

Support Vector Machines

Christos Dimitrakakis

November 21, 2023

Outline

Background

- Maximal margin

Support vectors

- Support Vector Machines

Hyperplane

If $\mathbf{x} \in \mathbb{R}^n$, then an affine subspace of dimension $n - 1$ is a hyperplane.

Definition

A hyperplane in \mathbb{R}^n is the set of points satisfying

$$\{\mathbf{x} : \beta_0 + \boldsymbol{\beta}^\top \mathbf{x} = 0\}$$

Separating hyperplane

Consider a dataset (\mathbf{x}_t, y_t) of points with $y_t \in \{-1, 1\}$. If

$$(\beta_0 + \boldsymbol{\beta}^\top \mathbf{x}_t)y_t > 0 \quad \forall t$$

then the hyperplane separates the dataset.

The maximal margin hyperplane

The margin

For any a dataset (\mathbf{x}_t, y_t) , and hyperplane we can define the margin

$$M = \min_t (\beta_0 + \boldsymbol{\beta}^\top \mathbf{x}_t) y_t$$

as the minimum distance between the hyperplane and a correctly classied point.

The maximal margin hyperplane

Similarly, there is some $\beta_0, \boldsymbol{\beta}$ that achieves the maximum separation:

$$\max_{\beta_0, \boldsymbol{\beta}} \min_t (\beta_0 + \boldsymbol{\beta}^\top \mathbf{x}_t) y_t$$

The maximum margin classifier

The optimisation problem

We can write the problem like this

$$\begin{aligned} \max_{\beta_0, \beta, M} M & \quad (\text{maximise the margin}) \\ \text{s.t. } \|\beta\| &= 1 & (\text{invariance}) \\ y_t(\beta_0 + \beta^\top \mathbf{x}_t) &\geq M \quad \forall t \in [T]. & (\text{margin for all examples}) \end{aligned}$$

And we can divide by $\|\beta\|$ to remove the norm constraint:

$$y_t(\beta_0 + \beta^\top \mathbf{x}_t) \geq M\|\beta\|, \quad \forall t \in [T]$$

Setting $\|\beta\| = 1/M$, we can rewrite this as

$$\begin{aligned} \min_{\beta_0, \beta} \|\beta\|^2 \\ \text{s.t. } y_t(\beta_0 + \beta^\top \mathbf{x}_t) &\geq 1 \quad \forall t \end{aligned}$$

Quadratic programming

A quadratic program has the form:

$$\begin{aligned} \min_{\beta} & \|\beta\|^2 \\ \text{s.t.} & \beta^\top x_t \geq 1 \forall t. \end{aligned}$$

A constrained optimisation problem

$$\begin{aligned} \min_{\beta} & f(\beta) \\ \text{s.t.} & g_i(\beta) = 0 \forall i \\ & h_i(\beta) \geq 0 \forall i. \end{aligned}$$

We can use the **Lagrange** method of multipliers to solve these problems.

Lagrange methods

Lagrange multipliers

For any local minimum β^* , there is a unique vector $\lambda^* \in \mathbb{R}^n$ so that

$$\nabla f(\beta^*) + \sum_i \lambda_i^* \nabla h_i(\beta^*) = 0.$$

The Lagrangian function

We first augment the original cost function to the **Lagrangian**

$$L(\beta, \lambda) = f(\beta) + \sum_{i=1}^n \lambda_i h_i(\beta)$$

For any local minimum $\beta^*, \lambda^*, \nabla_{\beta} L(\beta^*, \lambda^*) = 0, \nabla_{\lambda} L(\beta^*, \lambda^*) = 0$.

The Lagrange dual function

The dual function D is always concave:

$$D(\lambda) = \inf_{\beta} L(\beta, \lambda).$$

Support vector machines

- ▶ Hyperplanes do not always work
- ▶ How about a non-linear boundary?
- ▶ Instead of mapping the inputs through a non-linearity, map inner products to a kernel

Kernelised linear functions

We can rewrite

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^n \beta_i x_i$$

in terms of a kernel $K : X \times X \rightarrow \mathbb{R}$:

$$f(\mathbf{x}) = \beta_0 + \sum_{t=1}^T \alpha_t K(\mathbf{x}, \mathbf{x}_t), \quad K(\mathbf{x}, \mathbf{x}_t) = \mathbf{x}^\top \mathbf{x}_t$$

because we can find α so that

$$\sum_i \sum_t (\alpha_t x_{t,i}) x_i = \sum_i \beta_i x_i.$$

In fact it is sufficient to have: $\sum_t \alpha_t x_{t,i} = \beta_i$.

Kernels

Radial Basis Function

A simple type of non-linear layer in neural networks:

$$f(x) = \sum_i \alpha_i K(x, c_i), \quad K(x, c_i) = \exp(-\|x - c_i\|^2),$$

where c_t are **fixed centroids**

Kernels in SVMs

Instead of fixed kernels, use the **training data**:

$$f(x) = \sum_t \alpha_t K(x, x_t),$$

Some common kernel choices

- ▶ Linear: $K(x, x') = x^\top x$.
- ▶ RBFs: $K(x, x') = \exp(-\|x - x'\|^2)$
- ▶ Polynomial: $K(x, x') = (1 + x^\top x)^d$.

Kernels as features*

Some kernels can be rewritten in terms of a feature mapping $\phi : X \rightarrow Z$

$$K(\mathbf{x}, \mathbf{x}') = \phi^\top(\mathbf{x})\phi(\mathbf{x})$$

- ▶ The mapping ϕ is implicit, and never computed.
- ▶ The dimension of Z can be infinite.
- ▶ So-called Mercer kernels are symmetric: $K(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}', \mathbf{x})$.

Mercer kernels

$K : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ is a Mercer kernel, if for any $\{\mathbf{x}_t : t \in [T]\}$, the kernel matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$

$$\mathbf{K} \triangleq [K(\mathbf{x}_i, \mathbf{x}_j)]_{i,j \in [T]}$$

is symmetric positive semi-definite, i.e.

$$\mathbf{z}^\top \mathbf{K} \mathbf{z} \geq 0 \quad \forall \mathbf{z} \in \mathbb{R}^n.$$