

入門ベイズ統計

第5章 情報検索とベイズ決定

10月10日

田中洸一

項目

- 問題提起
- 5. 1 分類子と関連性
- 5. 2 文献からのサンプリング

問題提起

図書館などの情報検索のシステムにベイズ決定を用いることができる



情報の関連性の有無は、「確率」が絡む



情報検索システムに応用するための手法を提示

5.1 分類子と関連性

- ・情報の関連性の有無に二つの基準を定義
-

例：図書館での検索1

「オペレーションズ・リサーチ」の手法を「図書館」の「閲覧」の問題に応用するための文献を検索することを考える。

- ・以下の四つの概念用語を用意

I1: '閲覧' I2: 'OR' I3: '図書館' I4: '閲覧'および'OR'

- ・概念用語を文献の関連性の有無の基準とする。

- ・概念用語を分類子と表現し、 m 個の分類子 I_1, I_2, \dots, I_m を定義

5.1 分類子と関連性

- 文献 $s \in \mathcal{S}$ に対し、関数 $Z_j(s)$

$$Z_j(s) = \begin{cases} 1 & (s \text{ に分類子 } I_j \text{ が与えられるとき}) \\ 0 & (\text{与えられないとき}) \end{cases} \quad (5.1.1)$$

$$(j = 1, 2, \dots, m)$$

例： $(1, 0, 0, 1, 1, \dots, 1, 0)$ (m 桁)

文献 s に m 桁の 2 進数 $(Z_1(s), Z_2(s), \dots, Z_m(s))$ が与えられる

どんな文献 s も分類子 I_1, I_2, \dots, I_m によって $\varphi (= 2^m - 1)$ 個の

どれかの部分集合 (カテゴリー) に分類される

5. 1 分類子と関連性

- 関連性の段階的尺度 $\{0,1,2,L ,q\}$ (かなり主観的な基準)

例: $q=1 \rightarrow 0$ (関連なし), 1 (関連あり)

$q=2 \rightarrow 0$ (関連なし), 1 (やや関連あり), 2 (関連あり)

- 2つの基準により、各文献 s に対し、たとえば以下のようなベクトルを表現

$(1, 0, 0, 1, 1, L , 1, 0, 5)$ ($m+1$ 桁)

$$Z_{m+1}(s) = 0, 1, 2, L , q$$

全文献 S の集合 $S = \bigcup_i \bigcup_k S_{ik}$ (集合の直和) (5.1.3)

$$S_{ik} = \{s \mid (Z_1(s), Z_2(s), L , Z_m(s)) = i, Z_{m+1}(s) = k\}$$

目的は分類子によるカテゴリーを検索すべきか否かを決定する

$$S_i = \bigcup_k S_{ik} \quad (5.1.4)$$

5.2 文献からのサンプリング

- 集合Sをランダム・サンプリングの問題として考える

$$\text{各カテゴリ } S_{ik} \text{ の確率 } \theta_{ik} = P(S_{ik}) \quad (5.2.1)$$

$$i \text{ 番目の文献カテゴリ } S_i \text{ の確率 } P(S_i) = \sum_{k=0}^q \theta_{ik} \quad (5.2.2)$$

$$(\ast \theta_{ik} \geq 0, \sum_{i=0}^p \sum_{k=0}^q \theta_{ik} = 1)$$

- 標本(ランダムサンプル)n点の文献よりrを表現

$$r = (r_{00}, r_{01}, \dots, r_{0q}, r_{10}, \dots, r_{pq}) \quad (5.2.4)$$

$$r \text{ の確率分布 } f(r|\theta) = \prod_{i=0}^p \prod_{k=0}^q \frac{n!}{r_{ik}!} (\theta_{ik})^{r_{ik}} \quad (\text{二項分布}) \quad (5.2.5)$$

5.2 文献からのサンプリング

表 5.2 ある図書館の文献構成

文献 カテゴリー	閲覧	OR	図書館	閲覧・OR	関連性	確率	標本 データ
S_0 {	S_{00}	×	×	×	×	なし	θ_{00} r_{00}
	S_{01}	×	×	×	×	ややあり	θ_{01} r_{01}
	S_{02}	×	×	×	×	あり	θ_{02} r_{02}
S_1 {	S_{10}	×	×	×	○	なし	θ_{10} r_{10}
	S_{11}	×	×	×	○	ややあり	θ_{12} r_{11}
	S_{12}	×	×	×	○	あり	θ_{12} r_{12}
S_{20}	×	×	○	×	なし	θ_{20} r_{20}	
...
$S_{15,2}$	○	○	○	○	あり	$\theta_{15,2}$	$r_{15,2}$
S						1	n (点)

各カテゴリーの確率とランダム・サンプルの結果が与えられている。各カテゴリーの定義は該当 (○)、非該当 (×) で定められている。

5.2 文献からのサンプリング

- 関連性が2段階なら $q=1$ で

$$f(r|\theta) = \prod_{i=0}^p \frac{n!}{r_{i0}! r_{i1}!} (\theta_{i0})^{r_{i0}} (\theta_{i1})^{r_{i1}} \quad (5.2.6)$$

- 文献カテゴリ S_i の検索基準 (関連性あり θ_{i1} , なし θ_{i0})

$$\frac{\theta_{i1}}{\theta_{i0}} \geq c \text{なら, 文献カテゴリ } S_i \text{ を検索} \quad (5.2.7)$$

$$\frac{\theta_{i1}}{\theta_{i0}} < c \text{なら, 文献カテゴリ } S_i \text{ を検索しない}$$

(※ c : 無関係文献を検索してしまうコストによる定数)

すべての θ は値がわかっているとして計算、

そうでなければデータ r から推定