

# Introduction to statistical learning

## 3.7 Supervised learning: SVM

V. Lefieux

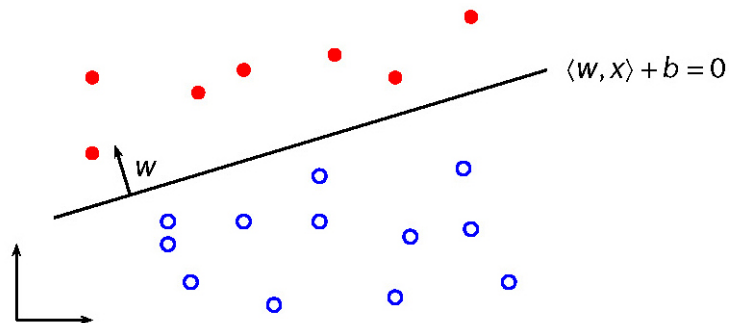
June 2018



- ▶ SVM (Support Vector Machine) are machine learning algorithms for classification (or regression) problems.
- ▶ First idea: find a plane that separates classes in the feature space.

# Linearly separable data I

References



# Linearly separable data II

Data are linearly separable if there exists  $w \in \mathbb{R}^p$  and  $b \in \mathbb{R}$  such that :

$$\begin{cases} y_i = 1 & \text{if } \langle w_i, x_i \rangle + b > 0 \\ y_i = -1 & \text{if } \langle w_i, x_i \rangle + b < 0 \end{cases}$$

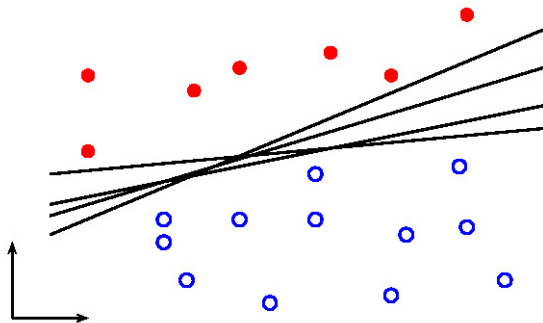
or in an equivalent manner:

$$y_i (\langle w_i, x_i \rangle + b) > 0 .$$

$\langle w_i, x_i \rangle + b = 0$  is the equation of a hyperplane.

# A lot of possible separating hyperplanes

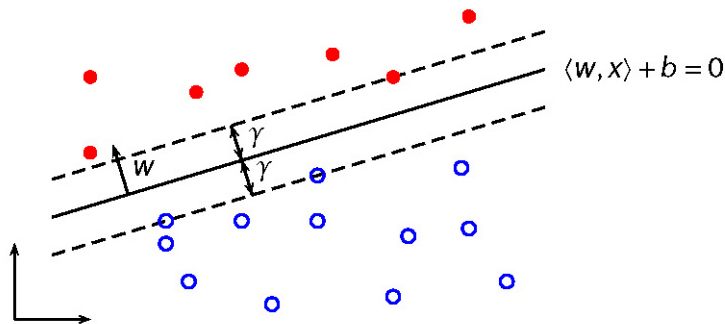
References



# Vapnik solution

Among all separating hyperplanes, find the one which maximizes margin  $\gamma$  between classes.

References



The margin is the distance between the separation boundary and the nearest samples which are called support vectors.

# Maximal margin classifier optimization problem I

References

Maximizing  $\gamma$  is equivalent to:

$$\begin{aligned} & \min \|w\|^2 \\ & \text{subject to } \forall i \in \{1, \dots, n\} : y_i (\langle w, x_i \rangle + b) \geq 1 . \end{aligned}$$

# Maximal margin classifier optimization problem II

References

$$\begin{array}{ll} \max & \gamma \\ \text{subject to} & \begin{cases} \sum_{j=1}^p w_j^2 = 1, \\ y_i (\langle w_i, x_i \rangle + b) \geq \gamma \end{cases} \end{array}$$

where  $C$  is a regularization parameter.



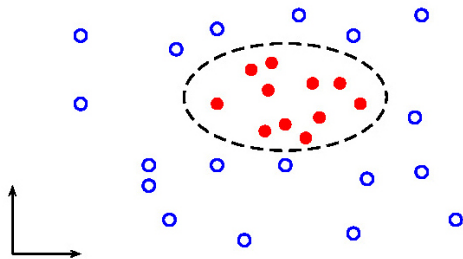
# Support vector classifier

The support vector classifier maximizes a soft margin.

$$\begin{array}{ll} \max & \gamma \\ \text{subject to} & \left\{ \begin{array}{l} \sum_{j=1}^p w_j^2 = 1 , \\ y_i (\langle w_i, x_i \rangle + b) \geq \gamma (1 - \xi_i) , \\ \xi_i \geq 0 , \\ \sum_{i=1}^n \xi_i \leq C \end{array} \right. \end{array}$$

# In the case of non-separable data I

References



# In the case of non-separable data II

References

- ▶ Sometimes doesn't work, whatever the value of  $C$ .
- ▶ Enlarge the space of features by including transformations, e.g.  $X_1, X_2, X_1^2, X_2^2$ .
- ▶ Fit a support-vector classifier in the enlarged space (non-linear decision boundaries in the original space).
- ▶ One can use polynomials, more generally one use kernels: polynomial, radial, gaussian, . . .

- Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The elements of statistical learning. Data Mining, inference, and prediction*. Springer Series in Statistics. Springer, 2 edition.
- James, G., Witten, D., Hastie, T., and Tibshirani, R. (2015). *An introduction to statistical learning with applications in R*. Springer Texts in Statistics. Springer.