強い」単語のベクトルはノルムが大きい傾向にあることが報告されています[93]。ここで単語 w の意味が「強い」とは、 w の周囲に出現する単語 c の確率分布 $\{p(c|w)\}$ と、単

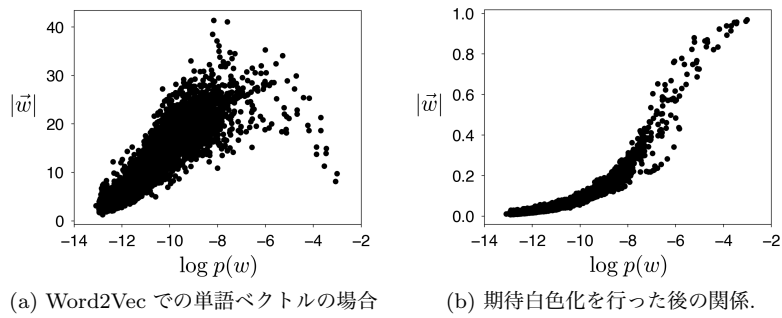


図 3.41: 単語の確率 (横軸) と単語ベクトルの長さ (縦軸) の関係.

語 c の平均的な確率分布 $\{p(c)\}$ との KL ダイバージェンス (式 (2.69)) が大きい, すなわち周辺語の確率分布が偏っていることを意味します.

なお, 単語ベクトルの期待白色化を行うと, 頻度と長さの関係は図 3.41(b) のように単純な比例関係に変わります. Word2Vec だけでなく, GloVe や PMI の場合にどうなるかなど, 単語ベクトルの長さについては統一的な理論が待たれています. (→演習 (7))

3章の演習問題

- (1) 実際のテキストで, Heaps の法則の γ を線形回帰で求めてみる. 線形回帰の係数の求め方については, GP 本の第 1 章を参照のこと.
- (2) 単語と異なり, 文字は言語では一般に有限です. 有限集合である文字の場合にテキストを読みながら文字の「語彙」の大きさを計算した場合, Heaps の法則は成り立つでしょうか? 文字の種類が少ない英語と多い日本語では, 違いがあるでしょうか?
- (3) Zipf の法則の α は, 実際のテキストではどの程度の値になるでしょうか. 図 3.4 の曲線を, 1 つの直線で完全に表すことはできるでしょうか.
- (4) 部分積分により, $\Gamma(x+1) = x\Gamma(x)$ を示せ.
- (5) 3.5.6 節節では, `ja.text8` について Word2Vec で学習された日本語の単語ベクトルについて考えてみましたが, GloVe や PMI で学習されたベクトルについては, 分布はどうなっているでしょうか. Word2Vec との間に, 何か違いがあるでしょうか?
また, 英語の単語ベクトルについてはどうでしょうか. 日本語の単語ベクトルの場合と, 何か違いがあるでしょうか.
- (6) 単語ベクトルの次元を変えたとき, 結果はどう変わるでしょうか. タスクについて.
- (7) GloVe を使って同様に単語類似度や, 単語の比例関係を計算してみる. C 言語の実装は公式ページにあります. 同じテキストから学習したとき, Word2Vec と何か違いがあるでしょうか.
- (8) A la Carte embedding について.

5 文書の統計モデル

これまで、文字・単語・文の順にテキストの単位と、その統計的なモデル化について考えてきました。実際に私たちが会おう場面では、単語だけ、文だけが問題となることは少なく、文が集まった**文書**を扱うことが多いでしょう。ここでいう文書とはいわゆる書類だけでなく、普段触れる Web ページや電子メール、アンケートの自由回答、SNS 上の記事なども含まれます。小説や論説、法案といったものも、時に非常に長くなりますが、文書の一つとっていいでしょう。

文書としてのテキストの特徴は、**意味的まとまり**を持っているということです*1。たとえば、図 5.1(a) のように完全に無関係な言葉の集まった文書は、まず存在しないでしょう*2。実際の文書は図 5.1(b) のように、何かのテーマ（ここでは「クラシック音楽」でしょう）に沿った内容が書かれているはずです。これはテキストが、何かの意図を持って情報を伝える媒体であることから明らかでしょう。それでは、文章からその「意味的まとまり」をどうやって数学的に取り出したらいいのでしょうか。

5.1 ナイーブベイズ法と単語集合表現

文書の意味的まとまりとして、最もわかりやすいのは文書の「ジャンル」*3、または分野を表すカテゴリといったものでしょう。新聞やニュース記事には「社会面」「家庭欄」といった区別があり、それぞれ特定の言葉を使って、一定のカテゴリの内容が書かれているのが普通です。身近なところでは、毎日届く迷惑メー

*1 これには本書で扱うもの以外に、文体（スタイル）上のまとまりも含まれます。一つの文書の中では一般にスタイルが統一されていることを利用して、単語の埋め込みベクトルを意味を表す成分とスタイルを表す成分に分けて学習する興味深い研究に、東北大学の赤間らによる[135]があります。

*2 3章での、意味を考えない n グラムモデルからの出力を思い出してみましょう。

*3 言語学では、レジスター（言語使用域）ともよばれています。

ルも一種のカテゴリで、この場合、メールは「通常のメール」と「迷惑メール」の二種類のカテゴリに分けられることとなります。また、Amazon のレビューなども付けられた星の数によって、「肯定的」および「否定的」なものに分けることができるでしょう。

Amazon レビューの 5 億個を超える多言語データセット *4 は公開されており [136], 日本語版だけ抜き出したものも Hugging Face で公開されています *5。一方で、新聞記事のようにカテゴリの付与された文書は有料のことが多く *6, 特に日本語で自由に使用できるコーパスは少ないのですが, [137] で公開されている Livedoor ニュースコーパスは, 表 5.1 に示した 9 個のカテゴリの記事からなる 7,367 文書の公開データです。たとえば「家電チャンネル」の記事には“売れ筋”, “加湿”といった言葉が, 「MOVIE ENTER」の記事には“公開”, “ヒーロー”といった言葉が多く出現するため, テキストを見れば, 9 つのカテゴリのどれに属する文書なのかが判定できそうです。なお以下では, 3.1 節で紹介した方法などで, 日本語の文書もあらかじめ単語に分けられていると仮定します。

単語集合表現 こうした文書の意味内容を表す最も簡単な方法は, 文書の中で単語の順番を考えず, 頻度だけを数えることです。というのは, 言語には一般に非常に多くの語彙があるため, どんな単語が何回出現したのかを見れば, 文書の大まかな意味はほぼ特定できるからです *7。図 5.1 のように, これは文書を単語の

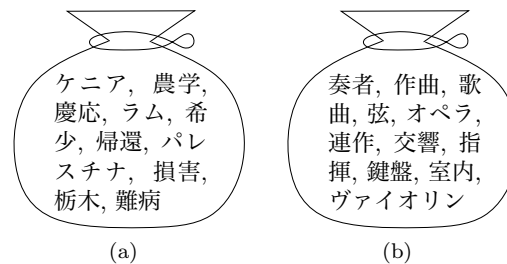


図 5.1: 単語集合表現 (Bag of Words) と意味的まとまりの様子. (a) 単語がまったくランダムに選ばれた文書, および (b) 内容に意味的まとまりがある文書の場合.

*4 <https://amazon-reviews-2023.github.io/>

*5 https://huggingface.co/datasets/SetFit/amazon_reviews_multi-ja

*6 新聞記事コーパスの入手については, 82 ページを参照してください.

*7 これに対して, たとえば DNA では文字が ATGC の 4 種類しかありませんから, 各文字の出現頻度よりも, その並び方の方がはるかに重要になります.

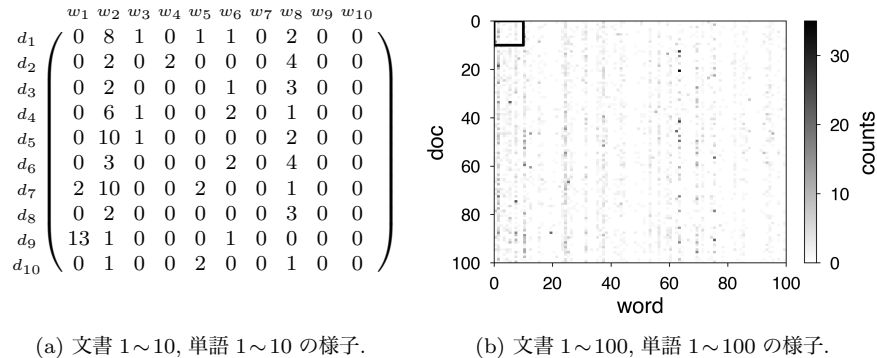


図 5.2: Livedoor コーパスでの文書-単語行列の一部. 横軸の単語は, コーパス全体での頻度順に並んでいます. この行列のうちほとんどは 0 (99.1%=123,721,815/124,833,815) で, 非常に疎な行列になっています. (b) の黒枠の中を拡大したものが (a) の行列です.

集合(袋詰め)にしたものと考えられることができるため, **単語集合** (Bag of Words) 表現とよばれています. 2章で学んだ用語を使えば, これは単語が(文書ごとに異なる)ユニグラム分布から生成された, と仮定していることとなります.*8

この場合, データは図 5.2 のように, 文書番号を縦に, 単語を横にとった**文書-単語行列**で表すことができます. この行列の (d, w) 要素が, d 番目の文書に単語 w が出現した回数 $n(d, w)$ になっています. 図からもわかるように, この行列はほとんどが 0 で, 非常に疎(スパース)な行列になっています. よって, 実際に表す際にはゼロでない要素だけを使って,

1:5 2:7 7:1 12:4 ...

のように, 3.4.3 節でも使った SVMlight 形式で表すとよいでしょう. $w:c$ のような表記は, 単語 w が c 回出現したことを表しています. つまりこの文書には, 単語 1 が 5 回, 単語 2 が 7 回, 単語 7 が 1 回, ...出現したことになります. これは単語集合表現ですから, 順番には特に意味はありません.

*8 「テキストを精密にモデル化する」という観点からは, これはずいぶん簡略化された表現です. 現在の深層学習を用いれば, 次の単語をより正確に予測することは可能ですが, テキストが何の意味を表しているかは完全にブラックボックスで, 分からなくなってしまいます. モデル化とは現実をそのまま再現することではなく, 「ある観点で見た」切り口を考えることで, 現象に関する理解を深め, 制御を可能にすることです[43]. これは無闇に複雑なモデルを使えばよいわけではない, というこの一つの現れといえるでしょう.

	w_1	w_2	w_3	w_4	w_5	w_6	w_7	ラベル	y
d_1	1		1		2		1	通常メール	0
d_2		1	2		1	1		迷惑メール	1
d_3	1	1		1		2			

図 5.3: ナイーブベイズ法での文書-単語行列の例. 空欄の部分は頻度が 0 であることを表しています.

たとえば, 最も簡単な例として, 文書-単語行列が図 5.3 のようになっているとしましょう. 文書 d_1 と d_2 は通常のメール ($y=0$), 文書 d_3 は迷惑メール ($y=1$) というラベル y がわかっていたとします. 単語 w_1 や w_3 は day や work などの普通の単語, w_4 や w_6 は cash や dollar といった迷惑メールによく現れる単語を表していると考えてください.

このとき, ラベル $y=0$ (通常メール) での単語の確率分布は, 最も簡単には d_1 と d_2 での頻度を合計して総和で割ればよく,

$$(5.1) \quad p(w|y=0) = \left(\frac{1}{10}, \frac{1}{10}, \frac{1+2}{10}, \frac{0}{10}, \frac{2+1}{10}, \frac{1}{10}, \frac{1}{10} \right) \\ = (0.1, 0.1, 0.3, 0, 0.3, 0.1, 0.1)$$

となります. 一方, $y=1$ (迷惑メール) の場合は文書は d_3 の 1 つだけなので,

$$p(w|y=1) = \left(\frac{1}{5}, \frac{1}{5}, \frac{0}{5}, \frac{1}{5}, \frac{0}{5}, \frac{2}{5}, \frac{0}{5} \right) = (0.2, 0.2, 0, 0.2, 0, 0.4, 0)$$

です. また, $y=0$ の文書は d_1 と d_2 の 2 個, $y=1$ は d_3 の 1 個ですから,

$$(5.2) \quad p(y) = \left(\frac{2}{3}, \frac{1}{3} \right) = (0.67, 0.33)$$

になります.

実際のテキストでは, $p(y)$ や $p(w|y)$ はどうなっているのでしょうか. サポートサイトにある `livedoor.py` を, Livedoor コーパス `ldcc-20140209.tar.gz` を展開したフォルダの下にある `text` フォルダに適用して

```
% livedoor.py 展開したフォルダ/text livedoor.txt
```

を実行すると、内部で MeCab で単語分割を行って URL 等を除き、結果を

```
dokujo-tsushin 友人 代表 の スピーチ、独女 は どう こなして いる？もうすぐ・ と呼ばれる 6 月。独女 の 中 には 自
分の 式 は まだ な のに ...
topic-news     園山 真希 絵 が 経営 する 料理 店 を 閉店、その 経
緯 に 非難 の 声 29 日、料理 研究 家 の 園山 真希 絵 が、自
身 が 経営 する 家庭 料理 「園山」を ...
```

のように、[ラベル]<TAB>単語列.. の形で1文書が1行のテキスト livedoor.txt に書き出します。この livedoor.txt に対して、 $p(y)$ を求めるスクリプト nblabels.py を実行してみると、

```
% nblabels.py livedoor.txt
⇒ dokujo-tsushin 0.1181
   it-life-hack   0.1181
   kaden-channel 0.1173
   :
   sports-watch  0.1222
   topic-news    0.1045
```

のようになりました。また、式(5.1)のように $p(w|y)$ を計算する nbprob.py を使って、“Peachy” カテゴリ ($y=6$) での単語分布 $p(w|y=6)$ を求めると、

```
% nbprob.py livedoor.txt peachy
⇒ 、      -> 0.047520
   の      -> 0.046050
   に      -> 0.029309
   :
   私      -> 0.000718
```

表 5.1: Livedoor コーパスのカテゴリ名と本書で用いるカテゴリ番号, および文書数.

y	カテゴリ	カテゴリ名	文書数
1	MOVIE ENTER	movie-enter	870
2	独女通信	dokujo-tsushin	870
3	Sports Watch	sports-watch	900
4	IT ライフハック	it-life-hack	870
5	livedoor HOMME	livedoor-homme	511
6	Peachy	peachy	842
7	家電チャンネル	kaden-channel	864
8	エスマックス	smax	870
9	トピックニュース	topic-news	770

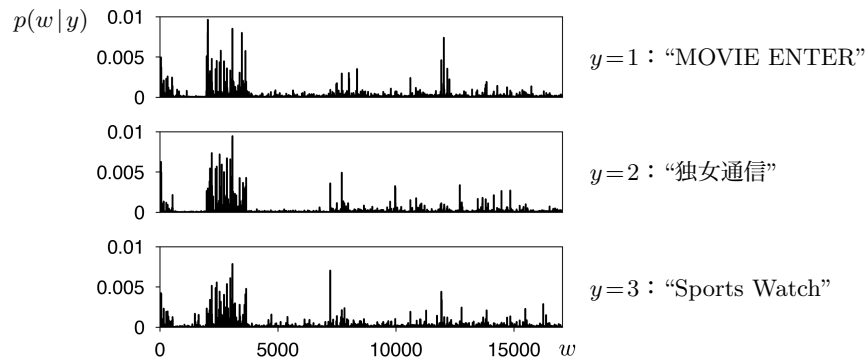


図 5.4: Livedoor コーパスにおけるカテゴリ別の単語確率分布 $p(w|y)$ の様子. 横軸の単語は, 辞書順に並んでいます. “か”, “の” など極端に確率の高い語があるため, 確率が 0.01 以上の単語はプロットから除外しています.

```

クリスマス -> 0.000716
恋愛       -> 0.000709
:
乗れ       -> 0.000002
考慮       -> 0.000002

```

のように計算することができます. このカテゴリ別の単語確率分布 $p(w|y)$ の様子を, 図 5.4 に示しました.

ノート：単語のカテゴリ所属確率の計算

上で計算したのはカテゴリ y の事前確率 $p(y)$ と, y から単語 w が出力される確率 $p(w|y)$ ですが, 実はこれから, w がカテゴリに所属する確率 $p(y|w)$ を求めることができます. というのは, 式(2.36)のベイズの定理から,

$$(5.3) \quad p(y|w) \propto p(y|w)p(y)$$

となるからです. 式(5.3)の右辺を和が1になるように正規化すれば, 単語 w の各カテゴリ y への所属分布 $p(y|w)$ が得られます.

Livedoor コーパスで実際に計算してみると, 表 5.2 のようになりました. 確かに, 各カテゴリの特徴がよく表れていることがわかります. 「劇場」のように, ほぼ特定のカテゴリにしか所属しない単語にどんなものがあるか, 調べてみると面白いでしょう (→演習 12).

カテゴリ	劇場	恋愛	Mac	携帯	ゴルフ	肌	(%)
movie-enter	87.9	8.4	0.0	1.5	0.3	1.5	
dokujo-tsushin	1.0	47.5	0.2	10.5	0.6	13.7	
sports-watch	0.0	2.3	0.0	1.9	11.7	1.7	
it-life-hack	0.0	0.3	83.4	15.4	0.9	0.3	
livedoor-homme	0.8	2.3	4.0	7.0	81.4	2.9	
peachy	3.6	29.9	0.0	3.1	1.2	75.0	
kaden-channel	1.7	2.4	8.9	32.1	2.1	4.2	
smax	1.5	0.2	3.5	22.3	1.1	0.1	
topic-news	3.5	6.7	0.0	6.2	0.8	0.6	

表 5.2: Livedoor コーパスのナイーブベイズ法で, 単語がカテゴリに所属する事後確率 $p(y|w)$ の例. わかりやすいよう, 確率を100倍して%で示しています.

5.1.1 文書の分類確率

これまでに求めた確率を使うと, 文書の確率を計算することができます. カテゴリ y のラベルを持つ文書 $d = w_1 w_2 \dots w_L$ の確率とは, d と y の同時確率のことですから, 式(2.20)の確率の連鎖則から

$$(5.4) \quad p(d, y) = p(y)p(d|y)$$

$$= p(y) \prod_{w \in d} p(w|y) \quad (\text{ナイーブベイズ法の文書確率})$$

となります。 $p(d|y)$ がユニグラム確率の積 $\prod_{w \in d} p(w|y)$ に分解されることが、単語集合の仮定に対応しています。たとえば図 5.3 のデータで、新しい文書 $d = w_2 w_2 w_5 w_7$ がカテゴリ $y=0$ から生成される確率は、

$$\begin{aligned} p(d, y=0) &= p(y=0) \prod_{w \in \{2,2,5,7\}} p(w|y=0) \\ &= 0.67 \times (0.1 \times 0.1 \times 0.3 \times 0.1) = 0.0002 \end{aligned}$$

になります。

このように文書の確率が計算できると、新しい文書 d がどのカテゴリに属するかを、確率的に予測できるようになります。これはつまり、確率分布

$$(5.5) \quad p(y|d) \quad (y \in \{1, \dots, K\})$$

を知りたいということです。式(2.30)のベイズの定理を用いれば

$$(5.6) \quad p(y|d) \propto p(y)p(d|y) = p(y) \prod_{w \in d} p(w|y)$$

となり、この確率は $y \in \{1, \dots, K\}$ について式(5.4)から求めることができます。

実際に計算してみましょう。図 5.3 のデータで $d = w_1 w_2 w_6$ のとき、 $p(y|d)$ は

$$(5.7) \quad p(y|d) \propto p(y)p(d|y) = \begin{cases} 0.67 \times 0.1 \times 0.1 \times 0.1 & (y=0) \\ 0.33 \times 0.2 \times 0.2 \times 0.4 & (y=1) \end{cases} = \begin{cases} 0.00067 & (y=0) \\ 0.00528 & (y=1) \end{cases}$$

となります。よって、これを和が1になるように正規化して、

$$(5.8) \quad p(y|d) = \left(\frac{0.00067}{0.00067+0.00528}, \frac{0.00528}{0.00067+0.00528} \right) = (0.113, 0.887)$$

が求める分布となります。つまり $d = w_1 w_2 w_6$ は $y=1$ 、すなわち迷惑メールである確率が0.887と非常に高い、ということがわかります。

このように、式(5.4)の文書確率をもとにベイズの定理から、式(5.6)を用いて

文書を属するカテゴリに確率的に分類する手法を、**ナイーブベイズ法**といいます。ナイーブ (素朴) とは、テキストの単語の順番を考えない単語集合の仮定のことをさしていますが、この単純な仮定にもかかわらず、この後でみるように、分類については高い性能を持つことがわかっており、文書分類の最も基本的なモデルになっています。

なお、式(5.7)のような確率は一般に単語数が増えると、確率が指数的に小さくなり、そのままでは計算機で表現できなくなってしまいます。よってこうした場合は、対数をとって

$$\begin{aligned}
 (5.9) \quad \log p(y=0|d) &\propto \log p(y) + \log p(d|y) \\
 &= \log(0.67) + \log(0.1) + \log(0.1) + \log(0.1) \\
 &= -0.405 - 2.303 - 2.303 - 2.303 = -7.314
 \end{aligned}$$

と計算するのがよいでしょう。このとき、 $\log p(y)$ および $\log p(w|y)$ のそれぞれを「スコア」とみなせば、式(5.9)は文書 d が $y=0$ に属する「スコア」を、ラベル y および d に含まれる各単語 w ごとに足し込んでいることになります。この合計スコアが最も大きいカテゴリが、分類の結果になります。^{*9}

ナイーブベイズ法の実験

実際に、Livedoor コーパスで実験してみましょう。先に作成したデータファイル `livedoor.txt` を、2章の2.6節で説明したように行をシャッフルして、ランダムに80%の学習データ、10%の開発データ、10%のテストデータに分割します。

```
% split.py livedoor.txt livedoor
```

これにより、`livedoor.train` (5888行)、`livedoor.dev` (736行)、`livedoor.test` (743行) が作成されます。ランダムに分割するため、以下の数値は人によって多少異なることに注意してください。

同じ場所にある `nb.py` を使って、学習データからナイーブベイズ法のモデルを計算します。この計算は数秒で終了します。

```
% nb.py livedoor.train nb.livedoor
N = 5888 docs, V = 17053 vocabs, K = 9 classes.
alpha = 0.01, threshold = 10.
saving model to nb.livedoor.. done.
```

^{*9} 一般に、確率的でないアルゴリズムで天下りに定義された「スコア」は、確率的に考え直すと、確率の対数 (=情報量) とみなせることが非常に多くみられます。

`nb.livedoor` はナイーブベイズ法のパラメータ, つまり $p(y)$ および $p(w|y)$ が格納された gzip 圧縮済みの pickle ファイルです. 紙面の都合でスクリプトは載せませんので, 必ず中身を読んでから使ってください.

学習したモデルを使って新しいテキストのカテゴリを予測するには, スクリプト `nbinf.py` を使って, 次のように実行します.

```
% cat nb.txt
誕生日にディナーでプレゼントをもらった！
% nbinf.py nb.livedoor nb.txt
[ 0.086 0.028 0.000 0.000 0.012 0.872 0.001 0.000 0.001]
```

このテキストは, 表 5.1 から 6 番目の “Peachy” カテゴリの確率が高い (0.872) と予測されることがわかります. なお, 次に確率の大きい (0.086) カテゴリは “MOVIE ENTER” でした.

特定のテキストではなく, モデル全体の性能を評価するには, スクリプト `nbeval.py` を使って, `livedoor.test` のカテゴリの予測精度を次のようにして測ります^{*10}.

```
% nbeval.py nb.livedoor livedoor.test
accuracy = 91.92%
```

このテストデータでの文書カテゴリの予測精度は, 91.9%になりました.

結果の検証と「正解」ラベル ということは, 8.2% (60/743 文書) はカテゴリの予測を「間違っただ」ということになります. ただし, この 8.2%がどんな場合なのかには注意が必要です. 実際に予測を「間違っただ」60件の文書を調べてみると, 表 5.3 のようになっており, 少なくとも人間の基準では, それほど間違いとは言えないことがわかります. 一方で, BERT のようなブラックボックスのニューラル手法は, このデータについて 97%程度の正解率を持っていますが[], これらは同じサッカーの記事でも, 「一般ニュース」と分類しているわけです. これはもしかするとニューラル手法が, Livedoor コーパスの各分野のライターの書き癖を学習しており, 内容に関係なく文体で分類しているのかもしれませんが (少なくとも, その可能性があります). したがって, 「誤り」とされたものが本当に誤りなのかは, こうして結果を確認して検討した方がよいでしょう. この

*10 ここでは開発データ `livedoor.dev` は使いません.

表 5.3: Livedoor コーパスでナイーブベイズ法が「間違えた」文書の例.

テキスト	予測ラベル	正解ラベル
川面を流れるの美しさに感動!散った桜が川面を埋め尽くす光景のあまりの美しさが話題に先週、花見客のマナーの悪さを紹介した「これからお花見の..	peachy	it-life-hack
五輪サッカー英韓戦を前に韓国では「最悪のシナリオ」の声 2 日、韓国のニュースサイト「」は、五輪サッカー男子で韓国がとの予選リーグ最終試合を..	sports-watch	topic-news
オトナ女子たちの圧倒的支持をうけ、ドラマ 10『』の一挙再放送が決定!昨年を出産した女優の木村佳乃さんが 2 年ぶりに連ドラ主演を務め、現在 NHK ..	movie-enter	dokujo-tsushin

結果は、人手による「正解」ラベルが本当に常に正しいのか、を教えてくれる例ともいえます。

テキストの感情極性分類 もう一つ、ナイーブベイズ法のわかりやすい応用に**感情分析**があります。感情分析には、テキストを肯定的・否定的などの**極性**に分類する感情極性分類 (sentiment analysis) と、“期待”、“驚き”などの基本感情を付与する基本感情分類 (emotion detection) があり、愛媛大学の梶原らは、日本語の SNS テキスト (ツイート) に対して上記の感情極性と基本感情の強度を付与したコーパス WRIME [138][139]を公開しています*11。

ここでは、より簡単な感情極性分類を行ってみることにしましょう。WRIME のデータをダウンロードし、含まれる `wrime-ver2.tsv` にサポートサイトにある `wrime.py` を実行すると、SNS のテキストを 90 ページの `neologdn` で正規化した後で MeCab で単語に分割し、`livedoor.txt` と同じ形式のデータ (8740 行) を作ることができます。neologdn をインストールしていない方は、先に `pip install neologdn` などでインストールしておいてください。

```
% git clone https://github.com/ids-cv/wrime
% cd wrime
```

*11 <https://github.com/ids-cv/wrime>

表 5.4: WRIME コーパスのツイート極性から計算した, positive と negative それぞれのカテゴリ y での単語確率 $p(w|y)$ の上位語.

(a) $y=positive$ の場合				(b) $y=negative$ の場合			
!	0.1060	♡	0.0730	ない	0.0930	怒ら	0.0800
♪	0.0890	好き	0.0720	“	0.0930	憂鬱	0.0800
最高	0.0850	フィンランド	0.0720	つらい	0.0880	無視	0.0800
楽しみ	0.0820	アニ	0.0710	嫌	0.0850	しんどい	0.0790
嬉しい	0.0810	ww	0.0710	しろ	0.0850	下痢	0.0790
曲	0.0800	~/	0.0710	イライラ	0.0850	怒り	0.0790
かわいい	0.0770	Saint	0.0710	悪い	0.0850	寂しい	0.0790
サマ	0.0750	Snow	0.0710	くさい	0.0830	—	0.0790
おいしい	0.0740	さん	0.0710	吐き	0.0820	注意	0.0780
歌	0.0740	楽しかつ	0.0710	迷惑	0.0820	生理	0.0780
o	0.0740	きれい	0.0700	むり	0.0820	やらかし	0.0780
可愛	0.0730	*	0.0700	無理	0.0810	地獄	0.0780
可愛い	0.0730	\(^	0.0700	なくなる	0.0810	吐き気	0.0780
良かつ	0.0730	しあわせ	0.0700	悲しい	0.0810	リスク	0.0780
ケーキ	0.0730	DVD	0.0700	めんど	0.0800	起き	0.0780

```
% wrime.py wrime-ver2.tsv > wrime.txt
% shuf wrime.txt | head -3
negative 掃除機が動かない(^q^)なんてこった..
positive おはようございます!今日は学びに使い..
negative 22時に閉まるすき家ってはじめてみた
```

このうち, ランダムに選んだ740ツイートをテストデータ, 残りの8000ツイートを学習データとして, 230ページと同様に `nb.py` でナイーブベイズ法のモデルを計算してみましょう. この計算は1秒未満で終わります.

```
% shuf wrime.txt > wrime.shuffled.txt
% head -8000 wrime.shuffled.txt > wrime.train
% tail -740 wrime.shuffled.txt > wrime.test
% nb.py wrime.train nb.wrime
N = 8000 docs, V = 1690 vocabs, K = 2 classes.
alpha = 0.01, threshold = 10.
saving model to nb.wrime.. done.
```

計算したナイーブベイズ法の単語確率 $p(w|y)$ (式(5.4)) の上位語を, y が positive と negative の場合のそれぞれについて表 5.4 に示しました. 単語の極性が, きわめて自然な形で学習されていることがわかります. このモデルをもとに, テストデータのツイートの感情極性を予測してみましょう.

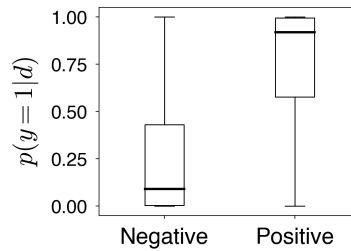


図 5.5: WRIME コーパスのテストデータのツイート極性の予測結果. 縦軸は、予測された極性が Positive である確率を表します. 正解が Positive な場合には確率は 1 に近く、Negative な場合には 0 に近くなっており、正しい学習が行えていることがわかります.

```
% nbinf.py nb.wrime wrime.test
⇒ 'negative': 0, 'positive': 1
negative [ 0.589 0.411] 手の痺れ一生治らん
positive [ 0.001 0.999] 先輩の息子が、まじ一瞬だけどつかまら..
positive [ 0.474 0.526] 私は、恐らく予定を組むのが好きな人な..
positive [ 0.464 0.536] チロルチョコにダイブしたい。
negative [ 1.000 0.000] わたしの伝え方がおかしいのか?なぜ、..
positive [ 0.739 0.261] 「落ちる」場面や、走っても走っても前..
positive [ 0.224 0.776] Shift+win キー+S キーで範囲指定の画面..
positive [ 0.091 0.909] おはようございます。
:
```

この場合、数字は順に negative, positive の確率を表します。短いツイートでは情報が少なく判断が難しくなりますが、ほぼ最初のカラムに示した「正解」のラベルに沿った予測ができていることがわかります。図 5.5 に、正解ラベルの Positive/Negative に分けた場合の Positive 確率の予測結果を示しました。次のように実行して、ランダムな 20 ツイートの極性を予測値でソートしてみると、**非常に簡単なモデルにもかかわらず**、確かにツイートの感情極性順に並んだ結果が得られることがわかります。

```
% nbinf.py nb.wrime wrime.test | shuf | head -20 | sort -k3 -n
⇒ positive [ 0.000 1.000] CD 工場のやつ今見た! 「どももっちーで..
positive [ 0.000 1.000] EXIT の DVD やっすううう 神谷浩史小野大..
positive [ 0.000 1.000] サスケとサクラちゃんが最高! って、 pi..
positive [ 0.000 1.000] 石川智晶姉さんきたあああ。大好きい..
positive [ 0.010 0.990] 飯塚昌明 ANNIVERSARYLIVE“e-XP02020..
```

```

positive [ 0.025 0.975] やる気出てきたぞー\ (´ω´)/
positive [ 0.026 0.974] 最近、肉蝮伝説という作品にハマって..
positive [ 0.029 0.971] 今期のロボットも二つとも面白いの?..
positive [ 0.071 0.929] RBC 入場!!
negative [ 0.445 0.555] お金無駄遣いなおしたい
negative [ 0.471 0.529] はーん
positive [ 0.519 0.481] はー今日のおもち本当に美味しかった..
positive [ 0.584 0.416] 大学会館前もある…
negative [ 0.611 0.389] パソコンがどんどんこわれてくひらが..
negative [ 0.743 0.257] 食費にいくら残る?
positive [ 0.752 0.248] 「京」の1000個のCPUの配布って最終回..
negative [ 0.824 0.176] 壊されてたまるかわたしは
negative [ 0.901 0.099] 君僕新刊もまだ買ってない(´Д`)天使も
negative [ 0.942 0.058] リボ払いコワイ(使ってないと言って..
negative [ 0.943 0.057] スペースもだけど夫のいびきうるさい..

```

5.2 ユニグラム混合モデル (UM)

前節で、ナイーブベイズ法が非常に簡単なモデル化と計算にもかかわらず、文書を高い性能で確率的に分類できること、また、人の付与した「正解」ラベルが必ずしも意味的な内容とは一致しないことをみてきました。

そもそも、テキストの内容が“一般ニュース”や“Peachy”といったラベルだけで語れるはずがなく、それぞれのカテゴリの中には、様々な話題が含まれているはずですが、しかし、それらの「話題」が何なのかテキストに書かれているわけではありません。また一般のテキストには、そもそもラベルすらないことがほとんどです。図 5.6 に示したのは、サポートサイトの本章のフォルダにある `jawiki.txt` の一部で、これは日本語 Wikipedia の記事全体 (約 62 万記事) からランダムに 10,000 記事を筆者が抜き出したテキストですが、これらの記事にはもちろんラベルはありません。

それでは、こうしたラベルのないテキストを扱うにはどうすればいいのでしょうか。基本的な戦略は、「ラベルがないのだから、自分で推定してしまえ」ということです。

コラム：対数確率と logsumexp

式(5.8)の確率は式(5.7)の事後確率の和が1になるように正規化して求めたものですが、文書の長さが大きくなると、この計算はそのまま行うことができません。というのは、一定以上小さな確率は計算機では表すことができず、0になってしまうからです。^{*12} よって、式(5.9)のように対数をとって計算するのがよいでしょう。

K 個のカテゴリについて、非常に小さな確率 (p_1, p_2, \dots, p_K) の対数 $(\ell_1, \ell_2, \dots, \ell_K)$ がわかっていたとき、 p を正規化して和を1にするには、

$$(5.10) \quad p'_k = \frac{p_k}{\sum_{k=1}^K p_k} = \exp\left(\ell_k - \log \sum_{k=1}^K \exp(\ell_k)\right)$$

で計算することができます。ただし上の式の $\log \sum_{k=1}^K \exp(\ell_k)$ は、そのまま計算すると、 $\exp(\ell_k)$ (たとえば $\exp(-1000)$) がすべて0になってしまい、求めることができなくなってしまう場合が多くあります。

このとき、 $\ell_1, \ell_2, \dots, \ell_K$ の中で最大のものを m とおけば、

$$(5.11) \quad \begin{aligned} \log \sum_{k=1}^K \exp(\ell_k) &= \log(e^{\ell_1} + e^{\ell_2} + \dots + e^{\ell_K}) \\ &= \log e^m (e^{\ell_1-m} + e^{\ell_2-m} + \dots + e^{\ell_K-m}) \\ &= m + \log \sum_{k=1}^K \exp(\ell_k - m) \quad (\text{logsumexp}) \end{aligned}$$

となります。こうすると ℓ_k の間の差だけが \exp の中に残るため、値がすべて0になるのを防ぐことができます。式(5.11)で $\log \sum \exp()$ の値を計算する関数を、logsumexp といいます [140, §2.5.4]。Python では `scipy.special.logsumexp()` で使えるほか、

```
from numpy import exp, log
def logsumexp(x):
    y = max(x)
    return y + log(sum(exp(x - y)))
```

で定義すれば計算することができます。

*12 執筆時の手元の環境 (64bit) では、倍精度実数の最小値は 2.385×10^{-277} でした。

```

<doc>
Mozilla Application Suite (モジラ・アプリケーション・スイート) また
は Mozilla Suite (モジラ・スイート) は Mozilla Foundation によりプロ
ジェクトを組んでオープンソースで開発されていたインターネットスイートで
あり、...
</doc>
<doc>
カート・ヴォネガット (Kurt Vonnegut、1922年11月11日 - 2007年4月11
日) は、アメリカの小説家、エッセイスト、劇作家。1976年の作品『スラップス
ティック』より以前の作品はカート・ヴォネガット・ジュニア ("Kurt Vonnegut
Jr.") の名で...
</doc>
<doc>
世界の放送方式 (せかいのほうそうほうしき) 高精細度テレビジョン放送
(HDTV: "High Definition Television") に対して従来のテレビ放送の画質は
標準テレビジョン放送 (SDTV: "Standard Definition Televisi on") と同
言われ、ここでは主に標準テレビに分類される方式について記述している。...
</doc>

```

図 5.6: 日本語 Wikipedia からランダムに抽出した記事を集めた `jawiki.txt` の一部.

いま図 5.7 のように、図 5.3 と同じデータでラベル y が未知だったとしましよ
う。この場合、推定するのは人がつけたラベル y ではなく、文書の潜在的なカテ
ゴリを表す**潜在変数**ですので、以後は y と区別して、 z と書くことにします。各
文書 d_1, d_2, d_3 についてカテゴリ $z \in \{0, 1\}$ がわからないのですから、(完全に対
称だと以下の計算が進まないため) $(0.5, 0.5)$ から少しずらした確率分布

$$p(z|d_1) = (0.4, 0.6), \quad p(z|d_2) = (0.6, 0.4), \quad p(z|d_3) = (0.4, 0.6)$$

を初期値としてみましょう *13。このとき $p(z)$ は、 $z \in \{0, 1\}$ について各文書が
もつ確率の総和として

$$(5.12) \quad p(z) \propto \sum_{n=1}^3 p(z|d_n) = \begin{matrix} & d_1 & d_2 & d_3 \\ \begin{cases} 0.4 + 0.6 + 0.4 = 1.4 \\ 0.6 + 0.4 + 0.6 = 1.6 \end{cases} & \propto & \begin{cases} 0.467 & (z=0) \\ 0.533 & (z=1) \end{cases} \end{matrix}$$

と求められます。これは、式(5.2)で文書がどちらかのカテゴリに所属するとし
て $0/1$ で数えていた頻度を、確率に拡張してソフトな和をとったと言っている
でしょう。また $p(w|z)$ も、同様に式(5.1)で頻度を z ごとに確率 $p(z|d_n)$ で重

*13 この確率をクラスタリングでは、文書の各クラスタへの**負担率** (responsibility) といいま
す。ここでは、負担率がほぼ等しい状態を初期値にして計算を始めています。

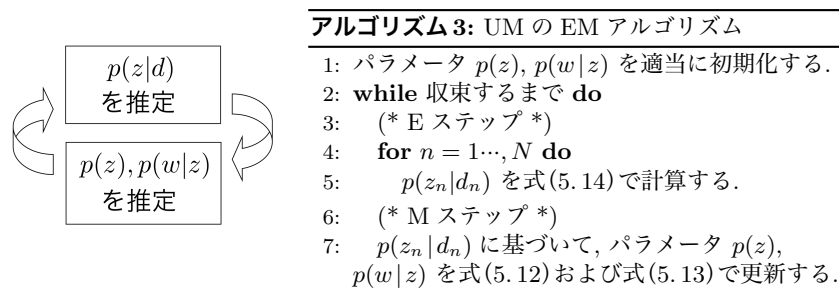


図 5.8: ユニグラム混合モデルの推定法の概要と EM アルゴリズム.

を使って, $p(z|d)$ を更新することができます. たとえば, 文書 d_2 については

$$\begin{cases} p(z=0|d_2) \propto 0.467 \times 0.143 \times (0.229)^2 \times 0.200 \times 0.200 = 0.00014008 \\ p(z=1|d_2) \propto 0.533 \times 0.125 \times (0.175)^2 \times 0.200 \times 0.200 = 0.00008162 \end{cases}$$

から,

$$p(z|d_2) = (0.632, 0.368)$$

となり, $p(z|d_2)$ が $(0.6, 0.4) \rightarrow (0.632, 0.368)$ に更新されることがわかります. このように, 適当な初期値から始めて

$$(5.15) \quad \begin{cases} p(z|d) \text{ を推定する} \\ p(z), p(w|z) \text{ を推定する} \end{cases}$$

ことを図 5.8 のアルゴリズムに従って繰り返すと, 表 5.5 に示したようにこの計算は 7 回程度の繰り返しで収束し, 結果として

$$p(z|d_1) = (1.000, 0.000), p(z|d_2) = (1.000, 0.000), p(z|d_3) = (0.000, 1.000)$$

が得られます. **何も教えていないのに, 教師データ y があつたときと同じ分類ができてしまいました!** $p(z|d_1)$ は最初, $(0.4, 0.6)$ と $z=1$ の方が確率が高かつたにもかかわらず, 式(5.15)の計算の繰り返しの途中で逆転し, 正しい確率分布が得られていることに注意してください.

このように, 式(5.15)の計算を繰り返すことで, 教師なしでテキストをクラス

表 5.5: ユニグラム混合モデルの学習過程. ここでは7ステップで学習が収束し, 教師なしで文書を2つの種類に分類できています.

文書 繰り返し	d_1		d_2		d_3	
	$z=0$	$z=1$	$z=0$	$z=1$	$z=0$	$z=1$
0(初期値)	0.400	0.600	0.600	0.400	0.400	0.600
1	0.399	0.601	0.630	0.370	0.367	0.633
2	0.423	0.577	0.668	0.332	0.303	0.697
3	0.519	0.481	0.715	0.285	0.183	0.817
4	0.769	0.231	0.776	0.224	0.047	0.953
5	0.982	0.018	0.873	0.127	0.002	0.998
6	1.000	0.000	0.988	0.012	0.000	1.000
7	1.000	0.000	1.000	0.000	0.000	1.000

タリングする方法を, **ユニグラム混合モデル** (Unigram Mixtures, UM) [141] といいます. これまでの説明でわかるように, UM は「教師なしナイーブベイズ法」ともいえ, 機械学習で K 平均法[23]として知られるクラスタリングを, テキストの多項分布の場合に適用したものになっています.

ユニグラム混合モデルの実験

実際の文書でも計算してみることにしましょう. 図 5.6 に示したテキストは, 日本語版 Wikipedia の全記事 (執筆時点で約 62 万記事) からランダムに 1 万記事の概要部分を筆者が抽出したもので, サポートサイトに `jawiki.txt` として置いてあります. これを SVMlight 形式のデータにするには, 同じフォルダにある `text2data.py` を次のように実行します. 最後の数字は, 単語頻度の閾値です. スクリプトの使い方は, 中身を読んでみてください.

```
% text2data.py jawiki.txt jawiki 10
⇒ writing dic to jawiki.lex.. done.
   writing data to jawiki.dat.. done.
```

これにより, 10000 行のデータ `jawiki.dat` と, 7509 個の語彙と ID の対照表 `jawiki.lex` が作られます.*15 このデータについて, UM の実装 `um.py` を実行してみましょう.

*15 UM の場合は「が」「の」のような機能語を入れてしまうと, それに引っぱられてクラスタリングがうまく働かないため, 平仮名だけからなる単語を除くといった前処理を行っています.

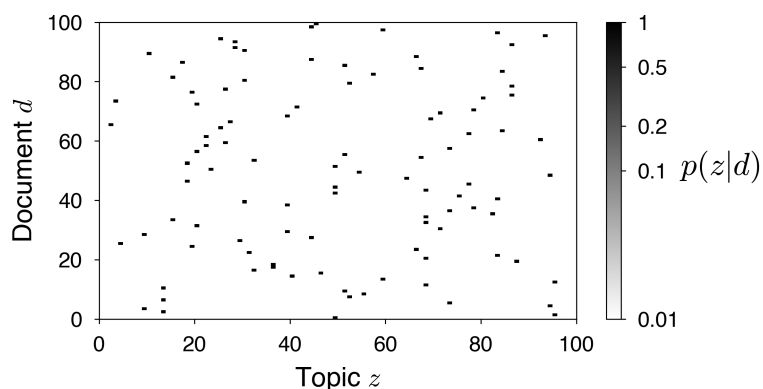


図 5.9: 日本語 Wikipedia コーパス `jawiki.txt` から学習された UM による, 各文書のトピックへの所属確率 $p(z|d)$ (最初の 100 文書). トピック数は $K=100$ としました. ほとんどの $p(z|d)$ はほぼ 1 になっており, 文書ごとに異なるトピックに所属していることがわかります. このプロットは `umplot.py` で作ることができます.

```
% um.py -K 100 jawiki um.jawiki.K100
⇒ UM: 10000 documents, 7509 words in vocabulary.
iter[ 1] : PPL = 2580.88
iter[ 2] : PPL = 1541.67
iter[ 3] : PPL = 684.78
iter[ 4] : PPL = 625.30
iter[ 5] : PPL = 617.85
iter[ 6] : PPL = 615.83
iter[ 7] : PPL = 614.67
iter[ 8] : PPL = 614.36
converged.
saving model to um.jawiki.K100.. done.
```

スクリプトの使い方は, `% um.py` をそのまま実行すれば表示されます. ここでは, $K=100$ 個のカテゴリを仮定しました. 図 5.9 に示したように, 学習されたモデルで $p(z|d)$ を表示してみると, 日本語 Wikipedia の文章がさまざまなカテゴリに所属する様子が, 自動的に学習されていることがわかります.

5.2.1 トピックの解釈と自己相互情報量

ラベルのない生テキストについてUMを学習することで、 K 個のカテゴリ z ごとの単語の出力分布 $p(w|z)$ と、カテゴリの事前確率 $p(z)$ を学習することができました。この場合の z はナイーブベイズ法のように人が与えたカテゴリではなく、データから統計的に学習されたものですから、一般にこれらを**トピック**（話題）とよびます。教師なしナイーブベイズ法であるUMは、最も簡単な**トピックモデル**の一つです。

それぞれのトピック z が何を表しているかは、単語分布 $p(w|z)$ の様子を見てみればわかるはずですが、ただし、この際には注意が必要です。各トピック z について単に確率 $p(w|z)$ の大きい単語を表示すると、表5.6(a)に示したように“年”や“日”といった語が上位に来てしまい、トピックの意味がほとんどわからなくなってしまいます。これは、考えると当然の現象です。UMでは、文書に含まれる語がすべて、あるトピック z から $p(w|z)$ に従ってサンプリングされたと考えています。すると、どの話題でも共通して現れる語は、どのトピックでも高い確率にならなければ、データをよく説明できないからです。ここでは“の”のような平仮名だけからなる語を前処理で除いていますが、除かない場合はこうした語がすべてのトピックで出現確率の上位を占めることになります。^{*16}

したがって重要なのは、単語 w のトピック z での生起確率 $p(w|z)$ 自体の大きさではなく、 w の平均的な出現確率

$$(5.16) \quad p(w) = \sum_{z=1}^K p(w, z) = \sum_{z=1}^K p(w|z)p(z)$$

との比をとった、

$$(5.17) \quad \frac{p(w|z)}{p(w)}$$

の値でしょう。式(5.17)は、単語 w がトピック z で「通常より何倍高い確率で

^{*16} この例からわかるように、言語では固定された「ストップワード」のリストを指定して、それらを除外すれば問題が解決するわけではありません。必ずこのように、リストには含まれない、同様の意味が薄い語が現れるからです。ある語がストップワードであるか否かは、ルールではなく、単語の振る舞いから統計的に判断すべき問題です。

表 5.6: 日本語 Wikipedia コーパスでの各トピックを表す特徴語の計算.

(a) 生成確率 $p(w z)$ を使った上位語							
Topic 1		Topic 10		Topic 20		Topic 30	
年	0.056	年	0.037	年	0.066	年	0.059
日	0.025	月	0.025	月	0.019	大学	0.027
月	0.024	日	0.017	日	0.019	月	0.026
音楽	0.022	システム	0.009	世	0.016	日本	0.021
作曲	0.019	社	0.009	王	0.012	日	0.019
曲	0.017	的	0.007	伯	0.010	会	0.013
家	0.016	法律	0.007	前	0.010	部	0.012
者	0.015	法	0.006	公	0.009	者	0.011
指揮	0.011	日本	0.006	紀元	0.008	長	0.010
作品	0.009	使用	0.006	語	0.007	委員	0.009

(b) PMI (w, z) を使った上位語							
Topic 1		Topic 10		Topic 20		Topic 30	
大刀	4.502	inotify	4.538	ロジスティック	4.528	法科	4.079
査証	4.343	カーネル	4.418	τ	4.516	学長	3.625
ジュリアン	4.313	ボトル	4.292	ネブカドネザル	4.508	ラトビア	3.612
グレコ	4.282	ボランティア	4.055	マウイ	4.499	民兵	3.545
管弦	4.170	上杉	3.989	巨星	4.244	医科	3.495
奏	4.066	宝石	3.787	諸侯	4.211	土木	3.491
石巻	4.051	長尾	3.778	写像	4.199	商科	3.416
ヴァイオリニスト	4.038	帯域	3.775	伯	4.126	法学	3.413
長調	4.029	GNU	3.753	司馬	4.068	志願	3.353
初演	3.913	ペット	3.684	ブランデンブルク	4.055	黒川	3.306

(c) NPMI (w, z) を使った上位語							
Topic 1		Topic 10		Topic 20		Topic 30	
作曲	0.462	カーネル	0.501	伯	0.476	工学	0.373
管弦	0.458	inotify	0.495	ロジスティック	0.457	大学	0.371
指揮	0.437	ボトル	0.490	τ	0.449	委員	0.368
ピアノ	0.437	ボランティア	0.463	ネブカドネザル	0.445	会長	0.353
交響	0.434	ペット	0.444	マウイ	0.441	理事	0.351
楽団	0.434	サーバ	0.438	ブランデンブルク	0.434	学会	0.347
大刀	0.434	帯域	0.430	諸侯	0.430	法学	0.345
グレコ	0.422	上杉	0.429	司馬	0.430	土木	0.332
音楽	0.422	Linux	0.427	辺境	0.428	長	0.327
協奏	0.421	宝石	0.425	写像	0.423	学長	0.327

現れるのか」を表しています。“帯域”のように、全体的な確率は低くても、特定のトピックでは高い確率で現れる語は、 $p(w|z)/p(w)$ の値は大きくなります。一方で“日”のように、出現確率は高くても、どのトピックでもほぼ同じ確率で現れる場合は、 $p(w|z)/p(w)$ はほぼ1になると考えられます。よって、式(5.17)の対数をとって、

$$(5.18) \quad \text{PMI}(w, z) = \log \frac{p(w|z)}{p(w)}$$

を計算すれば、“的”のような語ではほぼ $\log 1 = 0$ ，“帯域”のような語では > 0 になるでしょう。つまり、この形で単語の統計的な「重み」が得られます。表5.6(b)に、jawikiコーパスについて式(5.18)で計算した値が大きい単語を各トピックについて示しました。これで、だいぶトピックの差がみえてきました。

式(5.18)を、 w と z の**自己相互情報量** (Pointwise Mutual Information, PMI)とといいます。というのは、式(2.30)のベイズの定理を用いれば、式(5.18)は

$$(5.19) \quad \log \frac{p(w|z)}{p(w)} = \log \frac{p(w, z)}{p(w)p(z)}$$

と変形することができ、この値は2章で学んだ、確率変数 X と Y の間の**相互情報量**

$$(5.20) \quad I(X, Y) = \sum_x \sum_y p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

で、 Σ の中で期待値をとる対象の各要素 (pointwise) になっているからです。

正規化自己相互情報量 PMIを使うことで、各トピックの意味はかなり明らかになりました。ただ、PMIの上位語にはまだ、「ジュリアン」「ラトビア」などの特殊な単語が含まれてしまっています。これはPMIに一般にみられる現象で、式(5.18)で分母に $p(w)$ があるため、 $p(w)$ が小さい、すなわち稀な語であればあるほど、PMIの値が大きくなってしまいます。

この欠点を補うために、**正規化自己相互情報量** (Normalized PMI, NPMI)が2009年に提案されました [66]。式(5.18)のPMIは、 $p(w|z) = 1$ 、すなわち $p(w) = p(z)$ のときに最大値 $-\log p(w)$ をとることに注意しましょう。^{*17}

*17 これから、 w と z が完全に相関している、すなわち単語 w がカテゴリ z でしか現れない場

よって, PMI をその最大値で割った

$$(5.21) \quad \log \frac{p(w|z)}{p(w)} \Big/ (-\log p(w))$$

を考えることができ, これを正規化自己相互情報量 (NPMI) といいます. NPMI は $-1 \leq \text{NPMI} \leq 1$ の値をとり,

- w と z が完全に相関しているとき 1,
- w と z が独立なとき 0,
- w と z が完全に逆相関しているとき -1

となり, w と z の相関を表すために理想的な量となっています.

なお, 式(5.21)は

$$(5.22) \quad \frac{\log p(w|z) - \log p(w)}{-\log p(w)} = 1 - \frac{\log p(w|z)}{\log p(w)}$$

とも書き直すことができ, 式(2.63)で学んだトピック z での単語 w の自己情報量 $-\log p(w|z)$ の, w の平均的な情報量 $-\log p(w)$ に対する比を最大値 1 から引いた値になっています. NPMI は一般には, 式(5.19)で最大値 $-\log p(w, z)$ ($= -\log p(w) - \log p(z)$) で割って,

$$(5.23) \quad \text{NPMI}(w, z) = \log \frac{p(w, z)}{p(w)p(z)} \Big/ (-\log p(w, z))$$

(正規化自己相互情報量)

と定義されています[66].

表 5.6(c) に, NPMI で計算した UM の各トピックの特徴語を示しました. PMI と比べ, “太刀” や “ラトビア” といった低頻度語の順位が下がり, 各トピックの意味をより適切に表していることがわかります. こうした特徴から, NPMI はトピックモデルにおいて, トピックを説明する標準的な指標として広く使われています[]. なお, 表 5.6 ではまだ 1 つのトピックに複数の話題が混入していますが, この後で説明するように, これは UM をベイズ学習することで, 大きく改善することができます (表 5.7).

合は, PMI は $p(w)$ が小さい稀な語ほど大きな値をとってしまう, ということがわかります.

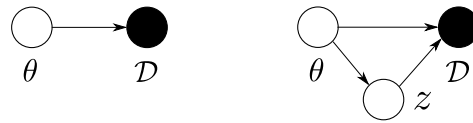
(a) 通常の統計モデル (b) 潜在変数 z がある統計モデル

図 5.10: 統計モデルとパラメータの推定. D はデータを, θ はパラメータを表します.

5.2.2 EM アルゴリズムによる学習

ところで, 図 5.8 の学習アルゴリズムはなぜ取束し, 教師なしで UM の学習ができるのでしょうか? それは, このアルゴリズムが **EM アルゴリズム** という一般的な学習法の一つになっているからです.

EM アルゴリズムは, 因果推論でも有名な Rubin らによって 1970 年代に提案され [142], 1990 年代には最先端の手法として機械学習で多く使われました [23]. EM アルゴリズムは, 最近のニューラルネットワークの VAE (変分自己符号化器) [143] の学習法の基礎ともなっている重要な方法ですので, 以下でみていくことにしましょう.

最尤推定とベイズ推定 まず, 統計モデルとは図 5.10 のように, データ D の背後にパラメータ θ があり, D を最もよく説明する θ を求めることです. このとき, θ から D が生成される確率 (尤度) を最大にする $\theta = \hat{\theta}$ を求める, すなわち

$$(5.24) \quad \hat{\theta} = \operatorname{argmax}_{\theta} p(D|\theta)$$

を点推定するのが**最尤推定**です. いっぽう, 有限のデータから $\hat{\theta}$ が一点で正確に求まるはずがありませんから, D が与えられた下での θ の確率分布

$$(5.25) \quad p(\theta|D) \propto p(D|\theta)p(\theta)$$

を求めるのが**ベイズ推定**です. 式(5.24)と式(5.25)を見比べると, 最尤推定とは θ の事前分布に一様分布 $p(\theta) \propto 1$ を仮定した上で, $p(\theta|D)$ をデルタ関数 $\delta(\theta = \hat{\theta})$ で一点近似することに対応することがわかります. EM アルゴリズムは最尤推定を行う方法ですので, ここではしばらく式(5.24)で考えることにしましょう.

EM アルゴリズム モデルにパラメータ θ だけでなく, UM の各文書のクラス

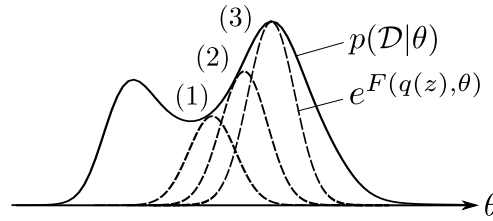


図 5.11: EM アルゴリズムによる学習. 実際には対数尤度 $\log p(\mathcal{D}|\theta)$ を最大化しますが, わかりやすさのため, もとの確率密度 $p(\mathcal{D}|\theta)$ の空間で示しています. EM アルゴリズムでは $\log p(\mathcal{D}|\theta)$ の下界 $F(q(z), \theta)$ を, 適当な初期値から始めて, E ステップと M ステップを繰り返すことで次々と最大化します. 初期値によっては, 左側の局所解に陥ることもあることに注意してください.

タ番号のように**潜在変数** z があるとき, 式(5.24)の尤度は

$$(5.26) \quad p(\mathcal{D}|\theta) = \int p(\mathcal{D}, z|\theta) dz$$

と書き換えることができます. z が離散の場合は上の積分は和になりますが, ここでは一般的に積分で表しました. 実際には, きわめて小さな値になる式(5.26)の確率の代わりに, その対数をとった

$$(5.27) \quad \log p(\mathcal{D}|\theta) = \log \int p(\mathcal{D}, z|\theta) dz$$

を最大化したいのですが, この式は \log の中に積ではなく和が入っているため, $\log p()$ の形に分解することができません. そこで, 一つ工夫をします. 分母・分子に同じ値 (補助分布) $q(z)$ をかけて, ??ページの Jensen の不等式を使えば, 式(5.27)の対数尤度は

$$(5.28) \quad \begin{aligned} \log p(\mathcal{D}|\theta) &= \log \int p(\mathcal{D}, z|\theta) dz = \log \int q(z) \frac{p(\mathcal{D}, z|\theta)}{q(z)} dz \\ &\geq \int q(z) \log \frac{p(\mathcal{D}, z|\theta)}{q(z)} dz \equiv F(q(z), \theta) \end{aligned}$$

と, 不等式で下から抑えることができます. 補助分布 $q(z)$ を用いることで, 無事, $\log p(\mathcal{D}, z|\theta)$ の形にすることができました!

式(5.28)は真の対数尤度 $\log p(\mathcal{D}|\theta)$ の下界ですから, 図 5.11 に示したよう

に式(5.28)を θ および $q(z)$ について最大化することで、真の対数尤度の下から近づくことができます。^{*18} これを行うのがEMアルゴリズムです。具体的には、

Eステップ : θ を固定して、 $F(q(z), \theta)$ を $q(z)$ について最大化

Mステップ : $q(z)$ を固定して、 $F(q(z), \theta)$ を θ について最大化

を交互に繰り返すことで下限を逐次最大化します。EはExpectation(期待値)、MはMaximization(最大化)を意味します。それぞれ、どんな計算になるでしょうか。やや抽象的ですが、先に一般論をみてみましょう。UMの場合の具体例は、その後で説明します。

Eステップ $q(z)$ については、下限 $F(q(z), \theta)$ は

(5.29)

$$\begin{aligned} F(q(z), \theta) &= \int q(z) \log \frac{p(\mathcal{D}, z|\theta)}{q(z)} dz = \int q(z) \log \frac{p(z|\mathcal{D}, \theta)p(\mathcal{D}|\theta)}{q(z)} dz \\ &= \int q(z) \log \frac{p(z|\mathcal{D}, \theta)}{q(z)} dz + \log p(\mathcal{D}|\theta) \\ &= \underbrace{-D(q(z)||p(z|\mathcal{D}, \theta))}_{(*)} + \log p(\mathcal{D}|\theta) \end{aligned}$$

となることに注意しましょう。*の部分は式(2.69)で学んだ、 $q(z)$ と $p(z|\mathcal{D}, \theta)$ のKLダイバージェンスですから、この値は

$$(5.30) \quad q(z) = p(z|\mathcal{D}, \theta)$$

のとき最小になり、よって F が最大になります。すなわち、 $q(z)$ としては実際には、現在のモデル(パラメータ θ)での z の予測分布 $p(z|\mathcal{D}, \theta)$ をとればよい、ということがわかります。

Mステップ θ の方はどうでしょうか。式(5.28)を θ について変形すると、

$$(5.31) \quad \begin{aligned} F(q(z), \theta) &= \int q(z) \log \frac{p(\mathcal{D}, z|\theta)}{q(z)} dz = \int q(z) \log p(\mathcal{D}, z|\theta) dz \\ &\quad - \int q(z) \log q(z) dz = \left\langle \log p(\mathcal{D}, z|\theta) \right\rangle_{q(z)} + H(q(z)) \end{aligned}$$

^{*18} この $F(q(z), \theta)$ を、物理学とのアナロジーで自由エネルギーともいいます。

です. $H()$ は式(2.65)のエントロピー関数で, $\langle \dots \rangle_p$ は, \dots を p について期待値を取ることを表しています. この式はこれ以上簡単になりませんので, F を θ について最大化するためには,

$$(5.32) \quad Q(\theta) = \left\langle \log p(\mathcal{D}, z | \theta) \right\rangle_{q(z)}$$

に対して, $\partial Q / \partial \theta = 0$ とおくなどして θ について最大化することになります. 式(5.32)は, 一般に **Q 関数** とよばれています. Q 関数は, 現在のモデルを使った潜在変数 z とデータの同時確率を, 未知の z について期待値をとったものであることに注意してください.

UM の EM アルゴリズムの導出

抽象的な話が続いたので, 具体的に計算してみましょう. UM で文書 d が生成される確率モデルは, 式(5.4) で y (ここでは z) について和をとった

$$(5.33) \quad \begin{aligned} p(d) &= \sum_{z=1}^K p(d, z) = \sum_{z=1}^K p(d|z)p(z) \\ &= \sum_{z=1}^K p(z) \prod_{w \in d} p(w|z) \end{aligned} \quad (\text{UM の文書確率})$$

となります. よって N 個の文書 $\mathcal{D} = \{d_1, d_2, \dots, d_N\}$ 全体の確率は,

$$(5.34) \quad p(\mathcal{D}) = \prod_{n=1}^N \left(\sum_{z=1}^K p(z) \prod_{w \in d_n} p(w|z) \right)$$

になり, 対数をとれば

$$(5.35) \quad \log p(\mathcal{D}) = \sum_{n=1}^N \log \left(\sum_{z=1}^K p(z) \prod_{w \in d_n} p(w|z) \right)$$

です. パラメータは $\theta = \{p(z), p(w|z)\}_{z=1}^K$ です.

E ステップ $p(z|\mathcal{D}, \theta)$ とはこの場合, 各文書 $d \in \mathcal{D}$ について, その負担率 $p(z|d)$ のことです. これは, 式(5.14)で計算できるのでした.

M ステップ Q 関数の中身である $\log p(\mathcal{D}, z|\theta)$ は, この場合は

$$\log p(\mathcal{D}, z|\theta) = \log \left(p(z) \prod_{w \in d} p(w|z) \right) = \log p(z) + \sum_{w \in d} \log p(w|z)$$

になります。Q関数はこの z に関する期待値なので、 $q(z) = p(z|d)$ でしたから、

$$(5.36) \quad \left\langle \log p(\mathcal{D}, z|\theta) \right\rangle_{q(z)} = \sum_{z=1}^K p(z|d) \log p(z) + \sum_{w \in d} \sum_{z=1}^K p(z|d) \log p(w|z)$$

となります。文書は N 個あるので、全体のQ関数は

$$(5.37) \quad Q(\theta) = \sum_{n=1}^N \left[\sum_{z=1}^K p(z|d_n) \log p(z) + \sum_{w \in d_n} \sum_{z=1}^K p(z|d) \log p(w|z) \right]$$

となります。

$p(z)$ について最大化： $Q(\theta)$ を $p(z)$ に関して微分すると *19, $(\log x)' = \frac{1}{x}$ ですから

$$(5.38) \quad \frac{\delta Q(\theta)}{\delta p(z)} = \sum_{n=1}^N \frac{p(z|d_n)}{p(z)}$$

です。 $\sum_z p(z) = 1$ なので、制約項 $(\sum_z p(z) - 1)$ を追加すると、ラグランジュの未定乗数法 *20 により、未定乗数を λ として

$$(5.39) \quad \frac{\delta}{\delta p(z)} \left(Q - \lambda \left(\sum_z p(z) - 1 \right) \right) = 0$$

を解けばよいことになります。よって式(5.38)から

$$\sum_{n=1}^N \frac{p(z|d_n)}{p(z)} - \lambda = 0$$

となり、これを解いて、

*19 $p(z)$ は関数なので、形式的に変分 $\delta Q / \delta p(z)$ を用いていますが、 Q の中に $p(z)$ の微分に関する項はありませんので、これは $p(z)$ に関する通常の偏微分になります[144]。 $p(z=k) = \mu_k$ のようにおいて、変数 μ_k について最大化すると考えてもよいでしょう。

*20 ラグランジュの未定乗数法について知らない方は、紙幅の問題で本書では解説しませんが、[10, §1.2.3]などを参照してください。

$$(5.40) \quad p(z) = \frac{1}{\lambda} \sum_{n=1}^N p(z|d_n)$$

となります。ここで λ についても微分すれば、当たり前ですが

$$\sum_z p(z) = 1$$

ですから、式(5.40)で λ は和がちやうど 1 になるように決めればよいことがわかり、式(5.12)が得られます。

$p(w|z)$ について最大化： $Q(\theta)$ を $p(w|z)$ について微分すると、同様に

$$(5.41) \quad \frac{\delta Q(\theta)}{\delta p(w|z)} = \sum_{n=1}^N \sum_{v \in d_n} \mathbb{I}(v=w) \frac{p(z|d_n)}{p(w|z)} = \frac{1}{p(w|z)} \sum_{n=1}^N p(z|d_n) n(d_n, w)$$

となります。よって、同様にラグランジュの未定乗数法により

$$(5.42) \quad \frac{1}{p(w|z)} \sum_{n=1}^N p(z|d_n) n(d_n, w) - \lambda = 0$$

から、

$$(5.43) \quad p(w|z) = \frac{1}{\lambda} \sum_{n=1}^N p(z|d_n) n(d_n, w)$$

となり、同様にして式(5.13)が得られました。□

EM アルゴリズムでは、E ステップと M ステップの繰り返しによって下限 $F(q(z), \theta)$ が単調に増加します。よって、240 ページの実装 `um.py` のように学習データのパープレキシティを計算して、減少が一定の閾値以下になったかどうかで収束を判定するとよいでしょう。

5.2.3* UM のベイズ学習

前節で、図 5.8 のアルゴリズムが EM アルゴリズムになっていること、およびその導出について理解することができました。EM アルゴリズムでは、E ステップと M ステップを交互に行うことで、真の対数尤度を図 5.11 のように下から近似して最大化します。これにより、式(5.28)の $F(q(z), \theta)$ はつねに増加して極

大値に近づきます。しかし、EM アルゴリズムは常に山登りを行うため、HMM の学習について図 4.20 でみたように、初期値に依存して一番近くの山 (局所解) につかまってしまうという大きな欠点があります。複雑な問題で対数尤度の曲面に凸凹が大きいほど、この問題は顕著になります。

これを解消するもっとも簡単な方法は、EM アルゴリズムの E ステップで負担率を計算して期待値をとるのではなく、この確率に従ってランダムにサンプリングを行うことでしょう。常に山を登るのではなく、時には下ることも許せば、より広い空間を探索できます。具体的には、UM の場合は E ステップでは負担率 $p(z|d_n)$ に従ってトピック z_n を文書ごとに毎回サンプリングし、M ステップではその z_n をナイーブベイズ法における「正解」ラベルと同様にみなして、パラメータ $\theta = \{p(z), p(w|z)\}$ を更新します。E ステップをサンプリングで置き換えるこの方法を、**モンテカルロ EM アルゴリズム** [145] といいます*21。

モンテカルロ EM アルゴリズムは局所解から脱出できる有用な方法ですが*22、M ステップでは θ について最大化を行っているため、全体にみると EM アルゴリズムの逐次最大化の一種となっており、真の最大値を探索することはできません。より正しくは、 θ についても最尤推定ではなく、事後分布全体からサンプリングを行うベイズ学習を行う必要があります。

幸いにして、UM の場合は 4 章で学習した **Gibbs サンプリング** を容易に行うことができます。各文書 d_1, \dots, d_N について、その潜在トピック $\mathbf{z} = z_1, \dots, z_N$ を $p(z_n|d_n)$ に従ってサンプリングしたとしましょう。このとき、 \mathbf{z} の中でトピックが k となった文書の数を $n(k) = \sum_{n=1}^N \mathbb{I}(z_n = k)$ とおけば、 $\boldsymbol{\mu} = \{p(z)\}$ ($z = 1, \dots, K$) が事前分布

$$(5.44) \quad \boldsymbol{\mu} \sim \text{Dir}(\alpha, \dots, \alpha)$$

に従っているとしたとき、その事後分布は式 (4.38) と同様に

$$(5.45) \quad \boldsymbol{\mu} | \mathbf{z} \sim \text{Dir}(\alpha + n(1), \dots, \alpha + n(K))$$

*21 モンテカルロとはカジノで有名な地中海の都市で、乱数を使って期待値を求める方法をこの名前にしたのは、第二次世界大戦中に原爆開発のために作られたアメリカのロスアラモス研究所で、数学者のウラムとフォン・ノイマンによるものです[146]。

*22 モンテカルロ EM アルゴリズムは、潜在変数 z に指数的な組み合わせの可能性があるなど、すべての場合を計算して期待値をとることが実質的に不可能な場合にも有益なアルゴリズムです。

アルゴリズム 4: UM の周辺化 Gibbs サンプルング

```

1: for  $n$  in  $1, \dots, N$  do           (* 初期化 *)
2:    $z_n = \text{randint}(K)$            (*  $1 \dots K$  の乱数で初期化 *)
3:    $n'(z_n)++$ 
4:   for  $w$  in  $d_n$  do
5:      $m'(z_n, w)++$ 
6: while 収束するまで do           (* Gibbs サンプルング *)
7:   for  $n$  in  $\text{randperm}(N)$  do     (*  $1 \dots N$  をシャッフル *)
8:      $n'(z_n)--$                    (* カウントを減らす *)
9:     for  $w$  in  $d_n$  do
10:       $m'(z_n, w)--$ 
11:       $z_n$  を式(5.14)および式(5.49)に従ってサンプルング
12:       $n'(z_n)++$                    (* カウントを増やす *)
13:      for  $w$  in  $d_n$  do
14:         $m'(z_n, w)++$ 
15: 式(5.49)を用いて,  $p(z)$ ,  $p(w|z)$  の期待値を求める

```

図 5.12: 周辺化 Gibbs サンプルングによる UM のベイズ学習のアルゴリズム.

となります。 $\phi_k = \{p(w|k)\}$ についても、トピック $z_n = k$ となった文書だけを集めて、その中での各単語 w の頻度 $m(k, w)$ を計算すれば、 ϕ_k が事前分布

$$(5.46) \quad \phi_k \sim \text{Dir}(\beta, \dots, \beta)$$

に従っているとき、事後分布は式(4.49)とまったく同様に、

$$(5.47) \quad \phi_k | \mathbf{z} \sim \text{Dir}(\beta + m(k, 1), \dots, \beta + m(k, V))$$

となります。 \mathbf{z} をサンプルングした後で、 μ と ϕ を式(5.45)および式(5.47)からサンプルングし、これを繰り返せば、正しい事後分布からサンプルングを行って学習することができます。

これでもよいのですが、4.4.2節で学習したように、**周辺化 Gibbs サンプルング**を行えば μ および ϕ を積分消去し、より効率的に \mathbf{z} だけをサンプルングして学習することができます。いま、 n 番目の文書 d_n のトピック z_n をサンプルングしたいとしましょう。周辺化 Gibbs サンプルングでは、 \mathbf{z} をまずランダムに初期化し、 n 番目以外の文書のトピック $\mathbf{z}_{-n} = z_1, \dots, z_{n-1}, z_{n+1}, \dots, z_N$ がわかっているとしたときに、 z_n を事後分布 $p(z_n | d_n, \mathbf{z}_{-n})$ からサンプルングします。 \mathbf{z}_{-n}

表 5.7: UM のベイズ学習で得られたトピックの特徴語 (NPMI). ここではトピック数 $K=100$ として学習しました. EM アルゴリズムによる学習と比べて, 局所解に陥らず, はるかによいトピックが学習されていることがわかります.

Topic 1	Topic 10	Topic 20	Topic 30
アキレウス 0.643	南極 0.686	党 0.431	オリンピック 0.663
アステロパイオス 0.638	度 0.651	議員 0.387	団 0.635
ヒポポトオーン 0.633	西経 0.648	民主 0.385	北京 0.583
母音 0.619	北極 0.647	選挙 0.380	アテネ 0.582
ギリシア 0.615	東経 0.645	内閣 0.375	ロンドン 0.581
トロイア 0.607	線 0.623	大臣 0.368	バルセロナ 0.573
省略 0.598	成す 0.621	政治 0.365	結果 0.530
神話 0.574	角度 0.617	衆議 0.356	名簿 0.523
ゼウス 0.560	点 0.591	議会 0.344	競技 0.521
槍 0.548	通過 0.585	政党 0.343	選手 0.519

がわかっているのですから, 式(5.45)および式(5.47)において, 文書 d_n の分を除いた頻度を $m'(k, v)$ および $n'(k)$ とおけば, それぞれの事後分布は, 同様に

$$(5.48) \quad \begin{cases} \boldsymbol{\mu} | \mathbf{z}_{-n} \sim \text{Dir}(\alpha + n'(1), \dots, \alpha + n'(K)) \\ \boldsymbol{\phi}_k | \mathbf{z}_{-n} \sim \text{Dir}(\beta + m'(k, 1), \dots, \beta + m'(k, V)) \end{cases}$$

となります. この期待値は, デイリクレ分布の期待値の公式(3.32)から

$$(5.49) \quad \begin{cases} \mathbb{E}[\mu_k | \mathbf{z}_{-n}] = \frac{\alpha + n'(k)}{\sum_{k=1}^K (\alpha + n'(k))} & (\mu_k = p(z_n = k)) \\ \mathbb{E}[\phi_{kv} | \mathbf{z}_{-n}] = \frac{\beta + m'(k, v)}{\sum_{v=1}^V (\beta + m'(k, v))} & (\phi_{kv} = p(v | z_n = k)) \end{cases}$$

となります. 式(5.49)を式(5.14)の $p(z_n | d_n)$ の計算に用いれば, 文書 d_n のトピック z_n をサンプリングすることができます. これを, すべての $n=1, \dots, N$ について繰り返します. 以上を行う UM のベイズ学習 (周辺化 Gibbs サンプリング) のアルゴリズムを, 図 5.12 に示しました.

式(5.49)で期待値をとっているため, パラメータ $\theta = \{p(z), p(v|z)\}$ の点推定値はもはや存在しないことに注意してください. \mathbf{z} だけをサンプリングすれば, θ の分布は式(5.47)および式(5.45)の形でデイリクレ事後分布として表されますので, その期待値を式(5.49)のように求めることができます. この方法は, パ

ラメータの事後分布からサンプリングを行う完全な MCMC 法となっているため、局所解の危険がなく、理論的には十分長くサンプリングを行えば、図 4.20 のように真の最適解 (を含む事後分布全体) からのサンプルを得ることができます。

表 5.7 に、サポートサイトにある UM のベイズ学習の実装 `bun.py` を用いて、同じ日本語 Wikipedia コーパスについて $K=100$ のトピック数で学習した結果を示しました。ここでは、100 回の Gibbs サンプリング (MCMC) の繰り返しでモデルを学習しています。EM アルゴリズムと比べて、さまざまな可能性を確率的に試すことで、意味的にはるかにまとまりの高い潜在トピックが学習されていることがわかります。

5.3 ディリクレ混合モデル (DM)

ここまで、テキストの単語が式 (5.4) のように多項分布に従って発生すると仮定してきました。この多項モデルでは、確率 $1/100=0.01$ の単語が 2 回発生する確率は $(1/100)^2=1/10000=0.0001$ になります。しかし実際には、これはあまり正しくありません。図 5.13 に、『シャーロック・ホームズの冒険』の本文の中で特定の単語が出現した場所を示しました。これを見ると、“lestrade” (レストレード警部) は一度出現すると続けて何度も出現しており、これはそんなに低い確率ではなさそうです。言葉の出現のこうした性質を、**バースト性**とといいます。ただし、バースト性の度合いは単語によって異なり、同じ頻度 (= 確率) の “interest” は、ほぼ一様に出現していることがわかります。^{*23}

バースト性の度合いは、次のようにしても調べることができます [147]。いま、コーパスの各文書を前半と後半に分け、2 章で行ったように前半を ‘train’ データ、後半を ‘test’ データとよぶことにしましょう。表 5.8 に、Livedoor コーパスの 7367 文書の中で、“ソフトバンク” が前半と後半に 1 回以上出現した文書の数をまとめました。^{*24} a は両方出現した、 b は前半のみ出現し後半にはなかった、 c は前半にはなかったが後半に出現した、 d は前半と後半どちらにも出現しなかった文書の数をそれぞれ表しています。

*23 バースト性のモデル化について。

*24 これは、サポートサイトにある `recur.py` で調べることができます。

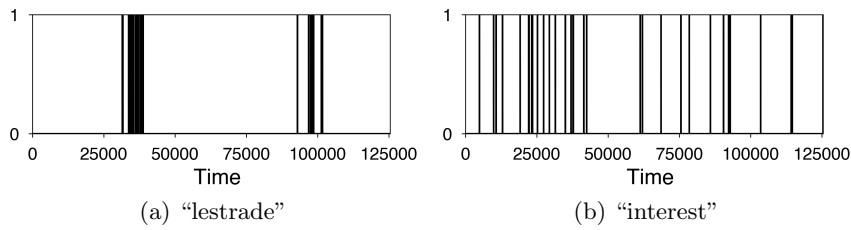


図 5.13: 『シャーロック・ホームズの冒険』での “lestrade” と “interest” の出現位置. この 2 つの単語の出現頻度はどちらも同じ 38 回ですが, “lestrade” (レストレード警部) のような言葉は, テキストの中で一部にバースト的に出現しています.

表 5.8: Livedoor コーパスの中で「ソフトバンク」が出現した文書の数. x はそこで単語が出現したことを, \bar{x} は出現しなかったことを表します.

	test	$\bar{\text{test}}$
train	$a = 108$	$b = 93$
$\bar{\text{train}}$	$c = 93$	$d = 7073$

“ソフトバンク” の文書レベルでの出現確率は $(a+c)/(a+b+c+d) = 0.027$ ですので, 前半に出現した中 $(a+b)$ で, 後半にも出現する a の確率は, 単語の出現が独立ならば 0.03 程度になるはずですが. しかし, この場合 $a/(a+b) = 0.537$ となっており, これは非常に高い確率です. つまり, “ソフトバンク” は 1 回出現すれば, 同じ文書でもう 1 回以上出現する確率は 0.537 もあるということです. これは明らかに, 単語のバースト性を示しています. [147]で (当時) ベル研究所の Church は, 出現確率 p の低い “Noriega” を例にとって 「“Noriega” が 2 回出現する (Two Noriegas) 確率は, p^2 ではなく $p/2$ だ」という言葉で, この現象を経験的に指摘しました. われわれの例でも, “ソフトバンク” が 2 回出現する確率は確かに, $p \cdot 0.537 \simeq p/2$ になっていることがわかります.

もちろん, バースト性は単語によって異なり, たとえば “順調” では $a/(a+b)$ は 0.065 と低い確率でした. テキストの確率モデルは, こうしたバースト性を表現できると望ましいでしょう.

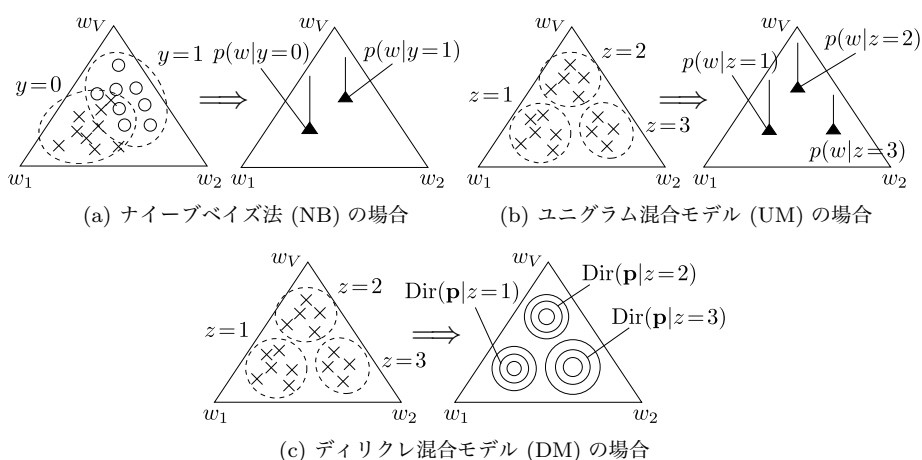
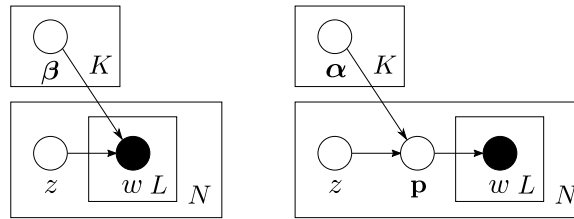


図 5.14: 文書モデルの幾何的表現. すべての多項分布は, $w_1 \dots w_V$ を頂点とする単語単体の中に存在しています. NB や UM では, それらをトピックの代表点 (▲) で針の立った δ 関数で表しますが, DM ではそれぞれのトピックをディリクレ分布で表現します.

5.3.1 単語単体と幾何的解釈

ここで, 少し違う視点からナイーブベイズ法と UM について振り返ってみましょう. 5.1 節の単語集合の仮定から, 各文書には含まれる単語を生成した V 次元の確率分布 $\mathbf{p} = (p_1, p_2, \dots, p_V)$ があることとなります. 3.4.1 節でみたように, この \mathbf{p} は図 5.14 の各図の左側のような, 各単語 w_1, w_2, \dots, w_V を頂点とする V 次元の単体の中に存在しています. これを**単語単体**といいます. ナイーブベイズ法や UM は, この単語単体上に点として, トピック分布 $\mathbf{p} = \{p(w|y)\}$ や $\mathbf{p} = \{p(w|z)\}$ を図 5.14(a) や (b) のように学習するものだと思います. ナイーブベイズ法は UM の教師あり版ですから, 以下では UM に統一して考えていくことにします.

しかし, よく考えると UM のモデルは, 制限が強すぎるのではないのでしょうか. 文書を生成したもともとの \mathbf{p} には様々な可能性があるにもかかわらず, UM では, トピック z ごとに特定の $\mathbf{p} = \{p(w|z)\}$ で代表してしまっています. 「単語の出やすさ」 \mathbf{p} のタイプが完全に固定されてしまっているので, むしろ, 図 5.14(c) のようにトピックごとに**分布**を考え, この分布から \mathbf{p} が生まれたと考える方が



(a) ユニグラム混合モデル (UM) (b) デリクレ混合モデル (DM)

図 5.15: UM と DM のグラフィカルモデル. ●は観測値, ○は確率変数を表します. 四角いプレートは繰り返しを, 右隅は繰り返しの数を表しています. DM では, 文書ごとに異なる \mathbf{p} があるのが特徴です. 実際には, この \mathbf{p} は積分消去することができます.

適切ではないでしょうか. \mathbf{p} 自体が多項分布なので, これは**確率分布の確率分布**となります.

この分布として最も簡単なのは, 3.4.1 節のデリクレ分布を考えることでしよう. 図 5.14(c) のように, 単語単体上にトピックを表現する複数のデリクレ分布を考えるこのモデルを, **デリクレ混合モデル** (Dirichlet Mixtures, **DM**) といいます.*25 DM は最初, バイオインフォマティクスの分野でアミノ酸配列の解析のために 1996 年に提案されました[148]. 筑波大学の山本らは 2003 年に, これがテキストにも適用でき, 5.4 節で説明する LDA よりも高い文書確率を与えることを示しました[149].*26

デリクレ混合モデルの生成モデル UM では, 文書 $d = w_1 w_2 \dots w_L$ は次のようにして生成されたと考えました.

1. $z \sim \lambda$ でトピック z を選択.
2. For $i = 1, \dots, L$,
 - a. $w_i \sim p(w|z)$ から単語 w_i を生成.

すなわち, 各単語はトピック z で決まるユニグラム分布 $\beta_z = \{p(w|z)\}$ から生成されると考えています. これをグラフィカルモデルで表すと, 図 5.15(a) のようになります. これに対して, DM では文書は

*25 これは, 通常のユークリッド空間において K 平均法の代わりに, ガウス混合モデルを考えると似ています. デリクレ分布は, 単体上のガウス分布のようなものと考えてよいでしょう.

*26 東京大学の佐藤らの研究[150]を用いれば, Pitman-Yor 過程を用いる洗練された方法で DM とこの後の LDA は統合でき, さらに高い文書確率を与えることが確かめられています.

アルゴリズム 5: DM の EM-Newton アルゴリズム.

-
- ```

1: $\lambda, \alpha_1, \dots, \alpha_K$ を適当に初期化する
2: while 収束するまで do
3: for n in $1, \dots, N$ do (* E ステップ *)
4: 式(5.51) で $p(z|d_n)$ を計算
5: 式(5.52) で $\lambda, \alpha_1, \dots, \alpha_K$ を更新 (* M ステップ *)

```
- 

図 5.16: DM の EM アルゴリズムによる学習.

1.  $z \sim \lambda$  でトピック  $z$  を選択.
2.  $\mathbf{p} \sim \text{Dir}(\alpha_z)$  でユニグラム分布  $\mathbf{p}$  を生成.
3. For  $i = 1, \dots, L$ ,
  - a.  $w_i \sim \mathbf{p}$  で単語  $w_t$  を生成.

のようにして生成されたと考えます. グラフィカルモデルで表すと, 図 5.15(b) のようになります. UM では決まったユニグラム分布  $\beta_1, \dots, \beta_K$  の中から選ぶのに対し, DM では文書ごとに異なるユニグラム分布  $\mathbf{p}$  が, トピックに従って毎回生成されるのが特徴です. これにより, DM ではたとえば同じ車の話題を扱う文書でも, 「トヨタ」が多い場合と「日産」が多い場合で異なる  $\mathbf{p}$  が使われていると考えて区別することができます. 後で数学的に示すように, このことから単語のバースト性を説明することができます.

幸いにして  $\mathbf{p}$  は直接求める必要はなく, 3章で説明したように, 期待値をとって積分消去することができます. いま, 文書  $d$  のトピックが  $z=k$  とわかっていたとしましょう. すると文書の確率は,  $\mathbf{p}$  が  $k$  番目のデイリクレ分布  $\text{Dir}(\alpha_k)$  ( $\alpha_k = (\alpha_{k1}, \alpha_{k2}, \dots, \alpha_{kV})$ ) から生成されたということなので, 期待値をとれば

$$\begin{aligned}
 (5.50) \quad p(d|z=k) &= \int p(d, \mathbf{p} | z=k) d\mathbf{p} = \int p(d|\mathbf{p})p(\mathbf{p}|z=k) d\mathbf{p} \\
 &= \int \prod_{i=1}^L \prod_{v=1}^V p_v^{\mathbb{I}(w_i=v)} \cdot \text{Dir}(\mathbf{p}|\alpha_k) d\mathbf{p} \\
 &= \int \prod_{v=1}^V p_v^{n_v} \cdot \frac{\Gamma(\sum_v \alpha_{kv})}{\prod_v \Gamma(\alpha_{kv})} \prod_{v=1}^V p_v^{\alpha_{kv}-1} d\mathbf{p} \\
 &= \frac{\Gamma(\sum_v \alpha_{kv})}{\Gamma(\sum_v \alpha_{kv} + L)} \prod_{v=1}^V \frac{\Gamma(\alpha_{kv} + n_v)}{\Gamma(\alpha_{kv})}
 \end{aligned}$$

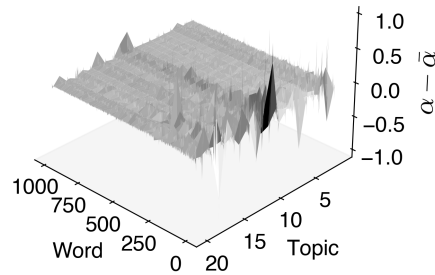


図 5.17: Livedoor コーパスで学習された DM のパラメータ  $\alpha_1, \dots, \alpha_K$  の様子. 全体の平均  $\bar{\alpha}$  との差分をプロットしています. (単語 1~1000, トピック数  $K=20$ )

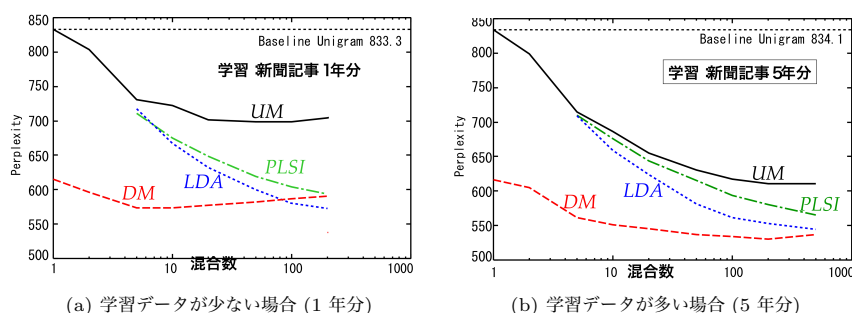
となり, 式(3.58)の**ポリア分布**で表すことができます. ここで, 1行目では同時確率の周辺化(公式(2.10))および確率の連鎖則(公式(2.20))を用いました.  $n_v = \sum_{i=1}^L \mathbb{I}(w_i=v)$  は,  $d$  の中で単語  $v$  が何回出たかを表します. 実際には  $z$  は未知ですから, DM での文書の確率は,  $z$  についても期待値をとり,

$$\begin{aligned}
 (5.51) \quad p(d) &= \sum_{k=1}^K p(d, z=k) = \sum_{k=1}^K p(d|z=k)p(z=k) \\
 &= \sum_{k=1}^K \lambda_k \frac{\Gamma(\sum_v \alpha_{kv})}{\Gamma(\sum_v \alpha_{kv} + L)} \prod_{v=1}^V \frac{\Gamma(\alpha_{kv} + n_v)}{\Gamma(\alpha_{kv})} \quad (\text{DM の文書確率})
 \end{aligned}$$

となります.\*27

**ディリクレ混合モデルの学習** DM のパラメータは,  $p(z)$  を表す  $\lambda = \{\lambda_1, \dots, \lambda_K\}$  およびトピック別のディリクレ分布のハイパーパラメータ  $\{\alpha_{kv} | k=1, \dots, K, v=1, \dots, V\}$  です. 他のトピックモデルも同様ですが, 一般に  $K=100$ , 語彙の大きさ  $V=10000$  と少なめするときでも,  $\alpha_{kv}$  は 100 万個程度のパラメータ数となることに注意してください.  $\alpha_{kv}$  は EM アルゴリズムの中で式(3.60)のニュートン法で最適化することもできますが, これは特殊関数であるダイガンマ関数  $\Psi(x) = d/dx \log \Gamma(x)$  が含まれているため, テキストデータについては非常に重い計算になります.

\*27 すなわち, 「ディリクレ分布の混合モデル」となっているのは  $\mathbf{p}$  に対してのもので,  $\mathbf{p}$  を積分消去すると, DM は観測値については「混合ポリア分布」になっています.



(a) 学習データが少ない場合 (1年分)

(b) 学習データが多い場合 (5年分)

図 5.18: 各トピックモデルの性能 ([151]から引用). 縦軸はテストデータに対するパープレキシティです. いずれの場合も, 単語のバースト性を捉える DM の性能が最も高くなっており, 混合数を増やすと LDA の性能と近づいてくることがわかります.

山本ら [149] はこれに代わり, [70] で示された高速な LOO (Leave-One-Out) 近似を使って, 図 5.16 のアルゴリズムおよび次式で DM のパラメータを高速かつ安定に学習できることを示しました.

$$(5.52) \quad \begin{cases} \lambda_k \propto \sum_{n=1}^N p(z_n = k | d_n) \\ \alpha'_{kv} = \alpha_{kv} \cdot \frac{\sum_{i=1}^N p(z_i = k | d_i) n_{iv} / (n_{iv} - 1 + \alpha_{kv})}{\sum_{i=1}^N p(z_i = k | d_i) n_i / (n_i - 1 + \sum_v \alpha_{kv})} \end{cases}$$

ここで  $n_{iv}$  は文書  $i$  における単語  $v$  の頻度,  $n_i = \sum_v n_{iv}$  は文書  $i$  の長さです. この式は特殊関数を含まないため計算が速く, 筆者の公開している C 言語による実装<sup>\*28</sup> を用いた場合, Livedoor コーパスで  $K = 100$  の場合は, 執筆時の環境では 3 分程度で EM アルゴリズムが収束します. Python では非常に遅くなるため, 注意してください. なお, この推定法は最尤推定のため過学習しやすく, 山本らはさらに, 階層ベイズの考え方で推定を安定化する Smoothed DM と, 変分ベイズ法の一つによる解法を示し [152], こちらの方が安定して高性能であることが報告されています.<sup>\*29</sup> 図 5.18 に, UM, DM および次節で説明する LDA を Livedoor コーパスで学習した場合の, テストデータのパープレキシティをト

\*28 <http://chasen.org/~daiti-m/dist/dm/>にある `dm-0.2.tar.gz` で C 言語による実装を公開しています.

\*29 <http://chasen.org/~daiti-m/dist/dm/>にある `sdm-0.2.tar.gz` で C 言語による実装を公開しています.

ピック数  $K$  ごとに示しました。これからわかるように、パープレキシティ(文書確率)の面では、DM は LDA よりも優れたモデルです。ただしこれは主に、文書内での単語のバースト性のモデル化から生じるもので、LDA の方が柔軟なモデルであることに注意してください。

なお、DM の各トピックを表すディリクレ分布  $\text{Dir}(\alpha_k)$  の期待値は、公式(3.32)で示したように  $\alpha_k$  の和を 1 に正規化すると求めることができ、これは UM におけるトピック分布  $p(w|z)$  と等価だとみなすことができます。DM ではどのようなトピックが推定されているか、5.2.1 節の UM の場合と同様にして調べてみると面白いでしょう。(→演習 5)

### 5.3.2 ポリア分布と単語のバースト性

いま、文書を生成したディリクレ分布  $\text{Dir}(\alpha)$  のパラメータを  $\alpha = (\alpha_1, \dots, \alpha_V)$  とすると、式(5.51)に示したように、文書  $d$  の確率はポリア分布

$$(5.53) \quad p(d) = \frac{L!}{\prod_{v \in d} n_v!} \frac{\Gamma(\sum_v \alpha_v)}{\Gamma(\sum_v \alpha_v + L)} \prod_{v \in d} \frac{\Gamma(\alpha_v + n_v)}{\Gamma(\alpha_v)}$$

となります。長さ  $L$  の文書で各単語の頻度が  $n_v$  になる組み合わせの数は、多項係数  $L!/\prod_v n_v!$  だけありますから、上では明示的に加えました。

ここで、3.4.3 節の図 3.17 で述べたように、言語の場合は  $\alpha_v$  は非常に小さく、ほとんどの場合は 0.01 あるいは 0.001 未満の値となることを思い出してください。このとき、式(5.53)で各単語に依存するファクター  $\Gamma(\alpha_v + n_v)/\Gamma(\alpha_v)$  は、式(3.39)の  $\Gamma$  関数の性質から

$$(5.54) \quad \begin{aligned} \frac{\Gamma(\alpha + n)}{\Gamma(\alpha)} &= \frac{(n + \alpha - 1)(n + \alpha - 2) \cdots \alpha \cancel{\Gamma(\alpha)}}{\cancel{\Gamma(\alpha)}} \\ &= \alpha(\alpha + 1) \cdots (n + \alpha - 1) \\ &\simeq \alpha(n - 1)! \quad (\because \alpha \ll 1) \end{aligned}$$

と近似することができます。図 5.19 に、この近似が実際に非常に正確であることを示しました。よって、この近似を式(5.53)に代入すれば

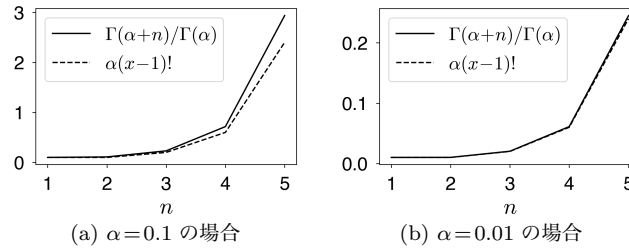


図 5.19: Pochhammer 関数  $\Gamma(\alpha+n)/\Gamma(\alpha)$  の近似.  $\alpha$  が小さい場合には, この値は近似  $\alpha(x-1)!$  と, ほとんど差がないことがわかります.

$$\begin{aligned}
 (5.55) \quad p(d) &\simeq \frac{L!}{\prod_v n_v!} \frac{\Gamma(\sum_v \alpha_v)}{\Gamma(\sum_v \alpha_v + L)} \prod_{v \in d} \alpha_v (n_v - 1)! \\
 &= L! \frac{\Gamma(\sum_v \alpha_v)}{\Gamma(\sum_v \alpha_v + L)} \prod_{v \in d} \frac{\alpha_v}{n_v} \quad (\text{指数分布族 DCM 分布})
 \end{aligned}$$

となります.  $n_v$  の項をくくり出せば, これは統計学で標準的な指数分布族の形で表すことができ, Elkan [153] はこれを指数分布族 DCM (EDCM) 分布と呼んでいます.

式(5.55)で, 文書中の単語  $v$  に依存するファクターは  $n \mid \begin{array}{cccc} 1 & 2 & 3 & \dots \\ p & \alpha_v & \frac{\alpha_v}{2} & \frac{\alpha_v}{3} & \dots \end{array}$   $\alpha_v/n_v$  になっており, これは右の表のように,  $n_v = 1$  のときは確率が  $\alpha_v$  倍,  $n_v = 2$  のときは  $\alpha_v/2$  倍, ... となることを意味しています. つまり, (頻度の低い) 単語  $v$  が 2 回出現しても, その確率は  $\alpha_v^2$  ではなく,  $\alpha_v/2$  倍にしかならないということです. これはまさに, 256 ページの “Two Noriegas” の発見そのものです. 単語が  $n+1$  回現れる確率は,  $n$  回現れる確率の  $n/(n+1)$  倍になっており, これから単語は一度現れれば, バースト的に出現することになります. ポリア分布を使うことで, われわれは実際にこの現象が成り立つことを数学的に説明することができました.

## 5.4 潜在ディリクレ配分法 (LDA)

これまで, UM や DM では各文書が 1 つのトピックから生成されたと仮定してきましたが, これには明らかに限界があると考えられます. たとえば, 「劇場



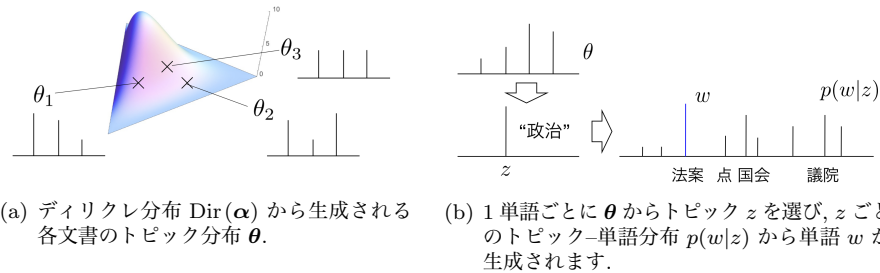


図 5.20: LDA の文書生成モデル. LDA では, 文書の各単語ごとに異なる潜在トピック  $z$  が存在します.

の改修」を伝えるニュースでは, 演劇についての言葉と建築についての言葉が両方現れるでしょう. また, 数理経済学の論文では, 数学の話と経済学の話が混じっていると考えられます. これらを単一のトピックで表現しようとする, すべての組み合わせに別々のトピックを用意しなければならないとなってしまいます. 組み合わせは2つとは限りませんから, これには膨大なトピック数が必要になります. また, 「数理経済学」と「開発経済学」のトピックが経済学という話題を共有していることもわからなくなってしまうでしょう.

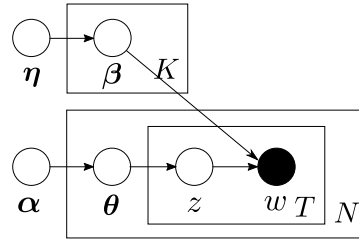
よって, 各文書にトピックが1つだけあるのではなく, 図 5.20(a) のようにトピックの**確率分布**  $\theta$  があると考えた方がよいでしょう. たとえばトピックが  $K=4$  種類あり, それぞれ「演劇」「経済」「建築」「数学」だったとき, 上の劇場の改修のニュースは  $\theta = (0.7, 0, 0.3, 0)$ , 数理経済学の論文は  $\theta = (0, 0.4, 0, 0.6)$  のようになると考えられます. 数学的には, 各文書  $d$  についてトピック  $z$  の確率分布  $\theta = \{p(z|d)\}$  を考えることになります. UM や DM は, これをあるトピック  $k$  についてだけ 1 をとるデルタ関数  $p(z|d) = \delta(z=k)$  で近似してしまうモデルとみなすことができます\*30.

この  $\theta$  は文書ごとに異なり, 最も簡単には 3 章で学習したディリクレ分布

$$(5.56) \quad \theta \sim \text{Dir}(\alpha) \quad (\alpha = (\alpha_1, \dots, \alpha_K))$$

\*30 なお, 単一トピックの UM や DM でも, 学習の際には  $\delta(z=k)$  の推定値として図 5.9 のように  $p(z|d)$  を考えることにはなりますが, 真の値はあくまで単一のトピックで, それを確率的に推定しているにすぎません. それに対して, ここでは真の値  $p(z|d)$  自体が多項分布であり, その推定値として, 多項分布の分布であるディリクレ分布が事後分布として得られることになります.

- For  $k = 1, \dots, K$ ,
  - Draw  $\beta_k \sim \text{Dir}(\eta)$ .
- For  $i = 1, \dots, N$ ,
  - Draw  $\theta_i \sim \text{Dir}(\alpha)$ .
  - For  $t = 1, \dots, T$ ,
    - \* Draw  $z_{it} \sim \theta_i$ .
    - \* Draw  $w_{it} \sim \beta_{z_{it}}$ .



(a) LDA の生成モデル

(b) LDA のグラフィカルモデル

図 5.21: LDA の生成モデルと、対応するグラフィカルモデル。●は観測された変数を、○は潜在変数を表しています。

から生成されたと考えるといいでしょう。図 5.20(b) に示したように、この  $\theta$  から、文書の各単語ごとに

- (1) まずトピック  $z \sim \theta$  をランダムに選び、
- (2)  $z$  ごとのトピック-単語分布  $p(w|z)$  から単語  $w$  を生成する

ことを繰り返して、文書が生成されたとするモデルを考えることができます。このトピック-単語分布  $p(w|z)$  も、 $z$  ごとにディリクレ事前分布  $\text{Dir}(\eta)$  から生成されたとしましょう。以上の生成モデルを、アルゴリズムの形で図 5.21(a) に示しました。学習の際には、文書ごとに潜在変数  $\theta$  のディリクレ事後分布を割り当てて推定するため、このモデルを**潜在ディリクレ配分法** (Latent Dirichlet Allocation, **LDA**) [141]とといいます。LDA は後で説明する PLSI のベイズ化として、Blei らにより 2001 年の論文[154]で最初に提案されました。

LDA のグラフィカルモデルを図 5.21(b) に示しました。図 5.15 の UM のグラフィカルモデルと比べると、 $z$  が内側の箱の中に入っており、文書に含まれる 1 つ 1 つの単語  $w$  ごとに潜在トピック  $z$  があることを示しています。100 万単語のコーパスがあれば、100 万個の  $z$  があるわけです。よって LDA は、UM よりずっと複雑、かつ柔軟な確率モデルになっています。

### Gibbs サンプリングによる学習

LDA の学習法には (1) 最初に提案された変分ベイズ法[141], (2) 周辺化の考え方を取り入れた周辺化変分ベイズ法[155], (3) Gibbs サンプリング[156]

[157] などがあります\*31。この中で、5.2節で説明したEMアルゴリズムのベイズ版である変分ベイズ法は数学的導出の難しさにもかかわらず、EMアルゴリズムの一種のため局所解の問題を逃れることができず、図5.22に示したように、Gibbsサンプリング(MCMC法)による学習が最終的に、最も性能が高いことがわかっています。導出や実装も容易なため、本書ではGibbsサンプリングによる学習について説明します。変分ベイズ法による学習について知りたい方は、教科書[71]を参照してください。

上の生成モデルによれば、LDAで文書  $\mathbf{w} = w_1 w_2 \cdots w_T$  と背後にある潜在変数

- (1) 各単語の潜在トピック  $\mathbf{z} = z_1 z_2 \cdots z_T$ ,
- (2) 文書の潜在トピック分布  $\boldsymbol{\theta}$ ,
- (3) トピック-単語分布  $\boldsymbol{\beta} = \{\beta_{kv} = p(v|k)\}$

の同時確率は、

$$(5.57) \quad p(\mathbf{w}, \mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\beta} | \boldsymbol{\alpha}, \boldsymbol{\eta}) = p(\mathbf{w} | \mathbf{z}, \boldsymbol{\beta}) p(\mathbf{z} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \boldsymbol{\alpha}) p(\boldsymbol{\beta} | \boldsymbol{\eta})$$

と分解することができます。具体的には、204ページで導入した指示関数  $\mathbb{I}()$  による記法を使うと、

$$(5.58) \quad p(\mathbf{w} | \mathbf{z}, \boldsymbol{\beta}) p(\mathbf{z} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \boldsymbol{\alpha}) p(\boldsymbol{\beta} | \boldsymbol{\eta}) \\ = \left( \prod_{t=1}^T p(w_t | z_t, \boldsymbol{\beta}) p(z_t | \boldsymbol{\theta}) \right) \cdot \text{Dir}(\boldsymbol{\theta} | \boldsymbol{\alpha}) \cdot \prod_{k=1}^K \text{Dir}(\boldsymbol{\beta}_k | \boldsymbol{\eta}) \\ \propto \underbrace{\prod_{t=1}^T \prod_{v=1}^V \prod_{k=1}^K \beta_{kv}^{\mathbb{I}(w_t=v) \mathbb{I}(z_t=k)}}_{= \beta_{z_t w_t}} \cdot \underbrace{\prod_{t=1}^T \prod_{k=1}^K \theta_k^{\mathbb{I}(z_t=k)}}_{= \theta_{z_t}} \cdot \prod_{k=1}^K \theta_k^{\alpha_k - 1} \cdot \prod_{k=1}^K \prod_{v=1}^V \beta_{kv}^{\eta - 1}$$

(LDAの文書同時確率)

と表すことができます。4章(204ページ)で説明したように、パラメータ  $\beta_{kv} = p(v|k)$  を表に出すために、 $p(w_t | z_t, \boldsymbol{\beta}) = \beta_{k v}$  (ただし  $k = z_t, v = w_t$ ) を指示関数  $\mathbb{I}()$  を使って、「ただし」の部分を実学的に

\*31 期待値伝搬法(EP)による学習も提案されていますが[158]、この場合はすべての単語についてそのトピック事後分布を保持しなければならないため、学習がスケールしないことが問題です。

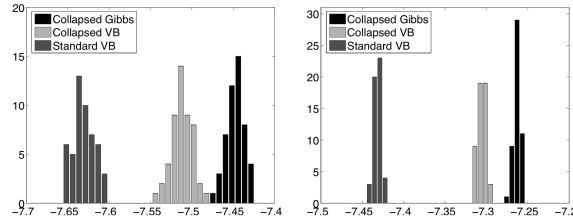


図 5.22: LDA の学習アルゴリズムの違いによる性能の比較 ([155]より引用). 左側は DailyKOS のニュース記事, 右側は NIPS の論文コーパスの結果で, 横軸はテストデータの 1 単語あたりの対数尤度を表します. いずれも, 変分ベイズ法に比べて周辺化 Gibbs サンプルング (黒棒) がもっとも高い対数尤度を与えていることがわかります.

$$(5.59) \quad p(w_t | z_t, \beta) = \prod_{v=1}^V \prod_{k=1}^K \beta_{kv}^{\mathbb{I}(z_t=k)\mathbb{I}(w_t=v)}$$

と表していることに注意してください.  $\theta_k$  についても同様です.

実際には, 文書は  $\mathbf{w}_1^N = \mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N$  の  $N$  個ありますから, データ全体の同時確率は, 対応する潜在トピックを  $\mathbf{z}_1^N = \mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N$ , 各文書のトピック分布を  $\theta_1^N = \theta_1, \theta_2, \dots, \theta_N$  とし

$$(5.60) \quad p(\mathbf{w}_1^N, \mathbf{z}_1^N, \theta_1^N, \beta | \alpha, \eta) \\ \propto \prod_{i=1}^N \left[ \prod_{t=1}^T \prod_{v=1}^V \prod_{k=1}^K \beta_{kv}^{\mathbb{I}(w_{it}=v)\mathbb{I}(z_{it}=k)} \cdot \prod_{t=1}^T \prod_{k=1}^K \theta_{nk}^{\mathbb{I}(z_{it}=k)} \prod_{k=1}^K \theta_{nk}^{\alpha_k - 1} \right] \prod_{k=1}^K \prod_{v=1}^V \beta_{kv}^{\eta - 1}$$

と表されます. 記法を簡単にするために, 文書の長さ  $T$  を上では同じにしていますが, 実際は文書によって異なる値であることに注意してください.

#### LDA の Gibbs サンプルング

このうち, わかっているのは観測された文書  $\mathbf{w}$  だけで, ほかの潜在変数はすべて未知の潜在変数で,  $\mathbf{w}$  から推定します. LDA のパラメータ  $\mathbf{z}, \theta, \beta$  は, 4 章で学んだ Gibbs サンプルングを使って, 順番にサンプルングすることを繰り返せば推定することができます.

**z のサンプルング** 文書  $n$  の  $t$  番目の単語  $w_{it}$  のトピックを表す潜在変数  $z_{it}$  は, 式(5.57)の 1 行目で  $z$  に関連する項は  $p(\mathbf{w} | \mathbf{z}, \beta)p(\mathbf{z} | \theta)$  だけですから, 式(5.60)

から対応する項を抜き出せば,

$$(5.61) \quad p(z_{it} | \mathbf{z}_{-it}, \boldsymbol{\theta}, \boldsymbol{\beta}) \propto p(w_{it} | z_{it}, \boldsymbol{\beta}) p(z_{it} | \mathbf{z}_{-it}, \boldsymbol{\theta}_i) \\ = \beta_{kv} \cdot \theta_{ik} \quad (v = w_{it})$$

となります. これは, 意味的には現在の文書を  $\mathbf{w}_i = d$ , 単語を  $w_{it} = v$ , サンプルングしたい  $z_{it} = k$  とおくと,

$$(5.62) \quad p(k | v, d) \propto p(v, k | d) = \underbrace{p(v | k)}_{\beta_{kv}} \underbrace{p(k | d)}_{\theta_{ik}}$$

を計算している, ということです. したがって,  $z_{it}$  は確率

$$(5.63) \quad p(z_{it} = k | \mathbf{z}_i \setminus z_{it}, \boldsymbol{\theta}, \boldsymbol{\beta}) = \frac{\theta_{ik} \beta_{kv}}{\sum_{k=1}^K \theta_{ik} \beta_{kv}}$$

に従ってサンプリングすればよいことがわかります.

**$\boldsymbol{\theta}$  のサンプリング** 同様に, 式(5.57)で  $\boldsymbol{\theta}$  について注目すると

$$(5.64) \quad p(\boldsymbol{\theta} | \mathbf{z}, \boldsymbol{\alpha}) \propto p(\mathbf{z} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \boldsymbol{\alpha})$$

ですから, 式(5.60)から対応する項を抜き出して

$$(5.65) \quad p(\boldsymbol{\theta}_i | \mathbf{w}_i, \mathbf{z}_i) \propto \prod_{k=1}^K \theta_{ik}^{\alpha_k - 1} \cdot \prod_{t=1}^T \prod_{k=1}^K \theta_{nk}^{\mathbb{I}(z_{it}=k)} = \prod_{k=1}^K \theta_{nk}^{\alpha_k - 1 + \sum_{t=1}^T \mathbb{I}(z_{it}=k)}$$

となり,  $n(i, k) = \alpha_k + \sum_{t=1}^T \mathbb{I}(z_{it}=k)$  とおけば,  $\theta_i$  はディリクレ事後分布

$$(5.66) \quad p(\boldsymbol{\theta}_i | \mathbf{w}_i, \mathbf{z}_i) = \text{Dir}(\alpha_1 + n(i, 1), \dots, \alpha_K + n(i, K))$$

からサンプリングすればよいことがわかります.

**$\boldsymbol{\beta}$  のサンプリング** 最後に, 最も重要なトピック-単語分布  $\boldsymbol{\beta}$  のサンプリングを考えましょう. 式(5.57)で  $\boldsymbol{\beta}$  に対応する項は

$$(5.67) \quad p(\boldsymbol{\beta} | \mathbf{w}, \mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\alpha}) \propto p(\mathbf{w} | \mathbf{z}, \boldsymbol{\beta}) p(\boldsymbol{\beta})$$

ですから,  $\boldsymbol{\beta}$  がトピック  $k$  ごとにディリクレ事前分布

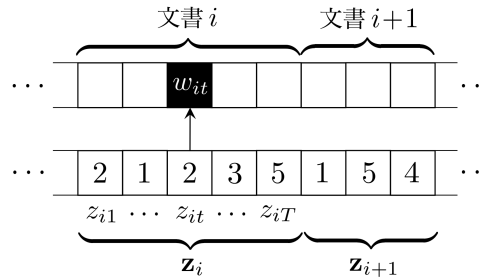


図 5.23: LDA の Gibbs サンプリングの様子. 各単語  $w_{it}$  を生成した潜在トピック  $z_{it}$  を, 事後分布からランダムにサンプリングして更新することを繰り返します.

$$(5.68) \quad p(\beta_k) = \text{Dir}(\eta, \dots, \eta) \propto \prod_{v=1}^V \beta_{kv}^{\eta-1}$$

に従うとすれば,  $\beta_k$  の事後分布は式(5.60)から

$$(5.69) \quad p(\beta_k | \mathbf{w}_1^N, \mathbf{z}_1^N) \propto \prod_{i=1}^N \prod_{t=1}^T \prod_{v=1}^V \beta_{kv}^{\mathbb{I}(w_{it}=v) \mathbb{I}(z_{it}=k)} \cdot \prod_{v=1}^V \beta_{kv}^{\eta-1}$$

$$= \prod_{v=1}^V \beta_{kv}^{\eta-1 + \sum_{i=1}^N \sum_{t=1}^T \mathbb{I}(w_{it}=v) \mathbb{I}(z_{it}=k)}$$

となり,  $m(k, v) = \eta + \sum_{i=1}^N \sum_{t=1}^T \mathbb{I}(w_{it}=v) \mathbb{I}(z_{it}=k)$  とおけば, こちらもディリクレ事後分布

$$(5.70) \quad p(\beta_k | \mathbf{w}_1^N, \mathbf{z}_1^N) = \text{Dir}(\eta + m(k, 1), \dots, \eta + m(k, V))$$

からサンプリングすることができます.

式(5.63), 式(5.66), 式(5.70) に従って  $\mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\beta}$  を次々とサンプリングすることを繰り返せば, LDA のパラメータを求めることができます. 実は, データ量やトピック数は圧倒的に少ないものの, 本質的に同じモデルが LDA より 1 年前にバイオインフォマティクスの分野で Pritchard ら [159] によって提案されており, 推論にはこの方法が使われていました.

---

**アルゴリズム 6:** LDA の周辺化 Gibbs サンプリング

---

```

1: for i in $1, \dots, N$ do (* 初期化 *)
2: for t in $1, \dots, T$ do
3: $z_{it} = \text{randint}(K)$ (* $1 \dots K$ の乱数で初期化 *)
4: $n'(i, z_{it})++$
5: $m'(w_{it}, z_{it})++$
6: while 収束するまで do (* Gibbs サンプリング *)
7: for i in $\text{randperm}(N)$ do (* 文書 i をランダムに選ぶ *)
8: for t in $1, \dots, T$ do
9: $n'(i, z_{it})--$ (* カウントを減らす *)
10: $m'(w_{it}, z_{it})--$
11: z_{it} を式 (5.75) に従って サンプリング
12: $n'(i, z_{it})++$ (* カウントを増やす *)
13: $m'(w_{it}, z_{it})++$
14: 式 (5.76) から θ_1^N, β の期待値を求める

```

---

図 5.24: 周辺化 Gibbs サンプリングによる LDA の学習アルゴリズム。

### LDA の周辺化 Gibbs サンプリング

しかし、パラメータ数が桁違いに大きい自然言語の場合は、HMM (210 ページ) や UM (251 ページ) でもわれわれが行ったように、パラメータについて期待値をとり、可能な限り積分消去することで、より効率的なサンプリングを行うことができます。実際には、こちらの方法で学習を行うのがよいでしょう。

式(5.61)で  $z_{it}$  のサンプリングを行う際に、 $\theta_{ik}$  が必要になります。ただし、 $z_{it}$  以外の  $\mathbf{z}_i \setminus z_{it}$  はわかっているのですから、 $\theta_i$  の事後分布は式(5.66)から、 $\mathbf{z}_i \setminus z_{it}$  の中でトピック  $k$  に割り当てられた単語の数を  $n'(i, k)$  とおくと、

$$(5.71) \quad p(\theta_i | \mathbf{z}_i \setminus z_{it}) = \text{Dir}(\alpha_1 + n'(i, 1), \dots, \alpha_K + n'(i, K))$$

です。よって、その期待値

$$(5.72) \quad \mathbb{E}[\theta_{ik} | \mathbf{z}_i \setminus z_{it}] = \frac{\alpha_k + n'(i, k)}{\sum_k (\alpha_k + n'(i, k))}$$

を使うことができます。

同様に  $\beta_{kv}$  についても、その期待値で置き換えることができます。 $\beta_k$  の

事後分布は式(5.70)から,

$$(5.73) \quad p(\beta_k | \mathbf{w}_1^N, \mathbf{z}_1^N \setminus z_{it}) = \text{Dir}(\eta + m'(k, 1), \dots, \eta + c'(k, V))$$

となることに注意しましょう. ここで  $m'(k, v)$  は  $z_{it}$  を除いたすべての潜在変数  $\mathbf{z}_1^N \setminus z_{it}$  の中で, 単語  $v$  にトピック  $k$  が割り当てられた回数を表しています. よって, 式(5.63)の  $\beta_{kv}$  の部分は

$$(5.74) \quad \mathbb{E}[\beta_{kv} | \mathbf{z}_1^N \setminus z_{it}] = \frac{\eta + m'(k, v)}{\sum_{v=1}^V (\eta + m'(k, v))}$$

となります.

以上より, 式(5.74)と式(5.72)から,  $z_{it}$  のサンプリングは式(5.63)に代えて, 次の式で行うことができます.\*32

$$(5.75) \quad p(z_{it} = k | w_{it} = v, \mathbf{z}_i \setminus z_{it}) \propto \frac{\alpha_k + n'(i, k)}{\sum_{k=1}^K (\alpha_k + n'(i, k))} \cdot \frac{\eta + m'(k, v)}{\sum_{v=1}^V (\eta + m'(k, v))}$$

#### (LDA の周辺化 Gibbs サンプリングの公式)

この場合,  $\theta$  や  $\beta$  のサンプリングは不要で,  $\mathbf{z}$  だけをサンプルすれば学習できることに注意してください.  $\mathbf{z}$  のサンプリングが取束すれば,  $\theta$  と  $\beta$  は式(5.66)および式(5.70)の期待値をとることで,

$$(5.76) \quad \mathbb{E}[\theta_{nk} | \mathbf{z}_i] = \frac{\alpha_k + n(i, k)}{\sum_{k=1}^K (\alpha_k + n(i, k))}, \quad \mathbb{E}[\beta_{kv} | \mathbf{z}_1^N] = \frac{\eta + m(k, v)}{\sum_{v=1}^V (\eta + m(k, v))}$$

と求めることができます. この学習アルゴリズムを図 5.24 に示しました.

**周辺化尤度の計算** なお, この場合は周辺化 Gibbs サンプリングの目的関数は, 式(5.60)で  $\theta_1^N$  と  $\beta$  を周辺化した

$$(5.77) \quad p(\mathbf{w}_1^N, \mathbf{z}_1^N | \alpha, \eta) = p(\mathbf{w}_1^N | \mathbf{z}_1^N, \eta) p(\mathbf{z}_1^N | \alpha)$$

\*32 実装の際には, 式(5.75)の第1項の分母の  $\sum_{k=1}^K (\alpha_k + n'(i, k))$  は  $k$  に依存しませんので, 比例式としては不要です. また, 第2項の分母の  $\sum_{v=1}^V (\eta + m'(k, v))$  も実際に  $V$  個の和を毎回計算する必要はなく,  $m'(k) = \sum_{v=1}^V m'(k, v)$  を保持しておいて更新すれば,  $V\eta + m'(k)$  で求められることに注意してください.



になります[157]. この各項は, 図 5.25 に示したように  $n(i, k)$  および  $m(k, v)$  を行列で表せば, 各行にそれぞれ潜在的な多項分布  $\theta_i, \beta_k$  が存在してディリクレ分布  $\text{Dir}(\boldsymbol{\alpha}), \text{Dir}(\boldsymbol{\eta})$  に従い, これから各行の頻度  $n(i, :)$  および  $m(k, :)$  が生まれたことを意味します. したがって, この尤度は, 3章の式(3.58)で学習したポリア分布になり, 式(5.77)の右辺の各項は  $m(k) = \sum_v m(k, v)$  とおくと,

$$(5.78) \quad \begin{cases} p(\mathbf{z}_1^N | \boldsymbol{\alpha}) &= \prod_{i=1}^N \left[ \frac{\Gamma(\sum_k \alpha_k)}{\Gamma(T_i + \sum_k \alpha_k)} \prod_{k=1}^K \frac{\Gamma(\alpha_k + n(i, k))}{\Gamma(\alpha_k)} \right] \\ p(\mathbf{w}_1^N | \mathbf{z}_1^N, \boldsymbol{\eta}) &= \prod_{k=1}^K \left[ \frac{\Gamma(V\eta)}{\Gamma(V\eta + m(k))} \prod_{v=1}^V \frac{\Gamma(\eta + m(k, v))}{\Gamma(\eta)} \right] \end{cases}$$

で計算できることに注意してください. これから式(2.68)を使って, 学習データの1単語あたりのパープレキシティを求めることができます.

### LDAの実験

LDAは単語ごとに潜在トピックが存在する確率モデルのため, Pythonのみで実装すると, 非常に計算が遅くなります. モジュールをC言語にコンパイルして高速化する Cython を使った実装を筆者が公開していますので, 公開ページ<sup>\*33</sup>, または本章のサポートページから `lda.py-0.2.tar.gz` をダウンロードした後,

```
% tar xvfz lda.py-0.2.tar.gz
% cd lda.py-0.2/
% make
```

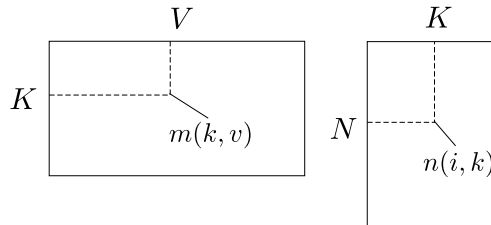


図 5.25: LDA の周辺化 Gibbs サンプリングで保持する統計量  $m(k, v), n(i, k)$  の行列. 行列の各行の頻度を生成した多項分布がそれぞれディリクレ分布に従っており, 全体の確率は, ポリア分布の積になります.

\*33 <http://chasen.org/~daiti-m/dist/lda-python/>

表 5.9: LDA ( $K=50$ ) で教師なし学習された Livedoor コーパスのトピックの例. 表 5.1 の教師ラベル (カテゴリ) では表されない, 細かい意味的トピックが学習されていることがわかります.

| Topic 1  |        | Topic 7 |        | Topic 43 |        | Topic 48 |        |
|----------|--------|---------|--------|----------|--------|----------|--------|
| 仕事       | 0.0357 | 応募      | 0.0435 | 料理       | 0.0276 | 結婚       | 0.0464 |
| 転職       | 0.0297 | 名       | 0.0431 | 食べ       | 0.0210 | 女        | 0.0389 |
| 会社       | 0.0246 | プレゼント   | 0.0370 | 店        | 0.0120 | 女性       | 0.0367 |
| 情報       | 0.0225 | 様       | 0.0338 | 味        | 0.0118 | 彼        | 0.0334 |
| 年収       | 0.0179 | 当選      | 0.0291 | 野菜       | 0.0105 | 男性       | 0.0303 |
| 人        | 0.0161 | 月       | 0.0226 | 食べる      | 0.0086 | 男        | 0.0258 |
| 自分       | 0.0160 | キャンペーン  | 0.0224 | 食        | 0.0080 | 恋愛       | 0.0244 |
| 万        | 0.0160 | 日       | 0.0193 | 酒        | 0.0074 | 相手       | 0.0216 |
| livedoor | 0.0158 | 終了      | 0.0177 | ビール      | 0.0071 | 彼女       | 0.0171 |
| 求人       | 0.0131 | 場合      | 0.0172 | 食事       | 0.0070 | 自分       | 0.0140 |

でモジュールをコンパイルし, 準備することができます. これには Cython が必要ですので, 事前に `% pip install cython` などを実行してインストールしておいてください.

この上で, Livedoor コーパスについて

```
% lda.py -K 50 -N 400 livedoor.dat model livedoor.lex
⇒ LDA: K = 50, iters = 400, alpha = 1, beta = 0.01
loading data.. documents = 7367, lexicon = 16946, nwords = 1838414
initializing..
Gibbs iteration [400/400] PPL = 4761.2851
saving model to model ..
including dictionary.
done.
```

を実行すると, トピック数  $K=50$  の LDA を学習することができます. この結果の一部を, 表 5.9 に示しました. 執筆時の環境では, この計算には約 6 分かかります. 上では MCMC の繰り返し数を 400 にしていますが, `model.log` に保存される学習データのパープレキシティの様子を見て, 収束するまで適宜増やすといいでしょう. 表 5.9 からわかるように, LDA により, 文書全体を 1 トピックとする UM や教師ラベルでは表現できない, 仕事, プレゼント, 食事, 恋愛などの細かいトピックが自動的に学習されていることがわかります.

表 5.10(a) に, 同様に `jawiki.txt` について表 5.7 の UM と同じ  $K=100$

表 5.10: LDA ( $K=100$ ) で学習した日本語 Wikipedia コーパスのトピック.(a) 生成確率  $p(w|z)$  を使った場合

| Topic 10 |        | Topic 20 |        | Topic 30 |        | Topic 100 |        |
|----------|--------|----------|--------|----------|--------|-----------|--------|
| 県        | 0.2955 | 世紀       | 0.0832 | 的        | 0.2476 | 月         | 0.3766 |
| 市        | 0.2018 | 時代       | 0.0777 | 化        | 0.0267 | 日         | 0.3252 |
| 福岡       | 0.0246 | 歴史       | 0.0341 | 基本       | 0.0255 | 年         | 0.2580 |
| 埼玉       | 0.0239 | 初期       | 0.0291 | 人間       | 0.0243 | テン        | 0.0030 |
| 広島       | 0.0229 | 頃        | 0.0286 | 方法       | 0.0243 | ルー        | 0.0017 |
| 兵庫       | 0.0229 | 年代       | 0.0224 | 主        | 0.0233 | 支え        | 0.0013 |
| 愛知       | 0.0215 | 説        | 0.0217 | 一般       | 0.0223 | 姓         | 0.0011 |
| 千葉       | 0.0194 | 当時       | 0.0215 | 表現       | 0.0206 | セントラル     | 0.0011 |
| 九州       | 0.0188 | 記        | 0.0172 | 用い       | 0.0200 | 政治        | 0.0008 |
| 名古屋      | 0.0172 | 古代       | 0.0167 | 手法       | 0.0176 | 会議        | 0.0008 |

(b) 正規化自己相互情報量  $\text{NPMI}(w, z)$  を使った場合

| Topic 10 |        | Topic 20 |        | Topic 30 |        | Topic 100 |        |
|----------|--------|----------|--------|----------|--------|-----------|--------|
| 県        | 0.7763 | 世紀       | 0.6493 | 的        | 0.7374 | 月         | 0.7545 |
| 市        | 0.6981 | 時代       | 0.5917 | 基本       | 0.5516 | 日         | 0.7281 |
| 福岡       | 0.5526 | 頃        | 0.5474 | 方法       | 0.5420 | 年         | 0.5677 |
| 埼玉       | 0.5511 | 初期       | 0.5386 | 人間       | 0.5392 | テン        | 0.4269 |
| 広島       | 0.5488 | 年代       | 0.5383 | 手法       | 0.5311 | ルー        | 0.3910 |
| 兵庫       | 0.5488 | 歴史       | 0.5379 | 主        | 0.5289 | 支え        | 0.3356 |
| 愛知       | 0.5442 | 説        | 0.5359 | 表現       | 0.5238 | モー        | 0.3236 |
| 千葉       | 0.5350 | 記        | 0.5305 | 分析       | 0.5201 | セントラル     | 0.3170 |
| 名古屋      | 0.5306 | 中世       | 0.5260 | 図        | 0.4977 | ウイング      | 0.2552 |
| 新潟       | 0.5306 | 不明       | 0.5195 | 成分       | 0.4944 | カルロス      | 0.2548 |

トピックの LDA を学習した場合のトピックの一部を示しました. LDA では, 文書全体を 1 トピックで表現する必要はないため, トピック 30 のように「的」「化」といった機能語に対応するトピックができることが多く<sup>\*34</sup>, 生成確率  $p(w|z)$  を使っても, 他のトピックで機能語が上位に来ることはあまりありません. 表 5.10(b) に示したように, この場合も NPMI を用いることで, 各トピックに相関する単語をより明確に取り出すことができます. 図 5.26 に, 各文書  $d$  に対する潜在トピック分布  $\theta = \{p(z|d)\}$  を示しました. 同じトピック数の図 5.9 の UM

\*34 ??ページで示すように, ハイパーパラメータ  $\alpha_k$  も推定して学習すると, こうした結果になりやすいことが確かめられています[160].

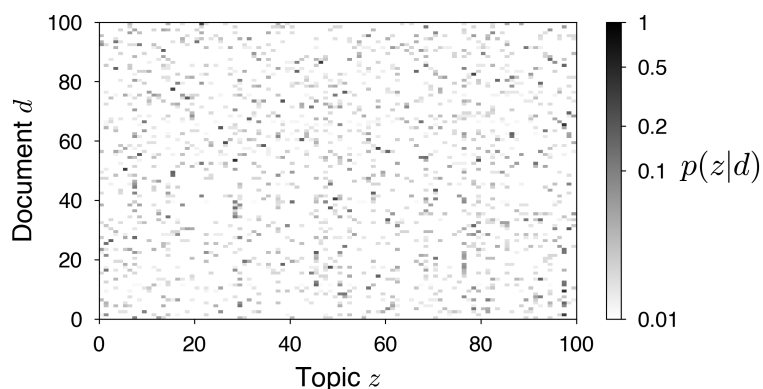


図 5.26: 日本語 Wikipedia コーパス `jawiki.txt` から学習された LDA による, 各文書のトピックへの所属確率  $p(z|d)$  (最初の 100 文書). トピック数は  $K=100$  としました. 図 5.9 の UM の場合と比較すると, 文書に対応する各行が, 複数の潜在トピックを確率的に持っていることがわかります.

の場合と比べると,  $p(z|d)$  がどれか 1 つのトピックに集中することがなく, 各文書が複数の潜在トピックからなる様子が学習されているのがわかります.

#### 5.4.1 LDA の幾何的解釈

LDA は, 幾何的には何をしていることに相当するのでしょうか. LDA では, トピック分布  $\theta$  が決まると, 文書の各単語は

1. トピック  $z \sim \theta$  をサンプルし,
2. 単語  $w \sim p(w|z)$  をサンプルする

という 2 段階で生成されます. 確率で書くと,  $\theta$  から  $w$  と  $z$  を生成する確率は

$$(5.79) \quad p(w, z|\theta) = p(w|z)p(z|\theta)$$

となっているわけです.

しかし, よく考えるとわれわれが観測しているのは  $w$  だけなので,  $z$  については和をとって周辺化してしまってもよいでしょう.\*35 つまり,  $w$  は確率分布

\*35 これを 216 ページの脚注のように, Rao-Blackwell 化ともいうのでした.

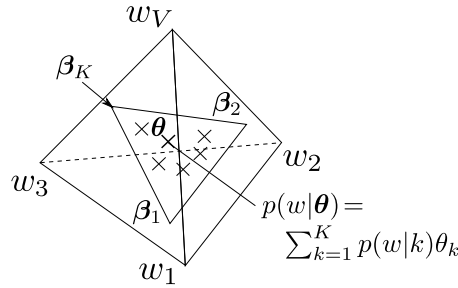


図 5.27: LDA の幾何的表現. LDA は×で示した単語の生成確率をモデル化できるよう、外側の  $V$  次元の単語単体の中で内側の  $K$  次元のトピック単体  $\beta_1, \dots, \beta_K$  の位置を最適化し、次元圧縮を行う統計モデルです.

$$(5.80) \quad p(w|\theta) = \sum_z p(w, z|\theta) = \sum_z p(w|z)p(z|\theta) = \sum_{k=1}^K p(w|k)\theta_k$$

に従う 1 段階で生成されたと考えても、数学的には等価になります.

トピック分布  $p(w|k)$  は、 $k$  ごとの単語の確率分布  $\beta_k = (p(w_1|k), \dots, p(w_V|k))$  のことですから、これは図 5.14 でみたように、単語単体の中の 1 点になります. 図 5.27 に示したように、式 (5.80) は、テキストが  $\beta_1, \dots, \beta_K$  を  $\theta_1, \dots, \theta_K$  の割合で内分した点  $p(w|\theta) = \sum_{k=1}^K p(w|k)\theta_k$  から生成されたことを表しています. このように、 $V$  次元の単語単体の中にあつて  $\beta_1, \dots, \beta_K$  によって張られる  $K$  次元の単体のことを、**トピック単体**とといいます.

$p(w|\theta)$  はこのトピック単体の中にあるのですが、一般に  $K \ll V$  ですから、トピック単体  $\beta_1, \dots, \beta_K$  をうまく選ばないと、文書を生成した確率分布  $\mathbf{p}$  を表現することができません. このように、

- (1) トピック、すなわちトピック単体の「角」 $\beta_1, \dots, \beta_K$  と、
- (2) 文書ごとの  $\theta$ 、すなわち上で定まるトピック単体内の各文書の位置

を同時に最適化しているのが LDA です. このとき  $\theta$  にはディリクレ事前分布を仮定して、各文書にディリクレ事後分布を割り当てる (allocation) ため、**潜在ディリクレ配分法**とよばれているわけです.

空間内に、データをうまく表現する低次元の部分空間を張るという意味では、これは実数値 (ユークリッド空間) における主成分分析 (PCA) に似ています. た

だし PCA と異なり, LDA は離散データのために単体上で定義された統計モデルなのが特徴で, 低頻度のデータでもうまく扱うことができます. Buntine ら [161] は, LDA およびその拡張を統一して離散 PCA (Discrete PCA) とよんでいます.

メモ：ディリクレ分布の  $\alpha$  のサンプリング\*

LDA のハイパーパラメータはトピック分布およびトピック-単語分布を生成するディリクレ分布のハイパーパラメータ  $\alpha, \eta$  ですが、これらはどうやって決めたらいいのでしょうか。

一般的にはトピック数を  $K$  として、 $\alpha = (50/K, \dots, 50/K)$ ,  $\eta = 0.01$  のようにヒューリスティックに決められることが多いのですが、これらのハイパーパラメータの影響を詳しく調べた研究[162]によると、これらもデータから学習した方がよいモデルとなることが示されています。とくに、非対称な  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_K)$  を推定することで、 $\alpha_k$  が大きいトピックに「が」「の」、英語なら “of”, “the” といった機能語<sup>\*36</sup> が集まり、他のトピックにこれらが混ざることが少なくなります。また、トピック数を増やしても不要なトピックは  $\alpha_k$  が小さく、自然に使われなくなることも報告されています。

$\alpha$  と  $\eta$  についての尤度関数は式(5.78)ですから、数学的にはこれは、このポリア分布のハイパーパラメータを推定することと等価です。ただし、このカウント  $n(i, k)$ ,  $m(k, v)$  はあくまで学習中の  $\mathbf{z}_1^N$  から計算される値ですから、学習中に式(3.60)を使って  $\alpha$  を最適化してしまうと、局所解に陥る危険があります。<sup>\*37</sup> よって、きちんと  $\alpha$  のベイズ推定を行うのが正しい方法でしょう。

ポリア分布を表す式

$$(5.81) \quad p(\mathbf{n}|\alpha) = \frac{\Gamma(\sum_k \alpha_k)}{\Gamma(N + \sum_k \alpha_k)} \prod_{k=1}^K \frac{\Gamma(\alpha_k + n_k)}{\Gamma(\alpha_k)}$$

は、 $\alpha_k$  がガンマ関数  $\Gamma()$  の中に含まれているため、そのままではサンプリングできる形になりません。しかし、巧妙な補助変数  $\theta, x$  を導入すると、 $\alpha_k$  はガンマ事後分布からサンプリングすることができます。導出はやや複雑なため、詳しくは、付録 E を参照してください。

\*36 もしこうした機能語を「ストップワード」として除いても、242 ページで説明したように、同様に意味の薄い語が必ず出現し、機能語とそれ以外は簡単に二分することができません。

\*37 すなわち、学習途中で  $\alpha$  を最適化してしまうと、学習アルゴリズム全体が Gibbs サンプリングではなく、252 ページのモンテカルロ EM アルゴリズムを行うことになってしまいます。

### 5.4.2 トピックモデルの評価

こうして学習したトピックモデルが、文書の「よいモデル」になっているかどうかは、どうやって測ればよいのでしょうか。LDA は先に説明したように、単語単体上の密度推定を行う統計モデルですから、テストデータに高い確率を与えられるかが、最も素直な評価方法といえるでしょう。LDA の評価方法について詳しく調べた研究[160]によると、これにはテストデータの各文書について、前半からトピック分布を学習し、それに基づいて後半の単語の予測パープレキシティを計算する方法が最も良かったことが報告されています。

サポートサイトの本章のフォルダに、この評価スクリプト `ldaeval.py` を示しました。LDA は順番を考慮しない単語集合のモデルですから、評価の際には各文書で単語をランダムにシャッフルし、前半  $\mathbf{w} = w_1 \dots w_{T/2-1}$  から計算したトピック分布  $\boldsymbol{\theta}$  を用いて、後半  $\mathbf{w}' = w_{T/2} \dots w_T$  の確率を

(5.82)

$$p(\mathbf{w}' | \boldsymbol{\theta}) = \prod_{t=T/2}^T p(w_t | \boldsymbol{\theta}) = \prod_{t=T/2}^T \sum_{z_t} p(w_t, z_t | \boldsymbol{\theta}) = \prod_{t=T/2}^T \sum_{k=1}^K p(w_t | z_t = k) \underbrace{p(z_t = k | \boldsymbol{\theta})}_{= \theta_k}$$

のように計算します。なお、 $\mathbf{w}$  からそれを生成した  $\boldsymbol{\theta}$  を求めるには式(5.75)を用いることもできますが、パラメータはすでに学習済みですので、変分ベイズ EM アルゴリズムを用いると高速に行えます。詳しくは、`ldaeval.py` のスクリプトの中身および[71]を参照してください。

図 5.28 に、Livedoor コーパスのうちランダムな 367 文書をテストデータ、残りの 7,000 文書を学習データとした場合の予測パープレキシティを、各トピック数  $K$  について示しました。LDA はベイズ的な手法のため、 $K$  を大きくしても過学習して性能が落ちることはなく、パープレキシティはゆっくり減少していくことがわかります。また、付録 E の方法でハイパーパラメータ  $\boldsymbol{\alpha}, \boldsymbol{\eta}$  を推定した場合、固定の場合と比べて予測精度が上がっていることも読みとれます。



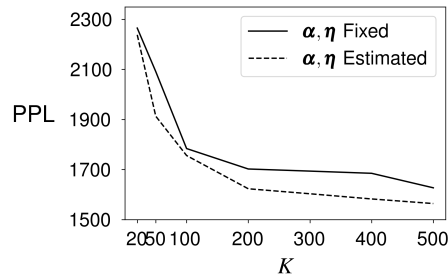


図 5.28: Livedoor コーパスにおけるテストデータの予測パープレキシティとトピック数  $K$  の関係. ハイパーパラメータ  $\alpha, \eta$  を固定した場合 (Fixed) に対して, それらもベイズ推定した場合 (Estimated) の方が, 一貫して低いパープレキシティを達成することがわかります.

### トピックの直接評価

上ではトピックモデルの性能を予測精度で評価しましたが, これは得られたトピックが意味的まとまりの高い, 「よいトピック」であることを保証しているわけではありません. 「よいトピック」であることは, どうやって測ればいいのでしょうか? トピックモデルは教師なし学習ですので, この問題にはもちろん「正解」はないのですが, いくつかの方法が提案されています[163][164]. 以下では, 各トピック  $k$  について, 生起確率  $p(w|k)$  の高い順に  $N$  語 (たとえば  $N=10$ ) とった単語集合を「トピック語」 $t(k)$  とよぶことにします.

**Word Intrusion** 各トピックが意味的によくまとまっていれば, トピック語に無関係な語 (侵略者, intruder) を混ぜてもたやすく検出できるはずですが, 逆にノイズが多いトピックなら, どれが侵略者なのかを判断できないでしょう. たとえば図 5.29 で, A のトピック語に混ぜた “京都” は簡単に発見できますが, B に混ぜた “大阪” はすぐには見分けられません. この判断は人間の被験者に行っても行うこともできますし, サポートベクトル回帰 (SVR) で自動化する方法も提案されています[164].<sup>\*38</sup>

**Topic Coherence** トピックが意味的にまとまっているかを, 直接測ることもできます. 5.2.1 節で議論したように, NPMI を使うとトピックと相関の高い単

<sup>\*38</sup> [https://github.com/jhlau/topic\\_interpretability](https://github.com/jhlau/topic_interpretability) で, 自動評価のための評価スクリプトが公開されています.

|        |  |                                       |
|--------|--|---------------------------------------|
| トピック A |  | 酸, 酵素, <b>京都</b> , 化合, 合成, 水素, 有機, 触媒 |
| トピック B |  | 湖, 不足, キル, PK, 通勤, <b>大阪</b> , テレビジョン |

図 5.29: Word intrusion の例. 意味的まとまりの高いトピック A では, 太字で示した侵略者 (intruder) はすぐに見分けることができますが, 意味的まとまりの低いトピック B では, どれが侵略者なのか判断が付きません.

語を取り出すことができますから, よいトピックであれば, それらの単語どうしの NPMI も高くなるはずです. トピック語のすべてのペア  $(v, w)$  について, 外部の大きなテキストでの共起から計算した  $\text{NPMI}(v, w)$  の平均値を Coherence と定義すると, これは上の Word Intrusion および人間の判断と高い相関を持つことが示されています[164].

**Topic Uniqueness** トピックが意味的にまとまっても, ほとんど同じトピックが複数あるのは望ましくないでしょう. 各トピックのトピック語を「文書」とみなしたとき, ある単語  $w$  が複数のトピック語に現れていれば, 293 ページで学ぶ文書頻度 (=何個のトピック語に現れたか)  $\text{df}(w)$  は 1 より大きくなります. よって, その逆数をとって平均し,

$$(5.83) \quad \text{TU} = \frac{1}{K} \sum_{k=1}^K \frac{1}{N} \sum_{w \in t(k)} \frac{1}{\text{df}(w)}$$

を計算すれば, トピック語がすべて異なっていれば 1, 重複があると  $< 1$  になるはずです. これを Uniqueness と定義して評価する方法が提案されています[165].

ただし, 式(5.83)の指標はトピック数  $K$  の違いについてロバストではありません.  $K$  が大きくなれば, ある単語が他のトピック語にも偶然含まれる可能性は, それだけ大きくなるからです. 式(5.83)は本質的に, すべてのトピック語をまとめた集合での各単語  $w$  の出現回数  $\text{df}(w)$  を計算しているだけですから, それをユニグラム分布  $p(w) = \text{df}(w)/(NK)$  に直せば, 2章で学習した  $p(\cdot)$  のエントロピー (式(2.65))  $H(p(\cdot))$  を指標とすればよいでしょう. この最大値は, すべての単語が 1 回ずつ現れる (=トピック語に重複がない) 場合で,  $\log(NK)$  になります. よって最大値との比をとり,

$$(5.84) \quad \text{TU}++ = H(p(\cdot))/\log(NK)$$

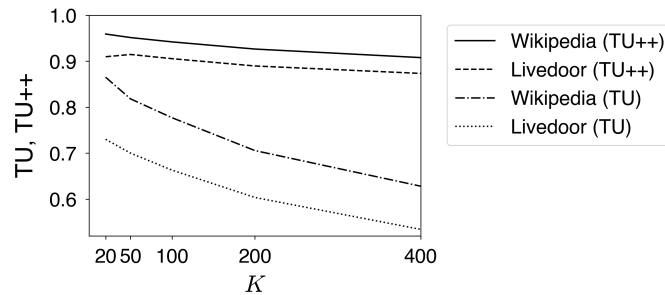


図 5.30: Livedoor コーパスおよび日本語 Wikipedia で計算した TU, TU++ の値とトピック数  $K$ .  $K$  の選択には使えませんが, 同じトピック数での比較のためには有効です.

を指標とすれば, これは  $K$  に依存しません. この指標 (本書のオリジナルです) を TU++ とよぶことにしましょう.

こうした評価は, LDA に限らずトピックモデル全般に用いることができます. 5.2 節の UM について, EM アルゴリズムで推定した表 5.6 のトピックと, Gibbs サンプリングで計算した表 5.7 のトピックを比べてみましょう. サポートサイトの `coherence.py` を使って計算すると,

```
% coherence.py model/um.jawiki.K100 data/ja.text8.txt
⇒ obtaining topic words..
 computing cooccurrences..
 calculating coherence..
 average = -0.0997
```

のようになります. ここでは, 日本語 text8 を共起計算のためのコーパスとして用いました.\*39 EM アルゴリズムの場合は上のように  $-0.0997$ , Gibbs サンプリングの場合は  $-0.0684$  となり, 確かに Gibbs サンプリングの方が意味的まとまりの高いトピックが学習されていることが定量化されました.

また Uniqueness については, Livedoor コーパスについて `uniqueness++.py` で TU および TU++ をさまざまな  $K$  について計算すると図 5.30 のようになり,  $K = 50$  程度でやや TU++ の値が高くなることがわかります. ただし, Wikipedia の場合は指標は  $K$  とともに一様に減少するようです.

\*39 この 10 倍のサイズがある 139 ページの `ja.text9` を使っても, 統計値はほとんど同じになりました.

**トピック数  $K$  の選択** トピックモデルにおいて、トピック数  $K$  の選択は重要な問題です。上でみたように、予測パープレキシティ、Uniqueness (TU++) のいずれも、それだけでトピック数を適切に決めることはできません。<sup>\*40</sup> 理論的には、LDA のトピック数は階層ディリクレ過程 (HDP) を用いた HDP-LDA で推定でき、この問題は解決されています[72]。ただし、HDP の理解と実装には測度論の知識が必要となり、本書の範囲を超えますので、簡便には (1) HDP-LDA が実装されている `gensim` などのパッケージ<sup>\*41</sup> を用いるか、または (2) 4 章で行ったように  $K$  を大きめにとって  $\alpha < 1$  と小さく設定するか、あるいは (3)  $\alpha$  をサンプリングして学習することで、実質的に必要なトピックだけを残すようにするとよいでしょう。HDP による無限モデルに興味のある方は、優れた導入である[166]を参照してください。

### LDA の拡張

- Supervised
- Embedding
- Tea party

---

<sup>\*40</sup> LDA で学習データの確率が最も高くなるのは、トピックの数と単語数  $V$  が等しく、各トピックで  $\beta_{kv}$  がある単語だけ確率 1、他は 0 となる場合です。しかし、式(5.76)の  $\beta$  の事後分布から、 $\eta > 0$  のときにそうなる確率は 0 です。したがって  $\eta$  を固定すれば、式(5.78)のエビデンス  $p(\mathbf{w}_1^N | \mathbf{z}_1^N, \eta)$  を最大にするトピック数  $K$  を求めることができます[157]。しかしこれは  $\eta$  に依存し、 $\eta = 0.1$  と  $\eta = 0.01$  ではまったく違うトピック数になってしまいます。式(0.139)から  $\eta$  をサンプリングする場合、 $K$  が大きいほど推定される  $\eta$  は小さくなり、この方法で最適な  $K$  を選ぶことはできません。

<sup>\*41</sup> <https://radimrehurek.com/gensim/models/hdpmodel.html>

## 5.5 ニューラル文書モデルと独立成分分析

LDA は、各文書にディリクレ分布に従う潜在的なトピック分布  $\theta$  があると考え、 $\theta$  とトピックの両方をデータから学習できる、解釈性に優れたモデルですが、LDA にもいくつかの問題があることには注意が必要です。

- 一つは、文書の内容を表現する  $\theta$  が多項分布に制限されていることです ( $\sum_k \theta_k = 1, \theta_k \geq 0$ )。こうすると、たとえば経済と数学の両方の内容を持つ文書は、経済トピックが強いほど数学トピックは弱くなり、その逆も成り立ちます。また、あまり内容がなかったり短い文書でも、 $\theta$  は内容の濃いテキストと同様に、和が1の確率分布になってしまいます。
- 二つ目は、単語を離散的にとらえているため、3章で学習した単語ベクトルのように単語間の関係を精密にモデル化できないことです。たとえば、“取る”と“取得”のようにほぼ同じ意味の言葉でも、LDA では「同じ潜在トピックから生成される確率が高い」という間接的な形でしか関係をとらえることができません。
- また、LDA は文書に含まれる各単語ごとにトピックを推定するため、精密なモデル化ができる一方で、大規模なコーパスでは計算量が大きくなってしまうという問題があります。

ここで、LDA で文書の  $\theta$  がわかっている場合、その文書での単語  $w$  の確率  $p(w|\theta)$  は、トピック  $z$  を考えて周辺化することで、式(5.80)でもみたように

$$(5.85) \quad p(w|\theta) = \sum_z p(w, z|\theta) = \sum_{k=1}^K p(w|z=k)p(z=k|\theta) = \sum_{k=1}^K p(w|k)\theta_k$$

と書けることに注意しましょう。以下ではわかりやすさのため、文書-単語行列の縦横を入れ替えて縦軸を単語、横軸を文書として説明します。<sup>\*42</sup> 上の確率を  $V$  個の単語について縦に並べれば、

<sup>\*42</sup> Python の Numpy では、行列の要素は標準では内部的に行方向が先に格納されるため (row-major といいます)、これまで説明した文書-単語行列の形が実装上は有利です。

$$(5.86) \quad \underbrace{\begin{pmatrix} p(w_1|\boldsymbol{\theta}) \\ p(w_2|\boldsymbol{\theta}) \\ p(w_3|\boldsymbol{\theta}) \\ \vdots \\ p(w_V|\boldsymbol{\theta}) \end{pmatrix}}_{\mathbf{p}} = \underbrace{\begin{pmatrix} p(w_1|1) & p(w_1|2) & \cdots & p(w_1|K) \\ p(w_2|1) & p(w_2|2) & \cdots & p(w_2|K) \\ p(w_3|1) & p(w_3|2) & \cdots & p(w_3|K) \\ \vdots & & & \vdots \\ p(w_V|1) & p(w_V|2) & \cdots & p(w_V|K) \end{pmatrix}}_{\mathbf{B}} \underbrace{\begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_K \end{pmatrix}}_{\boldsymbol{\theta}}$$

と書くことができ、 $\beta_k$  を並べた行列  $\mathbf{B}=(\beta_1, \beta_2, \dots, \beta_K)$  を使って  $\mathbf{p}=\mathbf{B}\boldsymbol{\theta}$  と表すことができます。式(5.86)を  $N$  個の文書について横に並べると、 $n$  番目の文書の  $\boldsymbol{\theta}$  を  $\boldsymbol{\theta}_n$  として

$$(5.87) \quad \underbrace{\begin{pmatrix} p(w_1|\boldsymbol{\theta}_1) & \cdots & p(w_1|\boldsymbol{\theta}_N) \\ p(w_2|\boldsymbol{\theta}_1) & \cdots & p(w_2|\boldsymbol{\theta}_N) \\ p(w_3|\boldsymbol{\theta}_1) & \cdots & p(w_3|\boldsymbol{\theta}_N) \\ \vdots & & \vdots \\ p(w_V|\boldsymbol{\theta}_1) & \cdots & p(w_V|\boldsymbol{\theta}_N) \end{pmatrix}}_{\mathbf{P}} = \underbrace{\begin{pmatrix} p(w_1|1) & \cdots & p(w_1|K) \\ p(w_2|1) & \cdots & p(w_2|K) \\ p(w_3|1) & \cdots & p(w_3|K) \\ \vdots & & \vdots \\ p(w_V|1) & \cdots & p(w_V|K) \end{pmatrix}}_{\mathbf{B}} \underbrace{\begin{pmatrix} \theta_{11} & \theta_{21} & \cdots & \theta_{N1} \\ \theta_{12} & \theta_{22} & \cdots & \theta_{N2} \\ \vdots & & & \vdots \\ \theta_{1K} & \theta_{2K} & \cdots & \theta_{NK} \end{pmatrix}}_{\boldsymbol{\Theta}}$$

となり、図 5.31(a) に示したように行列形式で

$$(5.88) \quad \mathbf{Y} \sim \mathbf{P}, \quad \mathbf{P}=\mathbf{B}\boldsymbol{\Theta}$$

と表すことができます。つまり LDA は、 $\mathbf{P}$  と同じ大きさの単語-文書行列  $\mathbf{Y}$  が多項分布に従って  $\mathbf{Y} \sim \mathbf{P}$  と生成されたと考え、この  $\mathbf{P}$  を式(5.88)のように分解する行列分解を教師なし学習している、と考えることができるわけです。ここで  $\sim$  とは、観測値  $\mathbf{Y}$  の各列が  $\mathbf{P}$  の同じ列をパラメータとして多項分布で生成されたことを表します。<sup>\*43</sup>

<sup>\*43</sup>  $\boldsymbol{\Theta}$  と  $\mathbf{B}$  はともに非負ですから、これは非負値行列因子分解 (NMF) の一種とみなすことができます。Lee らが 2000 年に提案した NMF [167] は、統計的には頻度がポアソン分布に従うことに対応しています。ポアソン分布はスケールが文書の長さに依存してしまい、正規化することで多項分布となるため、式(5.88)はその一般化といえます。NMF とテキストの Gamma-Poisson モデル (GaP) [168] については、筆者のメモ[169]も参照してください。

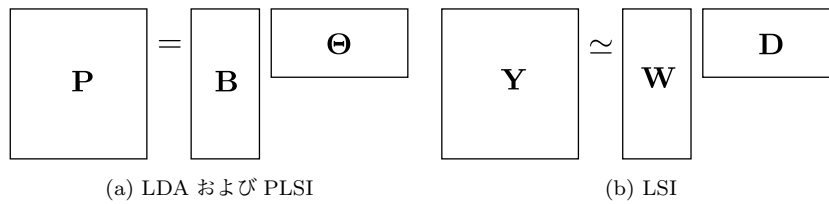


図 5.31: 行列分解による文書モデル. 左辺の行列を, 右辺の二つの行列の積で近似しています. このとき, 右辺の行列の各行と各列が「単語ベクトル」および「文書ベクトル」になります.

式(5.87)で, 確率分布の形で各文書の  $\theta_n$  は文書  $n$  の「埋め込み」,  $\mathbf{B}$  の各行  $\phi(w) = (p(w|1), p(w|2), \dots, p(w|K))$  も単語  $w$  の一種の「埋め込み」と考えることができます. しかし,  $\theta$  も  $\phi(w)$  も負になることができず, 縦横に和が1になる必要もあるため, この行列分解には大きな制約が伴います.

それならば, いっそのこと  $\mathbf{Y}$  を図 5.31(b) のように直接,

$$(5.89) \quad \mathbf{Y} \simeq \mathbf{W}\mathbf{D}$$

と, 一般の実行列  $\mathbf{D}, \mathbf{W}$  に行列分解 (150 ページ) すれば,  $\mathbf{W}$  と  $\mathbf{D}$  の各行は制約のない単語ベクトル, 文書ベクトルになるのではないのでしょうか. この考えに基づくのが, 情報検索の分野で 1990 年に提案された **LSI** (Latent Semantic Indexing) [170] でした. LSI では, 頻度  $\mathbf{Y}$  に含まれる頻度  $Y(w, d)$  をそのまま使うとほとんどを機能語が占めてしまうため,  $\tilde{Y}(w, d) = \text{tfidf}(Y(w, d))$  のように, この後で説明する **tf.idf** のような重みづけを行ってから, 式(5.89)のように行列分解を計算して文書ベクトル・単語ベクトルを求めます.

ただし, LSI では **tf.idf** のような単語の重みづけ法がアドホックで, 得られるベクトルの数学的な意味が弱いという大きな欠点がありました. そこで, LSI を多項分布による確率モデルとして式(5.88)のようにとらえ直した, **PLSI** (Probabilistic Latent Semantic Indexing) という革新的なモデルが 1999 年に提案され [171], それが 2001 年にベイズ化されて LDA になった [154] という歴史があります. 実験的にも, パープレキシティで測った性能は  $\text{LDA} > \text{PLSI} > \text{LSI}$  の順に高いことがわかっており [141], テキストをアドホックでなく, 確率的に考えることの重要性を示しています.

### 5.5.1 文書ベクトルと Doc2Vec

ただし、実行列による行列分解自体に意味がないわけではありません。皆さんは図??のような行列分解は、本書でこれまでに覚えがあるのではないのでしょうか。これは、3章の図 3.32 でみた、行列分解によるニューラル単語ベクトルの学習

$$(5.90) \quad \mathbf{X} \simeq \mathbf{W}\mathbf{C}^T$$

と、ほとんど同じ形をしていることに注意しましょう。単語ベクトルの場合、もとの共起行列  $\mathbf{X}$  の要素  $X(w, c)$  は、単語  $w$  とその周辺に出現した文脈語  $c$  との非負自己相互情報量 (PPMI)

$$X(w, c) = \text{PPMI}(w, c) = \max\left(\log \frac{p(w, c)}{p(w)p(c)}, 0\right)$$

でした (式(3.110))。そこで、単語-文書行列の場合は単語  $w$  とそれが出現した文書  $d$  の PPMI を用いて、

$$(5.91) \quad X(w, d) = \text{PPMI}(w, d) = \max\left(\log \frac{p(w, d)}{p(w)p(d)}, 0\right)$$

を並べた行列  $\mathbf{Y}$  を作り、これを図 5.32 のように

$$(5.92) \quad \mathbf{X} \simeq \mathbf{W}\mathbf{D}$$

と行列分解すれば、行列  $\mathbf{D}$  の各列  $\vec{d}$  は Word2Vec と数学的に等価な「ニューラル文書ベクトル」に、 $\mathbf{W}$  の各行  $\vec{w}$  は「ニューラル単語ベクトル」になるのではないのでしょうか。<sup>\*44</sup>

実際に、Mikolov らが Word2Vec の後に提案し、よく使われている **Doc2Vec** [172] は文書ベクトルから文書に含まれる単語を予測するモデルで、単語ベクトルのスキップグラム (3.5.3 節) と同じ目的関数をしており、式(5.92)と等価です。しかし、確率的勾配法による繰り返し計算で近似を行う Doc2Vec に対して、

<sup>\*44</sup> すなわち、LSI に足りなかったのは適切な単語の重みづけと、それを支える背後の理論だったということです。149 ページの脚注で述べたように、実験的には PPMI による単語重みづけが高性能であることは、深層学習以前に一部ではすでに知られていました。



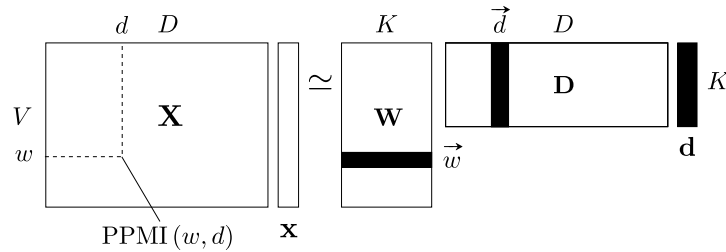


図 5.32: DocVec による文書ベクトルと単語ベクトルの計算. 右辺の  $\vec{d}$  と  $\vec{w}$  の内積で左辺の  $\text{PPMI}(d, w)$  を近似し, これは Doc2Vec (Word2Vec) の学習と数学的に等価です.

式 (5.92) から得られる文書ベクトル (本書ではこれを **DocVec** とよぶことにします) はこの後で説明するように線形代数で高速に解くことができ, 数学的な最適解が求まるために, 性能も高いことを筆者が確かめています [173].

なお, 式 (5.91) はベイズの定理から  $p(w|d) = p(w, d)/p(d)$  ですから,

$$(5.93) \quad \text{PPMI}(w, d) = \max \left( \log \frac{p(w, d)}{p(w)p(d)}, 0 \right) = \max \left( \log \frac{p(w|d)}{p(w)}, 0 \right)$$

で求めることができます. すなわち, 式 (5.93) の PPMI は, 「文書  $d$  の中で単語  $w$  が通常より何倍多く出現したのか」という値の対数となっています.  $d$  の中で  $w$  (たとえば “the”) の出現確率が平均的な確率とほぼ同じならば,  $p(w|d) \simeq p(w)$  ですから, 式 (5.93) の比は

$$\log \frac{p(w|d)}{p(w)} \simeq \log 1 = 0$$

となることに注意してください. 機能語はどの文書でも確率がほぼ一定なので, こうすれば「ストップワード」を準備しなくても, 機能語の PPMI はほぼ 0 になることがわかります.

DocVec の学習アルゴリズムを, 図 5.34 に示しました. この計算は, サポートサイトの `docvec.py` で行うことができます \*45. 図 5.32 からわかるように, このとき **D** の各行として文書ベクトル  $\vec{d}$  が, **W** の各行として単語ベクトル  $\vec{w}$  が

\*45 疎行列の特異値分解を行う SciPy の `svds()` を使う場合, 特異値が大きい順ではなく, 小さい順に返されますので注意してください. 標準では, 次元が重要な順の逆に並ぶことになります.

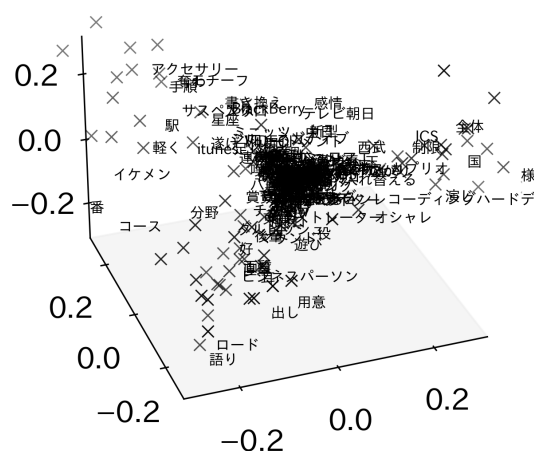


図 5.33: Livedoor コーパスから DocVec で計算した文書ベクトル (×印) と単語ベクトルの一部. 3.5.6 節で議論したように, 最大の特異値に対応する次元は全体のバイアスを表しているため, ここでは特異値の大きい方から 2,3,4 次元目の値でプロットしました. これから, Livedoor コーパスの文書ベクトル・単語ベクトルは大きく「女性軸」(左上), 「電気製品軸」(右), 「一般軸」(左下) に沿って分布していることがわかります. ただし, 座標軸が必ずしもそれらに沿ってとられているわけではありません.

得られることに注意してください. 図 5.33 に示したように, これらのベクトルは同じ  $K$  次元の潜在空間に存在しており,

$$(5.94) \quad \text{PPMI}(w, d) \simeq \vec{w} \cdot \vec{d}$$

となるように文書ベクトル  $\vec{d}$ , 単語ベクトル  $\vec{w}$  が学習されることとなります.

**文書ベクトルの計算** 実際に実験してみましょう. サポートサイトにある docvec.py を使うと, Livedoor コーパスのデータ livedoor.dat について次のように実行すれば,  $K=200$  次元の文書ベクトルと単語ベクトルが計算できます.

```
% docvec.py -K 200 -d livedoor.lex livedoor.dat model.docvec
⇒ using dictionary: livedoor.lex
livedoor.dat: K = 200, output = docvec.livedoor.K200
parsing data..
creating sparse matrix..
computing data vectors..
```

```
done.
writing model to model.docvec..
done.
```

この計算は、執筆時の環境では 14.6 秒で終了しました。これから、たとえば 200 番目の文書に似た文書ベクトルを検索したい場合は、コーパスの実際の中身が書かれているテキスト `livedoor.txt` (226 ページ) を同時に与えて、次のように実行します。

```
% docvec-similar.py docvec.model livedoor.txt 200
⇒ loading model from docvec.model.. done.
 1.0000 dokujo-tsushin 30 歳を過ぎた大人の肌へ、世界が認めた美容液「
0.6536 dokujo-tsushin 女性も驚愕!? スキンケアしていない人は 70% と男
0.6530 peachy あなたの見た目年齢を上げているのは “シミ” だった
0.6492 peachy 自宅でできるでふっくらお肌を目指そう! 毎朝のメイ
0.6315 peachy 日本初! DHC、“10 倍濃度” の Q10 シリーズを
0.6148 kaden-channel ライオン、毛と音波振動でくすみを落とす「プラチア
0.6104 peachy 男性の 4 割が “すっぴん” にながっかり! 最強のすっぴ
0.5819 peachy “究極のクリーム” で 10 年後も 20 年後もずっと綺
```

一番上の文書は自分自身 (類似度は  $\cos 0 = 1$ ) ですが、ジャンルを超えて内容的に近い文書が検索できていることがわかります。実装については、スクリプトの中身を読んでみてください。

**キーワードによる検索** それでは、あるキーワード集合に近い文書を探したい場合はどうすればいいのでしょうか。実は、これは自明ではありません。<sup>\*46</sup> というのは、いま探したい “映画 東京” のような検索語だけを含む文書は学習データには存在しないからです。

しかし、キーワード集合に対応する仮想的な「文書ベクトル」 $\mathbf{d}$  が計算できれば、学習した文書ベクトルと  $\mathbf{d}$  を上記のように比べればよいでしょう。ここで、図 5.32 では、左辺の観測データ  $\mathbf{Y}$  は常に式 (5.91) の PPMI の形で与えられることに注意してください。この値が 0 のとき、文書内での単語の確率は標準的な確率と等しいか小さいことを意味します。よって、検索したい語の PPMI をたとえば 1 にすれば、これは仮想的な「文書」の中で確率が  $e^1 \approx 2.72$  倍高いこと

<sup>\*46</sup> 簡易的には、キーワードに対応する単語ベクトルの平均を計算して文書ベクトルと比較するという方法が考えられますが、モデル上これが最適である保証はありません。実際、この後で計算して対応する列の和をとる行列  $(\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T$  は、 $\mathbf{W}$  自体とは異なったものです。

を意味します。そこで、キーワード集合に含まれる語について1, 他は0になる  $V$  次元の縦ベクトル  $\mathbf{x}$  を用意し、図 5.32 に表したように

$$(5.95) \quad \mathbf{x} \simeq \mathbf{W}\mathbf{d}$$

が成り立つ文書ベクトル  $\mathbf{d}$  を求めればよい、ということになります。

特異値分解は両辺の二乗誤差を最小化する方法ですので、最小二乗の意味で式(5.95)が成り立つ  $\mathbf{d}$  を求めるには、単純に最適化を行うこともできますが、式(5.95)はよく知られた線形回帰モデル (OLS) ですから、 $\mathbf{d}$  の最適解は

$$(5.96) \quad \mathbf{d}^* = (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T \mathbf{x}$$

で与えられます[8, §1.2]。あらかじめ回帰行列  $\mathbf{R} = (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T$  を計算しておけば<sup>\*47</sup>, 式(5.96)は

$$(5.97) \quad \mathbf{d}^* = \mathbf{R}\mathbf{x}$$

と、一瞬で最適解が求められます。<sup>\*48</sup> なお、新しい文書自体が与えられている場合は、式(5.93)から PPMI を計算して  $\mathbf{x}$  を作れば、同様に最適な文書ベクトル  $\mathbf{d}^*$  を求めることができます。

この方法で Livedoor コーパスの文書について、意味を考慮したキーワード検索を行った例は次のようになります。ニューラル単語ベクトルを介しているため、検索語自体が出現していなくても、意味が似ていれば (=単語ベクトルが近ければ) 文書が検索されることに注意してください。

```
% docvec.py -K 200 -R -d livedoor.lex livedoor.dat model.docvec-R
-R をつけて実行し、回帰行列を事前に計算しておく
% docvec-search.py model.docvec-R livedoor.txt 映画 東京
⇒ loading model from model.docvec-R.. done.
keyword: 映画 東京
0.0158 movie-enter この夏、東京スカイツリーが映画に 2008 年 7 月の
0.0152 movie-enter 小栗旬は“使えない若者”、映画『キツツキと雨』の
0.0152 movie-enter スパイダーマンが地上 75 メートルの通天閣を『アメ
```

<sup>\*47</sup>  $\mathbf{R}\mathbf{W} = (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T \mathbf{W} = \mathbf{I}$  ですから、この  $\mathbf{R}$  は  $\mathbf{W}^+$  とも書かれる Moore-Penrose の擬似逆行列です。

<sup>\*48</sup>  $\mathbf{y}$  の要素はほとんどが 0 ですから、式(5.97)の計算は実際には、 $\mathbf{R}$  のうち  $\mathbf{x}$  の要素が 1 の列を足すだけの操作になり、`mmap()` などで  $\mathbf{R}$  をディスクに保存すれば、空間計算量も最小限に抑えることができます。

---

**アルゴリズム 7:** 文書ベクトル (DocVec) の計算アルゴリズム.
 

---

- 1: 式 (5.91) に従って文書-単語行列  $\mathbf{X}$  を作成する. (疎行列フォーマットで)
  - 2:  $\mathbf{U}, \mathbf{S}, \mathbf{V} = \text{svds}(\mathbf{X}, K)$  と  $K$  次元に特異値分解する.
  - 3: 文書ベクトル行列  $\mathbf{D} = \mathbf{U}\mathbf{S}^{1/2}$ , 単語ベクトル行列  $\mathbf{W} = \mathbf{V}^T\mathbf{S}^{1/2}$  を計算する.
- 

図 5.34: 特異値分解による文書ベクトル (DocVec) の計算アルゴリズム.

|                    |                             |
|--------------------|-----------------------------|
| 0.0149 peachy      | 【スナップレポート】東京ガールズコレクション 2011 |
| 0.0149 movie-enter | 食べて、で、映画を観れる『東京ごほん映画祭』が今    |
| 0.0149 peachy      | 【スナップレポート】東京ガールズコレクション 2011 |
| 0.0148 movie-enter | 基礎から勉強しよう!初心者でもわかる「東京国際映    |
| 0.0144 movie-enter | 三池監督が映画『一命』について「満島ひかりがいろ    |
| 0.0141 movie-enter | 東京国際映画祭、『最強のふたり』がグランプリを受    |

### ノート：tf.idfと単語の重みづけ

自然言語処理の多くの場面で、“日”“方法”のように意味の弱い語には低い重みを、“ペプチド”“厳密”のような意味の強い語には高い重みを与えたい場合があります。このための古典的な方法として、**tf.idf**というヒューリスティックが知られています。

前者のような言葉の特徴は、「どんな文書にもまんべんなく現れる」ということでしょう。対して、後者のような言葉は特定のトピックと結びついており、全体のごく一部の文書にしか現れないのが特徴です。よって、 $N$  個の文書の中で、ある単語  $w$  が 1 回以上出現した文書の数を  $df(w)$  とおくと（これを**文書頻度** (document frequency) といいます）\*49、 $N$  に対する割合（文書確率）

$$p = \frac{df(w)}{N}$$

が 1 に近いほど  $w$  の意味は弱く、0 に近いほど意味が強いと考えられます。log  $1=0$  ですから、この情報量 (式(2.63)) をとった

| 単語 $w$ | $df(w)/N$ | $idf(w)$ |     |        |        |
|--------|-----------|----------|-----|--------|--------|
| 日      | 0.7003    | 0.3563   | ソニー | 0.0300 | 3.5075 |
| 人      | 0.5640    | 0.5726   | 知識  | 0.0272 | 3.6028 |
| 年      | 0.5014    | 0.6904   | 芸能  | 0.0245 | 3.7081 |
| 的      | 0.4632    | 0.7696   | 夏休み | 0.0191 | 3.9595 |
| 日本     | 0.3651    | 1.0075   | 手作り | 0.0163 | 4.1136 |
| 時間     | 0.2371    | 1.4395   | 画質  | 0.0136 | 4.2959 |
| 使用     | 0.1253    | 2.0767   | 新曲  | 0.0109 | 4.5191 |
| 姿      | 0.0899    | 2.4089   | 真っ白 | 0.0082 | 4.8067 |
| NTT    | 0.0790    | 2.5381   | 危機  | 0.0082 | 4.8067 |
| 説明     | 0.0654    | 2.7273   | 豪雨  | 0.0054 | 5.2122 |
| Web    | 0.0327    | 3.4205   | 愛着  | 0.0054 | 5.2122 |

表 5.11: Livedoor コーパスから計算した、単語  $w$  の文書確率  $df(w)/N$  と  $idf(w) = \log N/df(w)$  (抜粋)。“日”のような語は全体の 7 割の文書に出現しているため  $idf$  は小さく、一方で“豪雨”のような語は全体の 1%未満の文書にしか出現していないため、 $idf$  が大きいことがわかります。

\*49  $n(d, w)$  を使うと、数学的には  $df(w) = \sum_{n=1}^N \mathbb{I}(n(d, w) > 0)$  と表すことができます。

$$(5.98) \quad \text{idf}(w) = -\log p = \log \frac{N}{\text{df}(w)}$$

を単語  $w$  の意味的な重みと考えることができます。df が分母に現れることから、これを**逆文書頻度** (inverse document frequency, idf) といいます。表 5.11 に、Livedoor コーパスから計算した単語の  $p$  と idf の例を示しました。

一方で、特定の文書  $d$  の中では、単語  $w$  の出現回数  $n(d, w)$  が大きいほど意味は強いでしょう。ただし、その影響は線形ではなく、実質的な強さはその対数くらいだと考えられます (ウェーバー=フェヒナーの法則)。つまり、「ソニー」が 50 回出現したからといってその情報量が 50 倍あるわけではなく、効果は  $\log 50 \simeq 3.91$  倍に比例する程度だということです。このままだと頻度 1 のとき  $\log 1 = 0$  になってしまいますので、しばしば

$$(5.99) \quad \text{tf}(n) = 1 + \log n$$

のように定義されます。これを**用語頻度** (term frequency) とよびます。tf と idf の二つを組み合わせると、頻度  $n(d, w)$  を

$$(5.100) \quad \text{tf}(n(d, w)) \cdot \text{idf}(w) = (1 + \log n(d, w)) \cdot \log \frac{N}{\text{df}(w)}$$

と変換すれば、適切な重みづけになると考えられます。この単語重みづけを **tf.idf** といいます。

実際には対数の底に任意性があり、tf や idf の定義も上に示したものに限らず、いくつかのバリエーションがあります [38, §15.2.2]。tf.idf は有効に働くヒューリスティックですが、自然言語処理のタスク全体としてこの重みが最適であるという保証はなく、現代的には 5.2.1 節で説明した PMI や、4.2.2 節で説明した SIF による重みづけの方が数学的な最適性があり、より効果的です。

### 5.5.2 単語ベクトル/文書ベクトルの解釈

こうして得られた文書ベクトルは高速に計算でき、数学的に Word2Vec(Doc2Vec)と同じニューラル文書ベクトルとなっているため、高い性能を持っています。唯一の欠点は、 $K$  個の次元が LDA のようにトピックとして解釈ができないということでしょう。たとえば、上の実験で得られた単語ベクトルを並べた行列  $\mathbf{W}$  について、その 1 次元目、2 次元目、…の値が大きい単語を求めると表 5.12 のようになり、ここには強い規則性は見出せそうにありません。

考えてみるとこれは当然で、式(5.92)による行列分解は式(5.94)のように内積だけを問題にしているため、空間全体を任意に回転しても、図 5.35 のように 2 つのベクトルの間の内積は同じになるからです。数学的には、ある直交行列  $\mathbf{R}$  をとって  $\tilde{\mathbf{D}} = \mathbf{D}\mathbf{R}$ ,  $\tilde{\mathbf{W}} = \mathbf{W}\mathbf{R}$  とベクトル全体を回転したとき、式(5.92)の右辺は

$$(5.101) \quad \tilde{\mathbf{D}}\tilde{\mathbf{W}}^T = \mathbf{D}\mathbf{R}(\mathbf{W}\mathbf{R})^T = \mathbf{D}\underbrace{\mathbf{R}\mathbf{R}^T}_{=\mathbf{I}}\mathbf{W}^T = \mathbf{D}\mathbf{W}^T$$

となり、もとと等しくなります。

すなわち、解釈のためには図 5.35 に示したように、単語ベクトルや文書ベクトルがちょうど軸の近くに配置されるような回転  $\mathbf{R}$  を見つける必要があります。こうした方法として、心理学ではバリマックス回転のような方法が知られ

| 次元 1 | 次元 2  | 次元 3   | 次元 4    |
|------|-------|--------|---------|
| 級    | 自身    | Watch  | 再生      |
| 節電   | イベント  | Sports | ホラー     |
| 全    | クリスマス | 婚      | 恐怖      |
| 電力   | http  | 活      | 音楽      |
| シリーズ | 作っ    | 答え     | デビュー    |
| AKB  | 城     | 音楽     | PC      |
| 母親   | シーズン  | 佐      | 現象      |
| 通話   | テレビ   | 曲      | 娘       |
| スマ   | 婦     | 学校     | ロンドン    |
| 出展   | 超     | 枝      | YouTube |
| 家族   | 家政    | 調査     | 必ず      |
| 額    | www   | 選手     | 韓国      |

表 5.12: Livedoor コーパスから DocVec で計算した単語ベクトルの各次元の値が大きい単語 (一部)。もとのままでは、各次元に明確な意味は見出せそうにありません。



ています[174]. しかし、これは言語のように数百次元以上にもなる高次元の場合にはあまりうまくいかず[175], 最近、京都大学の平らにより、ICA (独立成分分析) を適用することで解釈が容易な軸が、しかも言語横断的に見つかることが示されました[176]. ICA はデータを線形変換して、各次元が可能な限り統計的に独立となるような座標軸を発見する方法です\*50. たとえば図 5.36 では、PCA(主成分分析)ではデータの分散を最大にするように真ん中のような座標軸がとられてしましますが、ICA ではデータの分布を独立な軸の積で説明する右のような座標軸を計算することができます.

「統計的に独立」とは、2章の式(2.21)でみたように、確率変数  $x$  と  $y$  の同時確率が  $p(x, y) = p(x)p(y)$  のように周辺確率の積に分解できることでした. われわれの場合、たとえば式(5.92)で得られた、単語ベクトルを並べた行列  $\mathbf{W}$  を行列  $\mathbf{A}$  を使って  $\mathbf{S} = \mathbf{WA}$  と線形変換したとき、 $\mathbf{S}$  の各行ベクトル  $\mathbf{s} = (s_1, s_2, \dots, s_K)$  について、分解

$$(5.102) \quad p(s_1, s_2, \dots, s_K) = p(s_1)p(s_2) \cdots p(s_K)$$

が可能な限り成り立つような行列  $\mathbf{A}$  を求めることになります\*51.

こうした  $\mathbf{A}$  は、機械学習の分野でICAを定式化したフィンランドの Hyvärinen

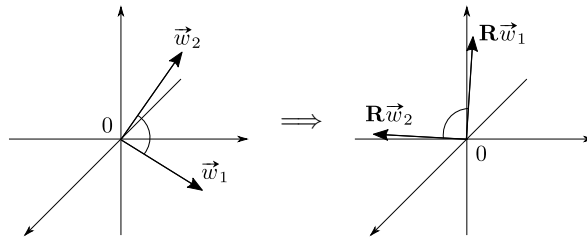


図 5.35: 単語ベクトルの回転. 単語ベクトルを直交行列  $\mathbf{R}$  で回転しても、二つの単語ベクトルのなす角度は不変です. 適切な回転  $\mathbf{R}$  を求めることで、単語ベクトルを軸に沿った形で、よりわかりやすく解釈することができます.

\*50 162 ページで説明した白色化はデータを無相関にする方法ですが、独立とは無相関を含んでおり、それより強い条件になります.

\*51 ICA は実際には、162 ページで説明したデータの白色化を行ってから  $\mathbf{R}$  による回転を行うことと等価で、白色化行列を  $\mathbf{B}$  とおくと  $\mathbf{A} = \mathbf{BR}$  になります.

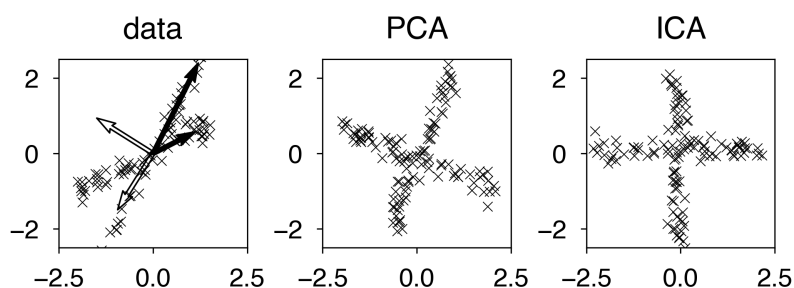


図 5.36: ICA の概要図. PCA(主成分分析) では白矢印のように、データ (×印) の分散を全体的に説明する軸がとられてしましますが、ICA(独立成分分析) では黒矢印のように、データの分布を独立な軸の積として説明する軸が求められているのがわかります。

によるパッケージ FastICA<sup>\*52</sup> で計算することができます。Python では、

```
from sklearn.decomposition import FastICA # ICA の計算
from scipy.stats import skew # 歪度の計算
import numpy as np
def ica (X):
 X = X - np.mean (X, axis=0)
 analyzer = FastICA (whiten="arbitrary-variance")
 S = analyzer.fit_transform (X)
 A = analyzer.components_
 # sort by skewness
 N,D = S.shape
 skews = np.abs (skew (S, axis=0)) # 0=Gaussian
 index = map (lambda x: x[1], sorted (zip(skews, np.arange(D)),
 key=lambda x: x[0], reverse=True))
 return S[:,list(index)], A
```

で  $\mathbf{S}$  および  $\mathbf{A}$  を求めることができます。多くの信号が混ざると、中心極限定理によって分布はガウス分布に近づくことから、逆に ICA で分解した軸では、各次元の周辺分布は非ガウ斯的になります。その度合いは、期待値  $\mu$ 、標準偏差  $\sigma$  をもつ確率変数  $X$  の歪度

$$(5.103) \quad \delta = \mathbb{E}[(X - \mu)^3] / \sigma^3$$

で表すことができます[177]。ガウス分布では  $\delta$  は 0 で、± になるほど左右の裾

\*52 <http://research.ics.aalto.fi/ica/fastica/>

が広くなることが知られています。ICA は PCA と異なり、独立性の高い軸から求まるとは限りませんので、上のコードでは  $\delta$  の絶対値を用いて、非ガウス性の高い順に次元を並び換えています。Livedoor コーパスから計算した単語ベクトルについて、ICA で変換した各次元の値の大きい単語を表 5.13 に示しました。ほとんどトピックモデルのような、解釈性の高い意味的な軸が求められていることがわかります。

ただしトピックモデルと異なり、これはベクトル空間の「軸」ですので、負の方向も存在します<sup>\*53</sup>。表 5.14 に、表 5.13 のいくつかの軸について値が負の単語を示しました。単なる文書-単語の共起行列から得られたにもかかわらず、iPhone ↔ Android、ビジネス ↔ 悩み相談のような興味深い対立軸が教師なしで得られていることがわかります。

## 5.6 確率的潜在意味スケーリング (PLSS)

前節の ICA の結果から、文書ベクトルおよび単語ベクトルの存在する埋め込み空間には、トピックモデルのトピックに対応するようなさまざまな「軸」が存在することがわかりました。実際に Word2Vec (これは前節で主成分分析で求めたものと、3.5.4 節で説明したように数学的には同じものでした) の空間には、ICA で見つかるものに限らず多くの意味的な軸が存在することがわかっています[178]。

特に重要なのは、表 5.14 でみたように多くの場合、この軸は正の方向と負の方向をもつ対立軸となっているということです<sup>\*54</sup>。すると、表 5.13 と表 5.14 で軸と向きが近い、または逆の単語を計算したように、図 5.37 のようにこの軸上に文書ベクトルを射影し、テキストの性質をある軸に沿って 1 次元の  $\pm$  の実数値として表すことができるでしょう。これを心理学の用語では、テキストの**尺度化** (scaling) といいます。

---

\*53 空間全体をある軸に関して折り返してもベクトル間の関係は変わりませんので、符号が正負どちらになるかに大きな意味はありません。

\*53 <https://news.livedoor.com/article/detail/6029862/>

\*54 ただし、Transformer で得られる埋め込み空間は注意機構のために複雑になり、こうした線形性があまり成り立たなくなっています。

| 次元 3   |        | 次元 5    |        | 次元 6   |        | 次元 10  |        |
|--------|--------|---------|--------|--------|--------|--------|--------|
| iPhone | 0.0224 | apps    | 0.0853 | ビジネス   | 0.0180 | 等      | 0.0272 |
| クリック   | 0.0185 | google  | 0.0851 | 経営     | 0.0174 | 由里子    | 0.0269 |
| ブック    | 0.0163 | 要件      | 0.0848 | 成功     | 0.0172 | 吉      | 0.0268 |
| 電子     | 0.0160 | store   | 0.0847 | キャリア   | 0.0164 | 愛      | 0.0265 |
| 術      | 0.0160 | play    | 0.0842 | オフィスエム | 0.0155 | サイト    | 0.0260 |
| アップ    | 0.0152 | ANDROID | 0.0821 | 笑      | 0.0147 | 幸せ     | 0.0246 |
| サイト    | 0.0150 | details | 0.0820 | 話し     | 0.0144 | 果たし    | 0.0243 |
| 携帯     | 0.0142 | Play    | 0.0779 | 管理     | 0.0139 | 公式     | 0.0242 |
| 既報     | 0.0135 | Google  | 0.0600 | スキル    | 0.0132 | 務め     | 0.0234 |
| 背面     | 0.0134 | Hisumi  | 0.0565 | 戦略     | 0.0129 | 篇      | 0.0233 |
| ライフ    | 0.0123 | 以上      | 0.0421 | 力      | 0.0129 | リアル    | 0.0224 |
| iPad   | 0.0123 | Android | 0.0420 | 重要     | 0.0124 | 女優     | 0.0223 |
| 次元 15  |        | 次元 19   |        | 次元 28  |        | 次元 39  |        |
| 高画質    | 0.0122 | サッカー    | 0.2559 | ロードショー | 0.0428 | スタイル   | 0.0218 |
| デジタル   | 0.0105 | 代表      | 0.2401 | 全国     | 0.0425 | オシヤレ   | 0.0197 |
| 実現     | 0.0095 | 戦       | 0.2193 | 女優     | 0.0237 | 着      | 0.0171 |
| 操作     | 0.0093 | 試合      | 0.2175 | 決意     | 0.0206 | ファッション | 0.0170 |
| ソニー    | 0.0084 | 杯       | 0.1897 | 土      | 0.0200 | 楽しむ    | 0.0166 |
| ズーム    | 0.0082 | W       | 0.1591 | 実力     | 0.0199 | シンプル   | 0.0164 |
| 進化     | 0.0081 | チーム     | 0.1532 | 最強     | 0.0198 | ダイエット  | 0.0153 |
| 保存     | 0.0080 | ゴール     | 0.1410 | 黄金     | 0.0194 | 味わい    | 0.0153 |
| シーン    | 0.0080 | 日本      | 0.1280 | 絆      | 0.0193 | 体重     | 0.0151 |
| 新      | 0.0079 | 予選      | 0.1277 | 祭      | 0.0187 | 食事     | 0.0151 |
| 軽量     | 0.0079 | 選手      | 0.1261 | ベルセルク  | 0.0186 | 誰      | 0.0147 |
| レス     | 0.0079 | リーグ     | 0.1239 | TOHO   | 0.0184 | 運動     | 0.0145 |

表 5.13: ICA で変換した単語ベクトルの各次元の値が大きい単語 (抜粋). 次元は歪度の絶対値の大きい順に並んでいます. 表 5.12 と比べて, ほとんどトピックモデルのような, 意味的な軸が学習されていることがわかります. 次元 10 の意味については, 表 5.14 を参照してください.

テキストを分類するのではなく, 連続値で測る尺度化は実際のさまざまな場面において有効です. たとえば, 法案を保守 (右翼) か革新 (左翼) のどちらかに分類するのは簡単でも, 実際の法案は日本ではほとんど保守 (自民党) の側から提出されており, 問題はむしろ, その法案がどれくらい保守的なのか, ということでしょう. また, 住民へのアンケートやホテルの評価に書かれたテキストを肯定・否定に分類するのは簡単ですが, 重要なのはどの意見が強い賛成・反対であり, どの意見が軽い文句なのかを見極めることでしょう. 小説家のテキストからその分裂症気質を測りたいといった場合でも, 分裂症かどうかという二値分類

| 次元 3    |         | 次元 6 |         | 次元 10 |         | 次元 19  |         |
|---------|---------|------|---------|-------|---------|--------|---------|
| apps    | -0.2364 | お答え  | -0.2046 | 妖     | -0.2595 | フィギュア  | -0.0405 |
| store   | -0.2363 | 辛口   | -0.2041 | ケ     | -0.2562 | スケート   | -0.0403 |
| details | -0.2362 | 悩め   | -0.2026 | 巻     | -0.2316 | 選手権    | -0.0364 |
| play    | -0.2300 | 説教   | -0.2018 | 劇場    | -0.1994 | 演技     | -0.0312 |
| google  | -0.2299 | 尽き   | -0.1938 | 勅使河原  | -0.1992 | 克也     | -0.0268 |
| 要件      | -0.2225 | 早    | -0.1856 | 妖怪    | -0.1967 | 浅田     | -0.0265 |
| ANDROID | -0.2187 | 面白く  | -0.1820 | 栄華    | -0.1961 | 真央     | -0.0262 |
| id      | -0.2119 | 姉妹   | -0.1689 | 仙人    | -0.1900 | メダリスト  | -0.0262 |
| Play    | -0.2006 | 悩み   | -0.1680 | お嬢様   | -0.1695 | 野村     | -0.0255 |
| com     | -0.1913 | type | -0.1663 | ネコ    | -0.1589 | ノム     | -0.0246 |
| Store   | -0.1798 | 若手   | -0.1650 | 話     | -0.1497 | バンクーバー | -0.0231 |
| カテゴリ    | -0.1682 | 瞬時   | -0.1596 | 次     | -0.1430 | 野球     | -0.0225 |

表 5.14: ICA で変換した単語ベクトルの各次元の値が小さい単語 (抜粋). 表 5.13 と見比べると, 次元 3 では iPhone の「反対の概念」として Android が, 次元 6 ではビジネスについてお悩み相談が, 次元 19 ではサッカーについてフィギュアスケートが得られていることがわかります. 次元 10 は元データの Livedoor ブログでの“いちおう妖ヶ劇場”という連載記事<sup>\*54</sup>に関連している軸で, 負の方にだけ意味を持っています.

にはあまり意味がなく (小説家の多くは分裂症気質のため), 分裂症気質の強さやその変化が主な興味の対象になると考えられます.<sup>\*55</sup>

こうした尺度化を行うもっとも簡単な方法は, 単語ベクトルや文書ベクトルの与えられた軸への近さを, 前節のように  $\cos$  距離を計算して求めることです. しかし, 3 章の図 3.37 で示したように, 埋め込みベクトルの間の  $\cos$  距離は 0 を中心には分布せず, 値も  $[-1, 1]$  の中の一部しかとりません. 別の言い方をすれば,  $\cos$  距離はあくまで後付けの値であり, ある軸に関係しないベクトルの値が 0 となる尺度としては設計されていない, ということです.

### 5.6.1 Wordfish

特に, 社会科学においてテキストの尺度化は重要な問題であるため, 政治学方法論 (Political methodology) の分野では, Slapin らが 2008 年に政党などの連続した極性の時間変化を求める Wordfish という方法を提案しました[180].

<sup>\*55</sup> 図 5.5 でもみたように, 分類モデルは尤度を上げるために極性を  $\pm 1$  のどちらかに寄せる傾向があり, SVM のようなベクトル空間の分類器でも, 分離平面からの距離がクラスに対応する尺度とは限りません[179].

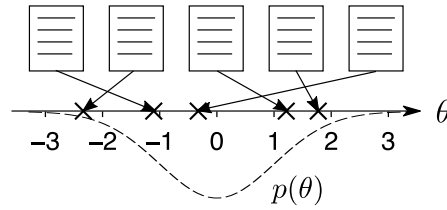


図 5.37: テキストの尺度化と潜在スケール (尺度)  $\theta$ . 各文書に対して、測りたい軸に沿った実数値  $\theta \in \mathbb{R}$  を推定します。項目反応理論に従い、 $\theta$  は標準正規分布  $\theta \sim \mathcal{N}(0, 1)$  に従う潜在変数だと考えます。

$y_{itv}$  を政党  $i$  が時刻  $t$  のテキスト (たとえば選挙時の公約) で単語  $v$  を使った回数とすると、Wordfish はデータ  $Y = \{y_{itv}\}$  ( $i = 1, \dots, I, t = 1, \dots, T, v = 1, \dots, V$ ) の確率を、次のようにポアソン分布  $\text{Po}(y|\lambda)$  でモデル化します。

$$(5.104) \quad \begin{cases} p(Y) = \prod_{i=1}^I \prod_{t=1}^T \prod_{v=1}^V \text{Po}(y_{itv} | \lambda_{itv}) \\ \lambda_{itv} = \exp(\alpha_{it} + \beta_v + \phi_v \cdot \theta_{it}) \end{cases}$$

ここで  $\alpha_{it}$  はテキスト  $it$  での固定効果 (ベースライン)、 $\beta_v$  は単語  $v$  の固定効果で、興味があるのは単語  $v$  の極性軸上での位置  $\phi_v \in \mathbb{R}$  と、政党  $i$  の時刻  $t$  での潜在位置 <sup>\*56</sup> $\theta_{it} \in \mathbb{R}$  です。

式(5.104)は、頻度  $y_{itv}$  はポアソン分布  $\text{Po}(\lambda)$  に従い、その期待値はテキスト  $it$  と単語  $v$  で決まるベースライン  $\alpha_{it} + \beta_v$  を、政党  $i$  の時刻  $t$  での極性  $\theta_{it}$  と単語の持つ極性  $\phi_v$  で上下して決まる、ということの意味をしています。 $\theta > 0$  が右翼、 $\theta < 0$  が左翼を表すとしたとき、政党の位置と単語の符号が一致する、すなわち  $\theta_{it} > 0$  かつ  $\phi_v > 0$  (たとえば  $v = \text{“軍備”}$ )、または  $\theta_{it} < 0$  かつ  $\phi_v < 0$  (たとえば  $v = \text{“社会保障”}$ ) のとき  $y_{itv}$  の期待値  $\lambda_{itv}$  は大きくなり、符号が逆ならば小さくなります。 $\{\alpha, \beta, \theta, \phi\}$  はすべて未知のため、推定には 246 ページで学習した EM アルゴリズムを用い、適切な初期値から始めて、 $(\alpha, \theta)$  の推定と  $(\beta, \phi)$  の推定を反復します。こうして求めたドイツの各政党の  $\theta_{it}$  の時間変化を、図 5.38 に示しました。

\*56 これを政治学の分野では、理想点 (ideal point) といいます。

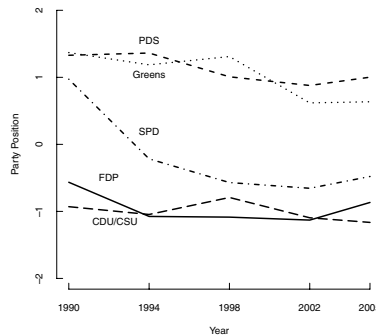


図 5.38: Wordfish によって推定された、ドイツの各政党の外交方針  $\theta_{it}$  の時間変化の例 (1990–2005) ([180]より引用). 縦軸の  $\theta$  は平均が 0 の潜在変数になっています.

### 5.6.2 確率的潜在意味スケーリング (PLSS)

Wordfish は与えられたテキストに対する系統的な統計モデルですが、完全な教師なし学習のため、これによって得られる  $\theta$  の極性が分析の目的と一致している、という保証はありません。適切な結果を得るためには、入力テキストを注意深く選ぶ必要がありますが、そのための客観的な方法は示されておらず、また学習データの単語頻度に対するポアソン分布はテキストの長さに暗黙に依存するため、新しいテキストでの  $\theta$  を計算できないという問題もあります。

そこで、心理統計学における**項目反応理論** (Item Response Theory, IRT) を参考に、潜在的な極性  $\theta \in \mathbb{R}$  をもつテキストで単語  $v \in \{1, \dots, V\}$  が出現する確率を、次のように多項分布でモデル化します。

$$(5.105) \quad \begin{aligned} p(v|\theta, \phi) &\propto p(v) \exp(\theta \cdot \phi_v) \\ &= \frac{\exp(\log p(v) + \theta \cdot \phi_v)}{\sum_{v=1}^V \exp(\log p(v) + \theta \cdot \phi_v)} \end{aligned}$$

ここで  $\phi_v \in \mathbb{R}$  は式(5.104)の Wordfish の場合と同様に、単語  $v$  の「極性」を表すパラメータで、 $\theta$  は 0 を中心とした標準正規分布

$$(5.106) \quad \theta \sim \mathcal{N}(0, 1)$$

に従うとします。

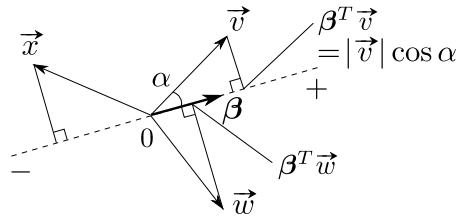


図 5.39: 方向ベクトル  $\beta$  による単語の極性の計算.  $\beta$  との内積, すなわち  $\beta$  上への正射影の長さで単語の極性  $\phi_v = \beta^T \vec{v}$  を計算します. ここで  $|\beta|=1$  で,  $\vec{v}, \vec{w}, \vec{x}$  は単語ベクトルを表します. 単語ベクトルが  $\beta$  と逆方向の場合は, 極性はマイナスになります.

式(5.105)より, このモデルでは Wordfish と同様に単語  $v$  の確率は  $\phi_v$  と  $\theta$  の正負と大きさが一致すれば事前確率  $p(v)$  より高くなり, 逆になれば低くなる, というモデルになっています.  $p(v)$  はコーパスから容易に計算できるため, このモデルは IRT の多項分布化, あるいは式(5.105)で表される多値ロジスティック回帰において, 説明変数  $\theta$  も回帰係数  $\phi$  も未知の場合の教師なし学習とみなすことができます.

このとき, テキスト  $d$  の確率は単語  $v$  のテキスト内での頻度を  $n_{dv}$  とおくと

$$(5.107) \quad p(d|\theta, \phi) = \prod_{v=1}^V p(v|\theta, \phi)^{n_{dv}}$$

と書けますから,  $D$  個のテキストからなるコーパス  $\mathcal{D}$  全体の確率は

$$(5.108) \quad p(\mathcal{D}|\Theta, \phi) = \prod_{d=1}^D \prod_{v=1}^V p(v|\theta_d, \phi)^{n_{dv}}$$

と表されます. Wordfish と同様に, 事前確率 (5.106) の下で式(5.108)を最大化するパラメータ  $\Theta = \{\theta_1, \dots, \theta_D\}$  および  $\phi_1, \dots, \phi_V$  を MCMC 法や EM アルゴリズムなどによって計算することができます.\*57

ただし, こうするとたとえ単語  $v$  と  $w$  が意味的に関係が深くても,  $\phi_v$  と  $\phi_w$  は別のパラメータとして推定しなければならないという問題があります. たとえば  $\phi_{\text{good}} > 0$  と推定できて, これは  $\phi_{\text{excellent}}$  や  $\phi_{\text{well}}$  とは無関係で, excellent

\*57 これは多項分布によるモデルのため, ポアソン分布による Wordfish と異なり, パラメータが学習テキストの長さに依存せず, 新しいテキストについても適用することができるという特徴があります.



や well がコーパスに現れなければ、まったく学習することができません。

そこで、 $\phi_1, \dots, \phi_V$  を独立に学習する代わりに、与えられたコーパスあるいは一般的なコーパスから事前に Word2Vec や GloVe など学習された  $K$  次元のニューラル単語ベクトル  $\vec{v}$  を用いて、 $\phi_v$  を

$$(5.109) \quad \phi_v = \beta^T \vec{v} \quad (\beta = (\beta_1, \beta_2, \dots, \beta_K)^T)$$

とモデル化します。この後で説明するように  $\beta$  の長さは 1 に正規化しますから、これは図 5.39 のように、各単語ベクトル  $\vec{v}$  と  $\beta$  のなす角を  $\alpha$  としたとき、単語の極性を  $\vec{v}$  の  $\beta$  上への正射影の長さ  $\beta^T \vec{v} = |\beta| |\vec{v}| \cos \alpha = |\vec{v}| \cos \alpha$  で「測つて」いることに相当します。これにより、 $V$  個の独立な  $\phi_1, \dots, \phi_V$  を求めるかわりに、 $K$  次元の方向ベクトル  $\beta$  を一つだけ推定すればよいことになります。このとき、式(5.105)は  $\ell_v = \log p(v)$  とおけば、

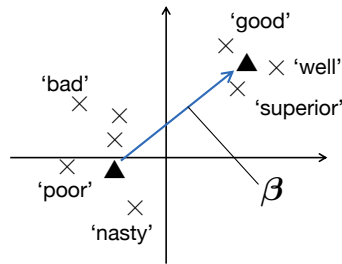
$$(5.110) \quad p(v|\theta, \beta) = \frac{\exp(\ell_v + \theta \cdot \beta^T \vec{v})}{\sum_{v=1}^V \exp(\ell_v + \theta \cdot \beta^T \vec{v})}$$

と表すことができます。筆者が開発したこの方法は、政治学分野で提案された、286 ページの LSI をベースにしたベクトル空間におけるアルゴリズムである LSS (Latent Semantic Scaling) [181] の確率化ともみなすことができるため、確率的 LSS, **PLSS** (Probabilistic Latent Semantic Scaling) [182] とよぶことにします。

この  $\beta$  は、単語埋め込みベクトルの空間において“良い-悪い”、“右翼-左翼”といった  $\theta$  の極性を与える「極性軸」、あるいは「意味方向」を表しています。 $\beta$  は完全に教師なしで学習することもできますが、5.6.1 節で述べたように、そうして得られた  $\beta$  が分析の目的と一致しているとは限りません。そこで、PLSS では表 5.15 のように、正例および負例として与える少数の極性語辞書から、単語ベクトルを用いて  $\beta$  を次のように計算します。なお、 $\beta$  のノルムは 1 に正規

表 5.15: 表 5.4 の WRIME コーパスの極性語を参考に作成した、標準的な極性辞書。+ が正例を、- が負例を表します。

|   |                                      |
|---|--------------------------------------|
| + | 最高 嬉しい 楽しい 楽しみ 可愛い きれい しあわせ 好き かわいい  |
| - | 悪い 悲しい 寂しい 嫌い 怒り ない つらい 無理 しんどい めんどく |

図 5.40: 極性辞書と単語ベクトルによる  $\beta$  の計算.

化します.

$$(5.111) \quad \beta \propto \left( \frac{1}{|S_+|} \sum_{v \in S_+} \vec{v} - \frac{1}{|S_-|} \sum_{w \in S_-} \vec{w} \right)$$

ここで  $S_+$  は極性辞書のうち正の語の集合,  $S_-$  は負の語の集合です. 式(5.111) は単純に,  $\beta$  として負の単語ベクトルの平均から, 正の単語ベクトルの平均へ向かう方向を取ることを表しています. この様子を図 5.40 に示しました. 極性辞書としては, 表 5.15 に示した標準的な極性辞書のように, ごく少数の単語のリストがあれば動きます. なお, 正例・負例の単語が「最も極端な語」である必要はなく, 意味方向のみが合っていればよいことに注意してください.

ここで「正例」「負例」とは必ずしも感情的な正負と関係していなくてもよく, 「右翼-左翼」「都会-田舎」「東洋-西洋」のように, 任意の意味的な軸を扱うことができます. ドイツの Schütze らは, 単語埋め込みの空間に実際にこうした, 与えられたタスクに関する低次元の部分空間が存在することを発見し, これを超密埋め込み (Ultradense embedding) と呼んでいます[178]. 超密埋め込みは, 最も単純には, この場合のように 1 次元の部分空間となります. ただし, 予備実験でこの超密埋め込みを行列の固有ベクトルとして求める DensRay [183] を使用したところ, 式(5.111)と比べて明らかにノイズの多い方向となったため, PLSS では単純な式(5.111)を採用することとしました.\*58

\*58 これは, DensRay の目的関数が, 極性辞書で単語  $v$  の属する極性を  $s(v)$  としたとき, 行列  $\mathbf{A} = \frac{1}{|S_+|} \sum_{\substack{(v,w): \\ s(v)=s(w)}} (\vec{v}-\vec{w})(\vec{v}-\vec{w})^T - \frac{1}{|S_-|} \sum_{\substack{(v,w): \\ s(v) \neq s(w)}} (\vec{v}-\vec{w})(\vec{v}-\vec{w})^T$  の固有ベクトルの計

表 5.16 に、表 5.15 の標準的な極性辞書と、後で説明する日本語 Wikipedia 記事から Word2Vec の方法で学習した  $K=100$  次元の単語ベクトルを用いて計算した単語の極性  $\phi_v = \beta^T \vec{v}$  を示しました。非常に少ない教師データにもかかわらず、肯定的な単語および否定的な単語が、その強さとともに連続的に取り出せている様子がわかります。

極性辞書を用いる場合は、PLSS は単語ベクトルの計算以外にパラメータの学習を必要としません。式(5.107)から、テキスト  $d$  と極性  $\theta$  の同時確率は

$$(5.112) \quad p(d, \theta) = \prod_{v=1}^V p(v|\theta, \beta)^{n_{dv}} \cdot p(\theta)$$

$$(5.113) \quad = \prod_{v=1}^V \left( \frac{\exp(\ell_v + \theta \cdot \beta^T \vec{v})}{\sum_{v=1}^V \exp(\ell_v + \theta \cdot \beta^T \vec{v})} \right)^{n_{dv}} \cdot \mathcal{N}(\theta | 0, 1)$$

となり、これを最大にするテキストの潜在的な極性  $\theta$  の MAP 解は、1次元の最適化で容易に計算することができます。

**PLSS の実験** 232 ページで使った WRIME コーパスのツイートの感情極性を、PLSS で分析してみましょう。本書では基本的に Python を使用していますが、政治学方法論分野では R が普及しているため、本節では R で実装を行っています。以下で使用しているスクリプトは、すべてサポートサイトからダウンロードすることができます。

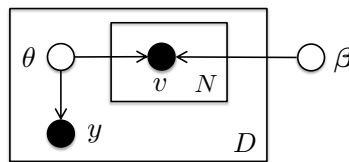


図 5.41:  $\beta$  を推定するための教師データのグラフィカルモデル。  $\beta$  だけでなく  $\theta$  も未知のため、  $\theta$  については周辺化を行って積分消去することで  $\beta$  を推定します。

算に帰着され、正例と負例の内部およびその間のペアを結ぶ個別のベクトルの和になっていることが原因だと考えられます。式(5.111)と異なり、この定式化では  $S_+$ ,  $S_-$  の中で和をとることでノイズを消す作用が働かず、  $v, w$  の選択に起因するノイズが  $\mathbf{A}$  に残ってしまうからです。

表 5.16: Wikipedia のコーパスから事前学習した単語ベクトルと、表 5.15 の極性辞書を用いて計算した WRIME コーパスでの単語の極性  $\phi_v = \beta^T \vec{v}$ .

| (a) 極性 $\phi_v > 0$ の上位語 |        |       |        | (b) 極性 $\phi_v < 0$ の下位語 |         |       |         |
|--------------------------|--------|-------|--------|--------------------------|---------|-------|---------|
| ミッフィー                    | 1.5639 | 星座    | 1.1436 | 苦しん                      | -1.8143 | 許せ    | -1.4058 |
| チェキ                      | 1.4430 | スク    | 1.1274 | 耐え                       | -1.7748 | 察し    | -1.3846 |
| キラキラ                     | 1.3147 | メゾン   | 1.1142 | 生じ                       | -1.7656 | しまっ   | -1.3837 |
| ♪                        | 1.2942 | プレゼント | 1.0828 | 治まっ                      | -1.7211 | デメリット | -1.3749 |
| かわいい                     | 1.2609 | リゾート  | 1.0742 | 予断                       | -1.6880 | しまう   | -1.3586 |
| アメニティ                    | 1.2367 | 展     | 1.0505 | 断ち切ら                     | -1.6177 | 起き    | -1.3475 |
| クロワッサン                   | 1.2283 | 土産    | 1.0492 | 原因                       | -1.6022 | 断ち切っ  | -1.3388 |
| たのしい                     | 1.2268 | ミント   | 1.0474 | ざる                       | -1.5514 | 後悔    | -1.3380 |
| ミルク                      | 1.1964 | マリオ   | 1.0255 | 容赦                       | -1.5478 | きれ    | -1.3336 |
| しあわせ                     | 1.1900 | 満喫    | 1.0134 | 断ち切れ                     | -1.5090 | 症状    | -1.3309 |
| ニベア                      | 1.1663 | 美容    | 1.0097 | 吐き                       | -1.4649 | 不安    | -1.3211 |
| 水着                       | 1.1637 | hello | 1.0028 | ひどく                      | -1.4551 | 許さ    | -1.3163 |
| 応募                       | 1.1579 | ほっと   | 0.9945 | 訴え                       | -1.4377 | 気付か   | -1.3146 |
| アニメイト                    | 1.1544 | 香菜    | 0.9894 | 隙                        | -1.4274 | 村八分   | -1.3140 |
| ♡                        | 1.1466 | コス    | 0.9867 | 困憊                       | -1.4211 | 雑言    | -1.3083 |

情報が少ないと精密な分析ができないため、ここでは `wrime.test` の中で 20 単語以上のツイートを対象にすることにして、次のようにしてテキストだけを抽出します。結果は、302 ツイートになりました。<sup>\*59</sup>

```
% awk 'NF>20' wrime.test | cut -f 2 > wrime.txt
% head wrime.txt
先輩の息子が、まじ一瞬だけどつかまらず立つ..
わたしの伝え方がおかしいのか？なぜ、全てや..
「落ちる」場面や、走っても走っても前に進..
Shift + win キー + S キーで範囲指定の画面キャプ..
```

`wrime.txt` はこのように、1 行に 1 文書 (1 ツイート) の単語が並んだ、単なるテキストファイルです。PLSS は正解ラベルを必要としませんが、検証のために取り出しておきましょう。

```
% awk 'NF>20' wrime.test | cut -f 1 > wrime.label
```

<sup>\*59</sup> `cut` は、タブで区切られたテキストの指定された番号のフィールドを表示する、Unix の標準コマンドです。`awk 'NF>20'` とは、空白で区切られたフィールド数が 20 より大きいとき標準アクション (その行を表示) を行う、という `awk` (71 ページ) のコマンドです。

単語ベクトルの種類は任意ですが<sup>\*60</sup>, ここでは Wikipedia の記事から作成された単語ベクトルの公開データ Wikipedia2Vec [184]<sup>\*61</sup> から, 100 次元の日本語単語ベクトルを使うことにします. サイトから “Japanese 100d(txt)” のベクトルを入手し, bzip2 で解凍します.

```
% bzip2 -dc jawiki_20180420_100d.txt.bz2 > jawiki.vec
```

分析に使う極性辞書は, 表 5.15 のものを用います. これは +, - の順に “ラベル<TAB>単語..” というフォーマットで, サポートサイトに posneg.ja という名前で置いてあります.

```
% cat posneg.ja
positive 最高 嬉しい 楽しい 楽しみ 可愛い きれ..
negative 悪い 悲しい 寂しい 嫌い 怒り ない つら..
```

この上で, 次のように plss を実行してツイートの極性を計算します. 使い方は plss の中を読むか, このスクリプトを単独で実行してみてください. 以下では R でテキストを扱うパッケージ quanteda<sup>\*62</sup> および関連ライブラリがインストールされていることを前提にしていますので, エラーが出た場合は install.packages でインストールしておいてください.

```
% plss wrime.txt posneg.ja jawiki.vec output
⇒ preparing word vectors..
total 2762 words selected.
running PLSS..
loading wordvectors from /tmp/wordvec-47752.vec.. done.
preparing data.. done.
documents = 302, vocabulary = 2762
computing theta..
computing 302/302.. done.
theta written to output.theta.
phi written to output.phi.
```

保存された output.theta が各ツイートの極性  $\theta \sim \mathcal{N}(0, 1)$ , output.phi は表 5.16 に示した, 内部で計算した単語の極性  $\phi_v$  です. (推定に用いなかった) 正

<sup>\*60</sup> 対象となるテキスト自体から, 3 章で紹介した方法で単語ベクトルを学習することも原理的には可能ですが, これには通常, かなり大量のテキストを必要とします.

<sup>\*61</sup> <https://wikipedia2vec.github.io/wikipedia2vec/>

<sup>\*62</sup> <http://quanteda.io/>; 執筆時は quanteda 4.0 の上で実装しています.

解のラベルおよびツイート本文とともに表示してみましょう。次のスクリプト `theta.sh` を実行すると、ランダムな 20 ツイートを極性  $\theta$  の値でソートして表示します。

```
% theta.sh output.theta wrime.label wrime.txt
⇒ positive 1.8362 あー(*´▽`)キセキたちが私を萌殺そうとす..
positive 0.9879 ぐっどこんでいしょん。心も頭もクリア。秋分..
positive 0.9121 やっぱ神席だった。キラート細胞さんってマチ..
positive 0.7518 3週間ぶりに卓球してきました♪高校生の時の..
positive 0.7212 ソフトバンク、楽天、元 zozo の前澤さんなど..
negative 0.6797 あっ、ぶんぐ博の荷物忘れた…。今日、先生に..
positive 0.6390 高級牛乳買ってから寝る前に牛乳を1杯飲むよ..
positive 0.5202 今日発売のじょーちゃん載ってる雑誌さすがに..
positive 0.3314 同席のご夫婦さんに JAW さんファン歴何年ですか..
positive 0.0837 早速遊びすぎて電池無くなったよ!家の中って被..
positive -0.0546 日々幸せだなあって思ってるんだけど、今日ほん..
negative -0.0575 あ、明日大阪に行ったら接触確認アプリの働きが..
negative -0.2291 ほわ〇ぶワナビとか借〇玉ワナビ、大抵頭(特に)..
positive -0.4468 うおー肩と腰が痛いわーごみ出しに行かないきゃー..
negative -0.6453 とある友人の言動が気になって仕方がない時があ..
negative -0.6546 気にかけてくれる素晴らしい上司に恵まれて、な..
negative -0.6990 どうしよう。来年行く旅行のことについて考えよ..
negative -0.7387 うちの母親でもここまでやったことはないかなあ..
negative -1.0035 体は疲労困憊なのに妙な緊張で眠れぬ今日もぐっ..
negative -1.1121 4時過ぎまで眠れず昨日打ち合わせた仕事モヤ..
```

このように、PLSS ではごく少数の極性語辞書だけで、テキストをその軸に沿って極性の強さに応じた連続値で並べることができ、陽性-陰性の場合には人手で付与した極性とも、ほぼ一致していることがわかります。

極性辞書は陽性-陰性だけでなく、任意の軸で作ることができます。図 5.42 では、右翼-左翼の政治的立場を表す (a) の極性語辞書を用いて、5.1 節の Livedoor コーパスの “topic-news” カテゴリーの記事を解析した様子を (b) に示しました。記事にはもちろん、右翼-左翼の軸に関するラベルはありませんが、確かにこの極性軸に沿って文書を並べることができることがわかります。

### 5.6.3\* PLSS の半教師あり学習

上では  $\theta$  の極性を表す基準として LSS と同様に少量の極性辞書を用いましたが、こうした辞書が分析対象について自明に作成できるとは限りません。例え

positive 保守 国家 伝統 経済 秩序 軍事 成長 資本  
 negative 平等 多様性 環境 平和 福祉 権利 労働

(a) 右翼-左翼の政治的立場を表す極性語辞書.

1.8022 「富士山をしろ」日本の軍事誌の内容に中国ネットユーザー..  
 1.6412 熊本の“政治家らしからぬ”経歴が話題 25 日に行われた熊..  
 1.4820 2011 年の日本の地震図に「すぎる」2011 年 1 月 1 日 00:00..  
 1.3375 富士山近郊です地震に不安の声相次ぐ 1 月 29 日に発生し..  
 1.2284 中華圏で大ブレイクした“日本作り” 19 日に「日経 NET」に..  
 1.1621 好きな女子アナランキングにネット騒然 9 日、STYLE が行っ..  
 1.1586 北朝鮮が『ミサイル』発射も 1 分で落下 13 日、TBS を含..  
 1.1239 世界で有名な日本人ベスト 3 に意外な人物が続出 16 日放送..  
 :  
 :  
 -0.9924 猫ひろし、日本国籍再取得にも「当たり前だ」の声法律相談..  
 -1.0196 死刑廃止論者がネット掲示板で 29 日、ダイヤモンドオン..  
 -1.0341 学校を休んでドイツニーランドへは、ありかしかり、費、読..  
 -1.0834 生活保護受給をめぐる物議を醸す河本親子に法的措置の可能..  
 -1.1430 見知らぬ団体が勝手にハロプロの YouTube 公式チャンネルを..  
 -1.2003 東京都がまたマンガ・アニメ規制の動きか東京都が男女平等..  
 -1.2050 河本の姉が片山さつき氏に「後で謝ることになる」ネットで..  
 -1.4021 生活保護のイメージ悪化に抗議も、反論の声多数 30 日、生活..

(b) 解析した  $\theta$  の上位および下位 8 個の記事.

図 5.42: PLSS による Livedoor ニュースコーパスの解析. “topic-news” のラベルをもつ 770 個の記事の内容を, (a) の極性語辞書を用いて分析した結果を (b) に示しました. 右寄りな記事にはより高い  $\theta$  が, 左寄りな記事にはより低い  $\theta$  が推定されています.

ば, 欧州において移民労働者についての賛成派と反対派を特徴づけるキーワードが, 分析前から明らかとは限らないからです.

しかし, そうした場合でも典型的な「正例」のテキストと「負例」のテキストは示せる場合が多いと考えられます. 直感的には, それぞれのテキストに共通して現れる単語 (単語ベクトル) から, 間接的に単語の極性が導かれるはずで, コーパスのうち, こうした極性が既知のテキストの集合を  $X_\ell$ , それらへの  $1/0$  のラベルを  $Y_\ell$  とすると, 各テキストとそのラベル  $(y, d) \in (Y_\ell, X_\ell)$  について,

$$(5.114) \quad p(y, d, \theta, \beta) = p(y|\theta) \prod_{v=1}^V p(v|\theta, \beta)^{n_{dv}} \cdot p(\theta)p(\beta)$$

を最大化する  $\beta$  を求めることを考えましょう. ここで第 1 項はロジスティック

回帰

$$(5.115) \quad p(y=1|\theta) = \sigma(\theta) = \frac{1}{1 + e^{-\theta}}$$

です。このグラフィカルモデルを図 5.41 に示しました。 $\theta$  が大きい、あるいは小さい方が第 1 項の教師データに対する識別モデルの尤度が高くなりますが、逆に第 2 項の単語の生成確率が下がる可能性があるため、両者のトレードオフで  $\theta$  が決まり、それによって式 (5.112) から  $\beta$  が決まることになります。

ここで問題なのは、回帰パラメータ  $\beta$  だけでなく、テキストの潜在的な極性  $\theta$  も未知なことです。 $\theta$  および  $\beta$  を同時に最適化することも可能ですが、心理統計学においてこうした同時推定は一致性を持たないことが知られています [185]。また、二値ラベルの  $y$  という弱い教師情報からは  $\theta$  は一意には決まらず、 $\theta$  を学習時に点推定することは教師データへの過学習をもたらす可能性があります。

**適応的ガウス-エルミート求積による解法** そこで、 $\theta$  を推定する代わりにモデルから積分消去し、

$$(5.116) \quad \begin{aligned} p(y, d, \beta) &= \int_{-\infty}^{\infty} p(y, d, \theta, \beta) d\theta \\ &= \int_{-\infty}^{\infty} p(y|\theta) \prod_{v=1}^V p(v|\theta, \beta)^{n_{dv}} p(\theta) d\theta \cdot p(\beta) \end{aligned}$$

を  $\beta$  について最適化することを考えます。<sup>\*63</sup>  $p(\theta)$  は標準正規分布  $\mathcal{N}(0, 1)$  ですから、式 (5.116) の形のガウス分布に関する積分はガウス-エルミート求積<sup>\*64</sup> と呼ばれる方法で、きわめて正確に数値的に求めることができます。

ガウス-エルミート求積では、関数空間での直交多項式を用いて、関数  $f(x)$  の  $e^{-x^2}$  に関する積分を次の形で高精度に近似します。

$$(5.117) \quad \int_{-\infty}^{\infty} f(x) e^{-x^2} dx \simeq \sum_{i=1}^H w_i f(x_i)$$

ここで  $\mathbf{x} = (x_1, \dots, x_H)$  は分点 (abscissa) と呼ばれる座標、 $\mathbf{w} = (w_1, \dots, w_H)$  は

<sup>\*63</sup> Hamiltonian MCMC 法 [186] を用いた  $\beta$  のベイズ推定も検討しましたが、MAP 推定を用いる方が安定した結果となりました。

<sup>\*64</sup> 求積 (quadrature) とは、定積分の値を数値的に求めることです。



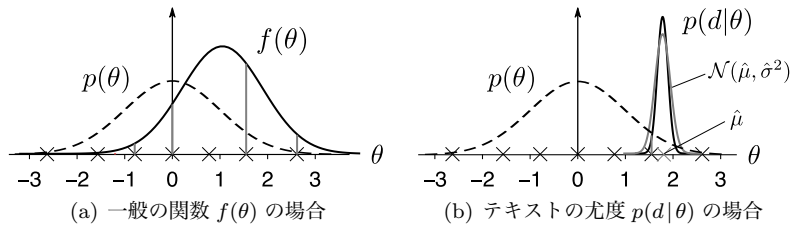


図 5.43: 適応的なガウス-エルミート求積. 一般の関数  $f(\theta)$  の期待値  $\int f(\theta)p(\theta)d\theta$  は,  $\times$ 印で示した分点  $x_i$  での  $f(\theta)$  と重み  $w_i$  から高精度に計算することができますが, テキストの場合は尤度  $p(d|\theta)$  が一部の  $\theta$  に集中するため, 変数変換を行って MAP 解  $\hat{\mu}$  の周辺で適応的なガウス-エルミート求積を行い, 期待値  $\int p(d|\theta)p(\theta)d\theta$  を計算します.

対応する重みで,  $H=9$  のとき (本書の例では  $H=20$  としました) は

$$\begin{aligned} \mathbf{x} &= (-3.191, -2.267, -1.469, -0.724, 0, 0.724, 1.469, 2.267, 3.191), \\ \mathbf{w} &= (0.00004, 0.005, 0.088, 0.433, 0.720, 0.433, 0.088, 0.005, 0.00004) \end{aligned}$$

です. これらの値は, 標準的なガウス-エルミート求積のパッケージで計算することができます.\*65 式(5.117)は  $e^{-x^2}$  についての積分なので,  $e^{-x^2} = e^{-\theta^2/2}$  すなわち  $\theta = \sqrt{2}x$  とおけば, 変数変換により

(5.118)

$$\int_{-\infty}^{\infty} f(\theta)\mathcal{N}(\theta|0,1)d\theta = \int_{-\infty}^{\infty} f(\theta)\frac{1}{\sqrt{2\pi}}e^{-\frac{\theta^2}{2}}d\theta \simeq \frac{1}{\sqrt{\pi}}\sum_{i=1}^H w_i f(\sqrt{2}x_i)$$

と計算することができます.

ただし, 式(5.118)による積分は, そのままでは非常に効率が悪くなります. 一般にテキストは多くの単語を含むため,  $\theta$  の事後分布はある値  $\hat{\theta}$  の近くに集中しており, 式(5.118)ではほとんどの分点での尤度が 0 になってしまうからです. そこで, [187]の方法を用いて, 積分を  $\hat{\theta}$  の周りで実行することにします.  $\theta$  の事後分布を近似する平均  $\mu = \hat{\theta}$  と分散  $\sigma^2$  は二分探索と二階差分により容易に求めることができますので, まず,  $\phi = \mu + \sigma\theta$  と変数変換すると, 簡単な計算により

\*65 Python では `numpy.polynomial.hermite.hermgauss` で, R では `gaussquad` パッケージの `hermite.h.quadrature.rules` で計算できます.

$$(5.119) \quad \int_{-\infty}^{\infty} f(\theta) \mathcal{N}(\theta|\mu, \sigma^2) d\theta \simeq \frac{1}{\sqrt{\pi}} \sum_{i=1}^H w_i f(\mu + \sqrt{2}\sigma x_i)$$

であることがわかります. このとき, 求める積分を次のように変形します.

$$(5.120) \quad I = \int_{-\infty}^{\infty} f(\theta) \mathcal{N}(\theta|0, 1) d\theta = \int_{-\infty}^{\infty} \underbrace{f(\theta) \frac{\mathcal{N}(\theta|0, 1)}{\mathcal{N}(\theta|\mu, \sigma^2)}}_{h(\theta)} \mathcal{N}(\theta|\mu, \sigma^2) d\theta$$

式(5.120)の最初の2項を  $h(\theta)$  とおけば, 式(5.119)より

$$(5.121) \quad I \simeq \frac{1}{\sqrt{\pi}} \sum_{i=1}^H w_i h(\mu + \sqrt{2}\sigma x_i)$$

と,  $\mathcal{N}(\theta|\mu, \sigma^2)$  についての適応的な積分で置き換えて求めることができます.

われわれの場合, 求めたい積分は式(5.116)でしたから, 対数で計算するために

$$(5.122) \quad \ell(\theta) = \log p(y|\theta) + \sum_{v=1}^V n_{dv} \log p(v|\theta, \beta)$$

と定義すれば,  $h(\theta)$  は

$$(5.123) \quad h(\theta) = e^{\ell(\theta)} \frac{\mathcal{N}(\theta|0, 1)}{\mathcal{N}(\theta|\mu, \sigma^2)} = \exp(\ell(\theta) + \log \mathcal{N}(\theta|0, 1) - \log \mathcal{N}(\theta|\mu, \sigma^2))$$

となります. 見やすくするために  $y_i = \mu + \sqrt{2}\sigma x_i$  とおけば,

$$(5.124) \quad p(y, d, \beta) = \int_{-\infty}^{\infty} h(\theta) \mathcal{N}(\theta|\mu, \sigma^2) d\theta \cdot p(\beta) \\ \simeq \frac{\sigma}{\sqrt{\pi}} \sum_{i=1}^H \exp \left[ \log w_i + \ell(y_i) + \frac{1}{2} \left( \frac{1}{\sigma^2} (y_i - \mu)^2 - y_i^2 \right) \right] \\ \cdot p(\beta)$$

となり,  $\ell_i$  の中に  $\beta$  が含まれるこの式の対数を  $\beta$  について偏微分し, L-BFGS法で最適化することで  $\beta$  を計算します. 式(5.124)の対数の偏微分は, 計算すると文書  $d$  の長さを  $L$  として,

$$(5.125) \quad \frac{\partial}{\partial \beta} \log p(y, d, \beta) = L\theta \left[ \frac{1}{L} \sum_{v \in d} \vec{v} - \sum_{v=1}^V p(v|\theta, \beta) \vec{v} \right]$$

という直感的にも妥当な結果となります。

**PLSS の実験 (続き)** それでは、実際のテキストを分析してみましょう。ここでは、国会の議事録を使ってみることにします。国会会議録検索システム<sup>\*66</sup>から、2023年度の第211回通常国会での衆議院・文部科学委員会のテキストをダウンロードします。これは16回の開催で、サポートサイトにあるスクリプトを用いて次のように実行すると、2,210個の発言がlivedoor.txtと同じ「発言者<TAB>発言内容..」の形式で得られます。

```
% diet-split.py *.txt | diet-format.py > all.txt
```

このうち、次のように実行して大臣および委員長の司会・答弁および50文字未満の短い発言を除外すると、984個の単語分割された発言とその発言者が得られました。<sup>\*67</sup>

```
% diet-sieve.awk all.txt | cut -f 2 | mecab -O wakati \
 -b 65536 > education.txt
% diet-sieve.awk all.txt | cut -f 1 > education.label
```

発言の多い上位順に宮本岳志委員(共産党)、<sup>ゆのき</sup>柚木道義委員(立憲民主党)、西岡秀子委員(国民民主党)のうち、違いの自明でない後者二名の発言から次のようにしてランダムに20個の発言を抽出し、これを正例-負例とみて違いの軸を抽出してみることにしましょう。なお、この二者はいずれも野党です。

```
% article-id.py education.label 柚木 20 > yunoki.id
% article-id.py education.label 西岡 20 > nishioka.id
% cat yunoki.id
4,13,149,156,160,162,165,167,234,240,243,250,251,310,311,..
```

この上で、次のようにしてPLSSの半教師あり学習を実行します。内部でplss-semi.Rを呼んでいますので、先にこのスクリプトを単体で実行して、必要なライブラリをインストールしておいてください。

<sup>\*66</sup> <https://kokkai.ndl.go.jp/>

<sup>\*67</sup> この場合のように1行が長いとMeCabの解析バッファが足りなくなるため、-b 65536のようにしてバッファを標準の4KBから64KBに増やしています。

```
% plss-semi education.txt yunoki.id nishioka.id jawiki.vec output
preparing word vectors..
total 8031 words selected.
running PLSS-semi..
loading wordvectors from /tmp/wordvec-38450.vec.. done.
preparing data of language 'ja'.. done.
optimizing beta from examples..
iter 1 value 15.745812
iter 2 value 14.242224
iter 3 value 13.447749
:
iter 29 value 10.210794
final value 10.210794
converged
number of docs = 984
word dimension = 100
norm of beta = 2.25
computing theta..
computing 984/984.. done.
theta written to output.theta.
phi written to output.phi.
```

output.theta および output.phi が推定した各発言の極性  $\theta$  と、学習した  $\beta$  に沿った単語の極性  $\phi_v = \beta^T \vec{v}$  です。次のように実行すると、 $\theta$  および発言者、発言内容を並べて表示し、 $\theta$  でソートすることができます。

```
% paste output.theta education.label education.txt \
| sed 's/ //g' | sort -nr
⇒ 1.400558 柚木委員 築副大臣、お考えについては答えてください。L..
1.369929 柚木委員 そんな、文部科学副大臣として、今おっしゃって..
1.292844 柚木委員 ということは、永岡文部科学大臣としては、LG..
:
-1.292683 西岡委員 大変様々な要件が課されておまして、また今後..
-1.305099 西岡委員 これまでも議論になっておりますけれども、やは..
-1.333402 西岡委員 関連いたしまして、先ほども申し上げましたよう..
```

この全体からランダムに抽出した 20 発言とその極性を図 5.44(a) に示しました。現場感覚を大事にするアクティブな柚木議員に対し、文教族の参議院議員を父親に持つ西岡議員は、より制度的なアプローチをとっていることが読み取れます。この違いは、 $\phi_v$  にも見てとることができます。図 5.44(b) では確かに、柚

| $\theta$ | 発言者  | 発言内容                                |
|----------|------|-------------------------------------|
| 0.9843   | 柚木委員 | まとめて答えていただいたので、最後、ちょっとだけ時間ができたので。.. |
| 0.7876   | 柚木委員 | ちょっと警察庁、この後、答弁いただきます。伊佐副大臣、最後の質問..  |
| 0.7387   | 梅谷委員 | そうですね。じゃ、この組織的あっせん事件の中心的人物が今大学で役..  |
| 0.7260   | 柚木委員 | これは驚きの答弁ですよ、大臣。私、一応、LGBTは種の保存に背く..  |
| 0.6559   | 宮本委員 | 学校給食法第十一条は障害にならない、当たり前です。しかも、自民党..  |
| 0.3464   | 吉川委員 | 指針ですけれども、先ほど触れたとおり、強制力を伴わない、望ましいだ.. |
| 0.3308   | 柚木委員 | もうこの項目は最後にしたいと思いますが、だとすると、まさに十三ペー.. |
| 0.2456   | 柚木委員 | 非常に、今後の方向性がある意味中間報告的に今整理して御答弁いただい.. |
| 0.1865   | 梅谷委員 | 検討中という話ですので、これ以上は申し上げませんが。ただ、私かも..  |
| -0.0026  | 菊田委員 | 大切なことは、やはり、現場の教職員の先生方がすごいプレッシャーとか.. |
| -0.0574  | 宮本委員 | 幼児教育無償化したからやれという話じゃなくて、されなくてももちろん.. |
| -0.1607  | 宮本委員 | おっしゃるとおりで、経緯があつて、だから多様になっているということ.. |
| -0.2457  | 吉川委員 | これから検討し、またパブコメ等ということですが、私は、審議会..    |
| -0.2807  | 牧委員  | 確におっしゃるように、ガバナンスは大切なんですけれども、逆に、萎..  |
| -0.3307  | 宮本委員 | 当然ですね。さらに、私に対応してくださった金沢大学の職員のお一人..  |
| -0.4056  | 鰐淵委員 | ありがとうございます。今御紹介もありましたけれども、当初は評議員..  |
| -0.4151  | 早坂委員 | 先ほど、私も私学出身で、両親が共働きで、どうか学校を出していただ..  |
| -0.5283  | 宮本委員 | 教育機会確保法附則の三には、「この法律の施行後三年以内にこの法律の.. |
| -0.6099  | 西岡委員 | 基準、要件については、これから法案が成立した後、審議会で議論をする.. |
| -1.0058  | 西岡委員 | やはり、支出の妥当性、透明性は大変重要だと思いますので、しっかりそ.. |

(a) 各発言に推定された  $\theta$  (一部)

|         |        |       |        |         |         |          |         |
|---------|--------|-------|--------|---------|---------|----------|---------|
| 勝目      | 3.0700 | 築     | 2.3292 | ワーキング.. | -2.6780 | 学寮       | -2.0325 |
| 辺野古     | 3.0331 | 浮島    | 2.3156 | 産学      | -2.3894 | 防災       | -2.0261 |
| ワールドカ.. | 2.6581 | 恥ずかしい | 2.2712 | 気象      | -2.2831 | 周期       | -2.0104 |
| 利府      | 2.6491 | 飲酒    | 2.2636 | か年      | -2.2318 | 高度       | -2.0074 |
| 有権者     | 2.5233 | わいせつ  | 2.2398 | 凝らし     | -2.2294 | 豊橋技術科学.. | -1.9977 |
| 中体連     | 2.4912 | 英弘    | 2.2019 | 耐震      | -2.2254 | 科目       | -1.9972 |
| 左折      | 2.4738 | 下線    | 2.1374 | 災       | -2.2171 | カリキュラム   | -1.9865 |
| 衆議院     | 2.4640 | 起点    | 2.1102 | 整い      | -2.1866 | 不断       | -1.9793 |
| 双葉      | 2.4262 | 知事    | 2.0699 | 屋根      | -2.1721 | 地震       | -1.9738 |
| 乗車      | 2.3836 | 野球    | 2.0688 | くらし     | -2.1421 | 築い       | -1.9514 |

 $\phi_v > 0$  (柚木議員側) $\phi_v < 0$  (西岡議員側)

(b) 推定された単語の極性

図 5.44: PLSS の半教師あり学習による極性の推定。極性辞書を使わなくても、与えた文書集合のうち、少数の正例と負例の文書の番号を与えれば、極性を表す軸  $\beta$  をそこから学習することができます。

木議員 (+ 側) では「左折」「乗車」「野球」といった単語の極性値が高く、西岡議員 (- 側) では「産学」「耐震」「防災」といった単語の極性値が高くなっており、同じ野党でありながらも、興味の違いや政治姿勢を反映していることがわかります。また他の議員の立場も同じ軸の上で解釈することができ、全体的な構造が一見不明な文部科学委員会での議論に、一定の計量的な視点を与えることができました。

## 5 章の演習問題

- (1) ナイーブベイズ法において、式(??)の平滑化パラメータ  $\alpha_w$  の値を変えると、性能にどのような影響があるでしょうか。これを 0 にすると、どんな問題が生じる可能性があるでしょうか。
- (2) ナイーブベイズ法で学習したモデルについて、各単語のラベル事後確率  $p(y|w)$  および対数オッズ比を計算してみましょう。また、単語の出現確率と対数オッズ比はどのような関係にあるか、プロットしてみましょう。
- (3) Livedoor コーパスの各カテゴリでの単語分布  $p(w|y)$  から、NPMI を使って、各カテゴリの特徴語を計算してみましょう。上位語、下位語はどのようなになっているでしょうか。また、UM で教師なしで得られたカテゴリの特徴語とは、どのような違いがあるでしょうか。
- (4) ナイーブベイズ法で、単語のカテゴリ所属確率  $p(y|w)$  が特定のカテゴリについて高い語には、どんなものがあるでしょうか。  $p(y|w)$  がどのくらい特定の  $y$  に集中しているかは、確率分布のエントロピー (??節) で測ることができます。  $p(y|w)$  のエントロピーの大小で単語をソートすると、どんな結果になるでしょうか。
- (5) ユニグラム混合モデルについて、学習したモデルから各単語  $w$  のトピック事後確率  $p(z|w)$  を計算してみましょう。意外なトピックに事後確率が高い単語はあるでしょうか。
- (6) 任意の単語  $w_1, w_2, w_3$  を 3 つ選んで、語彙がこれらの単語だけからなるとしたとき、適当なコーパスの各文書について最尤推定で計算した確率分布  $\mathbf{p}=(p_1, p_2, p_3)$  を、図 5.14 のように単体上にプロットしてみましょう。意味的に相関のある単語のペアを選んだとき、プロットにはどんな傾向が現れるでしょうか。
- (7) DM で学習されたディリクレ分布のハイパーパラメータ  $\alpha_k$  から、期待値を求めることで、各トピックがどんな意味を持っているかを調べてみましょう。UM と何か違いはあるでしょうか。
- (8) ~で行ったように、適当な小説などのテキストに対して、ある単語が出現

してから、次に出現するまでの時間 (=何単語後か) をプロットしてみましよう。地震のようなイベントがランダムな時間に起きる様子を記述するポアソン過程の言葉では、これは前方再起時間 (forward recurrence time) とよばれています。単語によって、どんな特徴があるでしょうか。この分布をモデル化するには、どうするとよいでしょうか。

- (9) 文書による単語の長さの分布の違い。
- (10) 文書の特徴を文字あるいは文字の種類で考えてみる。
- (11) Livedoor コーパスについて、LDA の  $\theta$  と教師ラベルとの対応。ラベルごとに、どんなトピックが多いのかをプロットしてみましよう。教師ラベルとしては、どんな話題が足りなかったといえるでしょうか。
- (12) ICA の各軸の小さい方を調べてみる。(iPhone<sub>j</sub>-<sub>i</sub>Android, 国内映画<sub>j</sub>-<sub>i</sub>海外映画等)
- (13) 3章で前後の共起窓から求めた単語ベクトルについても、同様に ICA で変換した独立成分軸を求めてみましよう。本章で文書内全部を共起の対象とした単語ベクトルと、何か違いがあるでしょうか。また、Word2Vec, GloVe, SVD による単語ベクトルの計算法で何か差がみられるでしょうか。
- (14) 同じコーパスでトピックモデルと文書ベクトルを学習したとき、学習された表 5.9 のようなトピックと、表 5.13 のような ICA の独立成分軸はどのように異なるでしょうか。どのトピックと軸が似ており、どれが似ていないかを定量化するには、どうすればいいでしょうか？
- (15) PLSS で、positive-negative 以外の辞書を使ってテキストを分析してみる。

## 5章の文献案内

- [156]:  $\Gamma$  関数の山は不要
- LDA のテクニカルレポート
- Embedded topic model.
- 統数研公開講座資料.
- 教師ありトピックモデル. (本文)
- BDA について. (もっと前?)



## コラム：テキストを処理するための言語

- Python, R, Julia
- R とそれ以外の言語の速度比較 (テキスト処理で)
- パッケージの存在.
  - RMeCab
  - MeCab-Python    Ginza, KyTea
  - Quanteda
- `awk`, `sed` の紹介 ← `awk` の簡単な~, `sedlec`

# 付録

## A ディリクレ分布の積分と期待値

## B ディリクレ分布の $\alpha$ のベイズ推定

ポリア分布の式

$$(0.126) \quad p(\mathbf{n}|\boldsymbol{\alpha}) = \underbrace{\frac{\Gamma(\sum_k \alpha_k)}{\Gamma(N + \sum_k \alpha_k)}}_{(A)} \underbrace{\prod_{k=1}^K \frac{\Gamma(\alpha_k + n_k)}{\Gamma(\alpha_k)}}_{(B)}$$

は  $\alpha_k$  がガンマ関数  $\Gamma(\cdot)$  の中に入っているため、そのままではサンプリングできる形になりません。しかし、巧妙な補助変数  $\theta, x$  を導入すると、ガンマ事後分布からサンプリングすることができます[73, Appendix C].\*68

(A) ディリクレ分布で2次元の場合を考えると、**ベータ関数**

$$(0.127) \quad B(a, b) = \int_0^1 \theta^{a-1} (1-\theta)^{b-1} d\theta = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$

が得られます。よって、

$$(0.128) \quad B(\sum_k \alpha_k, N) = \frac{\Gamma(\sum_k \alpha_k)\Gamma(N)}{\Gamma(\sum_k \alpha_k + N)}$$

ですから、(A) は補助変数  $\theta$  を積分消去した形で、

$$(0.129) \quad \frac{\Gamma(\sum_k \alpha_k)}{\Gamma(N + \sum_k \alpha_k)} = \frac{B(\sum_k \alpha_k, N)}{\Gamma(N)} = \frac{1}{\Gamma(N)} \int_0^1 \theta^{\sum_k \alpha_k - 1} (1-\theta)^{N-1} d\theta$$

\*68 これは MCMC 法の基本テクニックの一つで、**補助変数法**とよばれています[188].

と表すことができます.

(B) 115 ページのコラムでみたように, Pochhammer 関数  $\Gamma(\alpha+n)/\Gamma(\alpha)$  は

$$(0.130) \quad \frac{\Gamma(\alpha+n)}{\Gamma(\alpha)} = \alpha(\alpha+1)\cdots(\alpha+n-1) = \prod_{i=0}^{n-1} (\alpha+i)$$

と表すことができますが, これも  $\{0,1\}$  の補助変数  $x_{ki}$  を周辺化した形で

$$(0.131) \quad \frac{\Gamma(\alpha_k+n_k)}{\Gamma(\alpha_k)} = \prod_{i=0}^{n_k-1} (\alpha_k+i) = \prod_{i=0}^{n_k-1} \sum_{x_{ki} \in \{0,1\}} \alpha_k^{x_{ki}} i^{1-x_{ki}}$$

と書くことができることに注意しましょう.

よって,  $\alpha_k$  がガンマ事前分布

$$(0.132) \quad p(\alpha_k) = \text{Ga}(a, b) = \frac{b^a}{\Gamma(a)} \alpha_k^{a-1} e^{-b\alpha_k}$$

に従っているとすると, 式(0.129), 式(0.131)より,  $\alpha_k$  の事後分布は, ベイズの定理から

$$(0.133) \quad p(\boldsymbol{\alpha}|\mathbf{n}) \propto p(\mathbf{n}|\boldsymbol{\alpha})p(\boldsymbol{\alpha}) = \prod_{k=1}^K \frac{b^a}{\Gamma(a)} \alpha_k^{a-1} e^{-b\alpha_k} \times \frac{1}{\Gamma(N)} \int_0^1 \theta^{\sum_k \alpha_k - 1} (1-\theta)^{N-1} d\theta \times \prod_{k=1}^K \left( \prod_{i=0}^{n_k-1} \sum_{x_{ki} \in \{0,1\}} \alpha_k^{x_{ki}} i^{1-x_{ki}} \right)$$

となります. これは  $\theta$  と  $x_{ki}$  を周辺化した形ですから, 同時確率は

$$(0.134) \quad p(\boldsymbol{\alpha}, \theta, x|\mathbf{n}) \propto \prod_{k=1}^K \frac{b^a}{\Gamma(a)} \alpha_k^{a-1} e^{-b\alpha_k} \times \frac{1}{\Gamma(N)} \theta^{\sum_k \alpha_k - 1} (1-\theta)^{N-1} \times \prod_{k=1}^K \left( \prod_{i=0}^{n_k-1} \alpha_k^{x_{ki}} i^{1-x_{ki}} \right)$$

です. よって各  $\alpha_k$  の事後分布は, 式(0.134)から  $\alpha_k$  に関する項を抜き出せば, 補助変数  $\theta, x_{ki}$  が与えられた下では

$$(0.135) \quad p(\alpha_k|\mathbf{n}, \theta, x) \propto p(\alpha_k, \theta, x|\mathbf{n}) = \alpha_k^{a-1} \cdot e^{-b\alpha_k} \cdot \theta^{\alpha_k} \cdot \prod_{i=0}^{n_k-1} \alpha_k^{x_{ki}} \\ = \alpha_k^{a+\sum_{i=0}^{n_k-1} x_{ki}-1} e^{-(b-\log \theta)\alpha_k}$$

$$\sim \text{Ga}(a + \sum_{i=0}^{n_k-1} x_{ki}, b - \log \theta)$$

となることがわかります. 同様にして補助変数  $\theta, x_{ki}$  についても対応する項を抜き出すと, その分布は

$$(0.136) \quad p(\theta | \mathbf{n}, \alpha, y) \sim \text{Be}(\sum_k \alpha_k, N)$$

$$(0.137) \quad p(x_{ki} | \mathbf{n}, \alpha, \theta) \propto \alpha_k^{x_{ki}} i^{1-x_{ki}} \propto \left( \frac{\alpha_k}{\alpha_k + i} \right)^{x_{ki}} \left( \frac{i}{\alpha_k + i} \right)^{1-x_{ki}} \\ \sim \text{Bernoulli} \left( \frac{\alpha_k}{\alpha_k + i} \right)$$

となります. この分布から補助変数  $\theta, x_{ki}$  をサンプリングしてから, 式(0.135)を使って  $\alpha_k$  をサンプリングします. 実際には, 式(0.126)の  $p(\mathbf{n} | \alpha)$  は  $N$  個の文書について存在して, 尤度は式(5.78)になっていますから, 同様に計算すると,  $T_i$  を文書  $i$  の長さとして

$$(0.138) \quad \begin{cases} p(\theta_i | n, \alpha, x) \sim \text{Be}(\sum_k \alpha_k, T_i) \\ p(x_{ikj} | n, \alpha, \theta) \sim \text{Bernoulli} \left( \frac{\alpha_k}{\alpha_k + j} \right) \\ p(\alpha_k | n, \theta, x) \sim \text{Ga} \left( a + \sum_{i=1}^N \sum_{j=0}^{n(i,k)-1} x_{ikj}, b - \sum_{i=1}^N \log \theta_i \right) \end{cases}$$

#### (ディリクレ分布の $\alpha$ のサンプリングの公式)

で,  $\alpha$  をガンマ事後分布からサンプリングすることができます.  $\square$

$\eta$  についても, 式(5.78)から同様にして計算すれば, 補助変数  $\phi_k, y_{kvj}$  を導入することで

$$(0.139) \quad \begin{cases} p(\phi_k | m, \eta, y) \sim \text{Be}(V\eta, \sum_{v=1}^V m(k, v)) \\ p(y_{kvj} | m, \eta, \phi) \sim \text{Bernoulli} \left( \frac{\eta}{\eta + j} \right) \\ p(\eta | m, \phi, y) \sim \text{Ga} \left( a + \sum_{k=1}^K \sum_{v=1}^V \sum_{j=0}^{m(k,v)-1} y_{kvj}, b - \sum_{k=1}^K \log \phi_k \right) \end{cases}$$

#### (ディリクレ分布の $\eta$ のサンプリングの公式)

で, ガンマ事後分布からサンプリングを行えることがわかります.  $\square$

C Kneser-Ney 平滑化の導出

D Jensen の不等式

## アルゴリズムの一覧

|                                       |     |
|---------------------------------------|-----|
| 1 HMM の Gibbs サンプリング . . . . .        | 209 |
| 2 HMM の周辺化 Gibbs サンプリング . . . . .     | 215 |
| 3 UM の EM アルゴリズム . . . . .            | 239 |
| 4 UM の周辺化 Gibbs サンプリング . . . . .      | 253 |
| 5 DM の EM-Newton アルゴリズム . . . . .     | 259 |
| 6 LDA の周辺化 Gibbs サンプリング . . . . .     | 270 |
| 7 文書ベクトル (DocVec) の計算アルゴリズム . . . . . | 292 |

## 公式の一覧

|                                       |    |
|---------------------------------------|----|
| (2.0) 確率の最尤推定 . . . . .               | 18 |
| (2.9) 同時確率の周辺化の公式 . . . . .           | 24 |
| (2.11) 同時確率の周辺化の公式 (連続値の場合) . . . . . | 26 |
| (2.18) 条件つき確率の公式 . . . . .            | 28 |
| (2.19) 確率の連鎖則 . . . . .               | 30 |
| (2.20) 独立な事象の同時確率 . . . . .           | 30 |
| (2.34) ベイズの定理 (同時確率版) . . . . .       | 37 |

|                                 |     |
|---------------------------------|-----|
| (2.35) ベイズの定理 (比例版)             | 37  |
| (2.62) 自己情報量                    | 59  |
| (2.68) パープレキシティ                 | 62  |
| (2.69) KL ダイバージェンスの非負性          | 63  |
| (2.75) クロスエントロピー                | 65  |
| (2.76) クロスエントロピーとエントロピー         | 65  |
| (3.0) Heaps の法則                 | 84  |
| (3.10) Zipf の法則                 | 88  |
| (3.30) ディリクレ分布 (簡易版)            | 103 |
| (3.31) ディリクレ分布の期待値              | 104 |
| (3.34) ディリクレ分布 (正式版)            | 105 |
| (3.36) ガンマ関数                    | 105 |
| (3.46) ディリクレ事後分布                | 109 |
| (3.48) ディリクレ平滑化                 | 110 |
| (3.58) ポリア分布                    | 113 |
| (3.60) ポリア分布の最適化の公式             | 114 |
| (3.71) 絶対平滑化                    | 124 |
| (3.73) Kneser–Ney 平滑化           | 126 |
| (3.85) シグモイド関数                  | 135 |
| (3.86) $\tanh$ とシグモイド関数の関係      | 136 |
| (3.88) シグモイド関数の補関数              | 136 |
| (4.18) HMM の同時確率                | 190 |
| (4.67) HMM の周辺化 Gibbs サンプリングの公式 | 215 |
| (5.4) ナイーブベイズ法の文書確率             | 229 |
| (5.11) $\text{logsumexp}$       | 236 |
| (5.23) 正規化自己相互情報量               | 245 |
| (5.33) UM の文書確率                 | 249 |
| (5.51) DM の文書確率                 | 260 |
| (5.55) 指数分布族 DCM 分布             | 263 |
| (5.58) LDA の文書同時確率              | 266 |

## 公式の一覧

327

|                                                |     |
|------------------------------------------------|-----|
| (5.75) LDA の周辺化 Gibbs サンプルングの公式 . . . . .      | 271 |
| (0.138) ディリクレ分布の $\alpha$ のサンプルングの公式 . . . . . | 323 |
| (0.139) ディリクレ分布の $\eta$ のサンプルングの公式 . . . . .   | 323 |



# 索引

## Symbols

$\mathbb{I}()$ , 204

$\infty$  グラム言語モデル, 44

$p()$ , 19

$n$  グラム, 44

$n$  グラム言語モデル, 44

0 グラム, 40

1 グラム, 20

2 グラム, 20

3 グラム, 44

4 グラム, 128

## A

All-but-the-top, 177

awk, 41, 49, 67, 71

## B

bag of words, 224

Baum–Welch アルゴリズム, 198, 201

BERT, 231

BFRY 過程, 89

bit, 59

## C

CBOW, → 連続的単語集合, 142

CCG, 180

CKY アルゴリズム, 185

CRF, 13, 79, 186

Cython, 272

## D

DCM 分布, → デイリクレ複合多項分布

DM, → デイリクレ混合モデル

DNA, ii, 128, 223

Doc2Vec, 287

DocVec, 288

double power-law, 89

## E

EDCM 分布, 263

EM アルゴリズム, 187, 201, 220, 246, 266, 301, 303

## G

Gibbs サンプルング, 198, 201, 252, 266, 267

**H**

hapax legomenon, → 孤語

Heaps の法則, 84

HMM, 188

**I**

ICA, 296

**J**

Jensen-Shannon ダイバージェンス, 64

JUMAN, 80, 185, 193

**K**

KL ダイバージェンス, → Kullback-Leibler ダイバージェンス, 166

Kneser-Ney 平滑化, 121, 125

Kullback-Leibler ダイバージェンス, 63

K 平均法, 240

**L**

LDA, 258, → 潜在ディリクレ配分法

LSI, 286, 287

**M**

MAP 推定, 311

MCMC 法, 184, 303

MeCab, 79, 185, 187, 226, 232, 314

mecab-ipadic-NEologd, 91

Moore-Penrose の擬似逆行列, 291

**N**

nat, 59

Normalized PMI, 244

NPMI, → 正規化自己相互情報量, 280

NPYCRF, 188

NPYLM, 46, 81

**P**

Pólya 分布, → ポリア分布

PCFG, 180

Perl, 41, 71

Pitman-Yor 過程, 258

PLSI, 265, 286

PLSS, 304

PMI, → 自己相互情報量

Pochhammer 関数, 115

Pointwise Mutual Information, 244

**Q**

quanteda, v, 8, 90, 308

Q 関数, 249

**R**

Rao-Blackwell 化, 216, 275

RNN, → 再帰的ニューラルネットワーク

**S**

sed, 71, 90

SIF, 176

Softmax 関数, 134

spaCy, 90

SVD, → 特異値分解

- SVM, 300  
SVMlight 形式, 118, 224
- T**  
tf.idf, 286, 293, 294
- U**  
Unigram Mixtures, 240  
uSIF, 177
- V**  
VAE, 246  
Viterbi アルゴリズム, 194, 207, 216
- W**  
Word2Vec, 137
- Z**  
Zipf の法則, 85, 88, 145
- あ**  
圧縮, 121  
アノテーション, 189  
アンダーフィット, 68  
意味, 182  
ウェーバー＝フェヒナーの法則, 294  
埋め込み文, 180  
エントロピー, 60  
オーバーフィット, 67  
音楽, ii
- か**  
階層構造, 4  
階層ディリクレ過程, 121, 283  
階層ディリクレ言語モデル, 117  
階層 Pitman-Yor 過程, 121, 124  
階層ベイズ, 261  
開発データ, 52  
ガウス-エルミート求積, 311  
ガウス過程, 10  
過学習, 67, 217  
係り受け解析, 180, 182  
学習データ, 52  
確率の連鎖則, 30  
確率分布, 20, 101, 264  
確率変数, 23  
確率モデル, 101  
隠れマルコフモデル, 188  
加算平滑化, 46, 99  
過少適合, 68  
カルマンフィルタ, 10  
感情分析, 232  
機械学習, 4  
幾何平均, 58  
期待値, 210  
期待値伝搬法, 266  
機能語, 174  
逆温度, 131  
逆文書頻度, 294  
教師あり学習, i, 4, 14, 171, 181, 186,  
193  
教師なし学習, i, 5, 171, 181, 183, 185,  
187, 198

- 教師なし形態素解析, 79, 81, 188  
教師なし構文解析, 185  
教師なし談話構造解析, 189  
協調フィルタリング, ii  
共役, 109  
行列分解, 285, 286  
局所解, 201, 252  
極性, 232  
均衡コーパス, 82  
グラフィカルモデル, 112  
クロスエントロピー, 65  
クロスバリデーション, 53, 69  
訓練データ, 52  
経験ベイズ法, 114  
形態素解析, 79  
形態論, 3  
結合確率, 23  
言語学, 4, 71, 174, 185, 186, 222  
言語資源, 183  
検証データ, 52  
交差検証, 53  
構文, 180  
項目反応理論, 302  
コーパス, 82  
コーパス言語学, 86  
孤語, 86  
誤差逆伝搬法, 134  
固有表現認識, 91  
混合モデル, 100, 175
- 再帰的ニューラルネットワーク, 135  
最尤推定, 20, 27, 217, 246  
算術符号, 44  
次元の呪い, 132, 143, 185, 194  
事後確率, 37  
自己教師あり学習, 6  
自己情報量, 59  
自己相互情報量, 10, 94, 149, → Point-wise Mutual Information  
指示関数, 204, 266  
地震, 318  
指数分布族, 263  
事前確率, 38  
自然言語処理, 4  
尺度化, 298  
自由エネルギー, 248  
周辺化, 24  
周辺化 Gibbs サンプルング, 215, 253  
周辺確率, 22  
主成分分析, 7, 276  
出力確率, 190  
条件つき確率, 27  
状態遷移確率, 190  
情報量, 9, 58, 197  
情報理論, 44, 50, 58, 76, 79, 94, 121, 194  
深層学習, 97, 128, 143, 149  
心理学, 295, 298  
心理統計学, 302  
スキップグラム, 137, 144
- さ

スコア, 197  
 ストップワード, 242  
 スパース, 224  
 スムージング, → 平滑化  
 正規化自己相互情報量, 94, → Normalized PMI  
 政治学方法論, 8, 300, 306  
 生成モデル, 111, 188  
 絶対平滑化, 121, 124  
 絶対割引, 122  
 ゼロ頻度問題, 45, 99, 132  
 潜在ディリクレ配分法, 217, 265, 276  
 潜在変数, 5, 237, 247  
 測度論, iv, 117, 121

**た**

ダイガンマ関数, 114, 260  
 対数尤度比, 152  
 タガー, 186  
 多項分布, 102, 190, 240, 302  
 単語集合, 142, → bag of words, 229  
 単語単体, 257, 276  
 単語分割, 79  
 単語-文書行列, 287  
 単語ベクトル, 132  
 単体, 102, 257  
 談話構造, 189  
 中心化, 161  
 中心極限定理, 297  
 チューリング完全, 72  
 長大語, 80

チョムスキー標準形, 184  
 低資源, 183  
 デリクレ混合モデル, 258  
 デリクレ複合多項分布, → ポリア分布  
 デリクレ分布, 103, 258, 264  
 デリクレ平滑化, 110, 122  
 データスパースネス, 132  
 テキストマイニング, i  
 テストデータ, 56  
 統計力学, 131, 144  
 同時確率, 23  
 動的計画法, 194  
 特異値分解, 150  
 独立, 31  
 トピック, 242  
 トピック単体, 276  
 トピックモデル, 242  
 トライグラム, → 3 グラム, 47, 120  
 トレリス, 193

**な**

ナイーブベイズ法, 230  
 内容語, 174  
 二重分節, 4  
 ニューラル  $n$  グラム言語モデル, 132

**は**

パーザー, 181  
 バースト性, 113, 255  
 パープレキシティ, 61

- バイオインフォマティクス, 258, 269  
バイグラム, → 2 グラム, 40, 99, 116,  
124, 203  
白色化, 162, 164, 296  
パス, 193  
バスケット分析, ii  
バックオフ, 120  
バリマックス回転, 295  
汎化性能, 68  
半教師あり学習, 6, 183, 188  
非負値行列因子分解, 285  
品詞, 79, 185  
品詞分析, 5  
符号長, 60, 197  
負担率, 237, 252  
仏教学, 8  
物理学, 154, 201, 216, 220, 248  
負例サンプリング, 144  
分岐数, 58, 61  
文境界認識, 170  
分散, 210  
分散表現, 132  
文書, 222  
文書-単語行列, 224  
文書頻度, 281, 293  
分配関数, 144  
文分割, 170  
文ベクトル, 172  
文法, 180  
文脈, 186  
分類器, 300  
平滑化, 10, 46, 145  
平均分岐数, 61  
ベイジアン, 34  
ベイズ学習, 252  
ベイズ推定, 246, 311  
ベイズの定理, 34  
ベータ関数, 321  
巾乗則, 88  
ベン図, 31  
変分ベイズ法, 123, 266  
ポアソン過程, 318  
補助変数法, 321  
ポリア分布, 113, 260, 262, 272, 278,  
321
- ま**  
マスク化言語モデル, 143  
マルコフ確率場, 143  
マルコフ性, 50, 188, 208  
マルコフブランケット, 207  
マルコフモデル, 50  
無限隠れマルコフモデル, 217  
文字列, 16, 38  
モンテカルロ EM アルゴリズム, 252,  
278  
モンテカルロ積分, 65
- や**  
尤度関数, 38  
ユニグラム, → 1 グラム, 39, 162, 224

ユニグラム混合モデル, → Unigram

Mixtures

用語頻度, 294

予測, 55

## ら

ラグランジュの未定乗数法, 20, 250

ラベルスイッチング, 216

離散データ, ii

理想点, 301

レジスター, 222

連続的単語集合, 137, 142

## わ

歪度, 297

## 参考文献

- [1] Chihiro Shibata, Kei Uchiumi, and Daichi Mochihashi. How LSTM Encodes Syntax: Exploring Context Vectors and Semi-Quantization on Natural Text. In *COLING 2020*, pages 4033–4043, 2020.
- [2] André Martinet. *Éléments de linguistique générale, 5th ed.* Armand Colin: Paris, 2008.
- [3] Jun Suzuki, Akinori Fujino, and Hideki Isozaki. Semi-Supervised Structured Output Learning Based on a Hybrid Generative and Discriminative Approach. In *Proceedings of EMNLP-CoNLL 2007*, pages 791–800, 2007.
- [4] Jun Suzuki and Hideki Isozaki. Semi-Supervised Sequential Labeling and Segmentation Using Giga-Word Scale Unlabeled Data. In *ACL:HLT 2008*, pages 665–673, 2008.
- [5] Hinrich Schütze. Dimensions of Meaning. In *Proceedings of Supercomputing'92*, pages 787–796, 1992.
- [6] Christos H.Papadimitriou, Prabhakar Raghavan, Hisao Tamaki, and Santosh Vempalad. Latent Semantic Indexing: A Probabilistic Analysis. *Journal of Computer and System Sciences*, 61:217–235, 2000.
- [7] 石井公成. 仏教学における N-gram の活用. In **東京大学東洋文化研究所附属東洋学情報センター報「明日の東洋学」**, number 8, pages 2–4, 2002.
- [8] 持橋大地, 大羽成征. **ガウス過程と機械学習**. 機械学習プロフェッショナルシリーズ. 講談社, 2019.
- [9] 石野亜耶, 小早川健, 坂地泰紀, 嶋田和孝, 吉田光男. **Python ではじめるテキストアナリティクス入門**. 実践 Data Science シリーズ. 講談社サイエンティフィク, 2022.
- [10] 高村大也. **言語処理のための機械学習入門**. 自然言語処理シリーズ. コロナ社, 2010.
- [11] 佐藤坦. **はじめての確率論 測度から確率へ**. 共立出版, 1994.
- [12] 清水泰隆. **統計学への確率論, その先へ: ゼロからの測度論的理解と漸近理論への架け橋**. 内田老鶴圃, 2021.
- [13] Pierre-Simon Laplace. *Essai Philosophique sur les Probabilités*. 1814. 『確率の哲学的試論』内井惣七 (訳), 岩波文庫, 1997.
- [14] Alastair J. Walker. An Efficient Method for Generating Discrete Random Variables with General Distributions. *ACM Transactions on Mathematical Software*, 31(3):253–256, 1977.



- [15] Feras Saad, Cameron Freer, Martin Rinard, and Vikash Mansinghka. The Fast Loaded Dice Roller: A Near-Optimal Exact Sampler for Discrete Probability Distributions. In *AISTATS 2020*, pages 1036–1046, 2020.
- [16] David J Ward, Alan F Blackwell, and David J. C. MacKay. Dasher - a Data Entry Interface Using Continuous Gestures and Language Models. In *UIST 2000*, 2000.
- [17] Li Du, Lucas Torroba Hennigen, Tiago Pimentel, Clara Meister, Jason Eisner, and Ryan Cotterell. A Measure-Theoretic Characterization of Tight Language Models. In *ACL 2023*, pages 9744–9770, 2023.
- [18] Daichi Mochihashi and Eiichiro Sumita. The Infinite Markov Model. In *Advances in Neural Information Processing Systems 20 (NIPS 2007)*, pages 1017–1024, 2008.
- [19] 持橋大地, 隅田英一郎. 階層 Pitman-Yor 過程に基づく可変長 n-gram 言語モデル. *情報処理学会論文誌*, 48(12):4023–4032, 2007.
- [20] Daichi Mochihashi, Takeshi Yamada, and Naonori Ueda. Bayesian Unsupervised Word Segmentation with Nested Pitman-Yor Language Modeling. In *Proceedings of ACL-IJCNLP 2009*, pages 100–108, 2009.
- [21] Andrei A. Markov. An Example of Statistical Investigation of the Text Eugene Onegin Concerning the Connection of Samples in Chains. In *Proceedings of the Academy of Sciences, St. Petersburg*, volume 7, pages 153–162, 1913.
- [22] Roman Jakobson. *一般言語学*. みすず書房, 1973.
- [23] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Information Science and Statistics. Springer, 2007.
- [24] 金子弘昌. *Python で学ぶ実験計画法入門 ベイズ最適化によるデータ解析*. KS 情報科学専門書. 講談社, 2021.
- [25] 韓太舜, 小林欣吾. *情報と符号化の数理*. 培風館, 1999.
- [26] David J. C. MacKay. *Information Theory, Inference, and Learning Algorithms*. Cambridge University Press, 2003.
- [27] Jihyeon Roh, Sang-Hoon Oh, and Soo-Young Lee. Unigram-Normalized Perplexity as a Language Model Performance Measure with Different Vocabulary Sizes. In *arXiv preprint*, 2020. <https://arxiv.org/abs/2011.13220>.
- [28] S. Kullback and R. A. Leibler. On Information and Sufficiency. *Annals of Mathematical Statistics*, 22(1):79–86, 1951.
- [29] Christophe Andrieu, Nando de Freitas, Arnaud Doucet, and Michael I. Jordan. An Introduction to MCMC for Machine Learning. *Machine Learning*, 50:5–43, 2003.
- [30] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory, 2nd Edition*. Wiley Series in Telecommunications. Wiley-Interscience, 2013.
- [31] C. E. Shannon. A mathematical theory of communication. *Bell System Technical Journal*, 27:379–423, 623–656, 1948.
- [32] Tom Brown, Benjamin Mann, Nick Ryder, et al. Language Models are Few-

- Shot Learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901, 2020.
- [33] Larry Wall, Tom Christiansen, and Jon Orwant. **プログラミング Perl 第 3 版**. オライリー・ジャパン, 2002.
- [34] 山口和紀, 古瀬一隆. **新 The UNIX Super Text (上)(下) 改訂増補版**. 技術評論社, 2003.
- [35] Dale Dougherty and Arnold Robbins. **sed & awk プログラミング 改訂版**. A nutshell handbook. O'Reilly Japan, 1997.
- [36] A.V. エイホ, B.W. カーニハン, P.J. ワインバーガー. **プログラミング言語 AWK**. アジソン ウェスレイ・トッパン 情報科学シリーズ. トッパン, 1989.
- [37] Andrew Marc Greene. BASIX: An interpreter written in TEX. *TUGboat*, 11(3):381–392, 1990. <https://www.ctan.org/tex-archive/macros/generic/basix>.
- [38] Christopher D. Manning and Hinrich Schütze. *Foundations of Statistical Natural Language Processing*. MIT Press, 1999.
- [39] Christopher D. Manning and Hinrich Schütze. **統計的自然言語処理の基礎**. 共立出版, 2017.
- [40] 北研二. **確率的言語モデル**. 言語と計算 4. 東京大学出版会, 1999.
- [41] C.M. Bishop, 元田, 栗田, 樋口, 松本, and 村田 (監訳). **パターン認識と機械学習: ベイズ理論による統計的予測 (上)(下) (Pattern Recognition and Machine Learning)**. Springer/丸善出版, 2007, 2008.
- [42] Thomas M. Cover and Joy A. Thomas. **情報理論 — 基礎と広がり**. 共立出版, 2012.
- [43] 伊庭幸人. 「情報」に関する 13 章 – 私家版・情報学入門 –. *物性研究*, 78(2):172–193, 2002. <https://www.ism.ac.jp/~iba/a19.pdf>.
- [44] 近藤泰弘, 近藤みゆき. 平安時代古典語古典文学研究のための N-gram を用いた解析手法. In **言語情報処理学会第 7 回年次大会 発表論文集**, 2001.
- [45] Folgert Karsdorp, Mike Kestemont, and Allen Riddell. *Humanities Data Analysis: Case Studies with Python*. Princeton University Press, 2021.
- [46] 持橋大地, 山田武士, 上田修功. ベイズ階層言語モデルによる教師なし形態素解析. **情報処理学会研究報告 2009-NL-190**, 2009.
- [47] Linguistic Data Consortium. Web 1T 5-gram Version 1, 2006. <https://catalog.ldc.upenn.edu/LDC2006T13>.
- [48] 言語資源協会. Web 日本語 N グラム 第 1 版, 2007. <https://www.gsk.or.jp/catalog/gsk2007-c>.
- [49] Henry Kučera and W. Nelson Francis. *Computational Analysis of Present-Day American English*. Brown University Press, 1967.
- [50] Dirk Goldhahn, Thomas Eckart, and Uwe Quasthoff. Building Large Monolingual Dictionaries at the Leipzig Corpora Collection: From 100 to 200 Languages. In *LREC 2012*, pages 759–765, 2012.
- [51] Harold Stanley Heaps. *Information Retrieval, Computational and Theoretical Aspects*. Academic Press, 1978.

- [52] Gustav Herdan. *Type-Token Mathematics: A Textbook of Mathematical Linguistics*. Mouton en Company, 1960.
- [53] George Kingsley Zipf. *The Psycho-Biology of Language: An Introduction to Dynamic Philology*. The MIT Press Classics Series. The MIT Press, 1935.
- [54] Micheline Petruszewycz. L’histoire de la loi d’Estoup-Zipf : documents. *Mathématiques et Sciences humaines, Tome*, 44:41–56, 1973.
- [55] 水谷静夫. **数理言語学**. 現代数学レクチャーズ D-3. 培風館, 1982.
- [56] Alain Lelu. Jean-Baptiste Estoup and the origins of Zipf’s law: a stenographer with a scientific mind (1868-1950). *Boletín de Estadística e Investigación Operativa*, 30(1):66–77, 2014.
- [57] R. Albert and A.-L. Barabási. Statistical mechanics of complex networks. *Reviews of Modern Physics*, 74:47–97, 2002.
- [58] 田中久美子. **言語とフラクタル 使用の集積の中にある偶然と必然**. 東京大学出版会, 2021.
- [59] S. Naranan and V.K. Balasubrahmanyam. Models for Power Law Relations in Linguistics and Information Science. *Journal of Quantitative Linguistics*, 5(1-2):35–61, 1998.
- [60] Martin Gerlach and Eduardo G. Altmann. Stochastic model for the vocabulary growth in natural languages. *Physical Review X*, 3:021006, 2013.
- [61] Fadhel Ayed, Juho Lee, and Francois Caron. Beyond the Chinese Restaurant and Pitman-Yor processes: Statistical Models with double power-law behavior. In *ICML 2019*, pages 395–404, 2019.
- [62] 佐藤敏紀, 橋本泰一, 奥村学. 単語分かち書き用辞書生成システム NEologd の運用 一文書分類を例にして 一. In **情報処理学会 自然言語処理研究会研究報告**, pages NL-229–15, 2016.
- [63] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed Representations of Words and Phrases and their Compositionality. In *Advances in Neural Information Processing Systems 26 (NIPS 2013)*, pages 3111–3119, 2013.
- [64] Peter F. Brown, Vincent J. Della Pietra, Peter V. deSouza, Jenifer C. Lai, and Robert L. Mercer. Class-Based n-gram Models of Natural Language. *Computational Linguistics*, 18(4):467–479, 1992.
- [65] Kenneth Ward Church and Patrick Hanks. Word Association Norms, Mutual Information, and Lexicography. In *ACL 1989*, pages 76–83, 1989.
- [66] Gerlof Bouma. Normalized (Pointwise) Mutual Information in Collocation Extraction. In *Proceedings of GSCL*, pages 31–40, 2009.
- [67] Rameshwar D. Gupta and Donald St. P. Richards. The History of the Dirichlet and Liouville Distributions. *International Statistical Review*, 69(3):433–446, 2001.
- [68] Luc Devroye. *Non-Uniform Random Variate Generation*. Springer-Verlag, 1986. <http://www.nrbook.com/devroye/>.
- [69] David J. C. MacKay and L. Peto. A Hierarchical Dirichlet Language Model. *Natural Language Engineering*, 1(3):1–19, 1994.

- [70] Thomas P. Minka. Estimating a Dirichlet distribution, 2000. <http://research.microsoft.com/~minka/papers/dirichlet/>.
- [71] 佐藤一誠. トピックモデルによる統計的潜在意味解析. 自然言語処理シリーズ 8. コロナ社, 2015.
- [72] Yee Whye Teh, Michael I. Jordan, Matthew J. Beal, and David M. Blei. Hierarchical Dirichlet Processes. *JASA*, 101(476):1566–1581, 2006.
- [73] Yee Whye Teh. A Bayesian Interpretation of Interpolated Kneser-Ney. Technical Report TRA2/06, School of Computing, National University of Singapore, 2006.
- [74] Phil Cowans. *Probabilistic Document Modelling*. PhD thesis, University of Cambridge, 2006. <http://www.inference.phy.cam.ac.uk/pjc51/thesis/index.html>.
- [75] Reinhard Kneser and Hermann Ney. Improved backing-off for m-gram language modeling. In *Proceedings of ICASSP*, volume 1, pages 181–184, 1995.
- [76] Graham Neubig and Chris Dyer. Generalizing and Hybridizing Count-based and Neural Language Models. In *EMNLP 2016*, pages 1163–1172, 2016.
- [77] John Bridle. Training Stochastic Model Recognition Algorithms as Networks can Lead to Maximum Mutual Information Estimation of Parameters. In *NIPS 1989*, pages 211–217, 1989.
- [78] Yoshua Bengio, Réjean Ducharme, and Pascal Vincent. A Neural Probabilistic Language Model. In *Advances in Neural Information Processing Systems*, volume 13, 2000.
- [79] Y. Bengio, R. Ducharme, P. Vincent, and C. Jauvin. A Neural Probabilistic Language Model. *Journal of Machine Learning Research*, 3:1137–1155, 2003.
- [80] G. E. Hinton, J. L. McClelland, and D. E. Rumelhart. *Distributed Representations*, pages 77–109. MIT Press, Cambridge, MA., 1986.
- [81] Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. Linguistic Regularities in Continuous Space Word Representations. In *NAACL 2013*, pages 746–751, 2013.
- [82] Jun Suzuki and Masaaki Nagata. Right-truncatable Neural Word Embeddings. In *NAACL 2016*, pages 1145–1151, 2016.
- [83] Omer Levy and Yoav Goldberg. Dependency-Based Word Embeddings. In *ACL 2014 (short paper)*, pages 302–308, 2014.
- [84] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *AISTATS 2010*, pages 297–304, 2010.
- [85] Omer Levy, Yoav Goldberg, and Ido Dagan. Improving Distributional Similarity with Lessons Learned from Word Embeddings. *Transactions of the Association for Computational Linguistics*, 3:211–225, 2015.
- [86] Omer Levy and Yoav Goldberg. Neural Word Embedding as Implicit Matrix Factorization. In *Advances in Neural Information Processing Systems 27*, pages 2177–2185, 2014.
- [87] John A. Bullinaria and Joseph P. Levy. Extracting Semantic Representations

- from Word Co-occurrence Statistics: A Computational Study. *Behavior Research Methods*, 39:510–526, 2007.
- [88] John A. Bullinaria and Joseph P. Levy. Extracting Semantic Representations from Word Co-occurrence Statistics: Stop-lists, Stemming and SVD. *Behavior Research Methods*, 44:890–907, 2012.
- [89] D. R. Cox. Regression Models and Life-Tables. *Journal of the Royal Statistical Society. Series B (Methodological)*, 34(2):187–220, 1972.
- [90] Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global Vectors for Word Representation. In *EMNLP 2014*, pages 1532–1543, 2014.
- [91] 横井祥, 下平英寿. 単語埋め込みの確率的等方化. In **言語処理学会第27回年次大会**, pages A7–1, 2021. [https://www.anlp.jp/proceedings/annual\\_meeting/2021/pdf\\_dir/A7-1.pdf](https://www.anlp.jp/proceedings/annual_meeting/2021/pdf_dir/A7-1.pdf).
- [92] Adriaan M. J. Schakel and Benjamin J. Wilson. Measuring Word Significance using Distributed Representations of Words. In *arXiv preprint*, 2015. <https://arxiv.org/abs/1508.02297>.
- [93] 大山百々勢, 横井祥, 下平英寿. 単語ベクトルの長さは意味の強さを表す. In **言語処理学会第28回年次大会**, pages A5–1, 2022.
- [94] H. S. Terrace, L. A. Petitto, R. J. Sanders, and T. G. Bever. Can an Ape Create a Sentence? *Science*, 206(4421):891–902, 1979.
- [95] Toshitaka N. Suzuki. Animal linguistics: Exploring referentiality and compositionality in bird calls. *Ecological Research*, 36:221–231, 2021.
- [96] Nipun Sadvilkar and Mark Neumann. PySBD: Pragmatic Sentence Boundary Disambiguation. In *Proceedings of Second Workshop for NLP Open Source Software (NLP-OSS)*, pages 110–114, 2020.
- [97] Rachel Wicks and Matt Post. A unified approach to sentence segmentation of punctuated text in many languages. In *ACL 2021*, pages 3995–4007, 2021.
- [98] Ryan Kiros, Yukun Zhu, Russ R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Skip-Thought Vectors. In *NIPS 2015*, 2015.
- [99] Prakhar Gupta, Matteo Pagliardini, and Martin Jaggi. Better Word Embeddings by Disentangling Contextual n-Gram Information. In *NAACL-HLT 2019*, pages 933–939, 2019.
- [100] Sanjeev Arora, Yingyu Liang, and Tengyu Ma. A simple but tough-to-beat baseline for sentence embeddings. In *ICLR 2017*, 2017.
- [101] Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. A Latent Variable Model Approach to PMI-based Word Embeddings. *Transactions of the Association for Computational Linguistics*, 4:385–399, 2016.
- [102] 伊庭幸人. **ベイズ統計と統計物理**. 岩波講座 物理の世界. 岩波書店, 2003.
- [103] Jiaqi Mu and Pramod Viswanath. All-but-the-Top: Simple and Effective Post-processing for Word Representations. In *ICLR 2018*, 2018.
- [104] Kawin Ethayarajh. Unsupervised Random Walk Sentence Embeddings: A

- Strong but Simple Baseline. In *Proceedings of The Third Workshop on Representation Learning for NLP*, pages 91–100, 2018.
- [105] Christopher D. Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven J. Bethard, and David McClosky. The Stanford CoreNLP Natural Language Processing Toolkit. In *ACL System Demonstrations*, pages 55–60, 2014. <http://www.aclweb.org/anthology/P/P14/P14-5010>.
- [106] Nikita Kitaev and Dan Klein. Constituency Parsing with a Self-Attentive Encoder. In *ACL 2018*, pages 2676–2686, 2018.
- [107] Mai Omura and Masayuki Asahara. UD-Japanese BCCWJ: Universal Dependencies Annotation for the Balanced Corpus of Contemporary Written Japanese. In *Second Workshop on Universal Dependencies (UDW 2018)*, pages 117–125, 2018.
- [108] 松田寛. GiNZA — Universal Dependencies による実用的日本語解析 —. *自然言語処理*, 27(3):695–701, 2020.
- [109] Mark Johnson, Thomas L. Griffiths, and Sharon Goldwater. Bayesian Inference for PCFGs via Markov Chain Monte Carlo. In *Proceedings of HLT/NAACL 2007*, pages 139–146, 2007.
- [110] J. ホップクロフト, R. モトワニ, J. ウルマン. *オートマトン 言語理論 計算論 I*. サイエンス社, 2003.
- [111] 能地宏. 文に隠れた構文構造を発見する統計モデル (特集「統計的言語研究の現在」). *統計数理*, 64(2):145–160, 2016.
- [112] Sadao Kurohashi and Makoto Nagao. Building a Japanese Parsed Corpus while Improving the Parsing System. In *Proceedings of LREC 1998*, pages 719–724, 1998. <http://nlp.kuee.kyoto-u.ac.jp/nl-resource/corpus.html>.
- [113] 持橋大地, 能地宏. 無限木構造隠れ Markov モデルによる階層的品詞の教師なし学習. *情報処理学会研究報告 2016-NL-226*, 12:1–11, 2016.
- [114] Julian Kupiec. Robust part-of-speech tagging using a hidden Markov model. *Computer Speech & Language*, 6(3):225–242, 1992.
- [115] Bernard Merialdo. Tagging English Text with a Probabilistic Model. *Computational linguistics*, 20(2):155–171, 1994.
- [116] Sharon Goldwater and Tom Griffiths. A Fully Bayesian Approach to Unsupervised Part-of-Speech Tagging. In *Proceedings of ACL 2007*, pages 744–751, 2007.
- [117] Jurgen Van Gael, Andreas Vlachos, and Zoubin Ghahramani. The infinite HMM for unsupervised PoS tagging. In *EMNLP 2009*, pages 678–687, 2009.
- [118] Phil Blunsom and Trevor Cohn. A Hierarchical Pitman-Yor Process HMM for Unsupervised Part of Speech Induction. In *ACL 2011*, pages 865–874, 2011.
- [119] Ryo Fujii, Ryo Domoto, and Daichi Mochihashi. Nonparametric Bayesian Semi-supervised Word Segmentation. *Transactions of ACL*, 5:179–189, 2017.
- [120] 竹内孔一, 松本裕治. 隠れマルコフモデルによる日本語形態素解析のパラメータ推定. *情報処理学会論文誌*, 38(3):500–509, 1997.
- [121] Andrew J. Viterbi. Error Bounds for Convolutional Codes and an Asymptoti-

- cally Optimum Decoding Algorithm. *IEEE Transactions on Information Theory*, 13(2):260–269, 1967.
- [122] G. David Forney Jr. The Viterbi Algorithm: A Personal History. In *arXiv preprint*, 2005. arXiv:cs/0504020 [cs.IT].
- [123] Laurence R. Rabiner. A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition. *Proceedings of the IEEE*, 77(2):257–286, 1989.
- [124] Stuart Geman and Donald Geman. Stochastic Relaxation, Gibbs Distributions, and the Bayesian Restoration of Images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 6(6):721–741, 1984.
- [125] Julian Besag. On the Statistical Analysis of Dirty Pictures. *Journal of the Royal Statistical Society, Series B*, 48(3):259–302, 1986.
- [126] Daichi Mochihashi. Unbounded Slice Sampling. Technical Report Research Memorandum No.1209, The Institute of Statistical Mathematics, 2020. arXiv:2010.01760.
- [127] Christophe Andrieu. On random- and systematic-scan samplers. *Biometrika*, 103(3):719–726, 2016.
- [128] Jun S. Liu. The Collapsed Gibbs Sampler in Bayesian Computations with Applications to a Gene Regulation Problem. *Journal of the American Statistical Association*, 89(427):958–966, 1994.
- [129] George Casella and Christian P. Robert. Rao-Blackwellisation of Sampling Schemes. *Biometrika*, 83(1):81–94, 1996.
- [130] M. J. Beal, Z. Ghahramani, and C. E. Rasmussen. The Infinite Hidden Markov Model. In *NIPS 2001*, pages 577–585, 2001.
- [131] David J.C. MacKay. Ensemble Learning for Hidden Markov Models. Technical report, Cavendish Laboratory, University of Cambridge, 1997.
- [132] Matthew J. Beal. *Variational Algorithms for Approximate Bayesian Inference*. PhD thesis, Gatsby Computational Neuroscience Unit, University College London, 2003. <https://cse.buffalo.edu/faculty/mbeal/thesis/>.
- [133] Steven L. Scott. Bayesian Methods for Hidden Markov Models. *Journal of the American Statistical Association*, 97:337–351, 2002.
- [134] Kevin Knight. Bayesian Inference with Tears, 2009. <https://kevincrawfordknight.github.io/papers/bayes-with-tears.pdf>.
- [135] Reina Akama, Kento Watanabe, Sho Yokoi, Sosuke Kobayashi, and Kentaro Inui. Unsupervised Learning of Style-sensitive Word Vectors. In *ACL 2018*, pages 572–578, 2018.
- [136] Phillip Keung, Yichao Lu, György Szarvas, and Noah A. Smith. The Multilingual Amazon Reviews Corpus. In *EMNLP 2020*, pages 4563–4568, 2020.
- [137] 株式会社ロンウィット. Livedoor ニュースコーパス, 2014. <http://www.rondhuit.com/download.html#ldcc>.
- [138] 宮内裕人, 鈴木陽也, 秋山和輝, 梶原智之, 二宮崇, 武村紀子, 中島悠太, 長原一. 主観と客観の感情極性分類のための日本語データセット. In **言語処理学会第 28 回年次**

- 大会, pages 1495–1499, 2022.
- [139] Haruya Suzuki, Yuto Miyauchi, Kazuki Akiyama, Tomoyuki Kajiwara, Takashi Ninomiya, Noriko Takemura, Yuta Nakashima, and Hajime Nagahara. A Japanese Dataset for Subjective and Objective Sentiment Polarity Classification in Micro Blog Domain. In *Proceedings of the 13th International Conference on Language Resources and Evaluation (LREC 2022)*, pages 7022–7028, 2022.
- [140] Kevin P. Murphy. *Probabilistic Machine Learning: An Introduction*. MIT Press, 2022.
- [141] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent Dirichlet Allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.
- [142] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum Likelihood from Incomplete Data via the EM Algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1):1–38, 1977.
- [143] Diederik P Kingma and Max Welling. Auto-Encoding Variational Bayes. In *arXiv preprint*, 2013. arXiv:1312.6114.
- [144] 篠崎寿夫, 松森徳衛, 吉田正広. **現代工学のための変分学入門**. 現代工学社, 1991.
- [145] Greg C.G. Wei and Martin A. Tanner. A Monte Carlo Implementation of the EM Algorithm and the Poor Man’s Data Augmentation Algorithms. *Journal of the American Statistical Association*, 85(411):699–704, 1990.
- [146] ウラム. **数学のスーパースターたち: ウラムの自伝的回想**. 東京図書, 1979.
- [147] Kenneth W. Church. Empirical Estimates of Adaptation: The chance of Two Noriegas is closer to  $p/2$  than  $p^2$ . In *COLING 2000*, pages 173–179, 2000.
- [148] K. Sjölander, K. Karplus, M.P. Brown, R. Hughey, R. Krogh, I.S. Mian, and D. Haussler. Dirichlet Mixtures: A Method for Improved Detection of Weak but Significant Protein Sequence Homology. *Computing Applications in the Biosciences*, 12(4):327–345, 1996.
- [149] 山本 幹雄, 貞光 九月, 三品 拓也. 混合ディリクレ分布を用いた文脈のモデル化と言語モデルへの応用. **情報処理学会研究報告 2003-SLP-48**, pages 29–34, 2003.
- [150] Issei Sato and Hiroshi Nakagawa. Topic models with power-law using Pitman-Yor process. In *KDD 2010*, page 673–682, 2010.
- [151] 山本幹雄, 持橋大地. Topic に基づく統計的言語モデルの最前線 –PLSI から HDP まで-. In **言語処理学会年次大会 2006 チュートリアル**, pages 11–28, 2006.
- [152] 貞光九月, 待鳥裕介, 山本幹雄. 混合ディリクレ分布パラメータの階層ベイズモデルを用いたスムージング法. **情報処理学会研究報告 2004-SLP-53**, pages 1–6, 2004.
- [153] Charles Elkan. Clustering Documents with an Exponential-Family Approximation of the Dirichlet Compound Multinomial Distribution. In *ICML 2006*, pages 289–296, 2006.
- [154] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. Latent Dirichlet Allocation. In *Neural Information Processing Systems 14*, 2001.
- [155] Yee Whye Teh, David Newman, and Max Welling. A Collapsed Variational Bayesian Inference Algorithm for Latent Dirichlet Allocation. In *Advances in Neu-*



- ral Information Processing Systems 19 (NIPS 2006)*, 2006.
- [156] Tom Griffiths. Gibbs sampling in the generative model of Latent Dirichlet Allocation. Technical Report 11, Stanford University, 2002.
- [157] Thomas L. Griffiths and Mark Steyvers. Finding scientific topics. *PNAS*, 101:5228–5235, 2004.
- [158] Thomas Minka. Expectation-Propagation for the Generative Aspect Model. In *UAI 2002*, page 352–359, 2002.
- [159] J.K. Pritchard, M. Stephens, and P. J. Donnelly. Inference of Population Structure Using Multilocus Genotype Data. *Genetics*, 155:945–959, 2000.
- [160] Hanna M. Wallach, Iain Murray, Ruslan Salakhutdinov, and David Mimno. Evaluation Methods for Topic Models. In *ICML 2009*, pages 1105–1112, 2009.
- [161] Wray Buntine and Aleks Jakulin. Applying Discrete PCA in Data Analysis. In *UAI 2004*, pages 59–66, 2004.
- [162] Hanna M. Wallach, David Mimno, and Andrew McCallum. Rethinking LDA: Why Priors Matter. In *NIPS 2009*, pages 1973–1981, 2009.
- [163] Jonathan Chang, Sean Gerrish, Chong Wang, Jordan Boyd-graber, and David Blei. Reading Tea Leaves: How Humans Interpret Topic Models. In *NIPS 2009*, 2009.
- [164] Jey Han Lau, David Newman, and Timothy Baldwin. Machine Reading Tea Leaves: Automatically Evaluating Topic Coherence and Topic Model Quality. In *EACL 2014*, pages 530–539, 2014.
- [165] Feng Nan, Ran Ding, Ramesh Nallapati, and Bing Xiang. Topic Modeling with Wasserstein Autoencoders. In *ACL 2019*, pages 6345–6381, 2019.
- [166] 石井健一郎, 上田修功. 続・わかりやすいパターン認識—教師なし学習入門—. オーム社, 2014.
- [167] Daniel Lee and H. Sebastian Seung. Algorithms for Non-negative Matrix Factorization. In *Advances in Neural Information Processing Systems*, volume 13, 2000.
- [168] John Canny. GaP: A Factor Model for Discrete Data. In *SIGIR 2004*, pages 122–129, 2004.
- [169] 持橋大地. GaP, NMF, and more, 2006. <http://chasen.org/~daiti-m/paper/gap-nmf.pdf>.
- [170] S. Deerwester, Susan T. Dumais, and George W. Furnas. Indexing by Latent Semantic Analysis. *Journal of the American Society of Information Science*, 41(6):391–407, 1990.
- [171] Thomas Hofmann. Probabilistic Latent Semantic Indexing. In *SIGIR 1999*, pages 50–57, 1999.
- [172] Quoc Le and Tomas Mikolov. Distributed Representations of Sentences and Documents. In *ICML 2014*, pages 1188–1196, 2014.
- [173] Daichi Mochihashi. Researcher2Vec: Neural Linear Model of Scholar Recommendation for Funding Agency. In *International Society for Scientometrics and*

- Informatics (ISSI 2023)*, 2023.
- [174] 南風原朝和. **心理統計学の基礎**. 有斐閣アルマ. 有斐閣, 2002.
  - [175] Sungjoon Park, JinYeong Bak, and Alice Oh. Rotated Word Vector Representations and their Interpretability. In *EMNLP 2017*, pages 401–411, 2017.
  - [176] Hiroaki Yamagiwa, Momose Oyama, and Hidetoshi Shimodaira. Discovering Universal Geometry in Embeddings with ICA. In *EMNLP 2023*, pages 4647–4675, 2023.
  - [177] 東京大学教養学部統計学教室. **統計学入門**. 基礎統計学 I. 東京大学出版会, 1991.
  - [178] Sascha Rothe, Sebastian Ebert, and Hinrich Schütze. Ultradense Word Embeddings by Orthogonal Transformation. In *NAACL 2016*, pages 767–777, 2016.
  - [179] Bianca Zadrozny and Charles Elkan. Transforming Classifier Scores into Accurate Multiclass Probability Estimates. In *KDD 2002*, pages 694–699, 2002.
  - [180] Jonathan B. Slapin and Sven-Oliver Proksch. A Scaling Model for Estimating Time-Series Party Positions from Texts. *American Journal of Political Science*, 52(3):705–722, 2008.
  - [181] Kohei Watanabe. Latent Semantic Scaling: A Semisupervised Text Analysis Technique for New Domains and Languages. *Communication Methods and Measures*, pages 1–23, 2020.
  - [182] 持橋大地. 確率の潜在意味スケーリング. In **情報処理学会研究報告 2021-NL-249**, number 9, pages 1–16, 2021.
  - [183] Philipp Dufter and Hinrich Schütze. Analytical Methods for Interpretable Ultradense Word Embeddings. In *EMNLP-IJCNLP 2019*, pages 1185–1191, 2019.
  - [184] Ikuya Yamada, Akari Asai, Jin Sakuma, Hiroyuki Shindo, Hideaki Takeda, Yoshiyasu Takefuji, and Yuji Matsumoto. Wikipedia2Vec: An Efficient Toolkit for Learning and Visualizing the Embeddings of Words and Entities from Wikipedia. In *EMNLP 2020: System Demonstrations*, pages 23–30, 2020.
  - [185] Malay Ghosh. Inconsistent maximum likelihood estimators for the Rasch model. *Statistics & Probability Letters*, 23(2):165–170, 1995.
  - [186] Radford M. Neal. *MCMC Using Hamiltonian Dynamics*. Chapman and Hall/CRC, 2011.
  - [187] Qing Liu and Donald A. Pierce. A Note on Gauss-Hermite Quadrature. *Biometrika*, 81(3):624–629, 1994.
  - [188] Andrew Gelman, John B. Carlin, Hal S. Stern, and Donald B. Rubin. *Bayesian Data Analysis, Third Edition*. Chapman & Hall/CRC, 2013. <http://www.stat.columbia.edu/~gelman/book/>.