

# テキストデータの統計科学入門

## 第8章 テキストにおける確率モデリング

茨城大学  
國井慎也

# 概要

- テキスト分析のための確率の基本的な概念と主に使用される確率分布、その確率分布を用いた簡単なモデリングを紹介する。

# 確率

- 標本空間

-起こりうるすべての事象の集合

$$\Omega = \{A_1, A_2, \dots, A_n\}$$

- 確率

N: 試行の回数

$x_i$ : N回の試行で事象 $A_i$ が起こる回数

$$P\{A_i\} \doteq \frac{x_i}{N}$$

# 確率の主な定理

- 乗法定理

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

$P(A \cap B)$ : 事象A, Bが同時に起きる確率

$P(A|B)$ : Bが起きた条件の下でのAが起きる確率

- ベイズの定理

-乗法定理から導かれる定理

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

$P(A)$ : 事象Aの事前確率

$P(A|B)$ : 事象Aの事後確率

# 確率分布

- 確率変数

  - 試行の標本空間の事象の変数

- 確率分布

  - 確率変数とその変数に対応する確率のセット

$$P(X=x_i)=P(A_{xi})$$

$$\sum_i P(X=x_i)=1$$

確率分布の例(6面のサイコロ)

確率変数	1	2	3	4	5	6
確率	1/6	1/6	1/6	1/6	1/6	1/6

# 主な分布(1)

- 二項分布

-成功と失敗、男性と女性など試行の結果が2通りある問題をモデル化するとき用いられる

$$P(X=x) = {}_n C_x p^x (1-p)^{n-x}$$

$$E(X) = np, \quad \text{Var}(X) = np(1-p)$$

- 幾何分布

-独立な試行を繰り返して、x回失敗してx+1回目に成功する確率を考えた分布

$$P(X=x) = pq^x \quad (x=0,1,2,\dots, 0 < p < 1, q=1-p)$$

$$E(X) = \frac{1-p}{p}, \quad \text{Var}(X) = \frac{1-p}{p^2}$$

# 主な分布(2)

- ポアソン分布

-一定期間内で、平均 $\lambda$ で起きる事象が $x$ 回起きる確率を考えた分布

$$P(X=x) = e^{-\lambda} \frac{\lambda^x}{x!}$$

- 正規分布

-平均値に集積するようなデータの分布に関する確率分布

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

特に $\mu=0$ 、 $\sigma=1$ の場合を標準正規分布と呼ぶ

# 主な分布(3)

- 対数正規分布

-正規分布の確率変数に対数を取った $Y=\log(X)$ が $N(\mu, \sigma)$ に従うとき、確率変数 $X$ が従う分布

$$f(x) = \frac{1}{\sigma x \sqrt{2\pi}} e^{-\frac{(\log(x) - \mu)^2}{2\sigma^2}}$$

$$E(X) = e^{\mu + \frac{\sigma^2}{2}}, \quad Var(X) = e^{2\mu + \sigma^2} - e^{2\mu}$$

# モデリングの例(1)

- 英文の手紙における単語を、音節を単位として調べ、単語の長さの頻度分布を幾何分布に当てはめる

単語の長さの度数表

音節 $x_i$	度数 $f_i$	相対度数 $p_i$
1	3978	0.7613
2	831	0.1587
3	281	0.0537
4	121	0.0231
5	15	0.0029
6	6	0.0004
合計	5237	1

$$\text{幾何分布 } P(X=x) = pq^x$$

1. 実測値の平均を求める

$$\bar{x} = \sum_{i=1}^n x_i p_i = 1.3491$$

2. 幾何分布の $p, q$ を求める

$$p = \frac{1}{\bar{x}} = \frac{1}{1.3491} = 0.7412, \quad q = 1 - 0.7412 = 0.2588$$

# モデリングの例(2)

## 3. 音節の推測相対度数を求める

$$P_1 = p(1-p)^{1-1} = 0.7412 \times 0.2588^0 = 0.7412$$

$$P_2 = p(1-p)^{2-1} = 0.7412 \times 0.2588^1 \doteq 0.1918$$

## 4. 推測度数を求める

$$NP_1 = 5237 \times 0.7412 \doteq 3882$$

$$NP_2 = 5237 \times 0.1918 \doteq 1002$$

音節xi	実測度数fi	推測度数NPi	推測相対度数Pi
1	3987	3882	0.7412
2	831	1002	0.1918
3	281	260	0.0496
4	121	76	0.0128
5	15	17	0.0033
6	2	5	0.0009
合計	5237	5235	1