

Lecture 12: Kernel Methods

26/09/2022

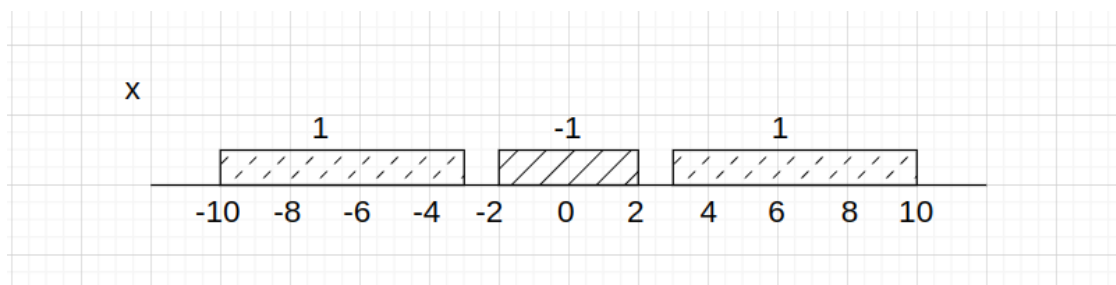
Lecturer: Abir De

Scribe: Group 23, Group 25, Group 26

In this Lecture we discuss the smooth transition from Linear to Non-Linear space of functions. We discuss treatments where $f(x)$ is non-linear in x that is, $f(x) \neq w^T x$, but we can construct it using our prior knowledge. We transform the data into a higher dimensional space and find out that data is separable when transformed into higher dimensions. We introduce "kernel method", a method of applying SVMs to problems with non linear classification boundaries. kernels are basically measure of similarity between two vectors.

1 Introduction

Remember we had a question in midsem (question 2.d) where the labels are +1 for all x such that $|x| > 2$ and -1 otherwise. We were asked to provide a 1D transformation $\phi : \mathbb{R} \rightarrow \mathbb{R}$ so that SVM applied to the given dataset.

Figure 1: Initial input space : x

As it can be seen that SVM can not be applied to the initial input space.

Now let us first define a mapping $\phi : \mathbb{R} \mapsto \mathbb{R}^2$ as follows:

$$\phi(x) = (x, x^2)$$

Now our prediction will be $h(\phi(x))$ instead of $h(x)$, where $\phi(x) = x^2$ and h is the learnt classifier, it is possible to apply SVM in the transformed space as seen in the Figure 2.

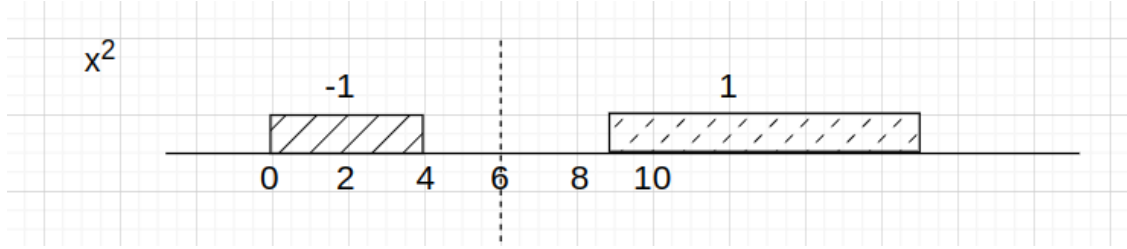


Figure 2: Transformed input space : x^2

The Main Concept behind Kernel Methods is that there is a possibility that, data-points that are not linearly-separable in lower dimensions are linearly-separable in higher dimensions.

2 Kernel Methods

The Basic Algorithm will be:-

- Given some Data Set $S = (\{\mathbf{x}_i\}_{i=0}^n, \{y_i\}_{i=0}^n)$, where $\mathbf{x}_i \in \mathbb{R}^d$
- Consider a function ϕ such that $\phi(\mathbf{x}) \in \mathbb{R}^{d'}$ where $d' > d$ (and can even be infinity)
- Create a new Data Set $\hat{S} = (\{\phi(\mathbf{x}_i)\}_{i=0}^n, \{y_i\}_{i=0}^n)$
- Train a linear Predictor h over \hat{S}
- And then the prediction of any point \mathbf{x}_{test} in the test dataset is given by $h(\phi(\mathbf{x}_{test}))$

Thus the prediction is given by $\mathbf{w}^T \phi(\mathbf{x})$ where both $\mathbf{w}, \phi(\mathbf{x}) \in \mathbb{R}^{d'}$

As d' can reach infinity, this is theoretically possible for non linear functions but calculating and storing \mathbf{w} and $\phi(\mathbf{x})$ becomes practically impossible, but the dot product $\mathbf{w}^T \phi(\mathbf{x})$ is a scalar.

Here we are enriching the expressive power of halfspaces by first mapping the data into a high dimensional feature space, and then learning a linear predictor in that space. While this approach greatly extends the expressiveness of half-space predictors, it raises both sample complexity and computational complexity challenges. We tackle this using the method of *kernels*.

A popular choice for the mapping ϕ is a polynomial mapping: $x \rightarrow (1, x, x^2, \dots)$

The Setup:

Given a mapping ϕ we are left to solve the optimization problem:

$$\min_w f(\{y_i\}_i^n, \{\mathbf{w}^T \phi(\mathbf{x}_i)\}_i^n) + \lambda R(\mathbf{w})$$

Where f is a loss function and R is a monotonic Regularization function.

2.1 Transformation

- In a finite case if we sample d dimension vector N times such that $d \ll N$. It is with high probability that there will be d vectors that are linearly independent of each other. But as the value of d itself tends to infinity it becomes less and less likely
- In the given case the vectors $\{\phi(\mathbf{x}_i)\}_i^n$ are of infinite dimension. So if the optimum solution is \mathbf{w}^* then there exists $\{\alpha_i\}_i^n$ such that \mathbf{w}^* can be represented as

$$\mathbf{w}^* = \sum_{i=0}^n \alpha_i * \phi(\mathbf{x}_i) + \mathbf{v}$$

where $\mathbf{v}^T \phi(\mathbf{x}_i) = 0$ for all i , that is, \mathbf{v} is orthogonal to the span of the vectors mapped by the function ϕ .

- Let us define $\mathbf{w} = \mathbf{w}^* - \mathbf{v}$

$$\|\mathbf{w}\|^2 = \|\mathbf{w}^* - \mathbf{v}\|^2$$

$$\|\mathbf{w}\|^2 = \|\mathbf{w}^*\|^2 - 2\mathbf{w}^* \cdot \mathbf{v} + \|\mathbf{v}\|^2 \quad \text{As the value } \mathbf{w}^* \cdot \mathbf{v} = \mathbf{v} \cdot \mathbf{v} = \|\mathbf{v}\|^2$$

$$\|\mathbf{w}\|^2 = \|\mathbf{w}^*\|^2 - 2\|\mathbf{v}\|^2 + \|\mathbf{v}\|^2$$

$$\|\mathbf{w}\|^2 = \|\mathbf{w}^*\|^2 - \|\mathbf{v}\|^2$$

- So as norm is a positive function we have $\|\mathbf{w}\| \leq \|\mathbf{w}^*\|$
- And since R is non-decreasing, we obtain $R(\mathbf{w}) \leq R(\mathbf{w}^*)$
(R is monotonic regularization function as defined above)

- Also as $\mathbf{w}^T \phi(\mathbf{x}) = (\mathbf{w}^* - \mathbf{v})^T \phi(\mathbf{x}) = \mathbf{w}^{*T} \phi(\mathbf{x})$
(since \mathbf{v} is orthogonal to $\phi(\mathbf{x})$, we have $\mathbf{v}^T \phi(\mathbf{x})=0$)
- And so,

$$f(\{y_i\}_i^n, \{\mathbf{w}^T \phi(\mathbf{x}_i)\}_i^n) = f(\{y_i\}_i^n, \{\mathbf{w}^{*T} \phi(\mathbf{x}_i)\}_i^n)$$

- We have shown that the loss function is same for both \mathbf{w}^* and \mathbf{w} and the regularization function is less for \mathbf{w} than \mathbf{w}^* ,
- Hence the objective function for \mathbf{w} is less than that for \mathbf{w}^* , but as \mathbf{w}^* is the optimum solution, we must have that \mathbf{w} is also an optimum solution .
- Hence we have proved that the value

$$\mathbf{w} = \sum_{i=0}^n \alpha_i * \phi(\mathbf{x}_i)$$

is an optimum solution for the objective function.¹

Theorem 2.1 (Representer Theorem). *Given a mapping from \mathbb{R}^d to $\mathbb{R}^{d'}$, there exists a vector $\alpha \in \mathbb{R}^{d'}$ such that $\mathbf{w} = \sum_{i=1}^{d'} \alpha_i \phi(\mathbf{x}_i)$ is an optimal solution.*

$$\min_{\mathbf{w}} f(\{y_i\}_i^n, \{\mathbf{w}^T \phi(\mathbf{x}_i)\}_i^n) + \lambda R(\mathbf{w})$$

2.2 Kernel

Substituting the fact that $\mathbf{w} = \sum_{i=1}^{d'} \alpha_i \phi(\mathbf{x}_i)$ is an optimum solution into the value $\mathbf{w}^T \phi(\mathbf{x})$

$$\begin{aligned} \mathbf{w}^T \phi(\mathbf{x}) &= \left(\sum_{i=1}^{d'} \alpha_i \phi(\mathbf{x}_i)^T \right) \phi(\mathbf{x}) \\ &= \sum_{i=1}^{d'} \alpha_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}) \quad \text{consider a function } K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \\ &= \sum_{i=1}^{d'} \alpha_i K(\mathbf{x}_i, \mathbf{x}) \end{aligned}$$

¹Also remember that this was the result that we got by applying the Lagrange multiplier on the objective functions

Here K is called a *Kernel Function* and it denotes *Similarity* between the data points \mathbf{x}_i and \mathbf{x}_j .

What does this expression physically mean? We want to make a prediction at a new test point \mathbf{x} , then $\mathbf{w}^T \phi(\mathbf{x})$ gives the **weighted mean** of labels $y(\mathbf{x}_i)$ with weights as similarity

$$S(\mathbf{x}, \mathbf{x}_i) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x})$$

thus giving more weight to similar neighbours. This essentially represents the idea of K-Means clustering.

We can note that this similarity function can be any form, not necessarily dot product, The term “kernels” is used in this context to describe inner product in the feature space

Implications of this result

- This representation takes us beyond SVM
- \mathbf{w} is linear combination of $\phi(\mathbf{x}_i)$ thus $\mathbf{w}^T \phi(\mathbf{x})$ can be easily calculated.

Some Popular Kernels are:

- **SVM** $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
- **Gaussian** $K(\mathbf{x}_i, \mathbf{x}_j) = e^{\frac{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}{\sigma^2}}$
- **k degree polynomial** $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \langle \mathbf{x}_i, \mathbf{x}_j \rangle)^k$

where $\langle ., . \rangle$ denotes the inner product.

Exercise

Que- Show that RBF can be written as inner product

It is easy to see that RBF(or Gaussian) Kernel satisfies the positive semidefinite condition, we can infact state and prove a possible inner product representation

Consider

$$\phi(\mathbf{x}) = \sum_{n=0}^{\infty} \frac{1}{\sqrt{n!}} e^{-\frac{\mathbf{x}^2}{2}} \mathbf{x}^n$$

Observe that

$$\langle \phi(x), \phi(x') \rangle = \sum_{n=0}^{\infty} \left(\frac{1}{\sqrt{n!}} e^{-\frac{x^2}{2}} x^n \right) \left(\frac{1}{\sqrt{n!}} e^{-\frac{x'^2}{2}} x'^n \right)$$

$$\langle \phi(x), \phi(x') \rangle = e^{-\frac{x^2+x'^2}{2}} \sum_{n=0}^{\infty} \left(\frac{(xx')^n}{n!} \right)$$

$$\langle \phi(x), \phi(x') \rangle = e^{-\frac{x^2+x'^2}{2} + xx'}$$

$$\langle \phi(x), \phi(x') \rangle = e^{-\frac{(x-x')^2}{2}}$$

Que- Do the same for polynomial kernel

Left for reader

Hint: Solution can be checked in book Understanding machine learning chapter Kernel Methods

3 SVM Kernel

Now going back to the dual formulation of Hard-SVM:

$$\begin{aligned} & \max_{\alpha \in R^D: \alpha \geq 0} \left(\sum_{i=1}^D \alpha_i - \frac{1}{2} \sum_{i=1}^D \sum_{j=1}^D \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \right) \\ &= \max_{\alpha \in R^D: \alpha \geq 0} \left(\sum_{i=1}^D \alpha_i - \frac{1}{2} \sum_{i=1}^D \sum_{j=1}^D \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \right) \\ &= \max_{\alpha \in R^D: \alpha \geq 0} \left(\sum_{i=1}^D \alpha_i - \frac{1}{2} (\alpha * Y)^T G (Y * \alpha) \right) \end{aligned}$$

* denotes element-wise multiplication

$G = \{K(\mathbf{x}_i, \mathbf{x}_j)\}_{i,j=1}^{D \times D}$ Known as the **Gram matrix**

Hence, we are optimising for α instead of w .

The advantage of working with kernels rather than directly optimizing w in the feature space is that in some situations the dimension of the feature space is extremely large while implementing the kernel function is very simple.

Theorem 3.1 (Mercers's Theorem). $K : [a, b] \times [a, b] \rightarrow \mathbb{R}$ is a kernel if and only if :

1. *Symmetric* : $K(x, y) = K(y, x)$ for all $x, y \in [a, b]$, and
2. *Positive Semi-definite* : $\sum_{j=1}^n \sum_{i=1}^n K(x_i, x_j) c_i c_j \geq 0$ for all finite sequences of points $\mathbf{x}_1, \dots, \mathbf{x}_n$ of $[a, b]$ and all choices of real numbers c_1, \dots, c_n

References

- [1] S. B.-D. Shai Shalev-Shwartz. *Understanding Machine Learning, Chapter 16*. Cambridge University Press, 2014.