

Minimalizacja ryzyka strukturalnego, podejście Vapnika

Wykład IV

Wista, grudzień 2009

Minimalizacja ryzyka strukturalnego

Problem klasyfikacji dla dwóch klas, $g = 2$.

$$L(f(\mathbf{x}), y) = I\{f(\mathbf{x}) \neq y\}.$$

Założmy, że dana jest rodzina \mathcal{C} potencjalnych klasyfikatorów

$$\phi \in R^p \quad \phi : R^p \rightarrow \{-1, 1\}.$$

(indeks klasy $Y = \pm 1$). Na przykład $\mathcal{C} = \{\phi_{\mathbf{a}} : \phi_{\mathbf{a}}(\mathbf{x}) = \text{sign}(\mathbf{a}'\mathbf{x})\}$

Problem: Ocena błędu predykcji dla reguły minimalizującej
ryzyko empiryczne (estymator błędu przez powtórne podstawienie).

Przykład

$p(x|-1) \sim U[0, 1]$ i $p(x|1) \sim U[1, 2]$

Minimalne ryzyko reguły (ryzyko bayesowskie równa) się 0.

A_1, \dots, A_n obserwacje z populacji pierwszej

B_1, \dots, B_n obserwacje z populacji drugiej.

LDA Fishera w tym przypadku redukuje się do:

klasyfikuj do drugiej populacji gdy $X > T = (\bar{A}_n + \bar{B}_n)/2$

$$Err_U = \frac{1}{2}(P(X > T|Y = -1) + P(X < T|Y = 1)) = |T - 1|/2$$

$$Err = E|T - 1|/2 \sim E|N(0, \frac{1}{24n})| = \frac{1}{4\sqrt{\pi \times n}}$$

Szacujemy błąd, jaki popełnimy, wybierając klasyfikator ϕ_n^* minimalizujący estymator błędu przez powtórne podstawienie (ryzyko empiryczne)

$$\bar{err}(\phi) = n^{-1} \sum_{i=1}^n I\{\phi(\mathbf{X}_i) \neq Y_i\}.$$

$$\phi_n^* = \operatorname{argmin}_{\phi \in \mathcal{C}} \bar{err}(\phi).$$

Ile różni się

$$Err_{\mathcal{U}}(\phi_n^*) = P(\phi_n^*(\mathbf{X}^0) \neq Y^0 | \mathcal{U})$$

od minimalnego błędu w klasie wszystkich klasyfikatorów $\inf_{\phi} Err(\phi)$.

$$Err_{\mathcal{U}}(\phi_n^*) - \inf_{\phi} Err(\phi) =$$

$$\underbrace{(Err_{\mathcal{U}}(\phi_n^*) - \inf_{\phi \in \mathcal{C}} Err(\phi))}_{\text{błąd estymacji}} + \underbrace{(\inf_{\phi \in \mathcal{C}} Err(\phi) - \inf_{\phi} Err(\phi))}_{\text{błąd aproksymacji}}$$

Zwiększając rodzinę \mathcal{C} zwiększamy (zmniejszamy) błąd estymacji (aproksymacji).

Zmniejszając rodzinę \mathcal{C} zmniejszamy (zwiększamy) błąd estymacji (aproksymacji).

Błąd estymacji $\tilde{\phi} = \operatorname{argmin}_{\phi \in \mathcal{C}} \operatorname{Err}(\phi)$

$$\operatorname{Err}_{\mathcal{U}}(\phi_n^*) - \operatorname{Err}(\tilde{\phi}) \leq \operatorname{Err}_{\mathcal{U}}(\phi_n^*) \underbrace{- \bar{e}rr(\phi_n^*) + \bar{e}rr(\tilde{\phi})}_{\geq 0} - \operatorname{Err}(\tilde{\phi}) \leq$$

$$\underbrace{\operatorname{Err}_{\mathcal{U}}(\phi_n^*) - \bar{e}rr(\phi_n^*)}_{\geq 0} + \underbrace{\bar{e}rr(\tilde{\phi}) - \operatorname{Err}(\tilde{\phi})}_{\leq 0} \leq 2 \sup_{\phi \in \mathcal{C}} |\operatorname{Err}(\phi) - \bar{e}rr(\phi)| \quad (*)$$

Ponadto

$$\operatorname{Err}_{\mathcal{U}}(\phi_n^*) \leq \bar{e}rr(\phi_n^*) + \sup_{\phi \in \mathcal{C}} (\operatorname{Err}(\phi) - \bar{e}rr(\phi)) \quad (**)$$

Problem (**): szacowanie błędu predykcji $\operatorname{Err}_{\mathcal{U}}(\phi_n^*)$ w terminach $\bar{e}rr(\phi_n^*)$.

$$\sup_{\phi \in \mathcal{C}} (\text{Err}(\phi) - \bar{e}r(\phi)) = \sup_{f \in \mathcal{F}} (Pf - P_n f)$$

$$\mathbf{Z} = (\mathbf{X}, Y), \mathbf{Z}_i = (\mathbf{X}_i, Y_i)$$

$$f(\mathbf{z}) = I\{\phi(\mathbf{x}) \neq y\}, Pf = Ef(\mathbf{Z}), P_n f = n^{-1} \sum_{i=1}^n f(\mathbf{Z}_i).$$

\mathcal{F} : rodzina funkcji f utworzona na podstawie \mathcal{C} .

Nierówność koncentracyjna McDiarmida

$g : X^n \rightarrow R$:

$$\sup_{x_1, \dots, x_n, x'_i \in X} |g(x_1, \dots, x_n) - g(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n)| \leq c \quad 1 \leq i \leq n$$

to dla $Z = g(X_1, \dots, X_n)$, gdzie X_1, \dots, X_n - niezależne,

$$P(Z - EZ > t) \leq \exp(-2t^2/nc^2) \quad (!)$$

Dla $Z = \sup_{f \in \mathcal{F}} Pf - P_n f$, $c = 2/n$ i (!) implikuje

Z prawdopodobieństwem $> 1 - \delta$

$$\sup_{f \in \mathcal{F}} [Pf - P_n f] \leq E \left(\sup_{f \in \mathcal{F}} [Pf - P_n f] \right) + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

Średnie Rademachera zbioru $A \subset R^n$

$$R_n(A) = E \left\{ \sup_{\mathbf{a}} \frac{1}{n} \sum_{i=1}^n \sigma_i a_i \right\},$$

$\mathbf{a} = (a_1, \dots, a_n)$ i σ_i nzl Rademachera.

Własności średnich Rademachera

- (i)

$$R_n(A) = R_n(\text{absconv}(A)),$$

gdzie $\text{absconv}(A) = \{\sum_{j=1}^N c_j \mathbf{a}^{(j)}, \sum_{j=1}^N |c_j| \leq 1, \mathbf{a}^{(j)} \in A\}$

- (ii)

$$R_n(\phi \circ A) \leq L_\phi R_n(A),$$

gdzie $\phi : R \rightarrow R$ Lipschitzowska ze stałą L_ϕ , a

$\phi \circ A = \{(\phi(\mathbf{a}_1), \dots, \phi(\mathbf{a}_n)), \mathbf{a} = (\mathbf{a}_1, \dots, \mathbf{a}_n) \in A\}$

- (iii) $A = \{\mathbf{a}_1, \dots, \mathbf{a}_m\}$ to $R_n(A) \leq \max_{i=1, \dots, m} \|\mathbf{a}_i\| \frac{\sqrt{2 \log m}}{n}$

$(\mathbf{Z}'_1, \dots, \mathbf{Z}'_n)$ niezależna kopia $(\mathbf{Z}_1, \dots, \mathbf{Z}_n)$ i $P'_n f = n^{-1} \sum_{i=1}^n f(\mathbf{Z}_i)$.

Technika sprzęgania daje

$$\begin{aligned} E\left(\sup_{f \in \mathcal{F}} [Pf - P_n f]\right) &= E\left(\sup_{f \in \mathcal{F}} E[P'_n f - P_n f | \mathbf{Z}_1, \dots, \mathbf{Z}_n]\right) \leq E \sup_{f \in \mathcal{F}} E[P'_n f - P_n f] \\ &= E \sup_{f \in \mathcal{F}} n^{-1} \sum_{i=1}^n \sigma_i(f(\mathbf{Z}'_i) - f(\mathbf{Z}_i)) \leq 2ER_n(\mathbf{Z}^n) \end{aligned}$$

Nasz zbiór A: $\mathcal{F}(\mathbf{z}^n) = \{(f(\mathbf{z}_1), \dots, f(\mathbf{z}_n))\}_{f \in \mathcal{F}}$, (iii) \Rightarrow

$$R_n(\mathbf{Z}^n) \leq \sqrt{\frac{2 \log S_{\mathcal{F}}(\mathbf{Z}^n)}{n}},$$

gdzie $S_{\mathcal{F}}(\mathbf{z}^n) = |\mathcal{F}(\mathbf{z}^n)|$, $\mathbf{z}^n = (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n)$.

Jeśli $V(\mathbf{z}^n)$ wymiar Vapnika-Czervonenkisa zbioru $S_{\mathcal{F}}(\mathbf{z}^n)$ to

$$\log S_{\mathcal{F}}(\mathbf{z}^n) \leq V(\mathbf{z}^n) \log(n+1)$$

i

$$R_n(\mathbf{Z}^n) \leq \sqrt{\frac{2V(\mathbf{Z}^n) \log(n+1)}{n}},$$

Jeśli $\sup_{\mathbf{z}^n} V(\mathbf{z}^n) \leq V$ to z prawdopodobieństwem $> 1 - \delta$

$$\sup_{f \in \mathcal{F}} [Pf - P_n f] \leq C \sqrt{\frac{V \log n}{n}} + \sqrt{\frac{2 \log(1/\delta)}{n}}$$

Czynnik $\log n$ można pominąć.

Minimalizacja ryzyka strukturalnego

Rodziny klasyfikatorów $\mathcal{C}_1 \subset \mathcal{C}_2 \dots \subset \mathcal{C}_N$

$$(1) \quad \text{Err}_{\mathcal{U}}(\phi_{n,i}^*) \leq \bar{e}rr(\phi_{n,i}^*) + B(n, \mathcal{F}_i) \quad i = 1, \dots, n$$

Minimalizacja prawej strony (1)

Problemy: (i) Klasy z małym wymiarem mogą mieć duży błąd aproksymacji (ii) minimalizacja $\bar{e}rr(\phi)$ często problem NP-trudny

(ii) nierówność

$$\text{Err}_{\mathcal{U}}(\phi_n^*) \leq \bar{e}rr(\phi_n^*) + \sup_{\phi \in \mathcal{C}} (\text{Err}(\phi) - \bar{e}rr(\phi)) \quad (**)$$

za gruba !

Oszacowania oparte na marginesach

Niech ϕ będzie regułą klasyfikacyjną opartą na funkcji klasyfikacyjnej d :
 $\phi(\mathbf{x}) = 1$ jeśli $d(\mathbf{x}) \geq 0$ i $\phi(\mathbf{x}) = -1$ w przeciwnym przypadku.

Wtedy

$$P(\phi(\mathbf{X}) \neq Y) = P(\text{sign}(d(\mathbf{X})) \neq Y) \leq E I\{d(\mathbf{X})Y \leq 0\}$$

Niech $\gamma : R \rightarrow R_+$ będzie funkcją kosztu taką, że $\gamma(x) \geq I\{x \geq 0\}$.

Typowe przykłady

$$\gamma(x) = e^x, \quad \gamma(x) = (1+x)_+, \quad \gamma(x) = \log_2(1+x)$$

Funkcjonał kosztu i jego wersja empiryczna

$$A(d) = E(\gamma(-d(\mathbf{X})Y)|\mathcal{U}), \quad A_n(d) = n^{-1} \sum_{i=1}^n \gamma(-d(\mathbf{X}_i)Y_i)$$

Mamy

$$Err(d) \leq A(d) \quad e\bar{r}_n(d) \leq A_n(d).$$

Niech $B = \|\gamma(-d(\mathbf{x})y)\|_\infty$. Wtedy

$$Err_{\mathcal{U}}(d_n) \leq A(d_n) = E(\gamma(-d_n(\mathbf{X})Y)|\mathcal{U}) \leq$$

$$A_n(d_n) + \sup_{d \in \mathcal{D}} (A(d) - A_n(d)) \leq$$

$$A_n(d_n) + 2L_\gamma ER_n(\mathcal{F}(\mathbf{Z}^n)) + B\sqrt{\frac{2 \log \frac{1}{\delta}}{n}},$$

L_γ - stała Lipschitza funkcji γ .

Systemy głosowania opierają się na podejmowaniu decyzji na podstawie kombinacji liniowych klasyfikatorów. \mathcal{C} - rodzina bazowych klasyfikatorów.

$$\mathcal{D}_\lambda = \left\{ d(\mathbf{x}) = \sum_{i=1}^m c_j g_j(\mathbf{x}) : m \in \mathbb{N}, \sum_{j=1}^n |c_j| \leq \lambda, g_1, \dots, g_n \in \mathcal{C} \right\}$$

Na podstawie własności średnich Rademachera zbioru

$$R_n(\mathcal{D}_\lambda(\mathbf{Z}^n)) = \lambda R_n(\mathcal{C}(\mathbf{Z}_n)) \leq \lambda \sqrt{\frac{2V_{\mathcal{C}} \log(n+1)}{n}}$$

Stąd

$$Err_{\mathcal{U}}(\text{sign}(d_n^*)) \leq A_n(d_n^*) + 2L_{\phi}\lambda\sqrt{\frac{2V_C \log(n+1)}{n}} + B\sqrt{\frac{2 \log \frac{1}{\delta}}{n}},$$

gdzie $\text{sign}(d_n^*)$ minimalizuje ryzyko empiryczne w klasie $\{\text{sign}(d)\}_{d \in \mathcal{D}}$
Oszacowanie błędu predykcji $\text{sign}(d_n^*)$ w terminach wymiaru
Vapnika-Chervonenkisa klasy bazowej.

Przykład

$$\phi_\eta(x) = 1 \quad x \leq -\eta$$

$$\phi_\eta(x) = 0 \quad x \geq 0$$

$$\phi_\eta(x) = 1 + x/\eta \quad \text{w p.p.}$$

W tym przypadku $B = 1$ i $L_\phi = 1/\eta$. Ponadto $A_n(d) \leq \bar{e}r_n^\gamma(d)$, gdzie

$$\bar{e}r_n^\gamma(d) = \frac{1}{n} \sum_{i=1}^n I\{d(\mathbf{X}_i)Y_i \leq \eta\}.$$

Liczymy nie tylko źle sklasyfikowane elementy, ale również te sklasyfikowane dobrze, ale z małym 'marginesem'

$$Err_{\mathcal{U}}(d_n) \leq \bar{err}_n(d_n) + \frac{2\lambda}{\eta} \sqrt{\frac{2V_C \log(n+1)}{n}} + \sqrt{\frac{2 \log \frac{1}{\delta}}{n}}.$$

- Własności skończeniopróbkowe estymatorów PMS i ich ryzyka (Barron, Birgé, Massart, Barraud ..)
- Metody rzadkiej regresji (Lasso, Least Angle Regression LAR,...). potencjalna liczba predyktorów w modelu liniowym $p \approx n$, $p > n$, jednakże $\#\{i : \beta_i \neq 0\} \ll p$
- Związki metod selekcji z wielokrotnym testowaniem, w szczególności z procedurami typu Benjaminiego-Hochberga ($M_1 \subset M_2 \Rightarrow \hat{M}_{AIC} = M_2 \iff LRT_{M_1, M_2} > 2(p_2 - p_1)$)
- metody selekcji w modelach regresyjnych oparte na wstępnym uporządkowaniu predyktorów.
- metody predykcji oparte na uśrednianiu modeli

Własności skończeni próbkowe

$$Y_i = f(x_i) + \varepsilon_i, \quad i = 1, \dots, n$$

$\varepsilon_i \sim (0, \sigma^2)$, σ^2 -znane, $f; [0, 1] \rightarrow R$, $x_i \in [0, 1]$.

M_m : funkcje stałe na odcinkach $[(j-1)/m, j/m)$.

$f \in \text{Holder}(L, \alpha)$, to dla $\hat{f} = \hat{f}_{MNK}(M_m)$

$$R_n = E\|f - \hat{f}\|_n^2 \leq L^2 M^{-2\alpha} + \frac{m\sigma^2}{n}$$

optymalny wybór m :

$$m(\alpha, L) = \left(\frac{L}{\sigma}\right)^{2/(1+\alpha)} m^{1/(2\alpha+1)}$$

Ryzyko $R_n = O(n^{-2\alpha/(2\alpha+1)})$. Konstrukcja estymatorów PMS o takim rządzie ryzyka, bez wykorzystywania wiedzy, że $f \in \text{Holder}(L, \alpha)$.
Kalibracja stałych w funkcji kryterialnej !!!

Referencje

- Bousquet, O., Boucheron, S., Lugosi, G. (2005) Theory of classification: a survey of recent results, ESAIM: Probability and Statistics, 323-375
- Claeskens, G., Hjort, N. (2008) Model selection and model averaging, Cambridge University Press
- Friedman, J. (1997) On bias, variance, 0/1-loss and the curse of dimensionality, Data mining and Knowledge Discovery, str. 55-77
- Hastie, T., Tibshirani, R., Friedman, J. (2009) The Elements of Statistical Learning, (wydanie drugie)
<http://www-stat.stanford.edu/~tibs/ElemStatLearn/>
- Konishi, S., Kitagawa, G. (2008) Information Criteria and Statistical Modelling, Springer
- Leeb, H., Pötscher, B. (2008) Sparse estimators and the oracle property, or the return of Hodges estimator, J. Econometrics, str. 201-211
- Leeb, H., Pötscher, B. (2008) Model selection, w Handbook of Financial Time Series (wyd Andersen, Davis, Kreiss, Mikosch), Springer, 2008

- Nishii, R. (1984) Asymptotic properties of criteria for selection of variables in multiple regression, An. Statist., str. 758-765.
- Schiavo, R., Hand, D.(2000) Ten more years of error rate research, International Statistical Review, str. 295-310.
- Shao, J. (1993) Linear model selection by cross-validation, JASA, str. 486-494.
- Shao, J. (1997) An asymptotic theory for linear model selection, Statistica Sinica, str. 221-264
- Yang, Y. (2005) Can the strengths of AIC and BIC be shared? A conflict between model identification and regression estimation, Biometrika, str. 937-950.