

# ESTYMACJA BŁĘDU PREDYKCJI I JEJ ZASTOSOWANIA

Jan Mielniczuk

Wiśła, grudzień 2009

## PLAN

- Błędy predykcji i ich podstawowe estymatory
- Estymacja błędu predykcji w modelu liniowym. Funkcje kryterialne
- Własności predykcyjne estymatorów postselekcyjnych
- Metoda minimalizacji ryzyka strukturalnego

# BŁĘDY PREDYKCJI I ICH PODSTAWOWE ESTYMATORY

## (i) Regresja wielokrotna z losowymi zmiennymi objaśniającymi

Wektor losowy  $(\mathbf{X}, Y) \in R^p \times R$  o rozkładzie  $P_{\mathbf{X}, Y}$  takim, że

$$Y = f(\mathbf{X}) + \varepsilon, \quad (*)$$

$f : R^p \rightarrow R$ ,  $\varepsilon$  zmienna niezależna od  $\mathbf{X}$ ,  $E(\varepsilon) = 0$ .

Uwaga. Założenia implikują

$$f(\mathbf{X}) = E(Y|\mathbf{X}) \quad i \quad \varepsilon = Y - E(Y|\mathbf{X}).$$

Dla takich  $f$  i  $\varepsilon$ ,  $(*)$  spełniona oczywiście zawsze, tu zakładamy dodatkowo niezależność  $\mathbf{X}$  i  $\varepsilon$ .

## (ii) Klasyfikacja pod nadzorem

Tak jak w (i), zakładamy dodatkowo, że

$$Y \in \{1, 2, \dots, K\}$$

( $Y$  - indeks przynależności do populacji).

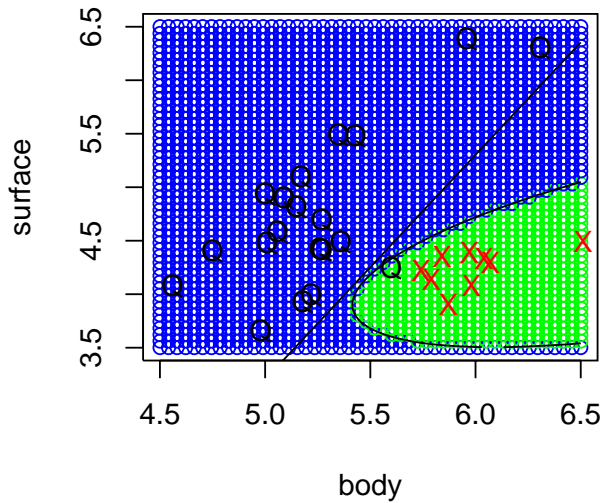
Nie zakładamy niezależności  $\mathbf{X}$  i  $\varepsilon$ .

Obserwowalna p.p.l.  $\mathcal{U} = \{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)\}$  z rozkładu  $P_{\mathbf{X}, Y}$ .

Na jej podstawie estymator  $\hat{f}(\mathbf{x})$  używany do prognozy zmiennej  $Y^0$ :

$$\hat{Y}^0 = \hat{f}(\mathbf{X}^0)$$

.



### (iii) Estymacja gęstości zmiennej losowej $Y$

$Y_1, \dots, Y_n$  p.p.l. z rozkładu  $P_Y$  o gęstości  $f$ . Estymacja  $f$  w oparciu o model parametryczny  $M = \{f_\theta\}_{\theta \in \Theta}$ .

### (iv) Estymacja gęstości warunkowej

$(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$  p.p.l. z rozkładu  $P_{\mathbf{X}, Y}$  na  $R^p \times R$ . Estymujemy gęstość warunkową  $Y$  pod warunkiem  $\mathbf{X} = \mathbf{x}$  w oparciu o model parametryczny

$$M = \{f_{\theta(\mathbf{x})}\}_{\theta \in \Theta, \mathbf{x} \in R^p}$$

$(\mathbf{X}^0, Y^0) \sim P_{\mathbf{X}, Y}$  niezależna od  $\mathcal{U}$ .

- Jak dalece  $\hat{f}(\mathbf{X}^0)$  różni się od  $Y^0$  ? (Problem (i) i (ii))
- Jak dalece gęstość  $f_{\hat{\theta}}$  różni się od gęstości  $f$ ? (problem (iii))
- Jak dalece gęstość  $f_{\hat{\theta}(\mathbf{x}^0)}$  różni się od gęstości  $f_{Y^0|\mathbf{x}^0}$ ? (problem (iv))

Błędy rozpatrywane w (iii) i (iv) będą związane z problemem prognozy.



Problemy (i) i (ii):  $L(y, \hat{f}(\mathbf{x})) \geq 0$  :

$(y - \hat{f}(\mathbf{x}))^2$ ,  $|y - \hat{f}(\mathbf{x})|$ ,  $I\{y \neq \hat{f}(\mathbf{x})\}$  (funkcja straty 0-1 dla problemu (ii)).

Problemy (iii) i (iv):  $L(y, \hat{\theta}(\mathbf{x}))$  :

$$L(y, \hat{\theta}) = -2 \log f_{\hat{\theta}}(y)$$

$$L(y, \hat{\theta}(\mathbf{x})) = -2 \log f_{\hat{\theta}(\mathbf{x})}(y)$$

# Błąd predykcji (uogólnienia)

$(\mathbf{X}^0, Y^0) \sim P_{\mathbf{X}, Y}$  niezależna od  $\mathcal{U}$ .

$$Err_{\mathcal{U}} = E(L(Y^0, \hat{f}(\mathbf{X}^0)) | \mathcal{U}) \quad (i), (ii)$$

$$Err_{\mathcal{U}} = E(L(Y^0, \hat{\theta}(\mathbf{X}^0)) | \mathcal{U}) \quad (iii), (iv)$$

Warunkowy błąd predykcji ( w (ii): prawdopodobieństwo błędnej klasyfikacji),  $Err_{\mathcal{U}} = G(\mathcal{U})$  jest funkcją próby  $\mathcal{U}$ .

Dla ustalonej próby uczącej  $\mathcal{U}$  konstruujemy predyktor  $\hat{f}$  i oceniamy jego działanie dla nowej obserwacji  $(\mathbf{X}^0, Y^0)$ .

Miara jakości odpowiadająca faktycznej funkcji predyktora.

## Bezwarunkowy błąd predykcji

$$Err = E(Err_{\mathcal{U}})$$

'Uśredniamy losowość próby uczącej, na podstawie której skonstruowano predyktor' (Hastie et al., ESL).

### Wewnątrzpróbkowy błąd predykcji (in-sample error)

$Y_i^0 \sim P_{Y|\mathbf{X}=\mathbf{x}_i}$  niezależne,  $i = 1, 2, \dots, n$ .

$$\mathbf{Y}^0 = (Y_1^0, \dots, Y_n^0)'$$

$$Err_{in} = \frac{1}{n} \sum_{i=1}^n E_{\mathbf{Y}^0}(L(Y_i^0, \hat{f}(\mathbf{X}_i)) | \mathcal{U})$$

Losujemy indeks  $l \in \{1, \dots, n\}$ , nowa obserwacja  $Y_i^0$  w punkcie  $\mathbf{X}_i$  dla  $i = l$ , dla niej oczekiwany błąd predykcji.

# Estymatory błędu predykcji

- estymator metodą próby testowej
- estymator metodą powtórnego podstawienia  $\bar{err}$
- estymator krosvalidacyjny  $\hat{Err}_{CV}$
- estymator metodą bootstrap  $\hat{Err}_{boot}$

# Estymator metodą próby testowej

Oprócz próby uczącej  $\mathcal{U}$  dysponujemy również w pełni obserwowalną próbą testową  $\mathcal{T}$  ( $\mathcal{D} = \mathcal{U} \cup \mathcal{T}$ ,  $\mathcal{U} \cap \mathcal{T} = \emptyset$ ),  $|\mathcal{T}| = m$

$$Err_{\mathcal{U}} = \frac{1}{m} \sum_{i=1}^m I\{(\mathbf{X}_i, Y_i) \in \mathcal{T} : \hat{f}(\mathbf{X}_i) \neq Y_i\}.$$

Mamy

$$E_{\mathcal{T}}(Err_{\mathcal{U}} | \mathcal{U}) = Err_{\mathcal{U}}$$

Nieobciążone oszacowanie prawdopodobieństwa błędnej klasyfikacji.  
Ale często w powyższym schemacie będziemy używać estymatora  $\hat{f}_{\mathcal{D}}$  a nie  $\hat{f}_{\mathcal{U}}$  - wtedy zbyt konserwatywnie z reguły oceniamy prawdopodobieństwo błędnej klasyfikacji tego klasyfikatora.

Trudny problem:  $Err_U = 0.05$ ,

$$2\sqrt{\frac{Err_U(1 - Err_U)}{m}} = 2\sqrt{\frac{0.05 \times 0.95}{m}} \leq 0.01$$

$$m \geq 1900 !!$$

Dwie możliwości zastąpienia próby testowej:

- metoda analityczna - uwzględnienie obciążenia estymatora metodą powtórnego podstawienia
- metoda repróbkiwania - efektywne ponowne użycie próby uczącej

Estymator metodą powtórnego podstawienia

$$\bar{err} = \frac{1}{n} \sum_{i=1}^n L(\hat{f}(\mathbf{X}_i), Y_i).$$

$\bar{err}$  jest naturalnym estymatorem  $Err_{\mathcal{U}}$ . **Ale:** elementy próby  $\mathcal{U}$  pełnią funkcję niezależnych od  $\mathcal{U}$  obserwacji  $(X, Y)$ . Jest estymatorem 'optymistycznym' w tym sensie, że z reguły możemy się spodziewać, że

$$\bar{err} < Err_{\mathcal{U}}.$$

Możemy porównać

$$E_{\mathbf{Y}}(\bar{err}|\mathbf{X}) \quad \text{z} \quad E_{\mathbf{Y}}(Err_{in}|\mathbf{X})$$

$$\mathbf{Y} = (Y_1, \dots, Y_n)', \quad \mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)'.$$

**Twierdzenie 1.** Dla kwadratowej funkcji straty w (i) i 0-1 w (ii)

$$E_{\mathbf{Y}}(e\bar{r}r|\mathbf{X}) = E_{\mathbf{Y}}(Err_{in}|\mathbf{X}) - \frac{2}{n} \sum_{i=1}^n \text{Cov}(\hat{Y}_i, Y_i|\mathbf{X}).$$

Im bardziej elastyczna metoda predykcji (silniej skorelowane  $\hat{Y}_i$  i  $Y_i$ ) tym większy optymizm  $e\bar{r}r$ . Człon

$$\frac{2}{n} \sum_{i=1}^n \text{Cov}(\hat{Y}_i, Y_i|\mathbf{X})$$

kara kowariancyjna.



Dowód Porównujemy przy ustalonym  $\mathbf{X}$  kwadraty wyrażeń

$$(Y_i - \hat{f}(\mathbf{X}_i)) \quad (1) \quad \text{z} \quad (Y_i^0 - \hat{f}(\mathbf{X}_i)) \quad (2)$$

$$(1) = (Y_i - f(\mathbf{X}_i)) - (\hat{f}(\mathbf{X}_i) - E(\hat{f}(\mathbf{X}_i)|\mathbf{X})) + (f(\mathbf{X}_i) - E(\hat{f}(\mathbf{X}_i)|\mathbf{X})) \\ =: A + B + C$$

$$(2) = A' + B + C \quad A' := Y_i^0 - f(\mathbf{X}_i)$$

Liczmy  $E((1)^2|\mathbf{X})$  i  $E((2)^2|\mathbf{X})$

$$2E(A \times B|\mathbf{X}) \neq 2E(A' \times B|\mathbf{X})$$

$$E(A'B|\mathbf{X}) = E_{\mathbf{Y}}(E(A'B|\mathbf{X}, \mathbf{Y})) = E_{\mathbf{Y}}(BE(A'|\mathbf{X}, \mathbf{Y})) = 0$$

$$E(AB|\mathbf{X}) = \text{Cov}(Y_i, \hat{Y}_i|\mathbf{X})$$

Przypadek szczególny (i):

$$Y_i = \mathbf{X}'_i \boldsymbol{\beta} + \varepsilon_i$$

$$\mathbf{X}_i = (X_{i1}, \dots, X_{ip})', \boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$$

$\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n)'$  macierz eksperymentu  $n \times p$  z wierszami =  $\mathbf{X}'_1, \mathbf{X}'_2, \dots, \mathbf{X}'_n$ .

Z definicji  $Err_{in} = n^{-1} \sum_{i=1}^n E_{\mathbf{Y}^0} (Y_i^0 - \hat{f}(\mathbf{X}_i))^2 | \mathcal{U}$

$$E_{\mathbf{Y}}(Err_{in} | \mathbf{X}) = \sigma^2 + \frac{1}{n} E(\|\mathbf{X} \hat{\boldsymbol{\beta}} - \mathbf{X} \boldsymbol{\beta}\|^2 | \mathbf{X})$$

Fakt Dla modelu liniowego

$$\sum_{i=1}^n \text{Cov}(\hat{Y}_i, Y_i | \mathbf{X}) = p\sigma^2,$$

jeśli rząd  $\mathbf{X} = p$  i  $\sigma^2 = \text{Var}(\varepsilon)$ .

Dowód  $\hat{\mathbf{Y}} = \mathbf{H}\mathbf{X}$ , gdzie  $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$

$$\begin{aligned} \sum_{i=1}^n \text{Cov}(\hat{Y}_i, Y_i | \mathbf{X}) &= E((\mathbf{H}\mathbf{Y} - \mathbf{X}\beta)' \varepsilon | \mathbf{X}) = \\ &= E(\varepsilon' \mathbf{H} \varepsilon) = p\sigma^2. \end{aligned}$$

Tylko w przybliżeniu prawdziwe dla innych modeli i metod !

Dla modelu liniowego i kwadratowej funkcji straty estymator  $E_Y(Err_{in}|\mathbf{X})$  ma postać

$$\bar{err} + 2\frac{p}{n}\hat{\sigma}^2 = \frac{1}{n}RSS + 2\frac{p}{n}\hat{\sigma}^2 \quad (*)$$

Jeśli  $E\hat{\sigma}^2 = \sigma^2$  to (\*) nieobciążony estymator  $E(Err_{in}|\mathbf{X})$ .

W takiej sytuacji nieobciążony estymator  $\frac{1}{n}E(\|\mathbf{X}\hat{\beta} - \mathbf{X}\beta\|^2|\mathbf{X})$  ma postać

$$C_p = \frac{1}{n}RSS + 2\frac{p}{n}\hat{\sigma}^2 - \hat{\sigma}^2.$$

Funkcja kryterialna Mallowsa.

# Estymator krosvalidacyjny (rotacyjny) oczekiwanego błędu prognozy $Err$

Modele (i) i (ii). Próbę  $\mathcal{U}$  dzielimy na  $K > 1$  w przybliżeniu równych części.  $\kappa : \{1, \dots, n\} \rightarrow \{1, \dots, K\}$  (funkcja przynależności do elementu podziału).  $\hat{Y}_j^{-i} = \hat{f}^{-i}(\mathbf{X}_j)$  prognoza  $Y_j$  na podstawie próby z usuniętą  $i$ -tą częścią.

$$CV = \frac{1}{n} \sum_{j=1}^n L(Y_j, \hat{Y}_j^{-\kappa(j)}).$$

$Y_j$  prognozowana na podstawie próby niezależnej od  $Y_j$  zgodnie z definicją  $Err$ .

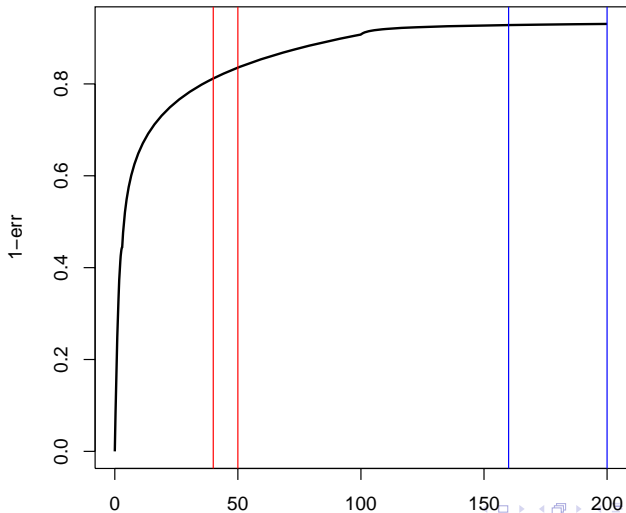
Co estymujemy?  $\mathcal{U}^{-i}$ : próba z usuniętą  $i$ -tą częścią.

Oszacowanie błędu

$$CV = \frac{1}{K} \sum_{i=1}^K \hat{E}(L(\mathbf{X}, Y) | \mathcal{U}^{-i}) = \hat{E}(Err_{\mathcal{U}^{-1}})$$

Nie estymujemy **warunkowego** błędu klasyfikatora  $\hat{d}_{\mathcal{U}}$  dla danej próby  $\mathcal{U}$  tylko **oczekiwany** błąd klasyfikatora  $\hat{d}_{\mathcal{U}-1} = Err_{(n(K-1)/K)}$

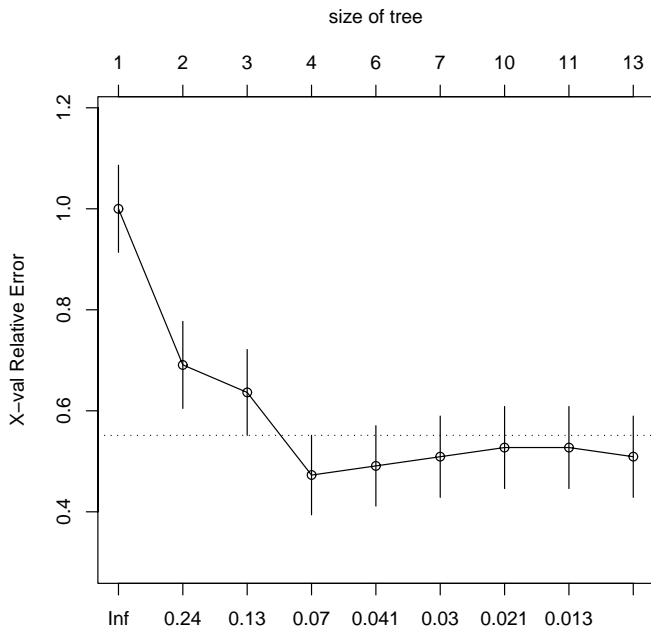
Hypothetical learning curve



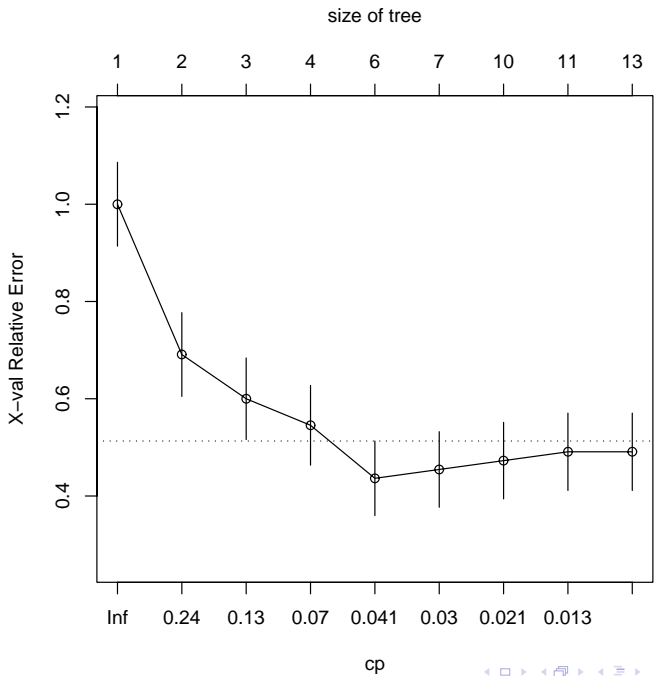
# Uwagi o estymatorze kroswalidacyjnym

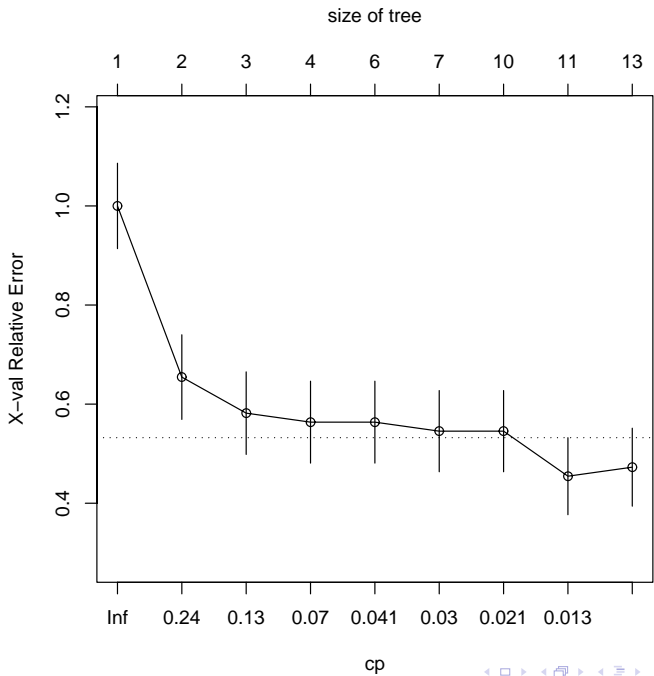
1. Liczba podziałów  $K = 5, 10$  i  $K = n$ . (leave-one out). W ostatnim przypadku najmniejsze obciążenie estymatora (jako oszacowania  $Err_n$  dla licznosci próby  $n$ ), największa wariancja (skorelowanie czynników w definicji CV).
2. Obciążenie estymatora kroswalidacyjnego jako estymatora  $Err_n$  dla  $K = 5, 10$  może być znaczne.

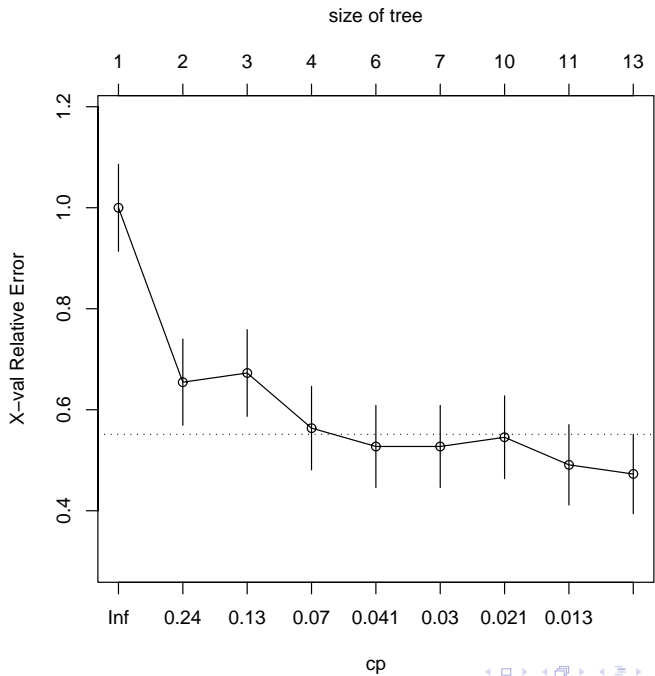
### 3. Niestabilność CV związana z losowością podziału na $K$ części.











## Przykład (ESL, rozdział VII)

$(Y, \mathbf{X}) \in \{0, 1\} \times R^{20}$ :  $Y = 1$  jeśli  $\sum_{i=1}^{10} X_i > 5$ , 0 w przeciwnym przypadku.

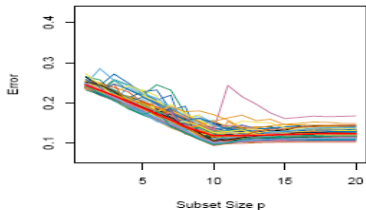
CV-10 i CV leave-one-out jak estymatory  $Err_{\mathcal{U}}$  i  $Err$  liczone dla 100 prób uczących ( $n = 80$ ).

Metoda estymacji: regresji liniowa dla najlepszego podzbioru predyktorów danej liczności ( $\operatorname{argmin}_{M:|M|=p} RSS_M$ )

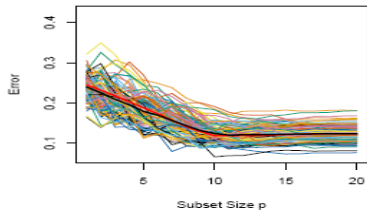
Czerwona linia -  $Err$

Czarna linia -  $ECV_{10}$  i  $ECV_n$

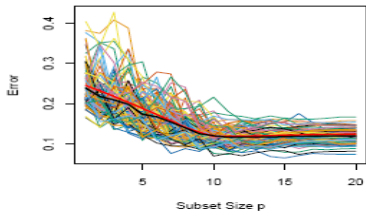
Prediction Error



10-Fold CV Error



Leave-One-Out CV Error



Approximation Error

