

# 潜在変数モデルに基づく $l_0$ 正則化付き最小二乗アルゴリズム $l_0$ -regularized least squares algorithm based on latent variable model

経営システム工学専攻 浅井 謙輔

## 1 はじめに

$p$  次元の説明変数  $\mathbf{x} := (x_1, \dots, x_p)^\top$  を用いて、従属変数  $y$  を以下の様に予測する線形モデルについて考える。

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

$\mathbf{y} := (y_1, \dots, y_N)^\top$ ,  $X := (\mathbf{x}_1, \dots, \mathbf{x}_N)^\top$  として、残差平方和  $\sum_{i=1}^N \varepsilon_i^2$  を最小にする  $\boldsymbol{\beta} := (\beta_0, \beta_1, \dots, \beta_p)^\top$  を求める問題である最小二乗法は以下の様に書ける。

$$\underset{\boldsymbol{\beta}}{\text{minimize}} \quad \|\mathbf{y} - X\boldsymbol{\beta}\|_2^2$$

実際のデータにおいては、得られたデータに対して、従属変数を説明する特徴量が非常に少ない場合が多い。そこで、制約条件を付けた以下の最適化問題

$$\begin{aligned} &\underset{\boldsymbol{\beta}}{\text{minimize}} \quad \|\mathbf{y} - X\boldsymbol{\beta}\|_2^2 \\ &\text{subject to} \quad \|\boldsymbol{\beta}\|_0 \leq K \end{aligned} \quad (1)$$

を考える。ここで、 $\|\boldsymbol{\beta}\|_0$  は  $\boldsymbol{\beta}$  のゼロでない要素の数を表し、 $l_0$  ノルムと呼ばれる。この問題は、NP 困難であることが知られている。そこで、貪欲法を用いて解を求める方法が一般的である。近年の研究では、(1) の制約条件である  $l_0$  ノルムを2つの凸関数の差で表し、最適化を行う DC アプローチ [1] が存在する。 $l_0$  ノルムを用いた関数は非凸であり、大域的最適解を得られるとは限らない。本論文では、(1) にラグランジュ乗数  $\lambda$  を導入した

$$\underset{\boldsymbol{\beta}}{\text{minimize}} \quad \|\mathbf{y} - X\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_0 \quad (2)$$

を考える。 $l_0$  正則化付き最小化問題である (2) の求解方法として、潜在変数モデルに基づく定式化を行い、EM アルゴリズムないし一般化 EM アルゴリズムに基づいたアルゴリズムによって、組み合わせ計算を行わずに良好な大域的収束性を得ることを目的とする。組み合わせ計算を行わないので、現実的に計算可能な時間で解が得られるという利点がある。

## 2 潜在変数モデルと EM アルゴリズム

パラメータ  $\theta$  とする観測データ  $X$  の確率密度関数は、確率変数  $X$  と観測されない確率変数  $Z$  の同時確率密度

関数の周辺分布であり、

$$f(x|\theta) = \sum_z f(x, z|\theta)$$

と表せる確率モデルについて考える。 $f(x, z|\theta)$  は  $X$  と  $Z$  の同時確率密度関数である。このような確率変数  $Z$  を潜在変数と呼ぶ。EM アルゴリズムは潜在変数を持つ尤度を最大化するアルゴリズムであり、以下の E ステップと M ステップを繰り返す。

1. パラメータの初期値  $\theta^{\text{old}}$  を選ぶ。
2. 【E ステップ】  
潜在変数の事後確率  $p(z|x; \theta^{\text{old}})$  を計算する。
3. 【M ステップ】  
 $Q$  関数を最大にする  $\theta$  を計算し、

$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}})$$

とする。ただし、 $Q$  関数

$$Q(\theta, \theta^{\text{old}}) := \sum_z p(z|x; \theta^{\text{old}}) \log f(x, z|\theta)$$

である。

4. 対数尤度関数またはパラメータの値が収束条件を満たしていれば終了。満たしていなければ、 $\theta^{\text{old}} \leftarrow \theta^{\text{new}}$  として、2. に戻る。

また、 $Q$  関数が増加するようにパラメータ  $\theta$  を更新するアルゴリズムを一般化 EM アルゴリズムという。

## 3 モデルの定式化

### 3.1 潜在変数モデルの導入

データセット  $d := ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N))$  ( $\mathbf{x}_i := (x_{i,1}, x_{i,2}, \dots, x_{i,p})^\top$ ) とする。第1章と同様、 $l_0$  正則化付き最小化問題

$$\underset{\boldsymbol{\beta}}{\text{minimize}} \quad \|\mathbf{y} - X\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_0 \quad (2)$$

を解くことを考える。(2) の解を得るために、0 と 1 を取る変数を追加した以下の問題を解く。

$$\begin{aligned} &\underset{\boldsymbol{\beta}, \mathbf{z}}{\text{minimize}} \quad \|\mathbf{y} - Xb(\mathbf{z})\|_2^2 + \lambda \|\mathbf{z}\|_0 \\ &\text{subject to} \quad \mathbf{z} \in \{0, 1\}^p, \boldsymbol{\beta} \in \mathbb{R}^{p+1} \end{aligned} \quad (3)$$

ただし、

$$b(\mathbf{z}) := (\beta_0, z_1 \beta_1, \dots, z_p \beta_p)^\top$$

である．潜在変数  $\mathbf{Z} = (Z_1, Z_2, \dots, Z_p)^\top$ , ( $\mathbf{Z} \in \{0, 1\}^p$ ) の条件付確率として，データセット  $d$  をとる確率を

$$P(\mathcal{D} = d | \mathbf{Z} = \mathbf{z}) = \exp \left\{ -\frac{1}{\|\mathbf{y}\|_2^2 + \lambda p} (\|\mathbf{y} - X\boldsymbol{\beta}(\mathbf{z})\|_2^2 + \lambda \|\boldsymbol{\beta}(\mathbf{z})\|_0) \right\}$$

とする．ただし， $\lambda (> 0)$  は定数である．潜在変数の確率分布を

$$g(\mathbf{z}) := P(\mathbf{Z} = \mathbf{z}) = \prod_{j=1}^p \pi_j^{z_j} (1 - \pi_j)^{1-z_j} \quad (\pi_j := P(Z_j = 1))$$

とする．すなわち， $Z_1, Z_2, \dots, Z_p$  は独立であり，それぞれベルヌーイ分布に従う．EM アルゴリズムによって， $\log P(\mathcal{D} = d)$  を最大にするパラメータ  $\boldsymbol{\beta}, \boldsymbol{\pi}$  を求める．この問題を定式化すると，以下の通り．

$$\begin{aligned} \max_{\boldsymbol{\beta}, \boldsymbol{\pi}} \quad & \log \sum_{\mathbf{z}} g(\mathbf{z}) \exp \left\{ -\frac{1}{t} (\|\mathbf{y} - Xb(\mathbf{z})\|_2^2 + \lambda \|\mathbf{z}\|_0) \right\} \\ \text{s.t.} \quad & g(\mathbf{z}) = \prod_{j=1}^p \pi_j^{z_j} (1 - \pi_j)^{1-z_j} \\ & \boldsymbol{\pi} \in [0, 1]^p, \boldsymbol{\beta} \in \mathbb{R}^{p+1} \end{aligned} \quad (4)$$

求められた  $\boldsymbol{\beta}$  と  $\boldsymbol{\pi}$  を用いて，係数の推定値を以下のように求める．

$$E[b(\mathbf{Z})] = (\beta_0, \pi_1 \beta_1, \dots, \pi_p \beta_p)^\top$$

ここで，以下の命題が成り立つ．

**命題 3.1** (2) の最適解  $\boldsymbol{\beta}^\bullet$  の集合を  $\mathcal{B}^\bullet$ ，(3) の最適解  $(\boldsymbol{\beta}^*, \mathbf{z}^*)$  の集合を  $(\mathcal{B}^*, \mathcal{Z}^*)$ ，(4) の最適解  $(\boldsymbol{\beta}^*, \boldsymbol{\pi}^*)$  の集合を  $(\mathcal{B}^*, \boldsymbol{\Pi}^*)$  とする．

$$\mathcal{B}^\bullet \subset \mathcal{B}^* = \mathcal{B}^*, \quad \mathcal{Z}^* \subset \boldsymbol{\Pi}^*$$

である．また，(2) の最適解  $\boldsymbol{\beta}^\bullet$  がただ一つするとき，(3) の最適解  $\mathbf{z}^*$  と (4) の最適解  $\boldsymbol{\pi}^*$  もただ一つであり，

$$\mathbf{z}^* = \boldsymbol{\pi}^*$$

を満たし，

$$\boldsymbol{\beta}^\bullet = b^*(\mathbf{z}^*) = b^*(\boldsymbol{\pi}^*)$$

が成り立つ．ただし，

$$\begin{aligned} b^*(\mathbf{z}^*) &= (\beta_0^*, z_1^* \beta_1^*, \dots, z_p^* \beta_p^*)^\top \\ b^*(\boldsymbol{\pi}^*) &= (\beta_0^*, \pi_1^* \beta_1^*, \dots, \pi_p^* \beta_p^*)^\top \end{aligned}$$

とする．

よって，(2) の最適解がただ一つするとき，EM アルゴリズムによって得られた (4) の最適解  $(\boldsymbol{\beta}^*, \boldsymbol{\pi}^*)$  を用いて，(2) の最適解が得られる．(2) の最適解がただ一つかどうかはデータセット  $d$  とハイパーパラメータ  $\lambda$  に依存する．

### 3.2 アルゴリズムの定式化

E ステップにおいて，

$$P(\mathcal{D} = d | \mathbf{Z} = \mathbf{z}) \doteq 1 - \frac{1}{t} (\|\mathbf{y} - X\boldsymbol{\beta}(\mathbf{z})\|_2^2 + \lambda \|\boldsymbol{\beta}(\mathbf{z})\|_0)$$

と近似し，M ステップにおいて，

$$P(\mathbf{Z} = \mathbf{z} | \mathcal{D} = d) \doteq \prod_{j=1}^p P(Z_j = z_j | \mathcal{D} = d)$$

と近似した値を用いた  $Q'$  関数を最大化する．提案するアルゴリズムは以下の通り．

1. パラメータの初期値  $\boldsymbol{\beta}^{\text{old}} = \mathbf{0}, \boldsymbol{\pi}^{\text{old}} = (0.5, \dots, 0.5)^\top$  とする．
2. 【E ステップ】 事後確率 ( $\gamma_j := P(Z_j = 1 | \mathcal{D} = d)$ ) の更新

$$\gamma_j = \pi_j \frac{(\|\mathbf{y}\|_2^2 + \lambda p) - (\|\mathbf{y} - X\boldsymbol{\beta}_j(\boldsymbol{\pi})\|_2^2 + R_j(\boldsymbol{\pi}) + \lambda \boldsymbol{\pi}_j^\top \mathbf{1})}{(\|\mathbf{y}\|_2^2 + \lambda p) - (\|\mathbf{y} - X\boldsymbol{\beta}(\boldsymbol{\pi})\|_2^2 + R(\boldsymbol{\pi}) + \lambda \boldsymbol{\pi}^\top \mathbf{1})} \quad (j = 1, \dots, m)$$

とする．ただし，

$$R(\boldsymbol{\pi}) := \sum_{k=1}^p g(\pi_k) \beta_k^2, \quad R_j(\boldsymbol{\pi}) := \sum_{k \neq j} g(\pi_k) \beta_k^2$$

$$g(\pi_k) := \pi_k (1 - \pi_k) \sum_{i=1}^N x_{i,k}^2, \quad \boldsymbol{\beta}(\boldsymbol{\pi}) := (\beta_0, \pi_1 \beta_1, \dots, \pi_p \beta_p)^\top$$

$$\boldsymbol{\beta}_j(\boldsymbol{\pi}) := (\beta_0, \pi_1 \beta_1, \dots, \beta_j, \dots, \pi_p \beta_p)^\top$$

$$\boldsymbol{\pi}_j := (\pi_1, \pi_2, \dots, 1, \dots, \pi_p)^\top$$

である．

3.  $\gamma_j = 0$  となった変数  $j$  があればそれを取り除く．
4. 【M ステップ】 パラメータ  $\boldsymbol{\beta}, \boldsymbol{\pi}$  の更新  $Q'$  関数

$$Q' = -\frac{1}{\|\mathbf{y}\|_2^2 + \lambda p} (\|\mathbf{y} - X^{(\gamma)} \boldsymbol{\beta}\|_2^2 + R(\boldsymbol{\gamma}) + \lambda \boldsymbol{\gamma}^\top \mathbf{1}) + \boldsymbol{\gamma}^\top \log \boldsymbol{\pi} + (\mathbf{1} - \boldsymbol{\gamma})^\top \log (\mathbf{1} - \boldsymbol{\pi})$$

を最大にする  $\boldsymbol{\beta}, \boldsymbol{\pi}$  を求める．ただし，

$$X^{(\gamma)} := \begin{pmatrix} 1 & \gamma_1 x_{1,1} & \gamma_2 x_{1,2} & \cdots & \gamma_m x_{1,p} \\ 1 & \gamma_1 x_{2,1} & \gamma_2 x_{2,2} & \cdots & \gamma_m x_{2,p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \gamma_1 x_{N,1} & \gamma_2 x_{N,2} & \cdots & \gamma_m x_{N,p} \end{pmatrix}$$

とする．これは解析的に求まり，

$$\begin{aligned} \boldsymbol{\beta}^{\text{new}} &= (X^{(\gamma)\top} X^{(\gamma)} + \text{diag}(g(\gamma_1), \dots, g(\gamma_p)))^{-1} X^{(\gamma)\top} \mathbf{y} \\ \boldsymbol{\pi}^{\text{new}} &= \boldsymbol{\gamma} \end{aligned}$$

である．逆行列の計算が困難であれば，

$$\boldsymbol{\beta}^{\text{new}} = \boldsymbol{\beta}^{\text{old}} + \alpha \frac{\partial Q'(\boldsymbol{\beta}^{\text{old}})}{\partial \boldsymbol{\beta}}$$

とする。ただし、 $\alpha$  は学習係数と呼ばれる定数である。

5.  $\|\pi^{\text{new}} - \pi^{\text{old}}\|_2^2 < \varepsilon$  を満たしていれば終了。満たしていなければ、

$$\beta^{\text{old}} \leftarrow \beta^{\text{new}}, \quad \pi^{\text{old}} \leftarrow \pi^{\text{new}}$$

として、2. に戻る。

## 4 数値実験

### 4.1 前立腺癌データ

前立腺癌のデータ [2] は 97 個のデータであり、8 個の変数がある。変数の詳しい説明については、Stamey et al. [2] を参照して頂きたい。 $\lambda$  を  $2^{-2}, 2^{-1}, 1, \dots, 2^7$  と変化させた時の推定された係数は図 1 の通り。収束を判定する  $\varepsilon$  は  $10^{-16}$  とした。

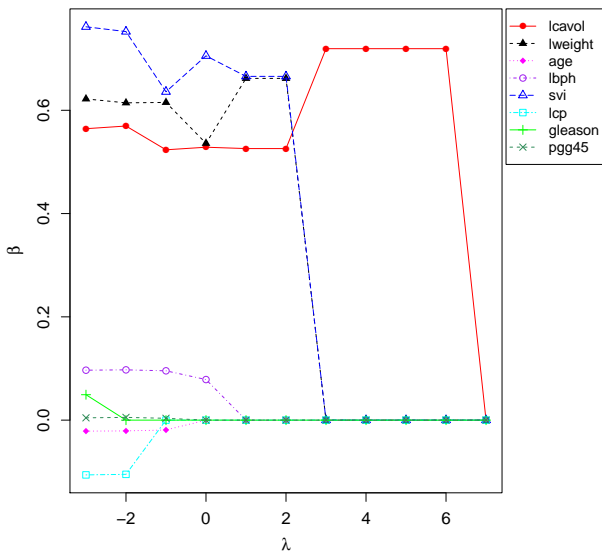


図 1  $\lambda$  を変化させたときの係数の推定値

一番左 ( $\log_2 \lambda = -3$  の位置に対応する) は全ての変数を用いた最小二乗法によって推定された推定値である。各  $\lambda$  において、係数が 0 にならなかった変数のみを用いて最小二乗法を行った場合の推定値と図 1 の値は一致した。 $\lambda = 2^3 \sim 2^6$  において、変数 lcaivol は全て同じ推定値であり、この変数を除いて推定値が 0 になっている。LASSO と違い、 $\lambda$  が違っていても選ばれている変数が同じであれば、推定値は同じ結果になる。これは、解いている問題が (3) である事を考えれば当然と言える。また、各  $\lambda$  について、ステップの経過と潜在変数の確率をプロットしたものが図 2 (凡例は図 1 と同じ) である。1 に収束する速さが遅い潜在変数の確率ほど、 $\lambda$  を大きくしていったとき、他の潜在変数よりも先に 0 に収束して

いる。例えば、 $\lambda = 2^{-1}$  において 1 に収束する速さは pgg45, age, lbph の順に遅く、pgg45 と age については  $\lambda = 1$  において 0 に収束している。また、どの  $\lambda$  においても潜在変数の確率は 0 か 1 付近に早く近づくが、値が変化しなくなるまでには時間が掛かっている。

(縦軸: 潜在変数の確率, 横軸: ステップ数)

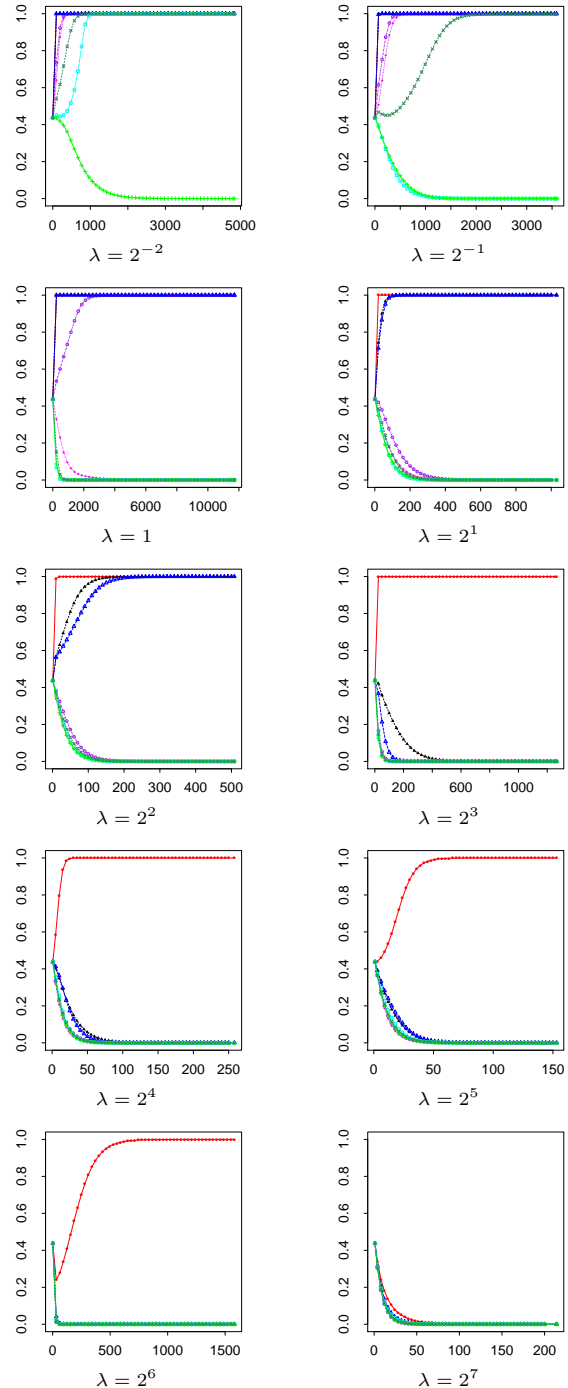


図 2 ステップ数と潜在変数の確率の変化

### 4.2 シミュレーション

人工データによるシミュレーションを行った。 $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$ ,  $\mathbf{x}_i \sim N(\mathbf{0}, \Sigma)$ ,  $\Sigma = (\sigma_{i,j}) = 0.5^{|i-j|}$

とする説明変数と

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \varepsilon_i \sim N(0,1) \quad (i = 1, \dots, n)$$

とする従属変数の標本を生成する。0でない係数  $\beta_j$  は  $p$  個の内 5 個であり, その 5 つの係数  $\beta = 1, 1.5, 2, 2.5, 3$  とする。  $j = 1, \dots, p$  の内どの 5 つかはランダムに決定されるものとする。上記のモデルに従う標本を訓練・検証・テスト用にそれぞれ 50 個ずつ生成する。訓練データで係数を推定し, 検証データを用いてハイパーパラメータ  $\lambda$  を決定する。その後, テスト用データを用いて MSE(平均二乗誤差) を計算しモデルの汎化誤差を評価する。

また, 潜在変数の確率は比較的早い繰り返し回数で 0,1 に近づき,  $\boldsymbol{\beta}$  は 0 でない係数のみを用いた最小二乗法の解となった 4.1 の結果を踏まえて, アルゴリズムの手順 5 を以下の様に変更する。

5. 以下の (a),(b) どちらかの条件を満たしていれば, 終了とする。

(a)  $\|\boldsymbol{\pi}^{\text{new}} - \boldsymbol{\pi}^{\text{old}}\|_2^2 < \varepsilon$

(b) E ステップと M ステップを  $I$  回以上繰り返している。

また, (b) で終了した場合,  $\boldsymbol{\pi}$  の値を小数点以下を四捨五入した 0,1 の値とし,  $\boldsymbol{\beta}$  の予測値を  $\boldsymbol{\pi} = 1$  となった変数のみで最小二乗法を行った予測値

$$\boldsymbol{\beta} = (X^T X)^{-1} X^T \mathbf{y}$$

とする。(  $\boldsymbol{\pi} = 1$  となった変数の数  $> n$  の場合, 終了せずに 2. に戻る。 )  
満たしていなければ,

$$\boldsymbol{\beta}^{\text{old}} \leftarrow \boldsymbol{\beta}^{\text{new}}, \quad \boldsymbol{\pi}^{\text{old}} \leftarrow \boldsymbol{\pi}^{\text{new}}$$

として, 2. に戻る。

150 個の標本を 200 回発生させて, 真の説明変数 5 つを過不足なく選んだ割合は表 1 の通りとなった。

表 1 正答率 (%)

$p$	$(l_0)$				$(l_1)$	$(l_2)$
	EM:500	EM:1000	EM:2000	GEM:5000	LASSO	Ridge
10	80.5	80.5	81.5	76.0	3.5	0.0
50	87.0	93.0	95.0	92.5	0.0	0.0
100	87.0	89.5	91.0	92.5	0.0	0.0

EM:500,1000,2000 はそれぞれ EM アルゴリズムに基づいたモデル (手順 4 において, 逆行列を計算する) において, 終了条件の繰り返し回数  $I = 500, 1000, 2000$  とした場合である。 GEM:5000 は一般化 EM アルゴリズムに基づいたモデル (手順 4 において, 逆行列を計算せず,  $\boldsymbol{\beta}$  を勾配方向に更新する) において, 終了条件の繰

り返し回数  $I = 5000$  とした場合である。  $\boldsymbol{\beta}$  の更新の際の学習係数  $\alpha = 1$  とした。 また,  $p = 10, 50, 100$  は真の説明変数 5 つとダミー変数を含めた説明変数の総数である。 Tibshirani [3] のシミュレーション結果からも分かる通り, LASSO は 0 でない説明変数を含む解を選ぶが過不足なく選ぶことは難しい。 Ridge 回帰については罰則化項の性質上, 完全に係数を 0 にする訳では無い。 また, どの説明変数を選ぶかは訓練データと検証データに依存しているため, アルゴリズムによって大域的最適解が得られている場合でも, 必ずしも真の説明変数 5 つを選ぶわけでは無いことに注意する。 しかし, 訓練データにおいて (2) の大域的最適解が得られ, 検証データにおいて適切な  $\lambda$  を選択出来ていれば, 真の説明変数を選ぶことが期待される。 表 1 より, EM アルゴリズムに基づいたモデルはどの繰り返し回数, 説明変数の数においても 80(%) を超える正答率を得られた。

## 5 おわりに

本論文では, 線形モデルに関する  $l_0$  正則化付き最小化問題に対して, 潜在変数モデルに基づいた定式化を行った。 また, 線形モデルに関する  $l_0$  正則化付き最小化問題の最適解が潜在変数モデルに基づいた定式化によって得られる最適解に含まれることを示し, 最適解を得るアルゴリズムを提案した。 数値実験の結果, 従来の LASSO と比べて高い割合で真の変数を過不足なく選び, 良好な結果を確認できた。

## 参考文献

- [1] Hoai An Le Thi, T Pham Dinh, Hoai Minh Le, and Xuan Thanh Vo. Dc approximation approaches for sparse optimization. *European Journal of Operational Research*, Vol. 244, No. 1, pp. 26–46, 2015.
- [2] Thomas A Stamey, John N Kabalin, John E McNeal, Iain M Johnstone, Fuad Freiha, Elise A Redwine, and Norman Yang. Prostate specific antigen in the diagnosis and treatment of adenocarcinoma of the prostate. ii. radical prostatectomy treated patients. *The Journal of urology*, Vol. 141, No. 5, pp. 1076–1083, 1989.
- [3] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, Vol. 58, No. 1, pp. 267–288, 1996.