

## Twierdzenie "No free lunch"

Jeśli uśrednimy jakość klasyfikacji po wszystkich możliwych problemach klasyfikacyjnych to wszystkie algorytmy dadzą ten sam błąd klasyfikacji na zbiorze testowym. To twierdzenie stanowi, że nie ma uniwersalnie najlepszego algorytmu. Zatem celem nie jest poszukiwanie takiego algorytmu, a znalezienie algorytmu optymalnego dla konkretnych danych.

### Dowód.

Rozważmy problem dwuklasowy. Jeśli Algorytm A wygrywa to zamieniamy etykiety i będzie wygrywał Algorytm B. ☐



Wolpert, D. (1996). *The Lack of a Priori Distinctions between Learning Algorithms*. Neural Computation 8(7): 341–1390.



Wolpert, D. & Macready, G. (1997). *No Free Lunch Theorems for Optimization*. IEEE Transactions on Evolutionary Computation 1(1):67–82.

## Kwadratowa analiza dyskryminacyjna

Najprostszym podejściem do zagadnienia klasyfikacji jest wykorzystanie estymacji funkcji gęstości i przyjęcie modelu parametrycznego dla gęstości, tj. przyjęcie założenia, że znana jest postać gęstości z wyjątkiem tkwiących w niej parametrów. Załóżmy, że  $\mathcal{Y} = \{1, 0\}$  i  $f_1(\mathbf{x}) = f(\mathbf{x}|Y = 1)$  oraz  $f_0(\mathbf{x}) = f(\mathbf{x}|Y = 0)$  są gęstościami  $p$ -wymiarowego rozkładu normalnego. Wówczas  $\mathbf{X}|Y = 1 \sim N_p(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$  i  $\mathbf{X}|Y = 0 \sim N_p(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ .

## Kwadratowa analiza dyskryminacyjna

### Twierdzenie

Jeżeli  $\mathbf{X} | Y = 1 \sim N_p(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$  i  $\mathbf{X} | Y = 0 \sim N_p(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ , to klasyfikator bayesowski ma postać

$$d_B(\mathbf{x}) = \begin{cases} 1, & \text{jeżeli } r_1^2 < r_0^2 + 2 \ln \left( \frac{\pi_1}{\pi_0} \right) + \ln \left( \frac{|\boldsymbol{\Sigma}_0|}{|\boldsymbol{\Sigma}_1|} \right), \\ 0, & \text{poza tym.} \end{cases}$$

gdzie

$$r_i^2 = (\mathbf{x} - \boldsymbol{\mu}_i)' \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i), \quad i = 1, 0$$

jest kwadratem *odległości Mahalanobisa*.

## Kwadratowa analiza dyskryminacyjna

Powyższy klasyfikator najczęściej zapisujemy w następującej postaci:

$$d_B(\mathbf{x}) = \arg \max_k \delta_k(\mathbf{x}),$$

gdzie

$$\delta_k(\mathbf{x}) = -\frac{1}{2} \ln |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) + \ln \pi_k.$$

Powyższa funkcja jest funkcją kwadratową argumentu  $\mathbf{x}$  i jest nazywana **kwadratową funkcją klasyfikującą** grupy  $k$ . Procedura klasyfikacji oparta na tej funkcji nosi nazwę **kwadratowej analizy dyskryminacyjnej (QDA)**.

## Kwadratowa analiza dyskryminacyjna

Parametry  $\pi_1$ ,  $\pi_0$ ,  $\mu_1$ ,  $\mu_0$ ,  $\Sigma_1$ ,  $\Sigma_0$  są w praktyce nieznane, zastępujemy je ich estymatorami z próby uczącej  $\mathcal{L}_n$  postaci:

$$\hat{\pi}_1 = \frac{1}{n} \sum_{i=1}^n Y_i, \quad \hat{\pi}_0 = 1 - \hat{\pi}_1,$$

$$\hat{\mu}_1 = \frac{1}{n_1} \sum_{Y_i=1} \mathbf{x}_i, \quad \hat{\mu}_0 = \frac{1}{n_0} \sum_{Y_i=0} \mathbf{x}_i,$$

$$\hat{\Sigma}_1 = \mathbf{S}_1 = \frac{1}{n_1 - 1} \sum_{Y_i=1} (\mathbf{x}_i - \hat{\mu}_1)(\mathbf{x}_i - \hat{\mu}_1)',$$

$$\hat{\Sigma}_0 = \mathbf{S}_0 = \frac{1}{n_0 - 1} \sum_{Y_i=0} (\mathbf{x}_i - \hat{\mu}_0)(\mathbf{x}_i - \hat{\mu}_0)',$$

gdzie  $n_1 = \sum_{i=1}^n Y_i$ ,  $n_0 = \sum_{i=1}^n (1 - Y_i)$ .

## Liniowa analiza dyskryminacyjna

Procedura ta znacznie się uprości, jeżeli dodatkowo założymy równość macierzy kowariancji w dwóch grupach  $\Sigma_0 = \Sigma_1 = \Sigma$ . W tym przypadku klasyfikator bayesowski ma postać

$$d_B(\mathbf{x}) = \arg \max_k \delta_k(\mathbf{x}),$$

gdzie

$$\delta_k(\mathbf{x}) = \mathbf{x}'\Sigma^{-1}\mu_k - \frac{1}{2}\mu_k'\Sigma^{-1}\mu_k + \ln \pi_k, \quad k = 1, 0.$$

Funkcja ta jest funkcją liniową  $\mathbf{x}$  i nosi nazwę **liniowej funkcji klasyfikującej** grupy  $k$ .

## Liniowa analiza dyskryminacyjna

Powierzchnia rozdzielająca

$$\{\mathbf{x} : \delta_1(\mathbf{x}) = \delta_0(\mathbf{x})\} = \{\mathbf{x} : \delta_{10}(\mathbf{x}) = \delta_1(\mathbf{x}) - \delta_0(\mathbf{x}) = 0\}$$

jest w tym przypadku hiperpłaszczyzną. Funkcja  $\delta_{10}(\mathbf{x}) = \delta_1(\mathbf{x}) - \delta_0(\mathbf{x})$  nazywa się **liniową funkcją dyskryminacyjną**. Ma ona postać:

$$\delta_{10}(\mathbf{x}) = (\mathbf{x} - \frac{1}{2}(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_0))' \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) + \ln \left( \frac{\pi_1}{\pi_0} \right).$$

Obserwację  $\mathbf{x}$  klasyfikujemy do grupy 1 wówczas, gdy  $\delta_{10}(\mathbf{x}) > 0$ . W przeciwnym razie klasyfikujemy do grupy 0. Procedura klasyfikacji bazująca na liniowej funkcji klasyfikującej lub liniowej funkcji dyskryminacyjnej nosi nazwę **liniowej analizy dyskryminacyjnej (LDA)**.

## Liniowa analiza dyskryminacyjna

Nieznane parametry estymujemy z próby uczącej tak jak poprzednio, przy czym:

$$\hat{\Sigma} = S = \frac{1}{n_1 + n_0 - 2} [(n_1 - 1)S_1 + (n_0 - 1)S_0].$$



## Przypadek większej liczby grup

### *Twierdzenie*

Założmy, że  $Y \in \{1, \dots, K\}$ . Jeżeli  $f_k(\mathbf{x}) = f(\mathbf{x}|Y = k)$  jest gęstością  $p$ -wymiarowego rozkładu normalnego (gaussowskiego)  $N_p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ , to klasyfikator bayesowski ma postać

$$d_B(\mathbf{x}) = \arg \max_k \delta_k(\mathbf{x}),$$

gdzie

$$\delta_k(\mathbf{x}) = -\frac{1}{2} \ln |\boldsymbol{\Sigma}_k| - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k) + \ln \pi_k.$$

Jeżeli ponadto wszystkie macierze kowariancji są sobie równe i równe macierzy  $\boldsymbol{\Sigma}$ , to

$$\delta_k(\mathbf{x}) = \mathbf{x}' \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k - \frac{1}{2} \boldsymbol{\mu}_k' \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k + \ln \pi_k.$$

## Przypadek większej liczby grup

Powierzchnia rozdzielająca grupy  $i$  oraz  $j$  ma postać

$$\{\mathbf{x} : \delta_i(\mathbf{x}) = \delta_j(\mathbf{x})\} = \{\mathbf{x} : \delta_{ij}(\mathbf{x}) = \delta_i(\mathbf{x}) - \delta_j(\mathbf{x}) = 0\},$$

gdzie

$$\begin{aligned}\delta_{ij}(\mathbf{x}) = & \frac{1}{2}\mathbf{x}'(\boldsymbol{\Sigma}_j^{-1} - \boldsymbol{\Sigma}_i^{-1})\mathbf{x} + (\boldsymbol{\mu}_i'\boldsymbol{\Sigma}_i^{-1} - \boldsymbol{\mu}_j'\boldsymbol{\Sigma}_j^{-1})\mathbf{x} \\ & + \frac{1}{2}\boldsymbol{\mu}_j'\boldsymbol{\Sigma}_j^{-1}\boldsymbol{\mu}_j - \frac{1}{2}\boldsymbol{\mu}_i'\boldsymbol{\Sigma}_i^{-1}\boldsymbol{\mu}_i \\ & + \frac{1}{2}\ln\left(\frac{|\boldsymbol{\Sigma}_j|}{|\boldsymbol{\Sigma}_i|}\right) + \ln\left(\frac{\pi_i}{\pi_j}\right), \quad i, j = 1, 2, \dots, K, \quad j \neq i\end{aligned}$$

w przypadku funkcji kwadratowej

## Przypadek większej liczby grup

oraz

$$\delta_{ij}(\mathbf{x}) = \mathbf{x}'\boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j) - \frac{1}{2}(\boldsymbol{\mu}_i + \boldsymbol{\mu}_j)'\boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j) \\ + \ln\left(\frac{\pi_i}{\pi_j}\right), \quad i, j = 1, 2, \dots, K, \quad j \neq i.$$

w przypadku funkcji liniowej.

## Przypadek większej liczby grup

Występujące w powyższych wzorach parametry nie są zazwyczaj znane i w praktyce należy zastąpić je ich estymatorami z próby uczącej. Jeżeli próba ucząca zawiera  $n_i$  obserwacji z  $i$ -tej grupy,  $n_1 + n_2 + \dots + n_K = n$  oraz  $\mathbf{X}_{ij}$  jest  $j$ -tą obserwacją z  $i$ -tej grupy, to estymatory nieznanych parametrów są równe:

$$\hat{\pi}_k = \frac{n_k}{n},$$

$$\hat{\mu}_k = \bar{\mathbf{X}}_k = \frac{1}{n_k} \sum_{j=1}^{n_k} \mathbf{X}_{kj},$$

$$\hat{\Sigma}_k = \mathbf{S}_k = \frac{1}{n_k - 1} \sum_{j=1}^{n_k} (\mathbf{X}_{kj} - \bar{\mathbf{X}}_k)(\mathbf{X}_{kj} - \bar{\mathbf{X}}_k)',$$

$$\hat{\Sigma} = \mathbf{S} = \frac{1}{n - K} \sum_{k=1}^K \sum_{j=1}^{n_k} (\mathbf{X}_{kj} - \bar{\mathbf{X}}_k)(\mathbf{X}_{kj} - \bar{\mathbf{X}}_k)'.$$

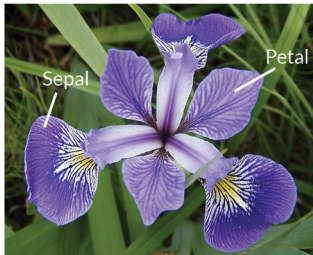
## Zbiór danych iris

The iris flower data set is a multivariate data set introduced by R. Fisher (1936). The data set consists of 50 samples from each of three species of iris (setosa, virginica and versicolor). Four features were measured: the length and the width of the sepals and petals, in centimeters. Based on it, Fisher developed a linear discriminant model to distinguish the species from each other.



Sir Ronald Aylmer Fisher  
(1890-1962)

## Zbiór danych iris



**Iris Versicolor**



**Iris Setosa**



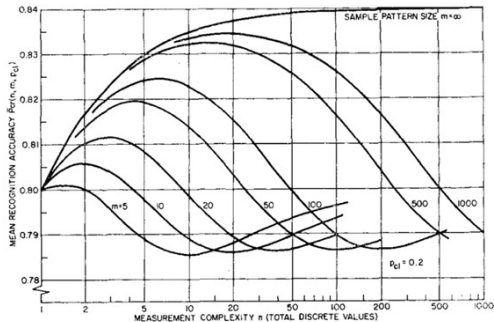
**Iris Virginica**



Fisher R. A. (1936). *The use of multiple measurements in taxonomic problems*. Annals of Eugenics 7(2):179–188.

## Przekleństwo wymiarowości

Dla ustalonej wielkości próby uczącej jakość klasyfikatora (regresji) na początku wzrasta wraz ze zwiększaniem liczby cech, a następnie spada. Zjawisko to nazywane jest paradoksem Hughes'a.



Hughes, G.F. (1968). *On the mean accuracy of statistical pattern recognizers*. IEEE Transactions on Information Theory 14(1):55–63.