

Recap:Primal:

$$\min \frac{1}{2} \|w\|^2$$

st.  $y_i(w^T x_i + b) \geq 1$   
 $\forall i=1, \dots, n$

Dual: (Easier problem to solve for  $d \gg n$ )

$$\max_{\lambda \geq 0} \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum \sum \lambda_i \lambda_j y_i y_j x_i^T x_j$$

st.  $\sum_{i=1}^n \lambda_i y_i = 0$

 $x_i \in \mathbb{R}^d$   
 $d \gg n$ 

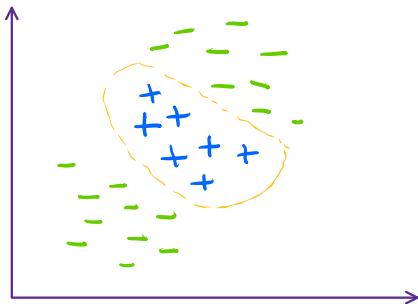
$$W^* = \sum_{i=1}^n \lambda_i^* y_i x_i$$

Why dual?

- less costly to solve
- Kernelization.

Recall that for a classification problem as below,

we need to use a basis  $\phi$  in,  $x_i \mapsto \phi(x_i)$ . But  $x_i$  are already in higher dimension, and  $\phi(x_i)$  is even higher dimension.



So to computationally manage this, we use Kernelization.

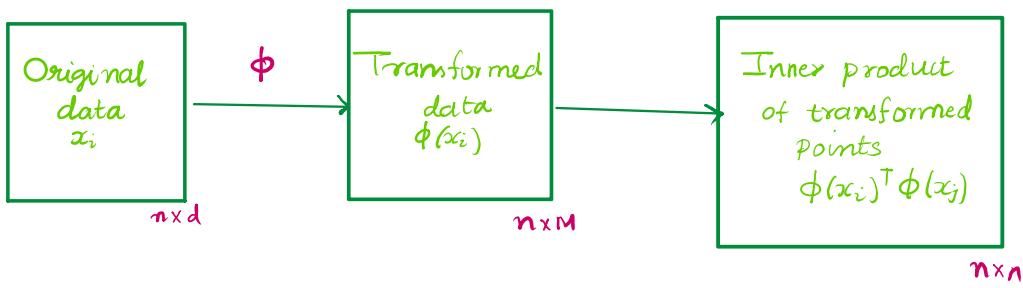
Given: linearly inseparable data

Goal: Project the data to a higher dimensional space  
 $\rightarrow$  solve using SVM  $\rightarrow$  find  $W, b$  in new dimension

Eg: Step 1:

$$x_i = \begin{pmatrix} x_i^{(1)} \\ x_i^{(2)} \end{pmatrix} \xrightarrow{\phi} \begin{pmatrix} 1 \\ x_i^{(1)} \\ x_i^{(2)} \\ x_i^{(1)} x_i^{(2)} \\ x_i^{(1)2} \\ x_i^{(2)2} \end{pmatrix} = \phi(x_i)$$

Step 2: Dual SVM: Objective fn has an inner product of the data in higher dimension.



Q: Do there exist  $f^C$ 's that calculate the inner product in the transformed space without explicitly computing the transformations?

A: Yes, such  $f^n$ 's are Kernel functions.

From the previous example, the terms involved =  $\{1, x_i^{(1)}x_j^{(1)}, x_i^{(2)}x_j^{(2)}, x_i^{(1)}x_i^{(2)}x_j^{(1)}x_j^{(2)}, x_i^{(1)}x_i^{(2)}x_j^{(1)}x_j^{(2)}, x_i^{(1)2}x_j^{(1)2}, x_i^{(2)2}x_j^{(2)2}\}$

Now consider the expression  $(1 + x_i^T x_j)^2$ . This has exactly same terms as before  
↳ KERNEL FUNCTION

Kernel trick: transformations are equivalent as long as we are calculating the dual of SVM.

### Different Kernels :

i. Linear :  $K(x, z) = x^T z$

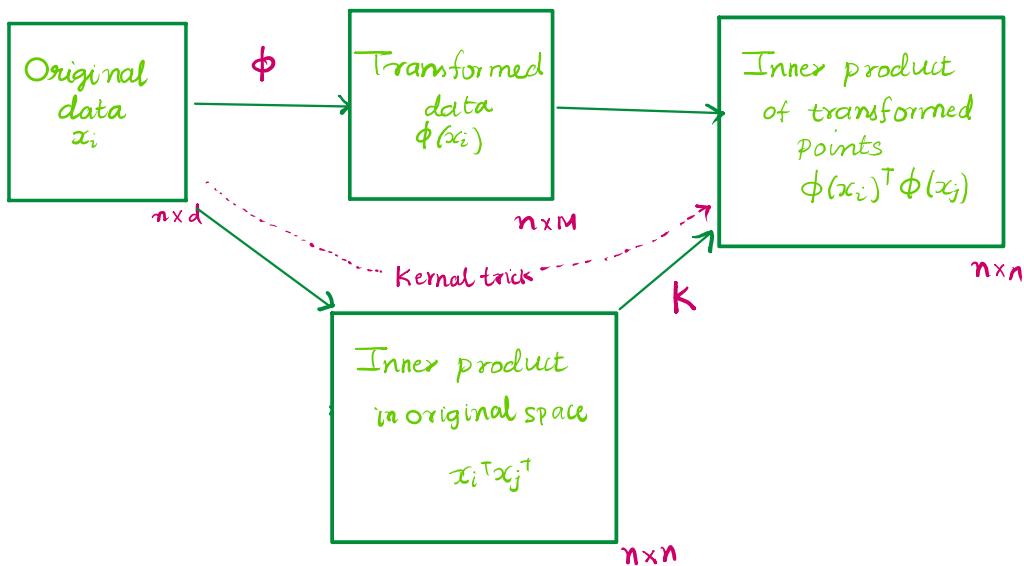
ii) Polynomial :  $K(x, z) = (1 + x^T z)^n$

iii) Gaussian Kernel :  $K(x, z) = e^{-\|x-z\|^2/2\sigma^2}$

iv) Laplace / Radial :  $K(x, z) = e^{-\|x-z\|/2\sigma}$

Uses of SVMs: Handwriting recognition, protein structure, medical image classification

A set of necessary and sufficient conditions govern the kernel functions.



The choice of Kernel to be used can be done using grid search.

### Limitations of SVM

- i) Binary Classification : One vs rest classifier
- ii) No probabilistic interpretation
- iii) Does not work well when data is noisy.

What we had done so far is **SUPERVISED LEARNING** where  $D = \{(x_i, y_i)\}_{i=1, \dots, n}$

We will now deal with **UNSUPERVISED LEARNING** where  $D = \{x_i\}_{i=1, \dots, n}$

### CLUSTERING: Unsupervised Learning

$$D = \{x_1, \dots, x_n\}, x_i \in \mathbb{R}^d$$

Goal: find a well-separated partition of data,  $D = D_1 \cup D_2 \dots \cup D_k$ ,  $D_i \cap D_j = \emptyset$  (hard clustering)

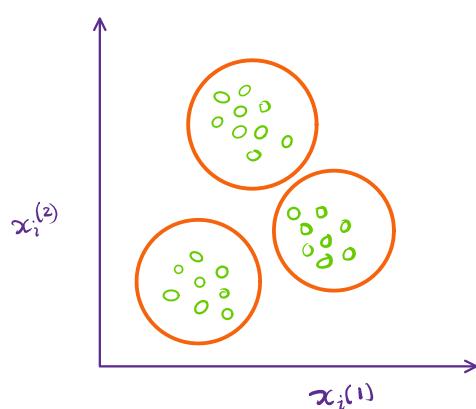
### K-means Clustering

Consider a function  $C$  as below:

$$C(i) = \hat{k}$$

Clustering function  
 $C: [n] \rightarrow [k]$   
 $[n] = \{1, 2, \dots, n\}$

*i<sup>th</sup> data point*      *cluster no.*



Algorithm:

- Input:  $D = \{x_1, \dots, x_n\}$
- Initialize:  $k$  cluster means  $m_1, \dots, m_k$  some arbitrary  $c^0$
- Repeat until convergence (until assignments don't change):
  - for every  $i = 1, 2, \dots, n$ 
    - if  $\exists k' \in C^{t+1}(i)$  such that  $\|m_{k'}^{t+1} - x_i\| < \|m_{C^t(i)}^{t+1} - x_i\|$ :
    - $C^{(t+1)}(i) \leftarrow k'$
    - Update cluster centres by average of its current associated datapoints:  $m_{k'}^{t+1}$

K-means clustering as an optimization problem:

$$\min_{C \in \wp} \sum_{i=1}^n \|x_i - m_{C(i)}\|^2 \rightarrow \text{NP-Hard!}$$

non-convex fn

$\wp$  Set of all possible clusters.

$|\wp| = k^n$

Convergence:

Lemma:  $\arg\min_x \sum_{i=1}^{n_k} \|x_i - x\|^2 = \bar{x} = \frac{1}{n_k} \sum_{i=1}^{n_k} x_i$

Theorem: K-means converges to a local minima

- consider  $t \rightarrow t+1$
- $SE(C^{t+1}, m^{t+1}) < SE(C^t, m^t)$  SE : squared error

Why is this sufficient? no clusters are repeated  
finite number of clustering

Proof:  $(m^t, C^t) \rightarrow (m^t, C^{t+1}) \rightarrow (m^{t+1}, C^{t+1})$

- first we will prove  $SE(C^{t+1}, m^t) < SE(C^t, m^t) \rightarrow$  from algorithm itself
- Now show that  $SE(C^{t+1}, m^{t+1}) \leq SE(C^{t+1}, m^t) \rightarrow$  from Lemma

Combining (i) and (ii)

Combining (i) and (ii)  $SE(c^{t+1}, m^{t+1}) < SE(c^t, m^t)$