

パターン認識 Pattern Recognition

佐藤真一
Shin'ichi Satoh

国立情報学研究所
National Institute of Informatics

May 16, 2023

宿題について

- 5月9-10日のITC LMSのシステム障害のため、5月9日締め切りであった4月25日出題の宿題の締め切りを一週間延長し5月16日とすることとした。
- よって直近の宿題の締め切りは以下の通り
 - 4月25日出題分、締め切り5月9日 → **5月16日に延長**
 - 5月2日出題分、締め切り5月16日
 - 5月9日宿題なし
 - 5月16日(本日)出題分、締め切り5月30日
- PDFフォーマットにてITC LMSに提出すること。
- Jupyter Notebook (ipynb) や Python (py) 形式の提出はおやめください。

今日の予定

- サンプルからの誤り率推定
- 再代入法 (R 法)
- 分割学習法 (H 法)
- 交差確認法 (CV 法)
- ブートストラップ法

識別誤り率推定

簡単のため、2クラス分類について考える。識別器は以下の関数により定義できる

$$h(X) \underset{\omega_2}{\overset{\omega_1}{\leq}} 0$$

ただし $h(X)$ はベクトル X の識別関数。この識別器の識別誤り率は

$$\varepsilon_1 = \int_{h(X) > 0} p_1(X) dX = \int u(h(X)) p_1(X) dX$$

ただし $u(\cdot)$ はステップ関数、 $p_i(X)$ はクラス i の確率密度関数。

識別誤り率推定

いまフーリエ変換とフーリエ逆変換を考える。

$$\mathcal{F}[x(t)] = X(\omega) = \int_{-\infty}^{\infty} x(t)e^{-j\omega t} dt$$

$$\mathcal{F}^{-1}[X(\omega)] = x(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} X(\omega)e^{j\omega t} d\omega$$

ステップ関数 $u(t)$ は

$$u(t) = \begin{cases} 0 & t < 0 \\ \frac{1}{2} & t = 0 \\ 1 & t > 0 \end{cases}$$

そのフーリエ変換

$$\begin{aligned}
 \mathcal{F}[u(t)] &= \frac{1}{2} \mathcal{F}[1] + \int_0^{\infty} e^{-j\omega t} dt \\
 &= \frac{1}{2} (2\pi\delta(\omega)) - \frac{1}{j\omega} [e^{-j\omega t}]_0^{\infty} \\
 &= \pi\delta(\omega) + \frac{1}{j\omega}
 \end{aligned}$$

そのフーリエ逆変換

$$u(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} [\pi\delta(\omega) + \frac{1}{j\omega}] e^{j\omega t} d\omega$$

識別誤り率推定

よって

$$\begin{aligned}\varepsilon_1 &= \int_{h(X) > 0} p_1(X) dX = \int u(h(X)) p_1(X) dX \\ &= \frac{1}{2\pi} \iint \left[\pi \delta(\omega) + \frac{1}{j\omega} \right] e^{j\omega h(X)} p_1(X) d\omega dX \\ &= \frac{1}{2} + \frac{1}{2\pi} \iint \frac{e^{j\omega h(X)}}{j\omega} p_1(X) d\omega dX\end{aligned}$$

識別誤り率推定

同様に

$$\varepsilon_2 = \int_{h(X) < 0} p_2(X) dX = \frac{1}{2} - \frac{1}{2\pi} \iint \frac{e^{j\omega h(X)}}{j\omega} p_2(X) d\omega dX.$$

全体の識別誤り率は

$$\begin{aligned} \varepsilon &= P_1 \varepsilon_1 + P_2 \varepsilon_2 \\ &= \frac{1}{2} + \frac{1}{2\pi} \iint \frac{e^{j\omega h(X)}}{j\omega} \tilde{p}(X) d\omega dX \end{aligned}$$

ただし

$$\tilde{p}(X) = P_1 p_1(X) - P_2 p_2(X).$$

標本からの誤り率推定

有限の標本 (サンプル) のみ利用可能な場合の「誤り数計数法」の正当性について考える。
その場合、 $p_i(X)$ は以下により置き換えられる

$$\hat{p}_i(X) = \frac{1}{N_i} \sum_{j=1}^{N_i} \delta(X - \mathbf{x}_j^{(i)})$$

ただし $\mathbf{x}_j^{(i)}$ は $p_i(X)$ から抽出した N_i 個の評価用の標本 (test sample) である。

標本からの誤り率推定

そうすると誤り率の推定は

$$\begin{aligned}\hat{\varepsilon} &= \frac{1}{2} + \frac{1}{2\pi} \iint \frac{e^{j\omega h(X)}}{j\omega} \left[\frac{P_1}{N_1} \sum_{j=1}^{N_1} \delta(X - \mathbf{x}_j^{(1)}) - \frac{P_2}{N_2} \sum_{j=1}^{N_2} \delta(X - \mathbf{x}_j^{(2)}) \right] d\omega dX \\ &= \frac{1}{2} + \frac{P_1}{N_1} \sum_{j=1}^{N_1} \alpha_j^{(1)} - \frac{P_2}{N_2} \sum_{j=1}^{N_2} \alpha_j^{(2)}\end{aligned}$$

ただし

$$\alpha_j^{(i)} = \frac{1}{2\pi} \int \frac{e^{j\omega h(\mathbf{x}_j^{(i)})}}{j\omega} d\omega = \frac{\text{sign}(h(\mathbf{x}_j^{(i)}))}{2}$$

標本からの誤り率推定

これにより

$$\begin{aligned}\frac{1}{N_1} \sum_{j=1}^{N_1} \alpha_j^{(1)} &= \frac{1}{2N_1} [(\# \omega_1\text{-errors}) - (\# \omega_1\text{-corrects})] \\ &= \frac{1}{N_1} (\# \omega_1\text{-errors}) - \frac{1}{2}.\end{aligned}$$

同様に

$$\frac{1}{N_2} \sum_{j=1}^{N_2} \alpha_j^{(2)} = -\frac{1}{N_2} (\# \omega_2\text{-errors}) + \frac{1}{2}.$$

まとめると

$$\hat{\epsilon} = P_1 \frac{(\# \omega_1\text{-errors})}{N_1} + P_2 \frac{(\# \omega_2\text{-errors})}{N_2}.$$

推定誤り率の平均と分散

推定誤りの評価用標本に対する平均と分散を求めよう。

$$\begin{aligned} E_t\{\alpha_j^{(i)}\} &= \bar{\alpha}_i = \frac{1}{2\pi} \iint \frac{e^{j\omega h(\mathbf{X})}}{j\omega} p_i(\mathbf{X}) d\omega d\mathbf{X} \\ &= \begin{cases} \varepsilon_1 - \frac{1}{2} & (i = 1) \\ \frac{1}{2} - \varepsilon_2 & (i = 2) \end{cases} \end{aligned}$$

$$\begin{aligned} E_t\{\alpha_i^{(i)2}\} &= E_t\left\{\left[\frac{1}{2\pi} \int \frac{e^{j\omega h(\mathbf{X})}}{j\omega} d\omega\right]^2\right\} = E_t\left\{\left[\frac{1}{2} \text{sign}(h(\mathbf{X}))\right]^2\right\} \\ &= \frac{1}{4} \end{aligned}$$

$$E_t\{\alpha_j^{(i)} \alpha_\ell^{(k)}\} = \bar{\alpha}_i \bar{\alpha}_k \quad (i \neq k \text{ or } j \neq \ell)$$

推定誤り率の平均と分散

まとめると

$$\begin{aligned} E_t\{\hat{\varepsilon}\} &= \frac{1}{2} + P_1\bar{\alpha}_1 - P_2\bar{\alpha}_2 \\ &= \frac{1}{2} + P_1\left(\varepsilon_1 - \frac{1}{2}\right) - P_2\left(\frac{1}{2} - \varepsilon_2\right) = \varepsilon \\ \text{Var}_t\{\hat{\varepsilon}\} &= \frac{P_1^2}{N_1} \text{Var}_t\{\alpha_j^{(1)}\} + \frac{P_2^2}{N_2} \text{Var}_t\{\alpha_j^{(2)}\} \\ &= \frac{P_1^2}{N_1} \left[\frac{1}{4} - \left(\varepsilon_1 - \frac{1}{2}\right)^2\right] + \frac{P_2^2}{N_2} \left[\frac{1}{4} - \left(\frac{1}{2} - \varepsilon_2\right)^2\right] \\ &= P_1^2 \frac{\varepsilon_1(1 - \varepsilon_1)}{N_1} + P_2^2 \frac{\varepsilon_2(1 - \varepsilon_2)}{N_2}. \end{aligned}$$

すなわち、 $\hat{\varepsilon}$ は $h(X)$ によらず不偏 (unbiased) で一致性のある (consistent) 誤り率推定である。

別解

$\hat{\tau}_i$ をクラス i について識別を誤った標本数とする。確率変数 $\hat{\tau}_1$ と $\hat{\tau}_2$ は独立であり、二項分布に従うので、

$$\begin{aligned} Pr\{\hat{\tau}_1 = \tau_1, \hat{\tau}_2 = \tau_2\} &= \prod_{i=1}^2 Pr\{\hat{\tau}_i = \tau_i\} \\ &= \prod_{i=1}^2 \binom{N_i}{\tau_i} \varepsilon_i^{\tau_i} (1 - \varepsilon_i)^{N_i - \tau_i}. \end{aligned}$$

ω_i 誤り、すなわち ε_i は、 $\frac{\hat{\tau}_i}{N_i}$ として推定できる

$$\hat{\varepsilon} = \sum_{i=1}^2 P_i \frac{\hat{\tau}_i}{N_i}.$$

二項分布の平均と分散から、

$$E\{\hat{\varepsilon}\} = P_1\varepsilon_1 + P_2\varepsilon_2 = \varepsilon$$
$$\text{Var}\{\hat{\varepsilon}\} = P_1^2 \frac{\varepsilon_1(1-\varepsilon_1)}{N_1} + P_2^2 \frac{\varepsilon_2(1-\varepsilon_2)}{N_2}$$

ベイズ誤り率の上界と下界

もし有限個の標本しか利用できない場合、われわれはどれを学習に用いるか、どれを評価に用いるかに関し、妥協を強いられる。

ここで識別誤りを学習セット (design/training set) と評価セット (test set) という二つの集合の関数として考える

$$\varepsilon(\mathcal{P}_D, \mathcal{P}_T)$$

ただし \mathcal{P} はクラスごとの分布を表す。例えば

$$\mathcal{P} = \{p_1(X), p_2(X)\}.$$

ベイズ誤り率の上界と下界

もし識別器が評価用のデータの分布に対するベイズ識別器ならば、識別誤りは最小となる

$$\varepsilon(\mathcal{P}_T, \mathcal{P}_T) \leq \varepsilon(\mathcal{P}_D, \mathcal{P}_T).$$

真の \mathcal{P} に対するベイズ誤り率は $\varepsilon(\mathcal{P}, \mathcal{P})$ となる。真の \mathcal{P} は得られないため、有限個の標本による推定 $\hat{\mathcal{P}} = (\hat{\mathbf{p}}_1(X), \hat{\mathbf{p}}_2(X))$ に基づく $\varepsilon(\mathcal{P}, \mathcal{P})$ の上界と下界について考えよう。

$$\varepsilon(\mathcal{P}, \mathcal{P}) \leq \varepsilon(\hat{\mathcal{P}}, \mathcal{P})$$

$$\varepsilon(\hat{\mathcal{P}}, \hat{\mathcal{P}}) \leq \varepsilon(\mathcal{P}, \hat{\mathcal{P}}).$$

誤り数計数法は不偏で一致性があることから、

$$\varepsilon(\hat{\mathcal{P}}, \mathcal{P}) = E_{\mathcal{P}_T} \{ \varepsilon(\hat{\mathcal{P}}, \mathcal{P}_T) \}$$

ただし $\hat{\mathcal{P}}_T$ は $\hat{\mathcal{P}}$ とは独立の別の集合。

ベイズ誤り率の上界と下界

再掲:

$$\varepsilon(\mathcal{P}, \mathcal{P}) \leq \varepsilon(\hat{\mathcal{P}}, \mathcal{P})$$

$$\varepsilon(\hat{\mathcal{P}}, \hat{\mathcal{P}}) \leq \varepsilon(\mathcal{P}, \hat{\mathcal{P}})$$

$$\varepsilon(\hat{\mathcal{P}}, \mathcal{P}) = E_{\hat{\mathcal{P}}_T} \{ \varepsilon(\hat{\mathcal{P}}, \hat{\mathcal{P}}_T) \}.$$

同様に、

$$E \{ \varepsilon(\mathcal{P}, \hat{\mathcal{P}}) \} = \varepsilon(\mathcal{P}, \mathcal{P}).$$

まとめると

$$E \{ \varepsilon(\hat{\mathcal{P}}, \hat{\mathcal{P}}) \} \leq \varepsilon(\mathcal{P}, \mathcal{P}) \leq E_{\hat{\mathcal{P}}_T} \{ \varepsilon(\hat{\mathcal{P}}, \hat{\mathcal{P}}_T) \}.$$

分割学習法 (Holdout/H 法)

\mathcal{D} から独立した二つの標本集合 $\hat{\mathcal{D}}$ と $\hat{\mathcal{D}}_T$ を得れば、 $\epsilon(\hat{\mathcal{D}}, \hat{\mathcal{D}}_T)$ は $\hat{\mathcal{D}}$ を学習セット、 $\hat{\mathcal{D}}_T$ を評価セットとして用いれば算出可能。

すでに見てきているように、 $E_{\hat{\mathcal{D}}_T}\{\epsilon(\hat{\mathcal{D}}, \hat{\mathcal{D}}_T)\}$ はベイズ誤り率の上界である。

同様に $E_{\hat{\mathcal{D}}}\{E_{\hat{\mathcal{D}}_T}\{\epsilon(\hat{\mathcal{D}}, \hat{\mathcal{D}}_T)\}\}$ も上界となる。

この手続きを分割学習法 (Holdout/H 法) と呼ぶ。

分割学習法 (Holdout/H 法)

- ① 与えられたデータセット S を二つの互いに素の集合、学習セット S_D と評価セット S_T に分割する
- ② S_D を用いて識別器を学習する
- ③ S_T を用いて誤り率を測定する

再代入法 (Resubstitution/R 法)

$E\{\varepsilon(\hat{\mathcal{P}}, \hat{\mathcal{P}})\}$ はベイズ誤り率の下界となる。

これは $\hat{\mathcal{P}}$ をベイズ識別器の学習に用い、同じ $\hat{\mathcal{P}}$ を評価に用いば得られる。

この手続きを再代入法 (Resubstitution/R 法) と呼ぶ。

再代入法 (Resubstitution/R 法)

- ① 与えられたデータセット S を用いて識別器を学習する
- ② 同じセットを用いて性能を評価する

交差確認法 (Cross-Validation/CV 法)

H 法はデータを学習用と評価用に分割する必要があるため、誤り推定に不利な点がある。推定の偏りは学習データ量の影響を受け、推定の分散は評価データ量の影響を受ける。この問題点を解決する方法として、交差確認法 (Cross-Validation/CV 法) が利用される。

交差確認法 (Cross-Validation/CV 法)

- ① 与えられたデータセット S を N の互いに素の集合 S_1, S_2, \dots, S_N に分割する
- ② 一つの集合を評価セットとし、残りの $N-1$ をあわせて学習セットとし、学習と誤り推定を行う
- ③ 評価セットを入れ替えながら上記を N 回繰り返す、得られる N 個の誤り率を平均する

交差確認法 (Cross-Validation/CV 法)

CV 法では、データセットの各要素を個別の集合として扱うこともできる。
この極端な事例は一つ抜き法 (leave-one-out/L 法) と呼ばれる。

ブートストラップ法

ベイズ誤り率 ε の推定において、まず R 法にて推定値 $\hat{\varepsilon}$ を得ることができる。
次いでその偏りを考えよう

$$b = \varepsilon - \hat{\varepsilon}.$$

もし b の適切な推定値を得ることができれば、以下によりベイズ誤り率を推定することができる

$$\varepsilon = \hat{\varepsilon} + b.$$

ブートストラップ法

ブートストラップ法 (bootstrap method) は、元の集合 S から、復元抽出により仮想集合 S^* を得る ($|S| = |S^*|$ とする)。

次いで偏り b を以下により推定する

$$b^* = \varepsilon^* - \hat{\varepsilon}^*$$

ただし ε^* は S^* を学習セット S を評価セットとして得た誤りの推定値であり、 $\hat{\varepsilon}^*$ は S^* を学習セット並びに評価セットとして得た誤りの推定値である。

これを B 回、 B 個の仮想集合を生成して繰り返し、その平均値を得ることで b^* を推定する。
ベイズ誤りは以下のように推定する

$$\varepsilon = \hat{\varepsilon} + b^*$$

ブートストラップ法

ブートストラップ法 (The Bootstrap Method)

- ① 与えられたデータセット S を学習用と評価用に用い (R 法) 識別誤り率の推定値 $\hat{\epsilon}$ を得る
- ② S からの $|S|$ 回の復元抽出により仮想セット S^* を得る
- ③ S^* を学習セット、 S を評価セットとして、識別誤り率の推定値 ϵ^* を得る
- ④ S^* を学習セット並びに評価セットとして用いて識別誤り率の推定値 $\hat{\epsilon}^*$ を得る
- ⑤ 偏りの推定値を算出する: $b^* = \epsilon^* - \hat{\epsilon}^*$
- ⑥ 仮想セットを都度生成しながらこの処理を B 回繰り返し、得られた B 個の偏りを平均して最終的な偏り b^* を得る
- ⑦ 誤り率の推定値を $\epsilon = \hat{\epsilon} + b^*$ で算出する

演習

二つのクラスの条件付確率分布として、二つの正規分布を考える

$$\text{Class 1 } \Sigma_1 = I, \mu_1 = [s, 0, \dots, 0]^T$$

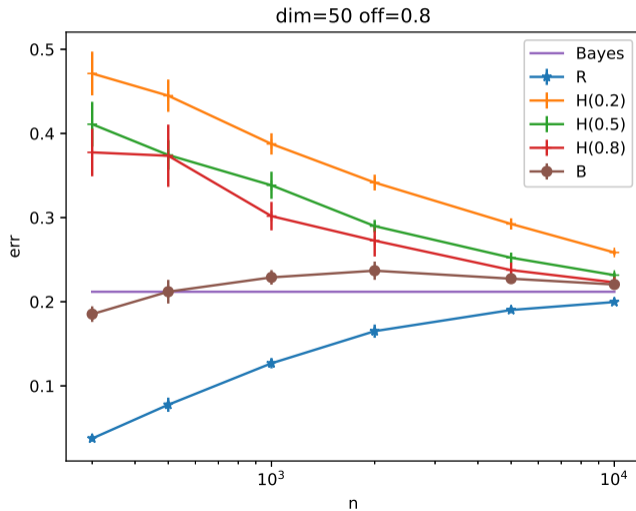
$$\text{Class 2 } \Sigma_2 = I, \mu_2 = [-s, 0, \dots, 0]^T$$

ベイズ誤り率は

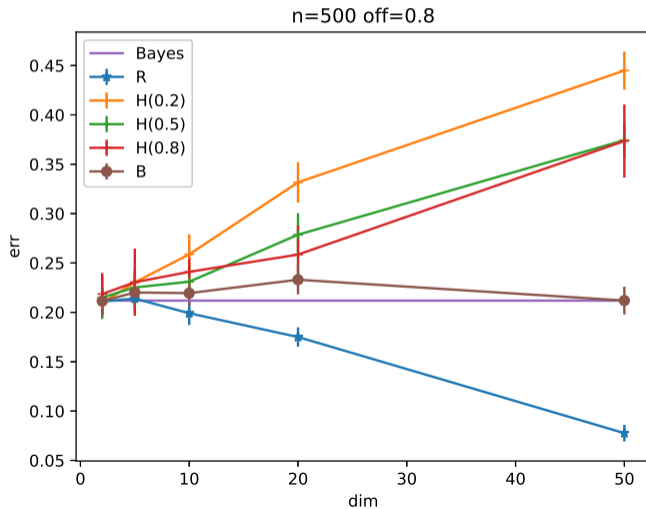
$$\begin{aligned} P^* &= \int_s^\infty \int_{-\infty}^\infty \cdots \int_{-\infty}^\infty N(0, I) dx_1 \cdots dx_n \\ &= \frac{1}{2} \left(1 - \operatorname{erf}\left(\frac{s}{\sqrt{2}}\right) \right) \\ \operatorname{erf}(x) &= \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt \end{aligned}$$

クラス1と2からそれぞれ n サンプル得たとして、二次識別器の誤り率を R 法と H 法で算出し、ベイズ誤り率と比較してみよう

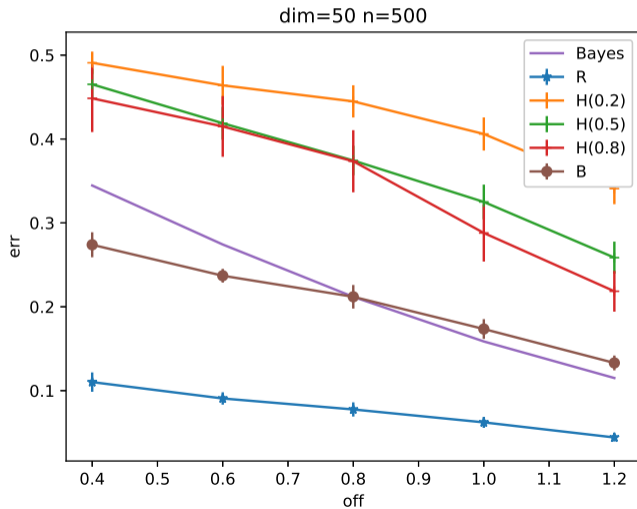
Exercise



Exercise



Exercise



宿題 (プログラミング課題)

- ブートストラップ法を実装せよ
- 演習の方法に従い、誤り率の推定結果を R 法と H 法並びにベイズ誤り率と比較せよ
- データの次元数、ベイズ誤り率、データセットの標本数を変えた時のふるまいを図示せよ
- 締切は 5 月 30 日
- PDF フォーマットにて ITC LMS に提出すること。
- Jupyter Notebook (ipynb) や Python (py) 形式の提出はおやめください。