

# Mathematics for Machine Learning

Prof Willie Brink

Applied Mathematics, Stellenbosch University

## Lecture 11: Classification with SVMs

# Contents of the module

Chapter 02: Linear Algebra

Chapter 03: Analytic Geometry

Chapter 04: Matrix Decompositions

Chapter 05: Vector Calculus

---

Chapter 06: Probability and Distributions

Chapter 07: Continuous Optimisation

Chapter 08: When Models Meet Data

Chapter 09: Linear Regression

---

Chapter 10: Dimensionality Reduction with Principal Component Analysis

Chapter 11: Density Estimation with Gaussian Mixture Models

**Chapter 12: Classification with Support Vector Machines**

## Introduction

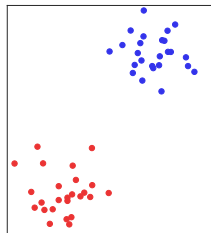
Consider the following **binary classification** problem.

Given training data  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$  consisting of input vectors  $\mathbf{x}_n \in \mathbb{R}^D$  and corresponding labels  $y_n \in \{-1, 1\}$ , learn a classifier  $f(\mathbf{x})$  such that

$$f(\mathbf{x}_n) \begin{cases} > 0, & \text{if } y_n = 1 \\ < 0, & \text{if } y_n = -1 \end{cases} \quad \text{or equivalently: } y_n f(\mathbf{x}_n) > 0$$

We'll consider **linear** classifiers of the form  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ .

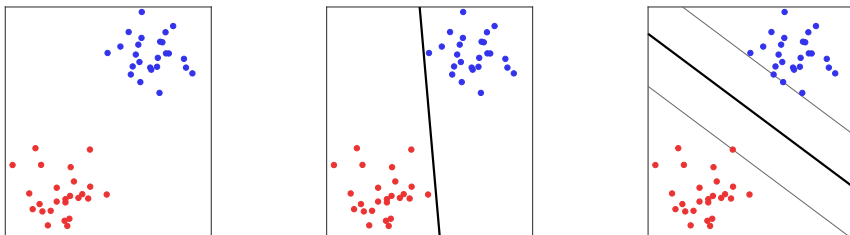
For  $D = 2$ ,  $f(\mathbf{x}) = 0$  represents a line that splits  $\mathbb{R}^2$  into two regions: where  $f(\mathbf{x}) > 0$ , and where  $f(\mathbf{x}) < 0$ .



## 12.1 Separating hyperplanes

A linear classifier  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$  leads to a linear decision boundary,  $f(\mathbf{x}) = 0$ , requiring the two classes to be **linearly separable**.

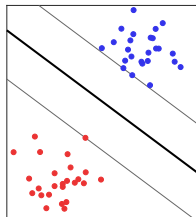
Parameters  $\mathbf{w}$  and  $b$  influence the position and orientation of the decision boundary. What would be the best choice for these parameters?



The **maximum margin solution** might be most stable under perturbations of the input.

We want to maximise the margin around the decision boundary.  
Training samples touching the margin are called **support vectors**.

Note:  $\mathbf{w}^T \mathbf{x} + b = 0$  and  $c(\mathbf{w}^T \mathbf{x} + b) = 0$  define the same line  
(or hyperplane in  $D$  dimensions).



We choose the scale of  $\mathbf{w}$  and  $b$  such that  $\mathbf{w}^T \mathbf{x}_+ + b = 1$  and  $\mathbf{w}^T \mathbf{x}_- + b = -1$  for positive and negative support vectors, respectively.

With this normalisation, the width of the margin turns out to be  $2/\|\mathbf{w}\|$ .

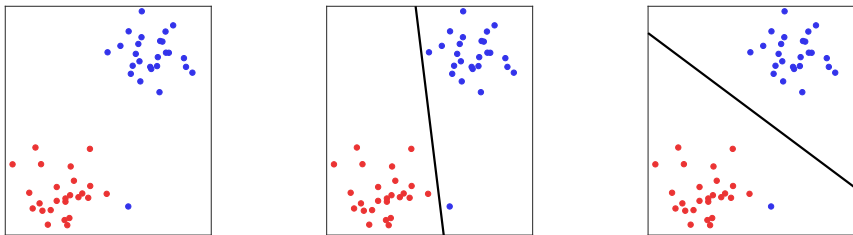
Learning the classifier is done through solving the following optimisation problem:

$$\max_{\mathbf{w}, b} \frac{2}{\|\mathbf{w}\|} \quad \text{subject to} \quad \mathbf{w}^T \mathbf{x}_n + b \begin{cases} \geq 1, & \text{if } y_n = 1 \\ \leq -1, & \text{if } y_n = -1 \end{cases} \quad n = 1, \dots, N$$

## 12.2 Primal support vector machine

Equivalent problem:  $\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$  subject to  $y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \quad n = 1, \dots, N$

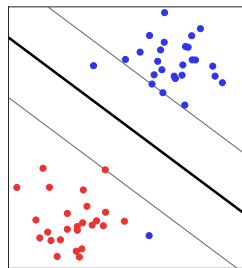
Now, should the classifier fit the training data *perfectly*?



There could be a trade-off between the **margin size** and the **number of mistakes** made on the training data.

We assign a **slack variable**  $\xi_n$  to every  $\mathbf{x}_n$  in the training set.

- if  $\xi_n \leq 0$ ,  $\mathbf{x}_n$  is classified correctly and is outside the margin
- if  $0 < \xi_n < 1$ ,  $\mathbf{x}_n$  is classified correctly but violates the margin
- if  $\xi_n \geq 1$ ,  $\mathbf{x}_n$  is misclassified



The optimisation problem becomes:

$$\min_{\mathbf{w}, b, \xi_1, \dots, \xi_N} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n \quad \text{subject to} \quad y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n \quad n = 1, \dots, N$$

The parameter  $C$  is a **regulariser**.

- a small  $C$  allows constraints to be easily ignored  $\implies$  large margin
- a large  $C$  makes constraints hard to ignore  $\implies$  narrow margin

## Multi-class SVM

SVMs find a maximum margin solution to separate **two** classes in  $\mathbb{R}^D$ , and there is no definitive multi-class SVM formulation.

### One-vs-rest

**Training:** learn an SVM for each class vs the others.

**Prediction:** apply each SVM to the test sample, and assign class according to the SVM that returns the highest decision score.

### One-vs-one

**Training:** learn an SVM for each pair of classes.

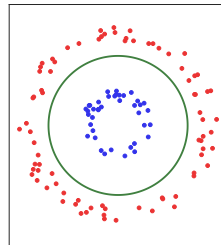
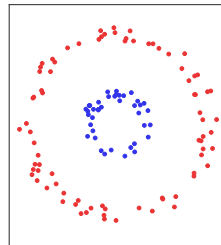
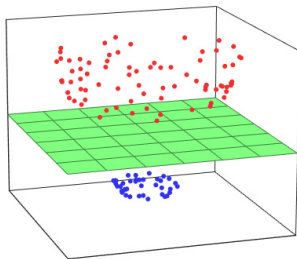
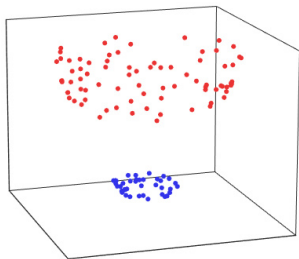
**Prediction:** each SVM votes on a class to assign to the test sample, and we pick the class with most votes (breaking ties with decision scores).



## 12.4 Kernels

What if the classes are not even almost linearly separable?

Map feature vectors to another space, where the classes are linearly separable!



Choose a **kernel function**  $k : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$  for which there exists a function  $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^{D'}$  such that  $k(\mathbf{a}, \mathbf{b}) = \phi(\mathbf{a})^\top \phi(\mathbf{b})$ .

*When we train an SVM, or predict with an SVM, we don't need to know what  $\phi$  is!*

We're never interested in the  $D'$ -dimensional version of a vector, only in dot products of two of those vectors, so we just apply the kernel function without ever specifying  $\phi$ .

Examples (where  $\gamma$ ,  $r$ ,  $d$  are hyperparameters):

RBF kernel:  $k(\mathbf{a}, \mathbf{b}) = e^{-\gamma \|\mathbf{a} - \mathbf{b}\|^2}$

Polynomial:  $k(\mathbf{a}, \mathbf{b}) = (\gamma(\mathbf{a}^\top \mathbf{b}) + r)^d$

Sigmoidal:  $k(\mathbf{a}, \mathbf{b}) = \tanh(\mathbf{a}^\top \mathbf{b} + r)$

The choice of a kernel and tuning of its hyperparameters can be done through validation.

