

## The kernel trick

We have seen the dual form of ~~the~~ SVM:

$$\max \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \begin{bmatrix} x_i^T x_j \\ 1 \end{bmatrix}$$

~~M2~~ result of  $\alpha_i$  &  $\alpha_j$

st.t.  $\sum_{i=1}^n \alpha_i y_i = 0, \alpha_i \geq 0$

$\begin{bmatrix} x_i^T x_j \\ 1 \end{bmatrix} \rightarrow$  is the function of  $x_i$  and  $x_j$   
we can use any similarity function  
of  $x_i$  and  $x_j$

$\rightarrow$  This is known as **Kernel Function**

$\rightarrow$  It is often replaced by  $K(x_i, x_j)$

Then we can write the given function as:

$$f(x_r) = \sum_{i=1}^n \alpha_i y_i K(x_i, x_r) + b$$

So, the replacement of  $x_i^T x_j$  with  $K(x_i, x_j)$

is called as **The Kernel Trick**

The most important idea of in SVM is the kernel trick.

If you don't apply the kernel trick, then it's called as Linear SVM.

→ There are many types of kernel functions available, and we are going to see few of them.

→ This is about the kernel trick, but the main reason why we do this?

We will see in a second.

Ok. so far what we have done is.

Linear SVM:  $x_i^T x_j$

Kernel SVM:  $K(x_i, x_j)$

In linear SVM, we are trying to form a hyper-plane, in space of  $x_i$ 's.

In logistic Regression, we are trying to find hyperplane that minimizes the logistic loss, in the space of  $x_i$ 's.

Most of the time, when you apply Logistic Reg. and SVM the result looks very-very similar. (not always)

So, margin maximizing is a nice idea, but not the world changing idea.

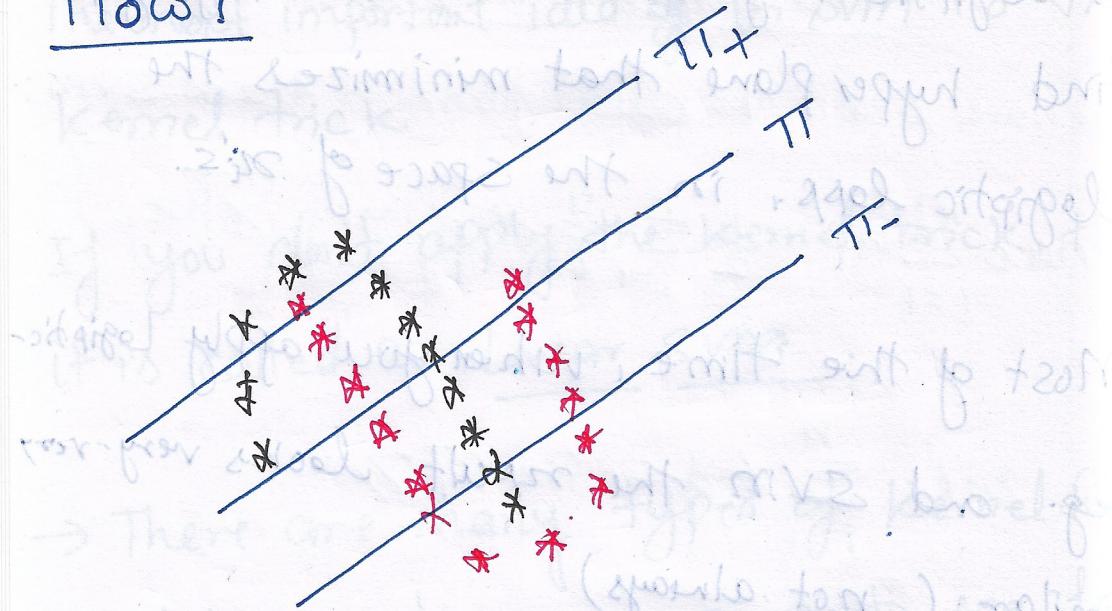
The world changing or the key idea of SVM is the kernel trick or kernelization.

Actually kernelization makes SVM to handle non-linearly separable datasets.  
(how?) (we will see.)

You can think of Kernelization kind of feature transformed Logistic Reg. (but they are different)

There is also Kernel-logistic Regression in Research

How?



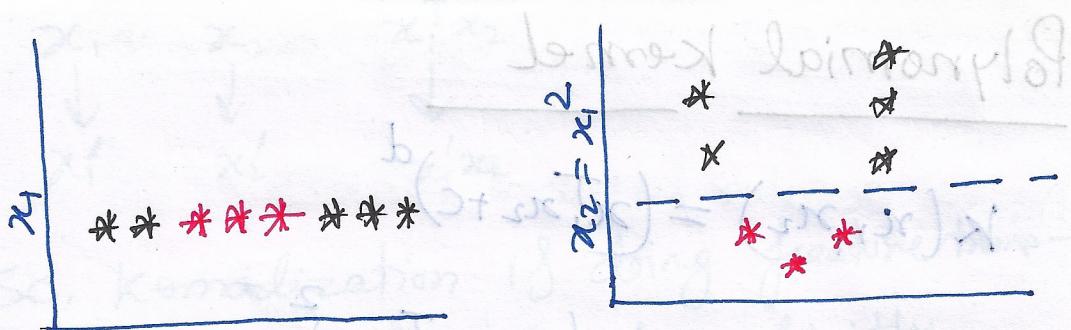
Most of the time, the real-life data sets are not nice, linearly separable.

Say I'm having moon-shaped data.

How to get the hyper plane?

Mapping to higher dimensions

→ The most fundamental principle of kernel trick is that, "In a high dimensional space, there is a bigger chance the data will become linearly separable".



From the above pair plot we can see that, as the dimension increases, we can separate the data.

When we've only one feature,  $x_1$  we can not separate the data with single line.

Adding another feature  $x_2 = x_1^2$ , makes it easy to separate the two classes.

So, Kernelization is basically a smart way of adding more features to the data in the hope of making linearly separable.

Some popular kernels are

- Polynomial ker.
- Gaussian Radial Basis Function RBF kernel

$\text{Cost}(w) = \frac{1}{2} w^T w$

as we saw below, following

## Polynomial kernel

$$k(x_1, x_2) = (x_1^T x_2 + c)^d$$

e.g.  $k(x_1, x_2) = (1 + x_1^T x_2)^2$

$c=1, d=2 \rightarrow$  Quadratic kernel

$$\begin{aligned} k(x_1, x_2) &= (1 + x_1^T x_2)^2 & x_1 &= \langle x_{11}, x_{12} \rangle \\ &= (1 + x_{11} x_{21} + x_{12} x_{22})^2 \end{aligned}$$

$$\begin{aligned} \Rightarrow k(x_1, x_2) &= 1 + x_{11}^2 x_{21}^2 + x_{12}^2 x_{22}^2 + 2x_{11} x_{21} \\ &\quad + 2x_{11} x_{22} + 2x_{12} x_{21} + 2x_{12} x_{22} \\ &= [1, x_{11}^2, x_{12}^2, \sqrt{2}x_{11}, \sqrt{2}x_{12}, \\ &\quad \sqrt{2}x_{11}x_{21}, \sqrt{2}x_{12}x_{21}, \sqrt{2}x_{11}x_{22}, \sqrt{2}x_{12}x_{22}] \end{aligned}$$

$$\therefore k(x_1, x_2) = (x'_1)^T (x'_2)$$

Basically, what we have done is:

$$\begin{matrix} x_1 & x_2 & x_1^T x_2 \\ \downarrow & \downarrow & \downarrow \\ x'_1 & x'_2 & x'_1 x'_2 \end{matrix}$$

So, kernelization is doing feature transformation internally and implicitly.

Kernalization  $d \xrightarrow[\text{internally}]{FT} d'$   $d' > d$   
 implicitly

But the challenge here is select the right kernel.

