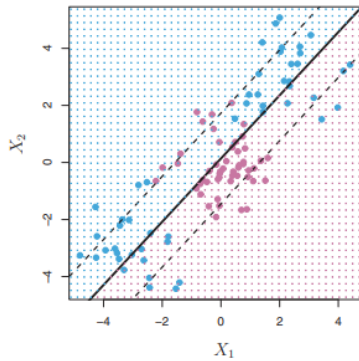
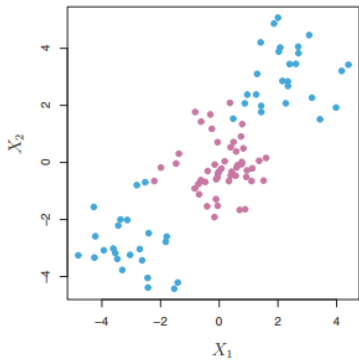


# Support Vector Machines

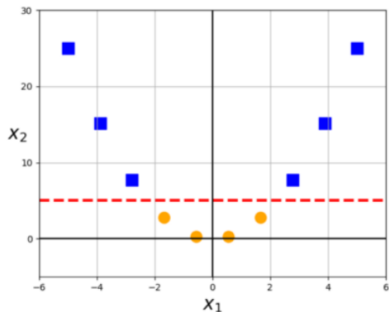
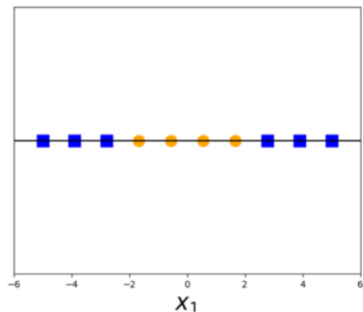
Zofia Dzedzic

Uniwersytet Wrocławski



**Rysunek:** Działanie klasyfikatora wektorów nośnych w przypadku, gdy obserwacje dzielą się na dwa zbiory rozdzielone nieliniową granicą. Źródło: „Introduction to Statistical Learning”

# Transformacja predyktorów – intuicja



Rysunek: Źródło: <https://towardsdatascience.com/the-kernel-trick-c98cdbcaeb3f>

Można powiększać przestrzeń cech, dodając funkcje obecnych predyktorów, np. ich kwadraty.

$$X_1, X_1^2, X_2, X_2^2, \dots, X_p, X_p^2$$

Wtedy problem optymalizacyjny wygląda następująco

$$\text{maksymalizuj } M \\ \beta_0, \beta_{11}, \beta_{12}, \dots, \beta_{p1}, \beta_{p2}, \epsilon_1, \dots, \epsilon_n$$

$$\text{pod warunkiem } y_i \left( \beta_0 + \sum_{j=1}^p \beta_{j1} x_{ij} + \sum_{j=1}^p \beta_{j2} x_{ij}^2 \right) \geq M(1 - \epsilon_i),$$

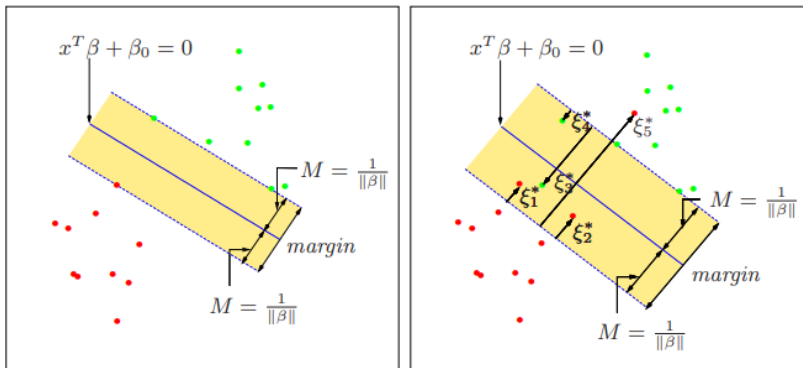
$$\sum_{i=1}^n \epsilon_i, \epsilon_i \geq 0, \sum_{j=1}^p \sum_{k=1}^2 \beta_{jk}^2 = 1$$

Takie podejście jest bardzo złożone obliczeniowo, a dla większej liczby cech niemal niewykonalne. Maszyny wektorów nośnych (SVM) pomagają poradzić sobie z tym problemem. Aby pokazać ich działanie, pokażemy najpierw, jak wygląda rozwiązanie problemu optymalizacyjnego klasyfikatora wektorów nośnych.

Użyjemy dla wygody obliczeń równoważnego zapisu tego zagadnienia

$$\min_{\beta_0, \beta} \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^N \xi_i$$

pod warunkiem  $y_i(x_i^T \beta + \beta_0) \geq 1 - \xi_i \quad \forall i$



Rysunek: Źródło: „The Elements of Statistical Learning”

$$L_P = \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i [y_i (x_i^T \beta + \beta_0) - (1 - \xi_i)] - \sum_{i=1}^N \mu_i \xi_i$$

$$\beta = \sum_{i=1}^N \alpha_i y_i x_i,$$

$$0 = \sum_{i=1}^N \alpha_i y_i,$$

$$\alpha_i = C - \mu_i, \forall i,$$

$$L_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{i'=1}^N \alpha_i \alpha_{i'} y_i y_{i'} x_i^T x_{i'},$$

$$\hat{\beta} = \sum_{i=1}^N \hat{\alpha}_i y_i x_i$$

Ostatecznie klasyfikator przyjmuje więc postać

$$\hat{f}(\mathbf{x}) = \mathbf{x}^T \hat{\beta} + \hat{\beta}_0 = \sum_{i=1}^N \hat{\alpha}_i y_i \langle \mathbf{x}, \mathbf{x}_i \rangle + \hat{\beta}_0.$$

Podobnie dla przypadku predyktorów przekształconych przez funkcję  $h$  będziemy mieli

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^N \hat{\alpha}_i y_i \langle h(\mathbf{x}), h(\mathbf{x}_i) \rangle + \hat{\beta}_0.$$

Jak widzimy, w tym rozwiązaniu istotne są więc jedynie iloczyny skalarne obserwacji testowej i obserwacji treningowych (a nawet więcej – tylko wektorów nośnych).



Możemy zastąpić ten iloczyn skalarny pewnym jego uogólnieniem. Taką funkcję nazwiemy funkcją jądra (kernel function). Określa ona podobieństwo między dwiema obserwacjami.

## Kernel function

Dla  $x, x' \in X$  i funkcji rzutu  $h : X \rightarrow \mathcal{R}^n$  funkcją jądra nazywamy funkcję  $K$ , taką że

$$K(x, x') = \langle h(x), h(x') \rangle$$

Jest to więc funkcja, która przyjmuje za argumenty wektory z oryginalnej przestrzeni, a zwraca iloczyn skalarny wektorów z przestrzeni cech (feature space).

$$\mathbf{x} = (x_1, x_2, x_3)^T$$

$$\mathbf{y} = (y_1, y_2, y_3)^T$$

$$\phi(\mathbf{x}) = (x_1^2, x_1x_2, x_1x_3, x_2x_1, x_2^2, x_2x_3, x_3x_1, x_3x_2, x_3^2)^T$$

$$\phi(\mathbf{y}) = (y_1^2, y_1y_2, y_1y_3, y_2y_1, y_2^2, y_2y_3, y_3y_1, y_3y_2, y_3^2)^T$$

$$\phi(\mathbf{x})^T \phi(\mathbf{y}) = \sum_{i,j=1}^3 x_i x_j y_i y_j$$

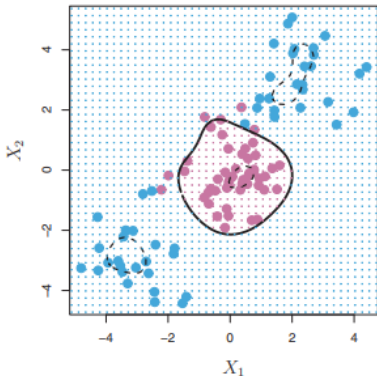
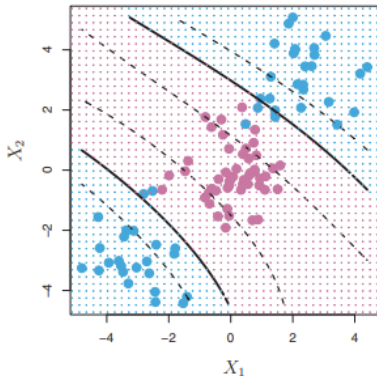
$$\begin{aligned} k(\mathbf{x}, \mathbf{y}) &= (\mathbf{x}^T \mathbf{y})^2 \\ &= (x_1 y_1 + x_2 y_2 + x_3 y_3)^2 \\ &= \sum_{i,j=1}^3 x_i x_j y_i y_j \end{aligned}$$

Najczęściej stosowane funkcje jądra:

*d*th-Degree polynomial:  $K(x, x') = (1 + \langle x, x' \rangle)^d$ ,

Radial basis:  $K(x, x') = \exp(-\gamma \|x - x'\|^2)$ ,

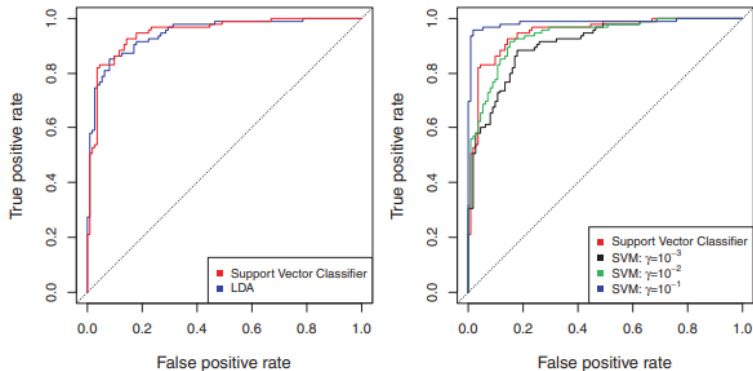
Neural network:  $K(x, x') = \tanh(\kappa_1 \langle x, x' \rangle + \kappa_2)$



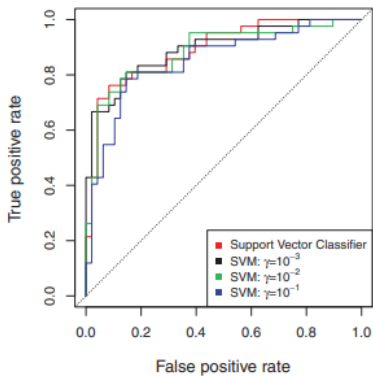
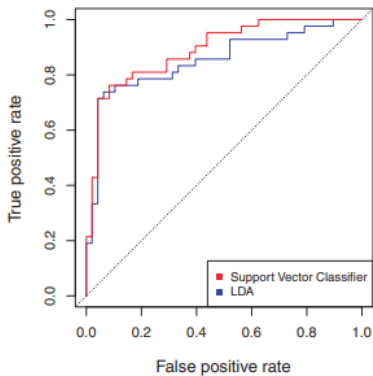
Rysunek: Po lewej: przykład zastosowania wielomianowego jądra stopnia 3 do danych z Rys. 1. Po prawej: jądro radialne. Źródło: „Introduction...”

$$\begin{aligned}
\exp\left(-\frac{1}{2}\|\mathbf{x} - \mathbf{x}'\|^2\right) &= \exp\left(\frac{2}{2}\mathbf{x}^\top \mathbf{x}' - \frac{1}{2}\|\mathbf{x}\|^2 - \frac{1}{2}\|\mathbf{x}'\|^2\right) \\
&= \exp(\mathbf{x}^\top \mathbf{x}') \exp\left(-\frac{1}{2}\|\mathbf{x}\|^2\right) \exp\left(-\frac{1}{2}\|\mathbf{x}'\|^2\right) \\
&= \sum_{j=0}^{\infty} \frac{(\mathbf{x}^\top \mathbf{x}')^j}{j!} \exp\left(-\frac{1}{2}\|\mathbf{x}\|^2\right) \exp\left(-\frac{1}{2}\|\mathbf{x}'\|^2\right) \\
&= \sum_{j=0}^{\infty} \sum_{\sum n_i=j} \exp\left(-\frac{1}{2}\|\mathbf{x}\|^2\right) \frac{x_1^{n_1} \cdots x_k^{n_k}}{\sqrt{n_1! \cdots n_k!}} \exp\left(-\frac{1}{2}\|\mathbf{x}'\|^2\right) \frac{x_1'^{n_1} \cdots x_k'^{n_k}}{\sqrt{n_1! \cdots n_k!}}
\end{aligned}$$

Rysunek: Źródło: <https://en.wikipedia.org>



**Rysunek:** Po lewej: porównanie krzywych ROC dla SVC i LDA. Po prawej: porównanie krzywych ROC dla różnych parametrów  $\gamma$  w radialnej funkcji bazowej. Źródło: „Introduction...”



Rysunek: Krzywe ROC jak w Rysunku powyżej, dla zbioru testowego. Źródło: „Introduction...”

# One-Versus-One Classification

Rozważmy teraz, jak zastosować SVM do więcej niż dwóch różnych klas. Pokażemy dwa najczęściej stosowane podejścia.

One-Versus-One Classification polega na dopasowaniu  $\binom{K}{2}$  modeli SVM, parami dla każdych dwóch z  $K$  klas, a następnie zaklasyfikowaniu obserwacji testowej do jednego zbioru z każdej pary, korzystając ze stworzonych modeli. Na koniec należy zliczyć częstotliwość klasyfikacji do każdej z klas i wybrać tę klasę, do której obserwacja testowa „wpadała” najczęściej.



# One-Versus-All Classification

W One-Versus-All Classification tworzymy  $K$  modeli, które pozwalają zbadać, czy obserwacja zalicza się do  $k$ -tej klasy, czy do którejś z  $K - 1$  pozostałych. Następnie dla testowej obserwacji  $x^*$  obliczamy  $\beta_{0k} + \beta_{1k}x^* + \dots + \beta_{pk}x^*$  i wybieramy tę klasę, dla której ta wartość jest największa.