

Lecture 12: Kernel Mode

26 / 09/ 2022

Lecturer: Abir De

Scribe: Dyuneesh, Harsha, Sanchit, Yasaswi

In this Lecture we discuss the smooth transition from Linear to Non-Linear space of functions. Since we are searching for non-linear functions, $f(x)$ is non-linear, that is, $f(x) \neq w^T x$. We discuss treatments where $f(x)$ is non-linear in x but we can construct it using our previous knowledge. We also look at tricks to reduce the order of computation required for the models.

1 Kernel Methods

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

The Midsem Question which stated that the points labeled '+1' are in the range $[-1, 1]$ and the points labeled '-1' are in the range $[-6, -2] \cup [2, 6]$. This is not linearly-separable but on converting to a higher dimension: $\phi(x) = (x, x^2)$, it converts to a form that we are more familiar with.

The Basic Algorithm will be:-

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

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

As d' can reach infinity, calculating and storing w and $\phi(x)$ becomes practically impossible, but the dot product $w^T \phi(x)$ 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 halfspace 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 that is $x \rightarrow (1, x, x^2, \dots)$

The Setup:

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

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

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

2 Solution to Setup

- 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 $\{\phi(x_i)\}_i^n$ are of infinite dimension. So if the optimum solution is w^* then there exists $\{\alpha_i\}_i^n$ such that w^* can be represented as

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

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

- Let us define $w = w^* - v$

$$\begin{aligned} ||w||^2 &= ||w^* - v||^2 \\ ||w||^2 &= ||w^*||^2 - 2w^*.v + ||v||^2 && \text{As the value } w^*.v = v.v = ||v||^2 \\ ||w||^2 &= ||w^*||^2 - 2||v||^2 + ||v||^2 \\ ||w||^2 &= ||w^*||^2 - ||v||^2 \end{aligned}$$

- So as norm is a positive function we have $||w|| \leq ||w^*||$
- And since R is non-decreasing, we obtain $R(w) \leq R(w^*)$
- Also as $w^T \phi(x) = (w^* - v)^T \phi(x) = w^{*T} \phi(x)$
- And so,

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

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

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

is an optimum solution for the objective function.¹

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

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 $w = \sum_{i=1}^{d'} \alpha_i \phi(x_i)$ is an optimal solution.*

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

3 Reducing Computation

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

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

- Similarly we have that,

$$\begin{aligned} ||w||^2 &= w^T w \\ &= \left(\sum_{i=1}^{d'} \alpha_i \phi(x_i)^T \right) \left(\sum_{i=1}^{d'} \alpha_i \phi(x_i) \right) \\ &= \sum_{i=1}^{d'} \sum_{j=1}^{d'} \alpha_i \alpha_j \phi(x_i)^T \phi(x_j) \\ &= \sum_{i=1}^{d'} \sum_{j=1}^{d'} \alpha_i \alpha_j K(x_i, x_j) \end{aligned}$$

- K is called a *Kernel Function* and it denotes *Similarity* between the data points x_i and x_j .
- Some Popular Kernels are:

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

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

4 Kernel Function

To show that a function is indeed a Kernel Function, We need to show that there exists a mapping ϕ to a higher dimension space such that

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

To show that K is a Kernel function, we first need it to be Positive Semi-Definite. That is,

for all real-valued c_i, c_j and for all $\mathbf{x}_i, \mathbf{x}_j$ in the domain of ϕ ,

$$\sum_{i=1}^D \sum_{j=1}^D c_i c_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) \geq 0$$

Since we have shown that K is Positive Semi-Definite, K implements an inner product in some Hilbert space and hence is a valid Kernel function.

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

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

* denotes element-wise multiplication

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

Hence, we are optimising for α instead of \mathbf{w} .

The advantage of working with kernels rather than directly optimizing \mathbf{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.

References

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