

## 確率分類ベクターマシンを用いた文書分類方式の検討

池上 裕之<sup>†</sup>, 阿部 洋志<sup>†</sup>, 小林 学<sup>†</sup>, 坂下 善彦<sup>†</sup>

<sup>†</sup>湘南工科大学大学院 工学研究科 電気情報工学専攻

### 1. はじめに

文書分類問題はベクトル空間モデルやサポートベクターマシンなど様々な手法により研究されてきた. 近年 Chen らにより確率分類ベクターマシンが提案されており, 基本的な分類問題に対して従来の種々の手法よりも優れた結果が示されている[3]. そこで本研究では確率分類ベクターマシンを拡張し, 文書分類問題に適用することを考える. このときカーネルにベクトル空間モデルの類似度を用いてその有効性の評価を行う.

### 2. 確率分類ベクターマシン

本稿では 2 値の文書分類問題を対象として, Chen らによる確率分類ベクターマシン (以下 PCVM と略す) に新しいパラメータ  $\gamma$  を導入した形で記述する. まず  $i$  番目の入力ベクトル  $\mathbf{x}_i$  とそれに対応するラベル  $t_i \in \{-1, +1\}$  の  $N$  個の組  $\{(\mathbf{x}_i, t_i)\}_{i=1}^N$  を学習データとする. また  $\mathbf{x}_i$  に依存するカーネル関数を  $\phi_i(\mathbf{x})$  と表し,

$$\boldsymbol{\phi}(\mathbf{x}) = (\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_N(\mathbf{x}))^T$$

と定義する. このときパラメータ

$$\mathbf{w} = (w_1, w_2, \dots, w_N)^T$$

と  $b$  を用いて線形識別関数  $y(\mathbf{x})$  を

$$y(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + b \quad (1)$$

と定義する. ここで観測できない隠れ関数  $h(\mathbf{x})$  を考え,

$$h(\mathbf{x}) = y(\mathbf{x}) + \varepsilon \quad (2)$$

とする. ただし観測できない雑音  $\varepsilon$  は平均 0, 精度  $\gamma$  の正規分布  $N(0, \frac{1}{\gamma})$  に従うものとする. このときラベル  $t$  は次式で決定されるモデルを考える.

$$t = \begin{cases} +1, & h(\mathbf{x}) \geq 0 \\ -1, & h(\mathbf{x}) < 0 \end{cases} \quad (3)$$

ここでプロビット関数  $\psi(x)$  を

$$\psi(x) = \int_{-\infty}^x N\left(t \mid 0, \frac{1}{\gamma}\right) dt \quad (4)$$

と定義すると,

Probabilistic Classification Vector Machines for Document Classification  
Hiroyuki IKEGAMI<sup>†</sup>, Hiroshi ABE<sup>†</sup>, Manabu KOBAYASHI<sup>†</sup>  
and Yoshihiko SAKASHITA<sup>†</sup>  
<sup>†</sup>Graduate School of Engineering, Shonan Institute of Technology

$$P(t = 1 | \mathbf{x}, \mathbf{w}, b) = P(y(\mathbf{x}) + \varepsilon \geq 0) = \psi(y(\mathbf{x})) \quad (5)$$

$$P(h(\mathbf{x}) | \mathbf{w}, b) = N\left(h(\mathbf{x}) \mid y(\mathbf{x}), \frac{1}{\gamma}\right) \quad (6)$$

が成り立つ. 学習データに対して

$$\boldsymbol{\Phi} = (\boldsymbol{\phi}(\mathbf{x}_1), \boldsymbol{\phi}(\mathbf{x}_2), \dots, \boldsymbol{\phi}(\mathbf{x}_N))^T \quad (7)$$

$$\mathbf{H} = (h(\mathbf{x}_1), h(\mathbf{x}_2), \dots, h(\mathbf{x}_N))^T \quad (8)$$

と定義すると, 各学習データは独立なため

$$P(\mathbf{H} | \mathbf{w}, b) = \prod_{i=1}^N N\{h(\mathbf{x}_i) | y(\mathbf{x}_i), \frac{1}{\gamma}\} \\ = \left(\frac{\gamma}{\sqrt{2\pi}}\right)^N \exp\left\{-\frac{\gamma}{2} \|\mathbf{H} - \boldsymbol{\Phi}\mathbf{w} - b\mathbf{I}\|^2\right\} \quad (9)$$

が成り立つ. ただし  $\mathbf{I}$  は要素が全て 1 のベクトルを表す. ここで  $\mathbf{w}$  および  $b$  に対する事前分布を

$$P(w_i | \alpha_i) = \begin{cases} 2N(w_i | 0, \alpha_i^{-1}), & y_i x_i \geq 0 \\ 0, & y_i x_i < 0 \end{cases} \quad (10)$$

$$P(b | \beta) = N(b | 0, \beta^{-1}) \quad (11)$$

と仮定する. このとき  $A = \text{diag}(\alpha_1, \alpha_2, \dots, \alpha_N)$  とすると,

$$\log P(\mathbf{w}, b | \mathbf{H}, A, \beta)$$

$$\propto \log P(\mathbf{H} | \mathbf{w}, b) + \log(\mathbf{w} | A) + \log(b | \beta) + \text{定数}$$

$$\propto -\frac{\gamma}{2} \|\mathbf{H} - (\boldsymbol{\Phi}\mathbf{w} + b\mathbf{I})\|^2 - \frac{1}{2} \mathbf{w}^T A \mathbf{w}$$

$$- \frac{1}{2} \beta b^2 + \text{定数}$$

$$= \gamma \mathbf{w}^T \boldsymbol{\Phi}^T (2\mathbf{H} - \boldsymbol{\Phi}\mathbf{w}) - 2\gamma b \mathbf{I}^T \boldsymbol{\Phi}\mathbf{w} + 2\gamma b \mathbf{I}^T \mathbf{H} \\ - \mathbf{w} A \mathbf{w} - \beta b^2 - \gamma b^2 N + \text{定数} \quad (12)$$

が成り立つ. 上式を最大とする  $\mathbf{w}$  および  $b$  を求めるために, EM アルゴリズムを用いる.

①E ステップ

$$Q(\mathbf{w}, b | \mathbf{w}^{old}, b^{old})$$

$$= E_{H, \alpha, \beta} [\log P(\mathbf{w}, b | t, H, \alpha, \beta) | t, \mathbf{w}^{old}, b^{old}]$$

$$= 2\gamma (\mathbf{w}^T \boldsymbol{\Phi}^T + b \mathbf{I}^T) E[H | t, \mathbf{w}^{old}, b^{old}]$$

$$- \mathbf{w}^T E[A | t, \mathbf{w}^{old}, b^{old}] \mathbf{w} - b^2 E[\beta | t, \mathbf{w}^{old}, b^{old}]$$

$$- \gamma \mathbf{w}^T \boldsymbol{\Phi}^T \boldsymbol{\Phi}\mathbf{w} - 2\gamma b \mathbf{I}^T \boldsymbol{\Phi}\mathbf{w} - \gamma b^2 N$$

$$= 2\gamma (\mathbf{w}^T \boldsymbol{\Phi}^T + b \mathbf{I}^T) \bar{H} - \mathbf{w}^T \bar{A} \mathbf{w} - b^2 \bar{\beta}$$

$$- \gamma \mathbf{w}^T \boldsymbol{\Phi}^T \boldsymbol{\Phi}\mathbf{w} - 2\gamma b \mathbf{I}^T \boldsymbol{\Phi}\mathbf{w} - \gamma b^2 N \quad (13)$$

ただし  $\bar{H} = (\bar{h}_1, \bar{h}_2, \dots, \bar{h}_N)^T$  は  $t_i = 1$  のとき

$$\bar{h}_i = E[h_i | t_i, \mathbf{w}^{old}, b^{old}]$$

$$\begin{aligned}
 &= \int h_i \cdot P(h_i|t_i, \mathbf{w}^{old}, b^{old}) dh_i \\
 &= \int_{y(x_i)+\varepsilon \geq 0} (y(x_i) + \varepsilon) N\left(\varepsilon|0, \frac{1}{\gamma}\right) d\varepsilon \\
 &= y(x_i)\psi(y(x_i)) + \frac{1}{\gamma} N\left(y|0, \frac{1}{\gamma}\right)
 \end{aligned} \tag{14}$$

となる。また  $t_i = -1$  のときは

$$\begin{aligned}
 \bar{h}_i &= \int_{y(x_i)+\varepsilon \leq 0} (y(x_i) + \varepsilon) N\left(\varepsilon|0, \frac{1}{\gamma}\right) d\varepsilon \\
 &= y(x_i)\psi(-y(x_i)) - \frac{1}{\gamma} N\left(y|0, \frac{1}{\gamma}\right)
 \end{aligned} \tag{15}$$

である。

②M ステップ

$$\frac{dQ}{d\mathbf{w}} = 2\gamma\Phi^T\bar{H} - 2\bar{A}\mathbf{w} - 2\gamma\Phi\Phi^T\mathbf{w} - 2\gamma b\Phi^T\mathbf{I} = 0 \tag{16}$$

として

$$\mathbf{w}^{new} = (2\bar{A} + 2\gamma\Phi\Phi^T)^{-1} \cdot \gamma\Phi^T(\bar{H} - b\mathbf{I}) \tag{17}$$

により  $\mathbf{w}$  を更新する。同様に

$$\frac{dQ}{db} = 2\gamma\mathbf{I}^T\bar{H} - 2b\bar{\beta} - 2\gamma\mathbf{I}^T\Phi\mathbf{w} - 2b\gamma N = 0 \tag{18}$$

として

$$b^{new} = (\bar{\beta} + \gamma N)^{-1} \gamma\mathbf{I}^T(\bar{H} - \Phi\mathbf{w}) \tag{19}$$

により  $b$  を更新する。

なお本節で  $\gamma = 1$  と置くと、Chen らの PCVM と等価となる。

3. 計算機による実験

本節では文書分類問題に PCVM を実際に用いて評価を行う。ここで  $i$  番目の学習文書  $\mathbf{x}_i$  を  $\mathbf{x}_i = (x_{i,1}, x_{i,2}, \dots, x_{i,M}) \in \{0,1\}^M$  と表現する。ただし  $x_{i,j}$  は  $i$  番目の学習文書に  $j$  番目の単語が存在すれば 1, そうでなければ 0 とする。新規文書  $\mathbf{x}$  についても同様に  $\mathbf{x} = (x_1, x_2, \dots, x_M) \in \{0,1\}^M$  と表記する。なお  $M$  は総単語数である。このとき本節で実験に用いるカーネル関数を

$$\phi_i(\mathbf{x}) = \frac{\mathbf{x}_i \cdot \mathbf{x}}{\|\mathbf{x}_i\| \|\mathbf{x}\|} \tag{20}$$

と定義する。すなわち  $\phi_i(\mathbf{x})$  は  $\mathbf{x}_i$  と  $\mathbf{x}$  の余弦とする。実験用のデータには CD-毎日新聞 94'データ集[4]を利用した。これには毎日新聞の1年分の記事データにラベルを付与したものが収録されている。ここで2カテゴリのデータを選択して2値分類を行った結果を表 1,2 に示す。表 1 は  $\gamma$  の値を変化させたときの2値分類の正分類率(%)を示している。ただし学習データ数を  $N = 400$  とし、評価のためのテストデータ数は 1600 とした。なお PCVM における EM アルゴリズムの繰り返し最大回数は 50 とした。表中カテゴリ Int., So., Sp., Cul., Ec., Ent はそれぞれ国際, 社会, スポーツ, 文化,

経済, 芸能を意味している。また非零の  $w_i$  の個数を調べるために、表 2 に最大繰り返し回数後に  $|w_i| > 10^{-8}$  である  $w_i$  の個数を示した。

表 1 : 各  $\gamma$  に対する 2 値分類の正分類率(%)

$\gamma$	0.05	0.1	0.5	1.0	2.0	4.0
Int./So.	84.3	86.5	86.2	86.1	85.8	85.6
So./Sp.	86.1	88.4	89.5	88.4	86.9	85.6
Cul./Ec.	94.9	95.8	96.9	96.9	97.1	97.4
Ec./Ent	95.3	97.9	97.3	96.1	95.0	94.1

表 2 : 各  $\gamma$  に対する  $|w_i| > 10^{-8}$  である  $w_i$  の個数

$\gamma$	0.05	0.1	0.5	1.0	2.0	4.0
Int./So.	6	10	33	42	62	88
So./Sp.	6	9	36	61	82	109
Cul./Ec.	7	10	24	36	52	82
Ec./Ent	8	10	30	40	59	81

表 1 を見ると、 $\gamma = 0.1$  あるいは  $0.5$  の正分類率が高いことが分かる。ただしカテゴリによって若干の差が見られ、文化と経済の分類(Cul./Ec.) では  $\gamma$  が大きいほど分類精度が良い結果となっている。このように与えられた問題によって最適な精度を与える  $\gamma$  に差があることが分かる。一方非零に収束する  $w_i$  の個数は  $\gamma$  が小さくなるにつれて単調に小さくなる。新規文書の分類は式(1)の  $y(\mathbf{x})$  の正負により判別するため、分類の計算量は非零の  $w_i$  の個数に比例する。従って計算量の点から見ると  $\gamma$  は小さい方が良いことが分かる。本実験による文書分類では  $\gamma$  の値は 0.1~0.5 程度がおおむね良い結果であることが分かる。

4. まとめと今後の課題

本研究では新しいパラメータ  $\gamma$  を導入した PCVM を用いて2値の文書分類の評価を行った。結果的に文書分類では  $\gamma$  の値は 0.1~0.5 程度がおおむね良い結果であった。多カテゴリに対する PCVM の適用と  $\gamma$  の自動決定法などが今後の課題である。

参考文献

[1] V. N. Vapnik, Statistical Learning Theory. New York: Wiley-Interscience, 1998.  
 [2] M. E. Tipping, "Sparse Bayesian learning and the relevance vector machine," J. Mach. Learn. Res., vol. 1, pp. 211-244, 2001.  
 [3] H. Chen, P. Ti and X. Yao, "Probabilistic Classification Vector Machines" IEEE Trans. On Neural Networks, vol. 20, no. 6, pp.902-914, JUNE 2009  
 [4] CD-毎日新聞 94'データ集, 日外アソシエーツ, 1995.