

# Spis treści

- 1 Obciążenie, wariancja i złożoność modelu
- 2 Bias-variance tradeoff
- 3 Optymizm błędu treningowego
- 4 Metody estymacji błędu in-sample
- 5 Krosvalidacja
- 6 Bootstrap

# Obciążenie, wariancja i złożoność modelu

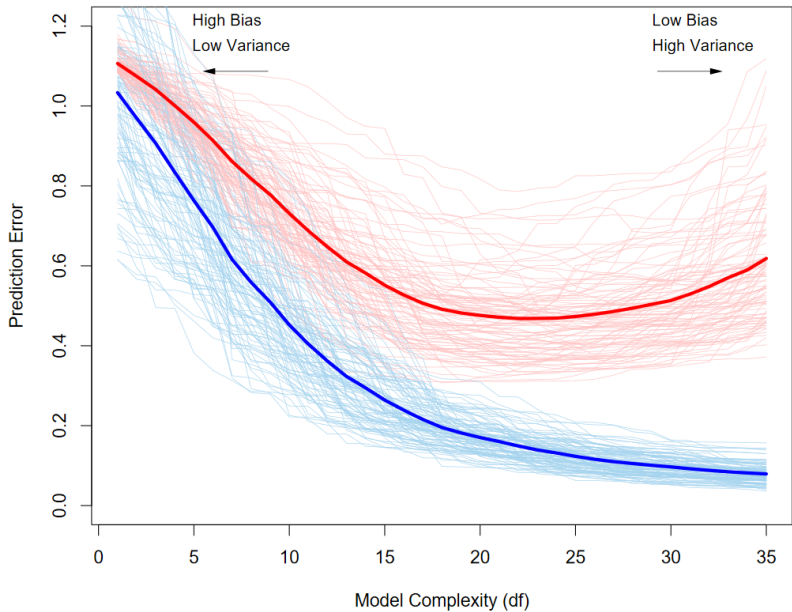
Mamy zmienną objaśnianą  $Y$ , macierz zmiennych objaśniających  $X$ , model predykcyjny  $\hat{f}(X)$ , zbiór treningowy  $\mathcal{T}_1$  oraz zbiór testowy  $\mathcal{T}_2$ .

Najpopularniejszymi funkcjami straty ( $L(Y, \hat{f}(X))$ ) są:

- 1  $(Y - \hat{f}(X))^2$  - błąd kwadratowy
- 2  $|Y - \hat{f}(X)|$  - błąd bezwzględny

Zdefiniujmy błąd treningowy:  $e\hat{r}r = \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{f}(x_i))$ , gdzie  $x_i$  pochodzą ze

zbioru treningowego. Błąd testowy:  $Err_{\mathcal{T}_2} = E[L(Y, \hat{f}(X)) | X \subset \mathcal{T}_2]$ . Oraz oczekiwany błąd predykcji  $Err = E[Err_{\mathcal{T}_2}]$



## Ta sama idea w modelach klasyfikacyjnych

Mamy kategorię zmienną objaśnianą  $G$ . Tutaj nasz model zwraca prawdopodobieństwa, że  $p_k(X) = P(G = k|X)$  dla  $k = 1, \dots, K$ .

Wynikiem modelu jest  $k$  o największej wartości z tych prawdopodobieństw:  $\operatorname{argmax}_k \hat{p}_k(X)$ . Najpopularniejszymi funkcjami straty ( $L(G, \hat{G}(X))$ ) są:

①  $I(G, \hat{G}(X))$  - strata 0-1

②  $-2 \sum_{k=1}^K I(G = k) \log \hat{p}_k(X) = -2 \log \hat{p}_G(X)$  -  $-2 \times$  logwarogodność

Błędy treningowe i testowe tak samo, tylko na innych funkcjach straty.

# Podział zbioru zmiennych objaśniających

Sytuacja idealna (duży zbiór danych):



## Bias-variance tradeoff

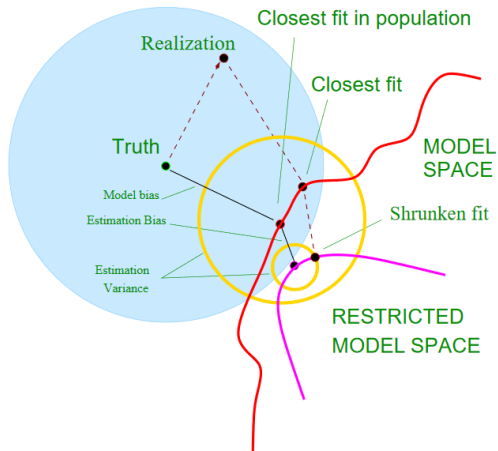
Zakładamy, że  $Y = f(X) + \epsilon$ , gdzie  $E[\epsilon] = 0$  i  $Var(\epsilon) = \sigma_\epsilon^2$ . Możemy teraz rozpisać oczekiwany błąd predykcji (na kwadratowej funkcji straty) jako:

$$Err(Y, \hat{f}(X)) = E[(Y - \hat{f}(X))^2] = \sigma_\epsilon^2 + (E\hat{f}(X) - f(X))^2 + E[(\hat{f}(X) - E\hat{f}(X))^2] = \sigma_\epsilon^2 + Bias^2 + Var$$

W szczególności dla modelu liniowego mamy:

$$Bias^2 = (f(X) - X^T(X^T X)^{-1}X^T Y)^2 \text{ i } Var = \|X(X^T X)^{-1}X\|^2 \sigma_\epsilon^2$$

# Bias-variance tradeoff na wykresie



## Optymizm błędu treningowego

Niech  $ERR_{\mathcal{T}} = E_{X^0, Y^0}[L(Y^0, \hat{f})(X^0)|\mathcal{T}]$ , punkt  $(X^0, Y^0)$  pochodzi z rozkładu  $(X, Y)$ . Niech  $ERR = E_{\mathcal{T}}[ERR_{\mathcal{T}}]$ . Zazwyczaj empiryczny błąd treningowy to  $\hat{err} = \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{f}(x_i)) < ERR_{\mathcal{T}}$ .

Zdefiniujmy błąd in-sample  $Err_{in} = \frac{1}{n} \sum_{i=1}^n E_{Y^0}[L(Y_i^0, \hat{f}(x_i))|\mathcal{T}]$

Definiujemy optymizm jako  $opt = Err_{in} - \hat{err}$

Oraz średni optymizm:  $\omega = E_{\mathcal{T}}[opt]$ .

W ogólności dla większości funkcji straty (między innymi dla straty

kwadratowej, 0-1 itd.)  $\omega = \frac{2}{n} \sum_{i=1}^n Cov(\hat{y}_i, y_i)$ .



I na koniec zyskujemy ważną równość:

$$E_Y[Err_{in}] = E_Y[e\hat{r}r] + \frac{2}{n} \sum_{i=1}^n Cov(\hat{y}_i, y_i)$$

Co w przypadku, gdy mamy zależność liniową pomiędzy objaśnianą i objaśniającymi upraszcza się do:

$$E_Y[Err_{in}] = E_Y[e\hat{r}r] + \frac{2d}{n} \sigma_\epsilon^2 \text{ (gdzie } d \text{ to liczba zmiennych użytych w modelu)}$$

# Metody estymacji błędu in-sample

Uogólniona formuła na estymator błędu in-sample:  $\hat{Err}_{in} = \hat{o}pt + \hat{e}rr$

Statystyka  $C_p$  Mallow'a:  $C_p = \hat{e}rr + \frac{2d}{n} \cdot \hat{\sigma}_\epsilon^2$ , np.  $\hat{\sigma}_\epsilon^2 = \frac{\|Y - X\hat{\beta}^{LS}\|^2}{n-d}$

AIC:  $-2E[\log P_{\hat{\theta}}(Y)] \approx \frac{-2}{n} E[\loglik] + \frac{2d}{n}$ , gdzie  $\loglik = \sum_{i=1}^n P_{\hat{\theta}}(y_i)$

AIC:  $\frac{-2}{n} \loglik + \frac{2d}{n}$

# Kryterium BIC

$$\text{BIC: } BIC = -2\loglik + \log(n) \cdot d \text{ i}$$

$$\text{AIC: } AIC = -2\loglik + 2d$$

Dla  $n > e^2 \approx 7.4$  BIC częściej odrzuca modele.

# MDL (minimum description length)

Message	$z_1$	$z_2$	$z_3$	$z_4$
Code	0	10	110	111

Twierdzenie Shannona:  $l_i = -\log_2 P(z_i)$

Wtedy zachodziłaby następująca zależność  $-\sum_{i=1}^n P(z_i) \log_2 P(z_i) \leq E[l]$

Dla  $P(z_1) = 1/2, P(z_2) = 1/4, P(z_3) = P(z_4) = 1/8$  optymalnym kodowaniem byłoby kodowanie z poprzedniego slajdu.

Dla ciągłych:  $-\int P(z) \log_2 P(z) dz$

Od teraz zmieniamy notację zamiast  $\log_2 P(Z)$  będzie po prostu  $\log P(Z)$

## MDL przy wyborze modelu

Mamy model  $M$  z parametrami  $\theta$ . Mamy  $P(Y|M, \theta, X)$ . Wtedy długością kodowania jest  $l = -\log P(Y|M, \theta, X) - \log P(\theta|M)$

Przykład: mamy  $Y$  z  $N(\theta, \sigma^2)$ ,  $\theta$  z  $N(0, 1)$  i nie mamy zmiennych wejściowych. Wtedy powyższy zapis uprości się do

$$l = c + \log \sigma + \frac{(Y-\theta)^2}{\sigma^2} + \frac{\theta^2}{2}$$

# Krosvalidacja

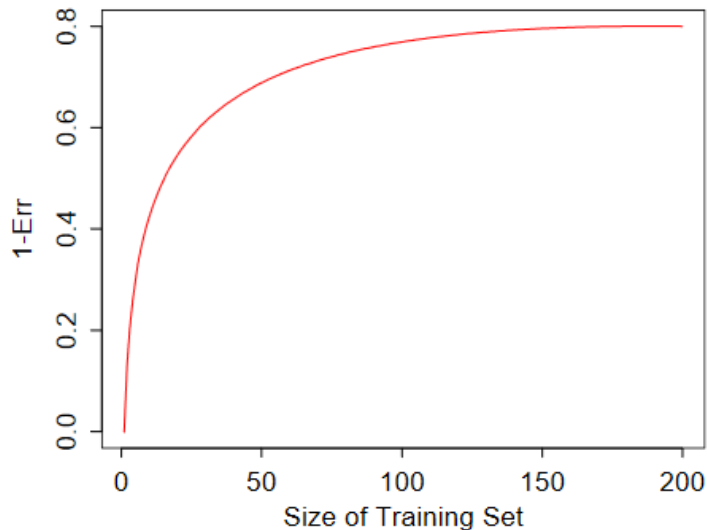
Niech  $\mathcal{K} : \{1, \dots, n\} \rightarrow \{1, \dots, k\}$ . Mamy  $\hat{f}^i(x)$

$$CV(\hat{f}) = \frac{1}{n} \sum_{j=1}^n L(y_j, \hat{f}^{\mathcal{K}(j)}(x_j))$$

Leave-one-out cv:  $\mathcal{K}(i) = i$ . Można ten zapis rozwinąć, ze względu na wybór modelu z wektorem parametrów  $\alpha$ , wtedy:

$$CV(\hat{f}, \alpha) = \frac{1}{n} \sum_{j=1}^n L(y_j, \hat{f}^{\mathcal{K}(j)}(x_j, \alpha))$$

## Jakość estymatora błędu predykcji





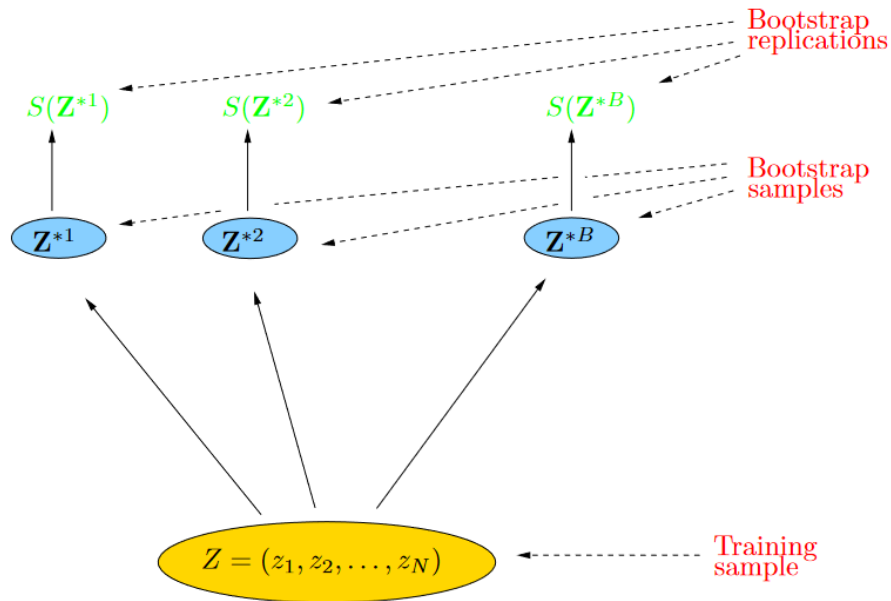
# Uogólniona krosvalidacja

Mamy jakieś liniowe dopasowanie:  $\hat{y} = Sy$ .

Wtedy

$$GCW(\hat{f}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}^{-i}(x_i))^2 = \frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - \hat{f}(x_i)}{1 - s_{ii}} \right)^2 = \frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - \hat{f}(x_i)}{1 - \text{trace}(S)/n} \right)^2$$

# Bootstrap



$$\hat{\text{var}}(\hat{S}(Z)) = \frac{1}{B-1} \sum_{b=1}^B (S(Z^{*b}) - \bar{S}^*)^2$$

Estymator błędu predykcji (gorszy):  $\hat{Err}_{boot} = \frac{1}{n} \frac{1}{B} \sum_{b=1}^B \sum_{i=1}^n L(y_i, f^{*b}(\hat{x}_i))$ .

Przykład: prawdziwy błąd predykcji 0.5. Za pomocą bootstrapu:

$P(\text{obserwacja } i \in b) = 1 - (1 - \frac{1}{n})^n \approx 1 - e^{-1} = 0.632$ , więc błąd predykcji z bootstrapu to  $1/2 * 0.368$  czyli zaniżony...

Estymator błędu predykcji (lepszy):  $\hat{Err}_{boot} = \frac{1}{n} \sum_{i=1}^n \frac{1}{\#C^{-i}} \sum_{b \in C^{-i}} L(y_i, \hat{f}^{*b}(x_i))$