

Support vector machine!

Meichen Lu (ml574@cam.ac.uk)

April 19, 2018

Contents

1	introduction	1
1.1	Motivation	1
1.2	Applications	1
2	Algorithms	1
2.1	Intuition	1
2.2	More formal walk-through	1

1 introduction

1.1 Motivation

1.2 Applications

- First applied for MNIST text recognition

2 Algorithms

2.1 Intuition

Begin with classification with two classes.

- Target is to find a boundary that is ‘best’ separates two classes. It means that we leave as big a **margin** as possible for both classes.
- To be able to compare different fitting, some kind of ‘normalisation’ is needed

2.2 More formal walk-through

For a linear decision boundary, we can express it as $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$. Thus, one class will be $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b > 0$ and the other class $f(\mathbf{x}) < 0$. Introduce the class label as a random variable

y being 1 and -1 for the two scenarios such that $\gamma^{(i)} = y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b)$ measures the distance of $\mathbf{x}^{(i)}$ from the decision boundary. It is called **functional margin** or **geometric margin**.

Support vector the vectors \mathbf{x} that lie on the boundary are called the support vectors. The following definition is introduced for computational convenience

$$yf(\mathbf{x}) = y(\mathbf{w}^T \mathbf{x} + b) = 1 \quad (1)$$

Next we need to find how to express the ‘maximum margin’ in maths. If we have two support vectors \mathbf{x}^+ and \mathbf{x}^- that belongs to the two classes, the width of the gap is thus $(\mathbf{x}^+ - \mathbf{x}^-) \cdot \frac{\mathbf{w}}{\|\mathbf{w}\|}$. Substituting eq. 1, we find that: width = $\frac{2}{\|\mathbf{w}\|}$

$$\max \frac{2}{\|\mathbf{w}\|} \equiv \min \|\mathbf{w}\| \equiv \min \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

Kernel method For data that is not linearly separable, we can increase the dimensionality. If we apply a fixed feature-space transformation $\phi(\mathbf{x})$, we can convert the original input space \mathbf{x} to a higher-dimensional space $\phi(\mathbf{x})$, where the contrast between two classes is exaggerated.

The functional margin becomes

$$y(\mathbf{w}^T \phi(\mathbf{x}) + b) = 1 \quad (2)$$

Lagrangian multiplier In order to solve this constrained optimization problem, we introduce Lagrange multipliers $a_n = 0$, with one multiplier an for each of the constraints.

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^N a_n \{y_n(\mathbf{w}^T \phi(\mathbf{x}_n) + b) - 1\} \quad (3)$$

Taking the derivative of $L(\mathbf{w}, b, \mathbf{a})$ w.r.t \mathbf{w} and b , we obtain

$$\mathbf{w} = \sum_{n=1}^N a_n y_n \phi(\mathbf{x}_n) \quad (4)$$

$$0 = \sum_{n=1}^N a_n y_n \quad (5)$$

Eliminating \mathbf{w} and b from $L(\mathbf{w}, b, \mathbf{a})$ gives the *dual representation*

$$\tilde{L}(\mathbf{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) \quad (6)$$

We define $k(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$ to be the kernel function. After solving the dual problem, we can predict y using the sign of the $f(\mathbf{x})$

$$f(\mathbf{x}) = \sum_{n=1}^N a_n y_n k(\mathbf{x}, \mathbf{x}_n) + b \quad (7)$$

The Karush-Kuhn-Tucker (KKT) conditions are satisfied in this case:

$$a_n \leq 0 \tag{8}$$

$$y_n f(\mathbf{x}_n) - 1 \leq 0 \tag{9}$$

$$a_n \{y_n f(\mathbf{x}_n) - 1\} = 0 \tag{10}$$

Therefore, either $a_n = 0$ or $y_n f(\mathbf{x}_n) = 1$. Any data point for which $a_n = 0$ (non-support vector) will not appear in the sum in Eq. 7 and hence plays no role in making predictions for new data points.

Regularization Cost function

$$\sum_{n=1}^N E_{\infty}(f(\mathbf{x}_n)y_n - 1) + \lambda \|\mathbf{w}\| \tag{11}$$

where

$$E_{\infty}(z) = \begin{cases} 0, & \text{if } z \leq 0, \\ \infty, & \text{otherwise} \end{cases} \tag{12}$$