

# 機械学習とデータマイニング レポート課題

28G23027 川原尚己

## 課題 A

[問 1] 以下の(1)~(3)式にて省略されている計算手順を明記すること.

$$\arg \min_{\theta \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n (\theta - X_i)^2 = \frac{1}{n} \sum_{i=1}^n X_i \quad (1)$$

$$\mathbf{E} \left( \frac{1}{n} \sum_{i=1}^n (X_i - \mathbf{E}_{\mu} X) \right)^2 = \frac{1}{n^2} \sum_{i=1}^n \mathbf{E} (X_i - \mathbf{E}_{\mu} X)^2 \quad (2)$$

$$\frac{1}{n^2} \sum_{i=1}^n \mathbf{E} (X_i - \mathbf{E}_{\mu} X)^2 = \frac{\text{var}_{\mu} X}{n} \quad (3)$$

[問 2] 式(4)の主張が任意の  $\epsilon > 0$  と  $\delta > 0$  に対して成り立つことを示すこと.

$$n \geq \frac{\text{var}_{\mu} X}{\epsilon} \left( \frac{1}{\delta} \right) \Rightarrow \mathbf{P}\{R(\bar{X}_n) - R^* > \epsilon\} \leq \delta \quad (4)$$

## 回答

[問 1]

$f(\theta)$  を式(A1)のように定義する.

$$f(\theta) := \frac{1}{n} \sum_{i=1}^n (\theta - X_i)^2 \quad (A1)$$

このとき,  $\frac{df}{d\theta}$  は式(A2)のように表される.

$$\frac{df}{d\theta} = \frac{1}{n} \left( 2n\theta - \sum_{i=1}^n X_i \right) \quad (A2)$$

$\frac{df}{d\theta} = 0$  となる  $\theta = \theta^*$  は式(A3)のようになる.

$$\theta^* = \frac{1}{n} \sum_{i=1}^n X_i \quad (A3)$$

$f(\theta)$  は下に凸な放物線であるから,  $\frac{df}{d\theta} = 0$  となる  $\theta$  において最小値をとる. 以上より, 式

(1)が導出された.

次に式(2)を示す.  $X_i \sim \mu$  であるから,  $\mathbf{E}[X_i] = \mathbf{E}_{\mu} X$  が成り立つ. 式(2)の左辺を直接変形することにより, 以下の等式を得る.

$$\begin{aligned}
\mathbf{E}\left(\frac{1}{n}\sum_{i=1}^n(X_i - \mathbf{E}_\mu X)\right)^2 &= \frac{1}{n^2}\mathbf{E}\left(\sum_{i=1}^n(X_i - \mathbf{E}_\mu X)\right)^2 \\
&= \mathbf{E}\left(\sum_{i=1}^n(X_i - \mathbf{E}_\mu X)^2 + 2\sum_{i=1}^n\sum_{j>i}^n(X_i - \mathbf{E}_\mu X)(X_j - \mathbf{E}_\mu X)\right) \tag{A4}
\end{aligned}$$

$X_i, X_j (i \neq j)$  は独立であるから、 $\mathbf{E}\left((X_i - \mathbf{E}_\mu X)(X_j - \mathbf{E}_\mu X)\right) = \mathbf{E}(X_i - \mathbf{E}_\mu X)\mathbf{E}(X_j - \mathbf{E}_\mu X) = 0$  が

成り立つ。よって、式(A4)の第二項は0に等しく、式(2)が得られる。

最後に式(3)を示す。

$$\begin{aligned}
\mathbf{E}(X_i - \mathbf{E}_\mu X)^2 &= \mathbf{E}\left(X_i^2 - 2X_i \mathbf{E}_\mu X + (\mathbf{E}_\mu X)^2\right) \\
&= \mathbf{E}_\mu(X_i^2) - (\mathbf{E}_\mu X)^2 \\
&= \mathbf{E}_\mu(X^2) - (\mathbf{E}_\mu X)^2
\end{aligned}$$

であり、 $\mathbf{E}(X_i - \mathbf{E}_\mu X)^2$  は分布 $\mu$ に対しての期待値で表せるから式(A5)が成り立つ。

$$\mathbf{E}(X_i - \mathbf{E}_\mu X)^2 = \mathbf{E}_\mu(X - \mathbf{E}_\mu X)^2 \tag{A5}$$

以上より、式(3)の左辺は式(A6)のように表される。

$$\frac{1}{n^2}\sum_{i=1}^n \mathbf{E}(X_i - \mathbf{E}_\mu X)^2 = \frac{1}{n}\mathbf{E}_\mu(X - \mathbf{E}_\mu X)^2 \tag{A6}$$

「学習問題の具体例」の式(4)より、 $\mathbf{E}_\mu(X - \mathbf{E}_\mu X)^2 = \text{var}_\mu X$ であるから、式(3)が得られる。

## [問2]

「学習問題の具体例」の8ページにおいて、 $R(\bar{X}_n) \geq 0$ であり、 $R^*$ は $R(\bar{X}_n)$ の下限であるから、 $R(\bar{X}_n) - R^* \geq 0$ が成り立つ。よって、Markovの不等式より式(A7)がなりたつ。

$$\mathbf{P}\{R(\bar{X}_n) - R^* > \epsilon\} \leq \frac{1}{\epsilon}\mathbf{E}[R(\bar{X}_n) - R^*] \tag{A7}$$

「学習問題の具体例」の7ページより $\mathbf{E}[R(\bar{X}_n) - R^*] = \frac{\text{var}_\mu X}{n}$ が成り立つから、式(A8)が得られる。

$$\mathbf{P}\{R(\bar{X}_n) - R^* > \epsilon\} \leq \frac{\text{var}_\mu X}{n\epsilon} \tag{A8}$$

このとき、任意の $\delta > 0$ に対し、 $n \geq \frac{\text{var}_\mu X}{\epsilon} \left(\frac{1}{\delta}\right)$ なる整数 $n$ が存在する。そのような $n$ をとると

き、式(A8)より式(A9)が成り立つ。

$$\mathbf{P}\{R(\bar{X}_n) - R^* > \epsilon\} \leq \delta \quad (A9)$$

以上より、式(4)が $\epsilon > 0$ と $\delta > 0$ に対して成り立つことを示せた。

## 課題 B

資料中の「分布の位置推定」の話に出てくる学習アルゴリズムはどのような一致性を満たすか。「学習法の一致性」の定義に出てくる $R, R_{\text{con}}^*, H_n$ はそれぞれ何に相当し、なぜ一致性が約束できるか示すこと。

## 回答

「分布の位置推定」の話での学習アルゴリズムにおいては $R, R_{\text{con}}^*, H_n$ はそれぞれ以下の事柄に相当する。

$R : X \sim \mu$ からの平均二乗誤差

$R_{\text{con}}^* : X \sim \mu$ からの平均二乗誤差の下限

$H_n : n$ この標本の平均

このとき、 $R, R_{\text{con}}^*, H_n$ は式(B1)~(B3)のように表される。

$$R(\theta) = \mathbf{E}_{\mu}(\theta - X)^2 \quad (B1)$$

$$R_{\text{con}}^* = \inf_{\theta \in \mathbb{R}} R(\theta) \quad (B2)$$

$$H_n = \frac{1}{n} \sum_{i=1}^n X_i \quad (B3)$$

この $R, R_{\text{con}}^*, H_n$ に対して式(A8)が成り立つから、 $n \rightarrow \infty$ に対する極限を考えることで一致性を満たすことが示せた。

## 課題 C

式(5)の導出過程を明記せよ。

$$N_{\epsilon, \delta}^* \leq \frac{\log(|\mathcal{H}|) + \log\left(\frac{1}{\delta}\right)}{\epsilon} \quad (5)$$

## 回答

式(C1)のように PAC 条件を考える。

$$\mathbf{P}(R(H_n) - R_{\text{con}}^* > \epsilon) \leq \delta \quad (C1)$$

$R_{\text{con}}^* = 0$ ,  $R(h) = \mathbf{P}\{h(X) \neq Y\}$ , 「サル真似」における $n+1$ 回目の更新の出力を $H_{\text{con}}^*$ として式(C1)に代入すると、式(C2)を得る。

$$\mathbf{P}\{\mathbf{P}\{H_n^{\text{con}}(X) \neq Y\} > \epsilon\} \leq \delta \quad (C2)$$

$\mathbf{P}\{\mathbf{P}\{H_n^{\text{con}}(X) \neq Y\} > \epsilon\}$ について評価を行うために、以下のような集合 $\mathcal{H}_{\text{bad}}$ を考える。

$$\mathcal{H}_{\text{bad}}(\epsilon) := \{h \in \mathcal{H} : \mathbf{P}\{h(X) \neq Y\} > \epsilon\} \quad (C3)$$

$H_n^{\text{con}}$ が $\mathcal{H}_{\text{bad}}(\epsilon)$ に含まれる確率は全ての $X_i$ を正しく分類できる $h \in \mathcal{H}_{\text{bad}}(\epsilon)$ よりも小さいこ

とから、式(C4)を得る.

$$P\{H_n^{con} \in \mathcal{H}_{bad}(\epsilon)\} \leq |\mathcal{H}_{bad}(\epsilon)|(1 - \epsilon)^n \quad (C4)$$

さらに、 $H_n^{con}(X) \neq Y$ をとる確率が $\epsilon$ より大きくなる確率が $P\{H_n^{con} \in \mathcal{H}_{bad}(\epsilon)\}$ より小さいことから、式(C5)を得る.

$$\begin{aligned} \mathbf{P}\{\mathbf{P}\{H_n^{con}(X) \neq Y\} > \epsilon\} &\leq P\{H_n^{con} \in \mathcal{H}_{bad}(\epsilon)\} \\ &\leq |\mathcal{H}| \exp(-n\epsilon) \end{aligned} \quad (C5)$$

式(C1)及び式(C5)より、式(C6)を満たす任意の $n$ で PAC 条件を満たすことがわかる.

$$\begin{aligned} |\mathcal{H}| \exp(-n\epsilon) &\leq \delta \\ \Leftrightarrow n &\geq \frac{\log(|\mathcal{H}|) + \log\left(\frac{1}{\delta}\right)}{\epsilon} \end{aligned} \quad (C6)$$

以上より、標本複雑度  $N_{\epsilon, \delta}^*$  は PAC 条件を満たす $n$ の最小値であるから、式(5)を満たす.

## 課題 D

二値分類の場合、99%の確率で $h^*(X) = Y$ 、残り1%の確率でラベルが反転して $h^*(X) \neq Y$ となるような確率ノイズ  $U$ を設計し、 $Y$ と $X$ と $U$ と $h^*$ を含む上記のような等式を明記した上で、その設計の妥当性を説明せよ.

## 回答

以下のような確率ノイズ $U$ を考える.

$$U(r) = \begin{cases} 1, & (0 \leq r < 0.99) \\ -1, & (0.99 \leq r < 1) \end{cases} \quad (D1)$$

ただし $r$ は区間 $[0,1]$ の一様分布に従う確率変数である.

このとき、 $Y = h^*(X)U(r)$ という $Y$ を考えると、99% の確率で正しいラベルが出力され、1% の確率でラベルが反転して出力されるため、題意を満たす設計であり、妥当性がある.

## 課題 E

以下の3つの不等式が成り立つことを図や数式などを用いて説明せよ.

1. 実数直線上の有限区間において

$$\text{shatter}(3) < 2^3 \quad (6)$$

2. 二次元平面における識別線において

$$\text{shatter}(4) < 2^4 \quad (7)$$

3.  $\mathbb{R}^d$ における長方形において

$$\text{shatter}(2d + 1) < 2^{2d+1} \quad (8)$$

## 回答

任意の $n$ に対して、式(E1)が成り立つ.

$$\text{shatter}(n) \leq 2^n \quad (E1)$$

式(E1)の統合が成立するのは、任意の二値ラベルの組み合わせでも入力データを分類可能な分類面が存在するときである. そのため、今回の課題においては分類不可能な入力データと二値ラベルの組を少なくとも一つ例示すればよい.

1.

実数数直線上において、図1のような3点を分類できる有限区間は存在しないため、式(6)は成立する.

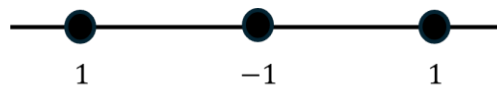


図1 実数数直線上において分類できない例

2.

二次元平面において、図2のような4点を分類できる識別線は存在しないため、式(7)は成立する.

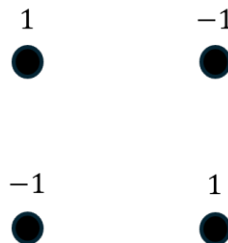


図2 二次元平面上において分類できない例

3.

どの $d+1$ 点も $(d-1)$ 次元長平面上にない場合を考えればよい.

$(2d+1)$ この点のうち、ある軸において最小値または最大値であるような点は丁度 $2d$ 個存在する. これら $2d$ 個の点を頂点とする領域 $D$ を考える.  $D$ は、 $d$ 個の軸に対する最小値または最大値をすべて含んでいるため、残りの一点は必ず $D$ に含まれることになる.  $D$ の頂点であるようなすべての点に同一のラベルを、残りの一点にもう一方のラベルを割り当てることを考えると、このような分類は実現できない. よって、式(8)は成立する.

## 課題 F

以下の等式 (9) と (10) がなぜ成立するか, 具体的に説明すること.

$$\mathbf{E}[R(H_T) - R(h^*)] = \mathbf{E}(\mathbf{E}[R(H_T) - R(h^*)|Z_T]) \quad (9)$$

$$\mathbf{E}(\mathbf{E}[R(H_T) - R(h^*)|Z_T]) = \frac{1}{T} \sum_{t=1}^T \mathbf{E}[R(H_t) - R(h^*)] \quad (10)$$

## 回答

式(9)の左辺は分布 $\mu$ から得られた独立標本 $Z_1, \dots, Z_T$ を用いて分類器 $H_1, \dots, H_T$ を学習する. その後,  $\text{Uniform}\{1, \dots, T\}$ よりランダムに選択した $T$ に対し,  $H_T$ を用いた場合の汎化性能の期待値を表している. 一方で, 式(9)の右辺は学習データを固定した場合において得られた $H_T$ の汎化性能の期待値に対し, さらに学習データ全体において期待値をとったものである. 任意の学習データは分布 $\mu$ に従って生成されることから, 学習データを固定して学習を行ったかどうかにかかわらず, 学習によって得られた分類器 $H_T$ の汎化性能の期待値は両者の間で等しくなる.

また, 式(10)の右辺の $Z_T$ は $R(H_T) - R(h^*)$ とは独立した事象であるから, 式(F1)のように表せる.

$$\begin{aligned} \mathbf{E}(\mathbf{E}[R(H_T) - R(h^*)|Z_T]) &= \sum_{Z_T} \mathbf{E}[R(H_T) - R(h^*)]P(Z_T) \\ &= \sum_{Z_T} \sum_{t=1}^T \mathbf{E}[R(H_t) - R(h^*)]P(Z_T) \end{aligned} \quad (F1)$$

学習データとして $Z_T$ が選ばれる確率は $\frac{1}{T}$ であるから, 最終的に(F2)式が得られる.

$$\mathbf{E}(\mathbf{E}[R(H_T) - R(h^*)|Z_T]) = \frac{1}{T} \sum_{t=1}^T \mathbf{E}[R(H_t) - R(h^*)] \quad (F2)$$

以上より, 式(9)及び式(10)が成立することが示された.