

Lecture 17 - 12-05-2020

1.1 Kernel functions

We use a notion of **feature expansion**. They are different but somehow they reach something similar. In fact Linear classifier have high bias. Linear predictor use hyper plane as basic brick to build prediction.

1.1.1 Feature expansion

Given $\phi : \mathbb{R}^d \rightarrow V$ V is typically \mathbb{R}^N $N \gg d$

For example:

$$d = 2 \quad N = 6 \quad \phi = \mathbb{R}^2 \rightarrow \mathbb{R}^6$$

$$\phi(x_1, x_2) = (1, x_1^2, x_2^2, x_1, x_2, x_1x_2)$$

We have a homogenous hyper plane.

$$w \in \mathbb{R}^6 \quad \{z \in \mathbb{R}^6 : w^T z = 0\} \quad z = \phi(x) \quad x \in \mathbb{R}^2$$

$$\forall x \in \mathbb{R}^2 \quad w^T \phi(x) = w_1 + w_2 x_1^2 + w_3 x_2^2 + w_4 x_1 + w_5 x_2 + w_6 x_1 x_2 = 0$$

$$w^T \phi(x) = 0$$

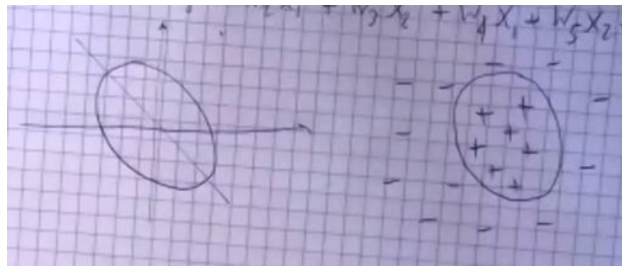


Figure 1.1:

$$\phi : \mathbb{R}^d \rightarrow \mathbb{R}^N \quad \prod_{s=1}^M x_{v_s} \quad v \in \{1, \dots, d\}^k \quad k = 0, \dots, n$$

$$h(x) = \text{sgn}(w^T \phi(x)) \quad w^T \phi(x) = \sum_{i=1}^N w_i \phi(x)_i$$

The problem of this feature expansion is the degree of the monomials!

$$N = \sum_{i=0}^n |\{1, \dots, d\}^k| = \sum_{k=0}^n d^k = \frac{d^{n+1} - 1}{d - 1} = \Theta(d^n)$$

So it's exponential! But this feature expansion can be implemented in a efficient way.

1.1.2 Kernels implements feature expansion (Efficiently)

$w^T \phi(x)$ Perception $w \leftrightarrow w + y_t x_t I\{y_t w^T x_t \leq 0\}$ MANCA quadcosa

$$w = \sum_{s \in S} y_s x_s \rightsquigarrow \sum_{s \in S} y_s \phi(x_s)$$

where S is a subset of training set where updates occurred.

Every time i make a mistake i add some of this product of data points.

If I run this using example that are images according to some feature expansion map (ϕ), I will get the perceptron after the mapping.

$$w^T \phi(x) = \sum_{s \in S} y_s \phi(x)^T \phi(x_s)$$

It's a inner problem and can have exponentially degree of the component.

Kernels help me compute this inner product $\phi(x)^T \phi(x_s)$ quickly

$$\phi : \mathbb{R}^2 \longrightarrow \mathbb{R}^6 \quad \phi(x_1, x_2) = (1, x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2)$$

$$\phi(x)^T \phi(z) = 1 + x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_2 z_2 + 2x_1 x_2 z_1 z_2 = (1 + x^T z)^2 = k(x, z)$$

$$w^T \phi(x) = \sum_{s \in S} y_s k(x, x_s)$$

$k(x, z)$ implements $\phi(x)^T \phi(z) \quad \forall x, z$ and ϕ defined as before

How to we generalise this?

$$k_n(x, x') = (1 + x^T x')^n$$

This is called polynomial kernel.

I want to check now what is the ϕ for K_n ?

I want to compute ϕ s.t. $\phi(x)^T \phi(x') = k_n(x, x') = (1 + x^T x')^n$

We can use Newtons binomial theorem:

$$(1 + x^T x')^n = \sum_{k=0}^n \binom{n}{k} (x^T x')^k$$

$$(x^T x')^k = \left(\sum_{i=1}^d x_i x'_i \right)^k = \sum_{v \in \{1, \dots, d\}^k} \left(\prod_{s=1}^k x_{V_s} x'_{V_s} \right)$$

$$\phi(x) = \left(\sqrt{\binom{n}{k}} \prod_{s=1}^k x_{V_s} \right) \quad k = 0, \dots, n \quad v \in \{1, \dots, d\}^k$$

When I am using polynomial kernel I am implicitly using the feature expansion ...

Can an algorithm work using kernel?

Perceptron works!

$S = 0$

For $t = 1, 2, \dots$

1) Get (x_t, y_t)

2) $\hat{y}_t = \text{sgn} \left(\sum_{s \in S} y_s K(x, x_s) \right)$

3) If $\hat{y}_t \neq y_t$ $S \leftarrow S \cup \{t\}$ $w \leftarrow w + y_t \phi(x_t)$

So I am representing y as a sum and not as a vector. In fact, $w = \sum_{s \in S} y_s \phi(x_s)$

ù

1.2 Gaussian Kernel

$$\gamma > 0 \quad k_\gamma(x, x') = \exp \left(-\frac{1}{2\gamma} \|x - x'\|^2 \right)$$

$$e^{-\frac{1}{2\gamma} (x-x')^2}$$

I can controll the distribution changing the value of γ

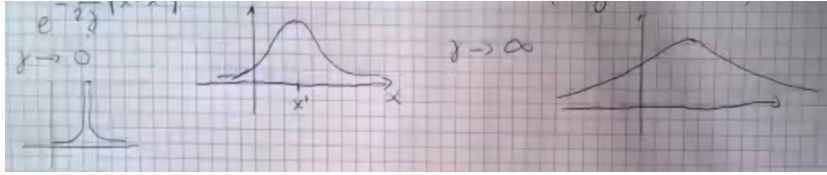


Figure 1.2:

$$\hat{y}_t = \text{sgn} \left(\sum_{s \in S} y_s K_\gamma(x, x_s) \right)$$

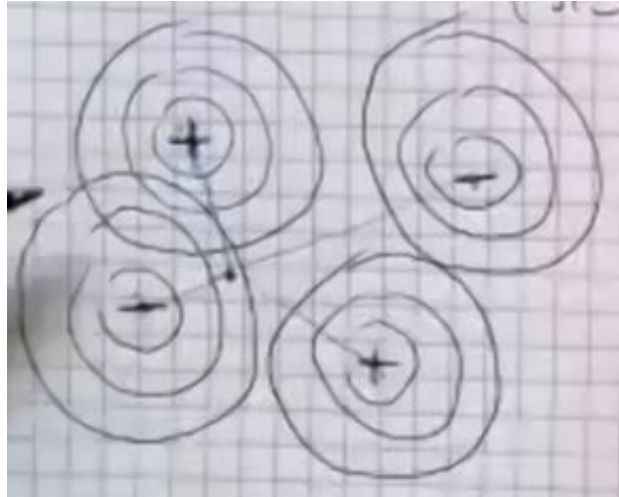


Figure 1.3:

Negative or positive gaussian component looking at the distance.

Now I want to compute: $\phi_\gamma : \mathbb{R}^n \rightarrow V$

$$\exp \left(-\frac{1}{2\gamma} \|x - x'\|^2 \right) = \exp \left(-\frac{1}{2\gamma} (\|x\|^2 + \|x'\|^2) \right) \cdot \exp \left(\frac{1}{\gamma} x^T x' \right) =$$

where $e = x + \frac{x^2}{12} \dots$

$$= \exp \left(-\frac{1}{2\gamma} \|x\|^2 \right) \cdot \exp \left(-\frac{1}{2\gamma} \|x'\|^2 \right) \cdot \sum_{n=0}^{\infty} \frac{1}{n!} \frac{(x^T x')^n}{\gamma^n}$$

Gaussian Kernel is a linear combination of infinitely many poly kernels.
 The higher I go the small is $n!$. Gaussian kernel mapping into a space that is very large. So large that it has infinitely many dimension. Why? Because each polynomial kernel maps to infinitely dimensions.

ϕ_γ maps \mathbb{R}^d into a space of infinitely many dimensions.

$$\phi_\gamma : \mathbb{R}^d \rightarrow V \quad k_\gamma(x, x') = \phi_\gamma(x)^T \phi_\gamma(x')$$

It maps to infinitely many dimension, so it maps to a function!

$\phi_\gamma(x)$ is a function.

In general, when I learn a linear predictor using k_γ

I learn $\sum_s \alpha_s k(x_s, \cdot) = f$
 $w^T \phi(x)$

$$H_\gamma \equiv \left\{ \sum_{i=1}^N \alpha_i k(x_i, \cdot) : x_1, \dots, x_N \in \mathbb{R}^d, \alpha_1, \dots, \alpha_N \in \mathbb{R}, N \in \mathbb{N} \right\}$$

Theorem

$\forall \gamma > 0 \forall f : \mathbb{R}^d \rightarrow \mathbb{R}$ continuous, $\forall \varepsilon > 0$

$\exists g \in H_\gamma$ that approximates f with error ε

We define a function with H . We see the \cdot and this tell us is a function. So we can evaluate every kind of x point in \cdot position.

We are able to get a super parametric algorithm and transform it in a non-parametric algorithm. Parametric algorithms is defined by an arbitrary number of parameter we cannot adapt it for every case.

Gaussian Kernels enable consistency by using feature expansion with infinitely many components.