

# Machine learning


---

Lect 10

Dr. Luana  
prokhry'ce

---

---



# # Non Linear SVM and Kernel functions

$$\phi: \mathcal{X} \rightarrow \phi(\mathcal{X})$$

original feature space can be mapped into new feature space where it is linearly separable  $\rightarrow$  higher dimensional

But what about computational cost??

Linear SVM:  $\sum L_i L_j y_i y_j x_i \cdot x_j$

$\nearrow d$ -dimensional  $\nearrow d$ -dimensional

we need to find pairwise dot products

if we have  $m$  training examples we need to do  $m^2$  computations like this

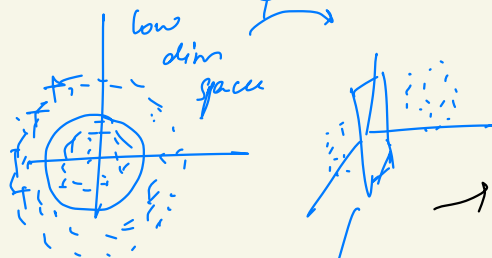
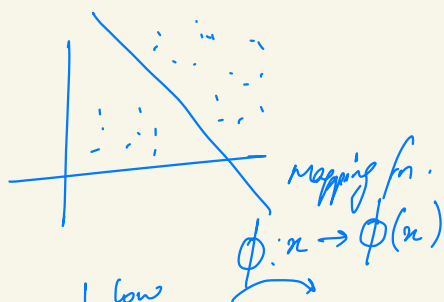
Computation time  $O(d^2)$

$\Downarrow$  Transform to high dimensional space  $D \gg d$

$$\sum L_i L_j y_i y_j \phi(x_i) \cdot \phi(x_j)$$

$O(D^2)$

$\Downarrow$  more computational cost



SVM inputs

- Gaussian RBF
- Linear kernel
- Polynomial kernel

$k(x_i, x_j)$

$\phi(x_i) = x_i$   
 $\phi(x_j) = x_j$

$\frac{1}{\sigma^2} e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$

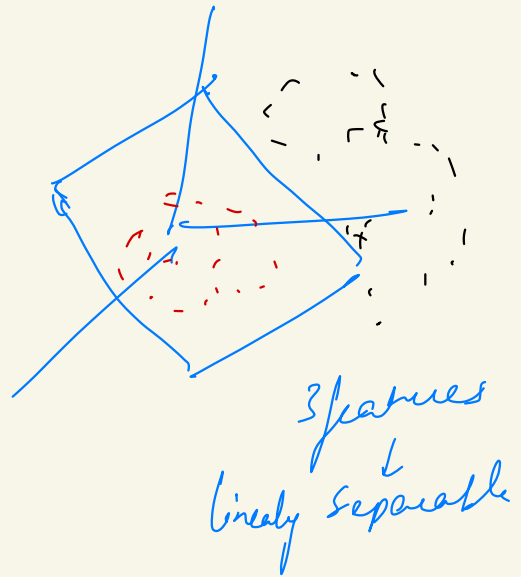
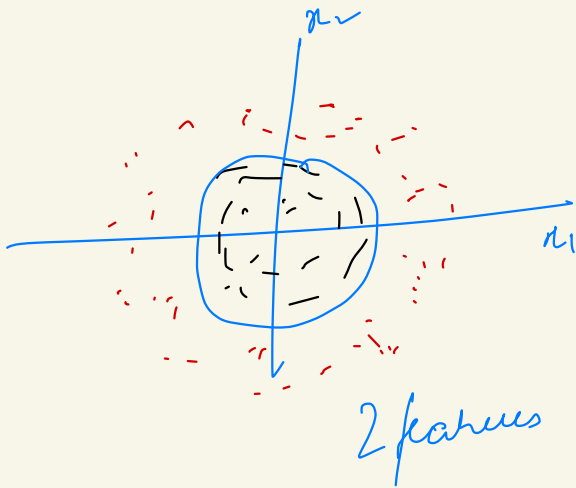
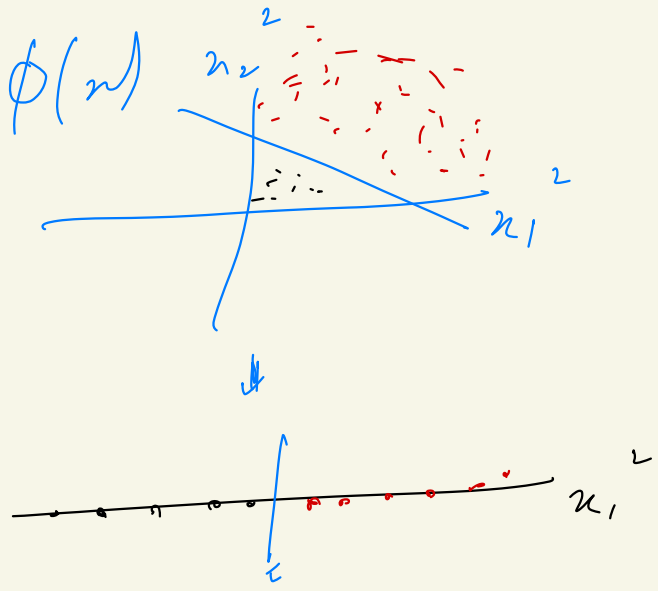
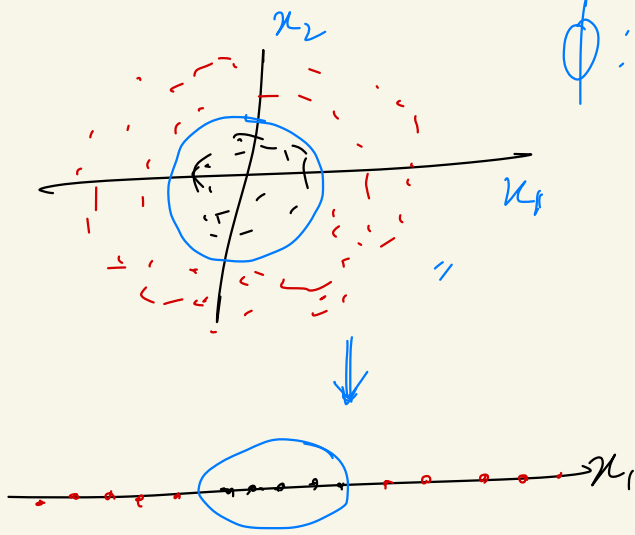
$(1 + x_i x_j)^{p-2}$

quad

Soln: Use kernel function to reduce computational costs

# # Visualize Feature Space

$$\phi: x \rightarrow \phi(x)$$



In certain feature transformation, SVM can be solved efficiently. where after transformation

kernel function	$x_a$	$x_b$	$x_a \cdot x_b$
	$\phi(x_a)$	$\phi(x_b)$	$\phi(x_a) \cdot \phi(x_b)$
			$[ \quad ]$

$$k(x_a, x_b) = \phi(x_a) \cdot \phi(x_b)$$

also called kernel trick  $\hookrightarrow$  may be easy to compute it as is a fn. of  $x_a$  &  $x_b$  without expanding  $\phi(x_a)$  or  $\phi(x_b)$

What is

Kernel: i) Original input attributes is mapped to a new set of input features via feature mapping  $\phi$

ii) Since the algo can be written in terms of the scalar product, we replace  $x_a \cdot x_b$  with  $\phi(x_a) \cdot \phi(x_b)$

iii) For certain  $\phi$ 's there is a simple operation on 2 vectors in the low-dim space that can be used to compute the scalar product of their 2 images

in the high dim space

dot product  
scalar or vector?

$$k(x_a, x_b) = \phi(x_a) \cdot \phi(x_b)$$

let the kernel do the work  
rather than do the scalar  
product in high dim space

# Non linear SVM: The kernel Trick  
→ With this mapping, our discriminant function is now:

$$g(x) = w^T \phi(x) + b$$

$$= \sum_{i \in SV} \alpha_i$$

$$\alpha_i \left[ \phi(x_i)^T \phi(x) \right] + b$$

$$\begin{aligned} &w^T x + b \\ &w^T \phi(x) + b \end{aligned}$$

Individual form  
 → We only use the dot product of feature vectors in both the training & test

→ A kernel function is defined as a function that corresponds to a dot product of two feature vectors in some expanded feature space:

$$k(x_a, x_b) = \phi(x_a) \cdot \phi(x_b)$$

This  $k(x_a, x_b)$  may be very inexpensive to compute even if  $\phi(x_a)$  may be extremely high dimensional

Kernel Example

2 dimensional vectors  $\vec{x} = [x_1, x_2]$

$$\text{let } k(x_i, x_j) = (1 + x_i \cdot x_j)^2$$

$$k(x_i, x_j) = (1 + x_i \cdot x_j)^2$$

$$= 1 + x_{i1}^2 x_{j1}^2 + 2 x_{i1} x_{j1} x_{i2} x_{j2} + x_{i2}^2 x_{j2}^2$$

$$= \begin{bmatrix} 1 & x_{i1}^2 & \sqrt{2} x_{i1} x_{i2} & x_{i2}^2 \\ 1 & x_{j1}^2 & \sqrt{2} x_{j1} x_{j2} & x_{j2}^2 \end{bmatrix} \cdot \begin{bmatrix} x_{i1} & x_{i2} \\ x_{j1} & x_{j2} \end{bmatrix}$$

$$= \phi(n_i) \cdot \phi(n_j)$$

where  $\phi(n) = \begin{bmatrix} 1 & n_1^2 & 2n_1 n_2 & n_2^2 & \sqrt{2} n_1 & \sqrt{2} n_2 \end{bmatrix}$

$$\left[ 1 + n_{i1} n_{j1} + n_{i2} n_{j2} \right]$$

$(1 + n_i \cdot n_j)^2$   
 $\begin{bmatrix} n_{i1} & n_{i2} \\ n_{j1} & n_{j2} \end{bmatrix}$

$$(a + b + c)^2 = a^2 + b^2 + c^2 + 2ab + 2bc + 2ac$$

$$= (1 + n_{i1}^2 n_{j1}^2 + n_{i2}^2 n_{j2}^2 + 2n_{i1} n_{j1} + 2n_{i2} n_{j2} + 2n_{i1} n_{j2} + 2n_{i2} n_{j1})$$

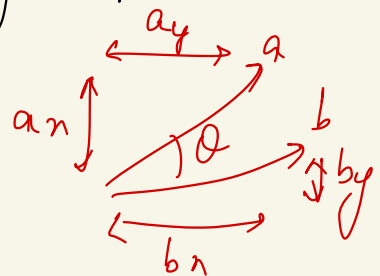
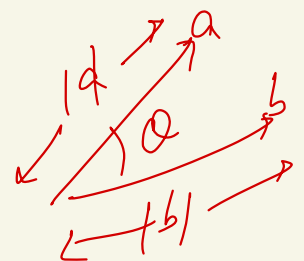
Dot product

$$a \cdot b = |a| |b| \cos \theta$$

or

$$a \cdot b = a_x b_x + a_y b_y$$

multiply the x-components  
then y-components



$$n_i = \begin{bmatrix} n_{i1} & n_{i2} \end{bmatrix}$$

$$n_j = \begin{bmatrix} n_{j1} & n_{j2} \end{bmatrix}$$

$$n_i \cdot n_j = n_{i1} n_{j1} + n_{i2} n_{j2}$$

$$(1 + n_i \cdot n_j)^2 = (1 + n_{i1} n_{j1} + n_{i2} n_{j2})^2 \Rightarrow (a + b + c)^2$$

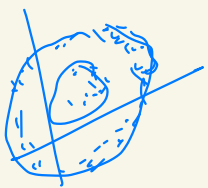
# Commonly used kernel functions

→ Linear SVM corresponds to linear kernel:  
 $k(x_i, x_j) = x_i \cdot x_j$  where  $\phi(x_i) = x_i$  &  $\phi(x_j) = x_j$

→ Polynomial of degree  $p$ :  
 $k(x_i, x_j) = (1 + x_i \cdot x_j)^p$  identity mapping

→ Gaussian (radial basis function):  
 $k(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$

→ Sigmoid:  
 $k(x_i, x_j) = \tanh(\beta_0 x_i \cdot x_j + \beta_1)$



In general, functions that satisfy Mercer's condition can be kernel functions

Eg: Text classification we can use kernels based on similarity of words of text

# Kernel Functions

→ Kernel function can be thought of as a similarity measure b/w the input objects

→ Not all similarity measure can be used as kernel function

→ Mercer's condition states that any positive semi-definite kernel  $k(x, y)$  is

$$\sum_i \sum_j k(x_i, x_j) c_i c_j \geq 0 \quad \text{any real no.}$$

positive semi-definite matrix

Generally this similarity measure is symmetric

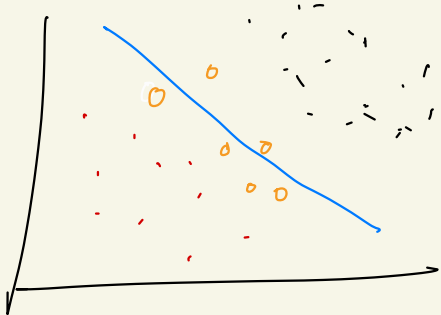
Any matrix  
is symmetric  
& all the  
eigen values  
are +ve

$$k(x_i, x_j) = k(x_j, x_i)$$

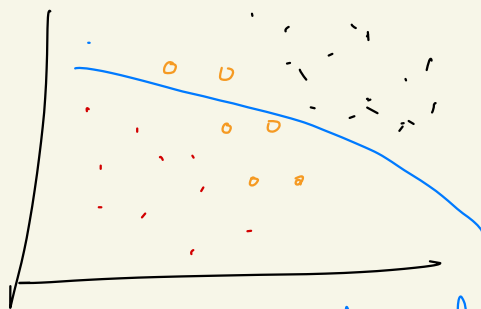
If you can write similarity b/w points  
as a matrix & this matrix is  
symmetric & semi definite then  
a kernel function will exist

→ Such functions can be expressed as a  
dot product in high dimensional space

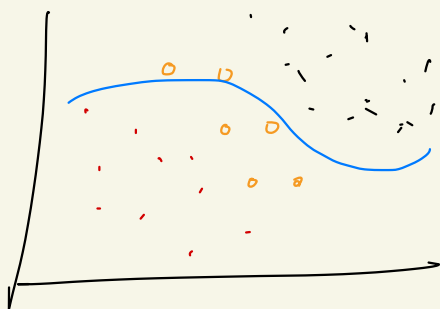
## SVM examples



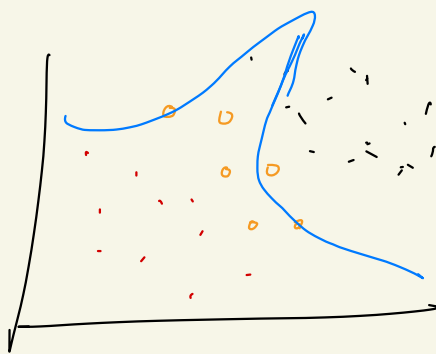
↓  
linear kernel  
with noise



↓ 2<sup>nd</sup> order  
polynomial



↓  
4<sup>th</sup> order  
polynomial



↓  
8<sup>th</sup> order  
polynomial



# Non linear SVM optimization formulation (Lagrangian Dual Problem)

$$\text{minimize} \quad \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j k(x_i, x_j)$$

$$\text{s.t.} \quad 0 \leq \alpha_i \leq C \rightarrow \text{hyperparameter}$$

$$C=0.001$$

KKT condition

$$\rightarrow \sum_{i=1}^n \alpha_i y_i = 0$$

The solution of the discriminant function is

$$g(n) = \sum_{i \in SV} \alpha_i k(x_i, n) + b$$

$\downarrow$   
 $\phi(x_i)^T \cdot \phi(n)$

## # Performance

→ SVM work very well in practice  
User must choose the kernel function & its parameters  $(C)$  appropriately

→ They can be expensive in time & space for big datasets

→ Computation of the maximum margin hyperplane depends on the square of number of training cases

→ We need to store all support vectors

→ The kernel trick can also be used to do PCA in a much higher-dimensional

9/21/2019

space, thus giving a non linear version of PCA  
in the original space

# multiclass classification

SVMs can only handle 2-class outputs

Learn N-SVMs

- SVM 1 learns Class 1 vs Rest  
- SVM 2 learns Class 2 vs Rest

SVM N learns Class N vs Rest

Then to predict off for a new input, just  
predict with each SVM & find out wh.  
which one put the prediction the furthest  
into the positive region

majority voting  
↓

Class 1    Class 2    Class 3 ... Class N  
                    OVO

— SVM 1    Class 1 vs Rest  
                    Class 2 vs Rest

OR  
one vs all

Class 1  
↑

negative

sk-learn

svm.predict  
svm.classify

python  
matlab

⇒ lib-svm library