

実践 Python による
ベイズ分析とトピックモデル
—— 先進的なデータ分析へのアプローチ ——

博士（工学） 藤野 巖【著】

コロナ社

【本書ご利用にあたって】

本書で解説している内容を実行・利用したことによる直接あるいは間接的な損害に対して、著作者およびコロナ社は一切の責任を負いかねます。利用についてはすべて読者個人の責任において行ってください。

本書に掲載されている情報は、本書執筆時点のもので、将来にわたって保証されるものではありません。特に、各社が提供しているソフトウェアパッケージは仕様やサービス提供に係る変更が頻繁にあり、Python のライブラリ群等も頻繁にバージョンアップがなされています。これらによっては本書で解説しているアプリケーション等が正常に動作しなくなることもあるので、あらかじめご了承ください。

本書の発行にあたって、読者の皆様に問題なく実践していただけるよう、できる限りの検証をしていますが、以下の環境以外では構築・動作を確認していないので、あらかじめご了承ください。

PC 本体：Windows11 Pro 64 bit (CPU：Intel(R) Core i7, メモリ：32 GB)

開発環境：Anaconda バージョン 2023.09 (Python バージョン 3.11)

また、上記環境を整えたいかなる状況においても動作が保証されるものではありません。ネットワークやメモリの使用状況および同一 PC 上にあるほかのソフトウェアの動作状況によって、本書のプログラムが正常に動作できなくなることがあります。併せてご了承ください。

なお、本書に記載している会社名、製品名は、それぞれ各社の商標または登録商標です。

本書の購入者に対する限定サービスとして、本書に掲載しているソースコードは、以下に示すコロナ社の Web ページからダウンロードできます。ぜひご利用ください。

<https://www.coronasha.co.jp/np/isbn/9784339029437/>

なお、本書に掲載しているソースコードについては、オープンソースソフトウェアの BSD ライセンス下で再利用も再配布も自由です。

まえがき

世の中にはデータがあふれかえっている。データと向き合うのは、もはや日常茶飯事である。データとは、現実社会の記録である。個々のデータを取り出してみることで、昔のことを具体的に正確に思い出すことができる。さらに、データをたくさん集めておけば、それらを分析することで、世の中の事柄のトレンド、事柄間の関連性のような価値のある情報が浮かび上がる。

データ分析の手法として、以前から平均値、標準偏差、四分位数などの統計量が用いられている。ただしこれらの手法では、収集したデータをそのまま利用しているため、いわゆる潜在的な情報、深層的な情報を発見することができない。

本書は、大量なデータから潜在的で深層的な情報を発掘できるトピックモデルという確率的なデータ分析手法の理論と実践の両方を解説するものである。トピックモデルは、自然言語処理（NLP）の手法の一つとして知られているが、その威力は文書データにとどまらず、画像データや軌跡データの分析にも応用できるような汎用的な技術であり、ディープラーニングと並んで人工知能の基本として認識されるべき技術である。

本書では、簡単な準備学習を経てから、ベイズ分析とトピックモデル、そしてその応用事例の順に、プログラムを作成しながら実践的に学習できるように章立てが用意されている。各章の内容を以下に簡単に紹介する。

- 1章：本書の学習に必要な確率と確率分布の知識およびそのプログラム実現を説明する。
- 2章：対比の位置付けとして、従来のデータ分析の基本手法を復習する。
- 3章：ベイズ分析の基本的な考え方を説明する。併せてベイズ分析のプログラム実現に使われる PyMC ライブラリの使い方を紹介する。
- 4章：対比の位置付けとして、従来の文書データ分析の基本手法を復習する。

- 5章：ユニグラムモデルを構成して、文書データの分析を行う。また、PyMC ライブラリにより、そのプログラム実現を示す。
- 6章：トピックの考え方を取り入れて、混合ユニグラムモデルを構成する。また、混合ユニグラムモデルを用いた文書解析プログラム例を示す。
- 7章：混合ユニグラムモデルをさらに発展させて、トピックモデルを構成する。また、トピックモデルを用いた文書解析プログラム例を示す。
- 8章：Scikit-learn ライブラリにあるトピックモデルのモジュールの使い方を説明する。それを利用して、20 ニュースグループデータセットの英語文書データからトピックを抽出する。
- 9章：Gensim というトピックモデルに特化したライブラリの使い方を説明する。それを利用して、Wikipedia の日本語文書データからトピックを抽出する。
- 10章：トピックモデルを拡張して、著者トピックモデルを構成する。そのうえで、Gensim ライブラリを利用して、Twitter（現、X）から収集した日本語の投稿データからトピックを抽出する。
- 11章：トピックモデルを画像データセットに応用する。Gensim ライブラリを利用して、Caltech 101 というデータセットから、小さく分割されたセルで表したトピックを抽出する。
- 12章：トピックモデルを軌跡データセットに応用する。Gensim ライブラリを利用して、船舶の AIS データから、航路（コース）となるようなトピックを抽出する。

読者の皆様におかれましては、本書を読むことで先進的なデータ分析スキルを身に付け、トピックモデルの技術を応用して実務の場で役立てて頂ければ幸いです。

最後に、本書の執筆にあたり、多方面から熱心に支えてくれた家族に感謝の意を表する。また、本書を出版する機会を与えてくださったコロナ社の皆さんに衷心より感謝申し上げます。

2024年2月

目 次

1. 確率と確率分布

1.1 確率と確率分布	1
1.1.1 確 率	1
1.1.2 確 率 分 布	2
1.2 条件確率と同時確率	3
1.2.1 条 件 確 率	3
1.2.2 同 時 確 率	4
1.3 乗法定理とベイズの定理	4
1.3.1 乗 法 定 理	4
1.3.2 ベイズの定理	5
1.4 各種確率分布	6
1.4.1 ベルヌーイ分布	6
1.4.2 二 項 分 布	7
1.4.3 カテゴリ分布	8
1.4.4 正規分布 (1次元ガウス分布)	8
1.4.5 指 数 分 布	9
1.4.6 ベ ー タ 分 布	10
1.4.7 ガ ン マ 分 布	10
1.4.8 デイリクレ分布	11
1.5 Scipy の確率統計モジュール stats	12
1.6 Python による各種確率分布のプログラム	15
演 習 問 題	19

2. データの統計分析の基本

2.1 統計分析	20
2.2 Pandas とは	23
2.2.1 シリズ	23
2.2.2 データフレーム	24
2.3 Pandas の基本的な使い方	26
2.4 Pandas による統計分析のプログラム例	31
演習問題	37

3. ベイズ分析の基本

3.1 ベイズ分析の基本的な考え方	38
3.2 PyMC の使い方 (その1)	39
3.3 Arviz について	42
3.4 PyMC によるサンプリングプログラム	43
3.4.1 ベルヌーイ分布からのサンプリング	43
3.4.2 正規分布からのサンプリング	46
3.5 PyMC によるパラメータ推定のプログラム	49
3.5.1 ベルヌーイ分布のパラメータ推定	49
3.5.2 正規分布のパラメータ推定	53
演習問題	57

4. 文書データ分析の基本

4.1 形態素解析と MeCab	58
4.2 Bag of Words	64
4.3 文書データの数値化	65
4.3.1 単語の出現回数	66
4.3.2 TFIDF	67
4.3.3 Scikit-learn ライブラリによる TFIDF 単語文書行列の実現法	69

4.3.4 TFIDF 単語文書行列を計算するプログラム例	71
4.4 コサイン類似度	74
4.4.1 コサイン類似度の計算式	74
4.4.2 Scikit-learn ライブラリによるコサイン類似度の実現法	74
4.4.3 コサイン類似度を計算するプログラム例	75
演習問題	77

5. ユニグラムモデル

5.1 文書生成の確率モデル	78
5.2 グラフィカルモデル表現	80
5.3 ユニグラムモデル	80
5.3.1 ユニグラムモデルとは	80
5.3.2 文書集合を生成するプログラム	82
5.3.3 文書集合の単語出現頻度のプログラム	85
5.4 ユニグラムモデルのパラメータ推定	87
5.4.1 データ分析とモデル	87
5.4.2 PyMC によるユニグラムモデルのパラメータ推定	89
5.4.3 ユニグラムモデルにおけるカテゴリ分布の発生確率推定	91
演習問題	96

6. 混合ユニグラムモデル

6.1 混合ユニグラムモデル	97
6.2 文書集合（複数の文書）の生成	101
6.3 PyMC の使い方（その2）	105
6.4 混合ユニグラムモデルにおけるトピック別単語分布の推定	107
演習問題	116

7. トピックモデル

7.1 トピックモデル	117
-------------------	-----

7.2 文書集合（複数の文書）の生成	120
7.3 トピックモデルにおけるカテゴリ分布のパラメータ推定	124
7.4 保存済み推定結果の利用	130
7.4.1 トピック別の単語の出現確率順表示	130
7.4.2 文書別トピック割合の表示	131
演習問題	135

8. Scikit-learn ライブラリによるトピックモデル

8.1 20 ニュースグループデータセット	137
8.2 英語文書の形態素解析	141
8.3 Scikit-learn のトピックモデルライブラリ	145
8.4 20 ニュースグループデータセットから単語集合を用意する	147
8.5 20 ニュースグループの単語集合からトピックを抽出する	150
8.6 トピック別の上位単語を取り出す	153
演習問題	158

9. Gensim ライブラリによるトピックモデル

9.1 Wikipedia 記事の単語集合の作成	159
9.2 Gensim を用いたトピック解析	164
9.3 Wikipedia データセットへのトピックモデルの適用	170
9.4 トピック別の上位単語と文書別のトピックの取り出し	173
演習問題	176

10. 著者トピックモデル

10.1 著者トピックモデル	177
10.2 PyMC による著者トピックモデルの実現	179
10.2.1 著者トピックモデルに基づいたデータセットの作成	179
10.2.2 PyMC による著者トピックモデルの実現	183

10.3 Gensim による著者トピックモデル	187
10.3.1 Gensim の著者トピックモデルクラス	187
10.3.2 Twitter データからのデータセットの準備	189
10.3.3 Twitter データセットへの著者トピックモデルの適用	192
演習問題	196

11. 画像データセットからのトピック抽出

11.1 画像データにトピックモデルを適用するための予備知識	197
11.1.1 画像データの仕組み	197
11.1.2 Matplotlib による画像の読み込みと表示	198
11.1.3 ベクトル量子化とコードの作成	202
11.2 画像データセットを用いた画像の文書集合の作成	205
11.3 画像の文書集合へのトピックモデルの適用	211
11.4 トピックモデルによる処理結果の可視化	212
演習問題	216

12. 船舶の航跡データからのトピック抽出

12.1 AIS データにトピックモデルを適用するための予備知識	219
12.1.1 AIS データ	219
12.1.2 Basemap による地図情報の表示	219
12.1.3 Basemap による航跡の表示	222
12.2 AIS データセットを用いた航跡の文書集合の作成	229
12.3 航跡のコード文書集合へのトピックモデルの適用	233
12.4 トピックをコードから航跡に復元	234
演習問題	240

付 録	241
参 考 文 献	243
索 引	244

1

確率と確率分布

本章では、今後のベイズ分析とトピックモデルの学習で必要となる各種確率分布について解説する。まず、確率の基本的な考え方を述べてから、条件確率、同時確率、ベイズの定理の順に基本的な概念を説明する。そのうえで、各種確率分布の計算式とグラフも示す。最後に、各種確率分布のグラフを作成する Python プログラムについて説明する。

1.1 確率と確率分布

1.1.1 確率

世の中のことは、確実に起こることもあれば、不確実に起こることもある。コインを高く投げたら、確実に落ちる。けれども、落ちたとき、上を向くのが、表なのか裏なのかは不確実である。不確実に起こることは、一度やってみるだけでは、起こることもあれば、起こらないこともある。けれども、何度もやってみると、だいたいそのうちの何回かが起こることになる。

具体例を使ってもう少し考えてみよう。例えば、コインを 100 回投げたら、そのうち表が出たのは 50 回であったとする。対して、サイコロを 100 回振ったら、そのうち 1 の目が出たのは 16 回であったとする。ここで、コインを投げたりサイコロを振ったりするような行為を**試行** (trial) と言う。また、コインの表やサイコロの 1 の目のような試行の結果を**事象** (event) と言う。つまり、この例では、コインが表という事象 A が 100 回試行のうち 50 回発生した。サイコロの目が 1 という事象 B は 100 回試行のうち 16 回発生した。このような状況を踏まえて、事象 A は事象 B より起こりやすいと考える。このように、不確実なことについては、その起こりやすさに違いがある。ある事象の起こりやすさを使えば、そのような事象の起こる確実性の度合いを表すことができる。

2 1. 確率と確率分布

この事象が起こりやすさのことを**確率** (probability) と言う。

1.1.2 確率分布

話をもう一步進めてみよう。まず、**確率変数** (random variable) c を使ってコインを投げる試行の結果を表すものとする。コインを投げる試行の結果、上に向く面は、表と裏の二つしかないので、確率変数 c の値が表 (1) と裏 (0) の二つしかない。確率変数 c のすべての値に対する確率が以下ようになる。

$$P(c=1) = \frac{1}{2}, \quad P(c=0) = \frac{1}{2} \quad (1.1)$$

このような確率変数のすべての値の確率を全体的に表したものを**確率分布** (probability distribution) と言う。

また、確率変数 d を使ってサイコロを振る試行の結果を表すものとする。サイコロの目は、1~6 まであるので、確率変数 c のすべての値に対する確率は以下ようになる。

$$\begin{aligned} P(d=1) &= \frac{1}{6}, & P(d=2) &= \frac{1}{6}, & P(d=3) &= \frac{1}{6} \\ P(d=4) &= \frac{1}{6}, & P(d=5) &= \frac{1}{6}, & P(d=6) &= \frac{1}{6} \end{aligned} \quad (1.2)$$

数式だけでは、全体的なイメージがつかみにくいので、グラフを使ってすべての確率を示すことがより一般的に行われている。また、コインやサイコロのように、確率変数が有限個整数から値を取るとき、**離散型** (discrete) **確率変数** と言い、離散型の確率分布の全体を表した関数を**確率質量関数** (**PMF** : probability mass function) と言う。対して、確率変数が実数から値を取るとき、**連続型** (continuous) **確率変数** と言い、連続型の確率分布の全体を表した関数を**確率密度関数** (**PDF** : probability density function) と言う。

1.2 条件確率と同時確率

簡単な例題を使って説明しよう。図 1.1 のように、箱の中には、黒いボールが一つと白いボールが二つ入っており、箱からボールを取り出すときの確率を考える。条件確率と同時確率を説明する前に、まずは、普通確率（無条件確率）を計算しておく。箱の中からボール b を一つ取り出す試行をしたとき、その結果が黒いボール、または白いボールである確率は以下ようになる。

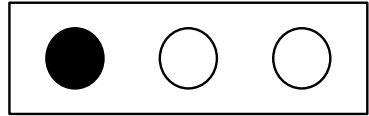


図 1.1 条件確率と同時確率の例題

$$P(b = \text{黒}) = \frac{1}{3}, \quad P(b = \text{白}) = \frac{2}{3} \quad (1.3)$$

1.2.1 条件確率

これからは、図 1.1 の箱の中からボールを一つ取り出すことを 2 回行うとする。まず、1 回目の試行のボール b_1 を取り出したあと、2 回目の試行のボール b_2 が黒または白となる確率を計算する。すべての組合せの計算結果を以下に示す。

$$\begin{aligned} P(b_2 = \text{黒} | b_1 = \text{黒}) &= 0, & P(b_2 = \text{白} | b_1 = \text{黒}) &= 1 \\ P(b_2 = \text{黒} | b_1 = \text{白}) &= \frac{1}{2}, & P(b_2 = \text{白} | b_1 = \text{白}) &= \frac{1}{2} \end{aligned} \quad (1.4)$$

ここでは、1 回目のボールが特定の色となったことを条件として、2 回目のボールが特定の色となる確率を計算している。このような確率のことを**条件確率** (conditional probability) と言う。一般的に言うと、条件確率は、ある事象 A が起こることを仮定したときの別の事象 B が起こる確率で、これを $P(B|A)$ と記す。

4 1. 確率と確率分布

1.2.2 同時確率

つぎに、1回目のボールが特定の色となり、かつ2回目のボールが特定の色となる確率を計算する。すべての組合せの計算結果を以下に示す。

$$\begin{aligned}P(b_1 = \text{黒かつ } b_2 = \text{黒}) &= \frac{1}{3} \times 0 = 0 \\P(b_1 = \text{黒かつ } b_2 = \text{白}) &= \frac{1}{3} \times 1 = \frac{1}{3} \\P(b_1 = \text{白かつ } b_2 = \text{黒}) &= \frac{2}{3} \times \frac{1}{2} = \frac{1}{3} \\P(b_1 = \text{白かつ } b_2 = \text{白}) &= \frac{2}{3} \times \frac{1}{2} = \frac{1}{3}\end{aligned}\tag{1.5}$$

ここでは、1回目のボールが特定の色となるかつ2回目のボールが特定の色となる確率を計算している。このような二つの事象を同時に起きる確率のことを**同時確率** (joint probability) と言う。一般的に言うと、同時確率は、ある事象 A と事象 B が同時に起こる確率で、これを $P(A \cap B)$ [†] と記す。

1.3 乗法定理とベイズの定理

1.3.1 乗法定理

ここで、1.2節の計算結果をよく観察してみると、同時確率が

$$P(b_1 = \text{白かつ } b_2 = \text{黒}) = \frac{2}{3} \times \frac{1}{2}\tag{1.6}$$

となる一方、一回目の試行では、無条件確率が

$$P(b_1 = \text{白}) = \frac{2}{3}\tag{1.7}$$

となり、2回目の試行では、条件確率が

$$P(b_2 = \text{黒} | b_1 = \text{白}) = \frac{1}{2}\tag{1.8}$$

[†] 本によっては、同時確率を $P(AB)$ または $P(A, B)$ と記すことがある。

となるので、同時確率が以下のように、無条件確率と条件確率の積で表せることがわかる。

$$P(b_1 = \text{白かつ } b_2 = \text{黒}) = P(b_1 = \text{白})P(b_2 = \text{黒} | b_1 = \text{白}) \quad (1.9)$$

これを一般的に言うと、事象 A と事象 B があるとすれば、同時確率

$$P(A \cap B) = P(A|B)P(B) \quad (1.10)$$

と表すことができる。この公式は**乗法定理** (multiplication theorem) と呼ばれる。これを利用すれば次式より条件確率 $P(A|B)$ を計算できるようになる。

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad (1.11)$$

さらに、ここで事象 A の確率が事象 B の影響を受けなければ、つまり

$$P(A|B) = P(A)$$

ならば、同時確率 $P(A \cap B)$ は、以下のように、それぞれの事象の無条件確率の積で計算できるようになる。

$$P(A \cap B) = P(A)P(B) \quad (1.12)$$

この場合は、事象 A と事象 B が**独立** (independent) であると言う。

1.3.2 ベイズの定理

乗法定理から、事象 A と事象 B の同時確率は

$$P(A \cap B) = P(A|B)P(B) \quad (1.13)$$

となり、事象 B と事象 A の同時確率は

$$P(B \cap A) = P(B|A)P(A)$$

となる。しかし、この二つの確率が等しいので、 $P(A \cap B) = P(B \cap A)$ から以下の**ベイズの定理** (Bayes' theorem) を簡単に導出することができる。

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (1.14)$$

索引

【あ行】	重み	66	指数分布	9	20 ニュースグループ	137
【か行】	確率	2	事前分布	39	【は行】	
	確率質量関数	2	四分位数	22	ヒストグラム	22
	確率分布	2	四分位数範囲	22	標準偏差	21
	確率変数	2	シリーズ	23	文書集合	65
	確率密度関数	2	条件確率	3	文書類度	68
	確率モデル	78	乗法定理	5	平均値	21
	カテゴリ分布	8	正規分布	8	ベイズ推定	39
	ガンマ分布	10	成功確率	7	ベイズの定理	5
	逆文書類度	68	【た行】		ベイズ分析	39
	グラフィカルモデル表現	80	単語頻度	67	ベータ分布	10
	形態素	58, 141	単語文書行列	66	バルヌーイ分布	6
	形態素解析	58, 141	中央値	21	変分推定法	40
	コサイン類似度	74	著者トピックモデル	177	【ま行】	
	混合ユニグラムモデル	97	ディリクレ分布	11	マルコフチェーン	
【さ行】			データフレーム	24	モンテカルロ法	39
	最小値	21	点推定	38	【や行】	
	最大値	21	統計	20	尤度関数	39
	サンプラー	105	統計量	20	ユニグラムモデル	80
	試行	1	同時確率	4	【ら行】	
	事後分布	39	独立	5	離散型確率変数	2
	事象	1	度数分布表	22	連続型確率変数	2
			トピック	97		
			トピックモデル	117		
			【な行】			
			二項分布	7		

【A~D, G】	AIS	219	Gensim	159, 164	PMF	2
	Arviz	42	【I, M, P】		PyMC	39, 89, 105, 179
	Basemap	219	IDF	68	【S, T, V】	
	Bag of Words (BoW)	65	MeCab	58	Scikit-learn	137
	Caltech 101	205	MCMC	40	TF	67
	DF	68	Pandas	23	TFIDF	69
			PDF	2	VI	40

— 著者略歴 —

1991年 東海大学大学院工学研究科博士課程修了
博士（工学）
1994年 東海大学短期大学部講師
2001年 東海大学短期大学部助教授
2004年 イギリス・サウサンプトン大学客員教授
～05年
2007年 東海大学短期大学部教授
2008年 東海大学情報通信学部教授
現在に至る
2016年 フランス海軍アカデミー招聘研究員
～17年

実践 Python によるベイズ分析とトピックモデル
—先進的なデータ分析へのアプローチ—

Perfect Practice Bayesian Analysis and Topic Modeling with Python
—An Advanced Approach to Data Analysis—

© Iwao Fujino 2024

2024年4月25日 初版第1刷発行



検印省略

著者	藤野 巖
発行者	株式会社 コロナ社
代表者	牛来真也
印刷所	壮光舎印刷株式会社
製本所	株式会社 グリーン

112-0011 東京都文京区千石 4-46-10
発行所 株式会社 コロナ社
CORONA PUBLISHING CO., LTD.
Tokyo Japan
振替00140-8-14844・電話(03)3941-3131(代)
ホームページ <https://www.coronasha.co.jp>

ISBN 978-4-339-02943-7 C3055 Printed in Japan

(松岡)



©COPY < 出版者著作権管理機構 委託出版物 >

本書の無断複製は著作権法上での例外を除き禁じられています。複製される場合は、そのつど事前に、出版者著作権管理機構（電話 03-5244-5088, FAX 03-5244-5089, e-mail: info@jcopy.or.jp）の許諾を得てください。

本書のコピー、スキャン、デジタル化等の無断複製・転載は著作権法上での例外を除き禁じられています。購入者以外の第三者による本書の電子データ化及び電子書籍化は、いかなる場合も認めていません。落丁・乱丁はお取替えいたします。