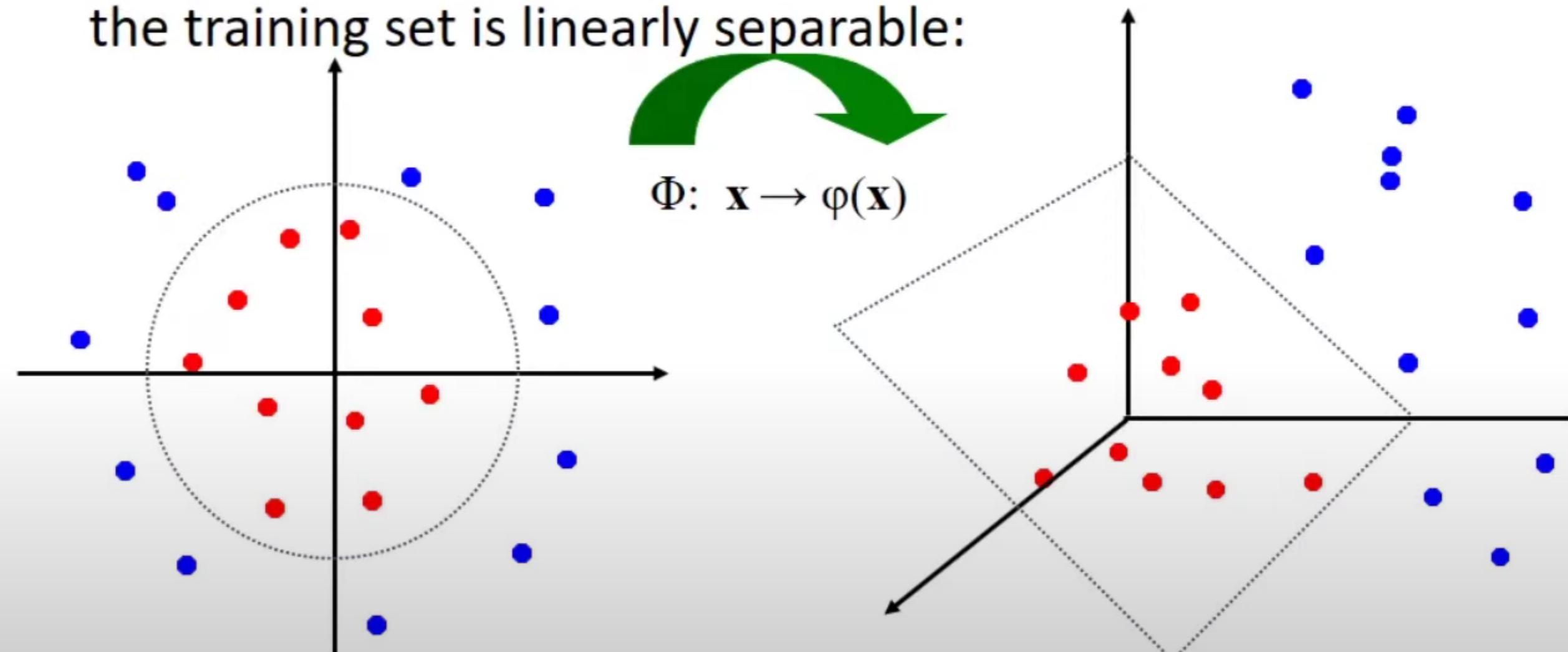




Making SVM capable of solving nonlinear classification task

- Idea: The original nonlinear input feature space can always be mapped to some higher-dimensional feature space where the training set is linearly separable:

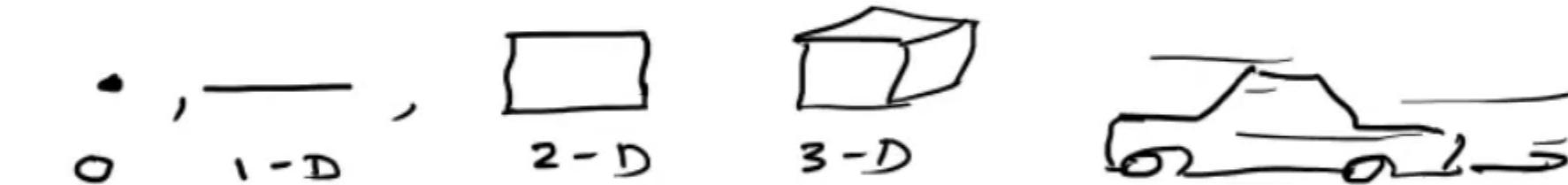
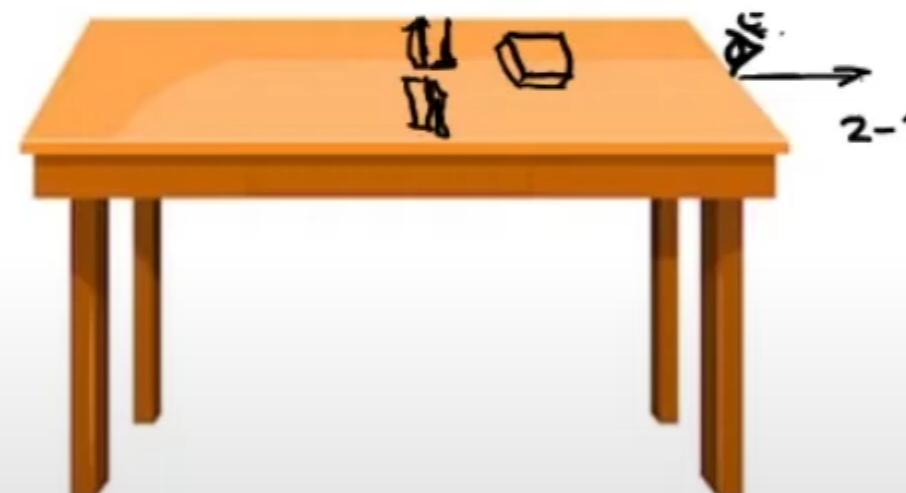
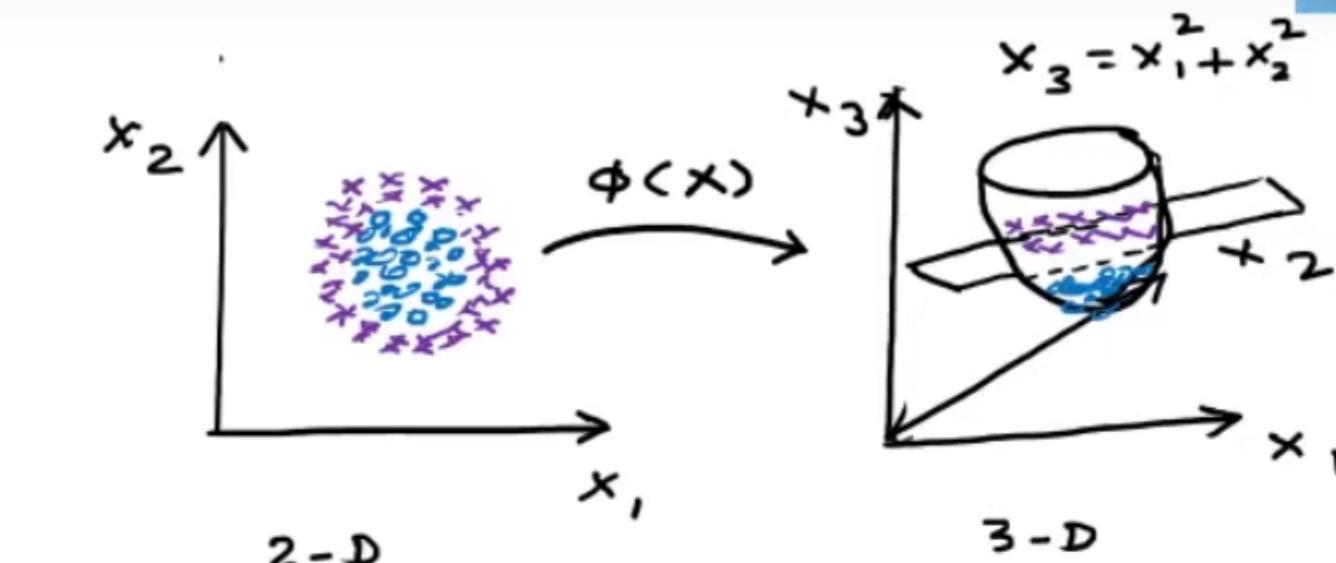
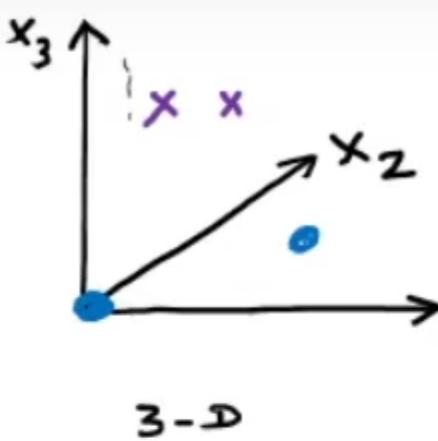
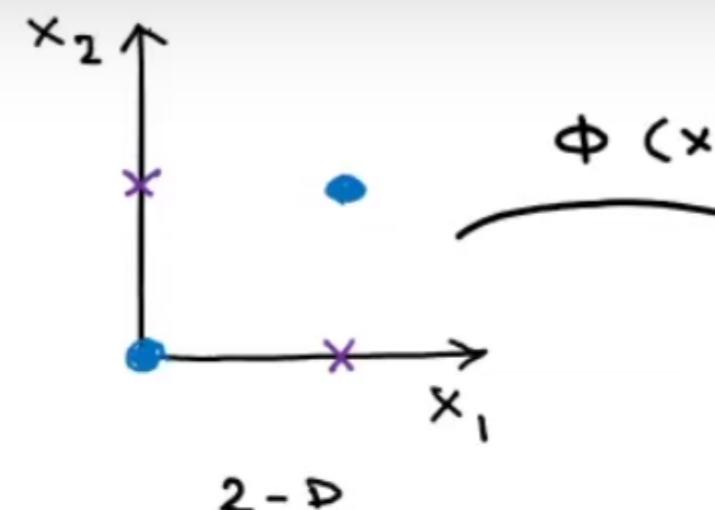


$$L = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y^{(i)} y^{(j)} \underbrace{\langle \phi(\bar{s}^{(i)}), \phi(\bar{s}^{(j)}) \rangle}_{\text{inner product}}$$

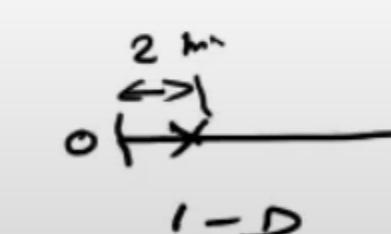


1.00

Kernel method in SVM



Physical Dimension.



$$\begin{aligned} & x_2 \uparrow \rightarrow \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in \mathbb{R}^2 \\ & x_3 \uparrow \quad x_2 \uparrow \quad x_1 \uparrow \rightarrow \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \in \mathbb{R}^3 \end{aligned}$$

Representational dimension.





Making SVM capable of solving nonlinear classification task

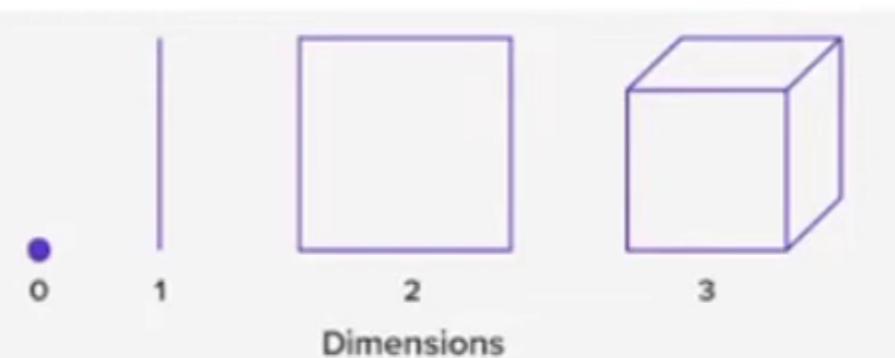
- Concept of Physical and Representational dimensions

- Concept of higher dimensional feature space

- Transformation function $\phi: x \rightarrow \phi(x)$; say we have samples:

$$\bar{S}_+^{(i)} = \left\{ \begin{pmatrix} 3 \\ 3 \end{pmatrix}, \begin{pmatrix} 3 \\ -3 \end{pmatrix}, \begin{pmatrix} -3 \\ 3 \end{pmatrix}, \begin{pmatrix} -3 \\ -3 \end{pmatrix} \right\} \rightarrow Y^{(i)} = 1$$

$$\bar{S}_-^{(i)} = \left\{ \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \begin{pmatrix} 2 \\ -2 \end{pmatrix}, \begin{pmatrix} -2 \\ 2 \end{pmatrix}, \begin{pmatrix} -2 \\ -2 \end{pmatrix} \right\} \rightarrow Y^{(i)} = -1$$



$$64 \times 64 \times 3 = 12288$$

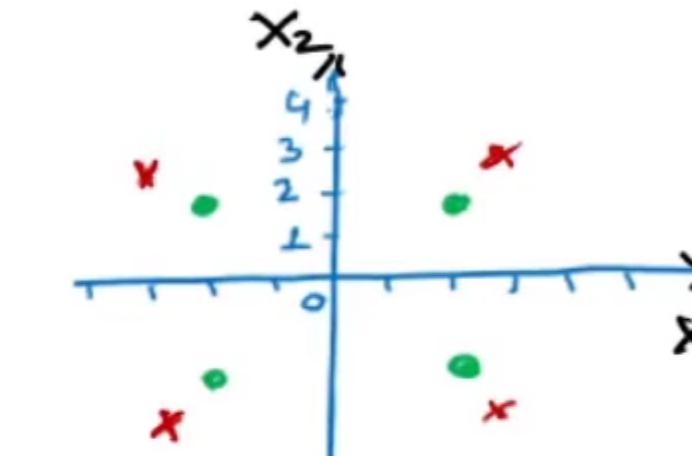
$$\begin{cases} 6 - \bar{s}^{(2)} + |\bar{s}^{(1)} - \bar{s}^{(2)}| & \text{if } \sqrt{\bar{s}^{(1)2} + \bar{s}^{(2)2}} > 3 \\ 6 - \bar{s}^{(1)} + |\bar{s}^{(1)} - \bar{s}^{(2)}| & \\ \bar{s}^{(2)} & \text{otherwise.} \end{cases}$$



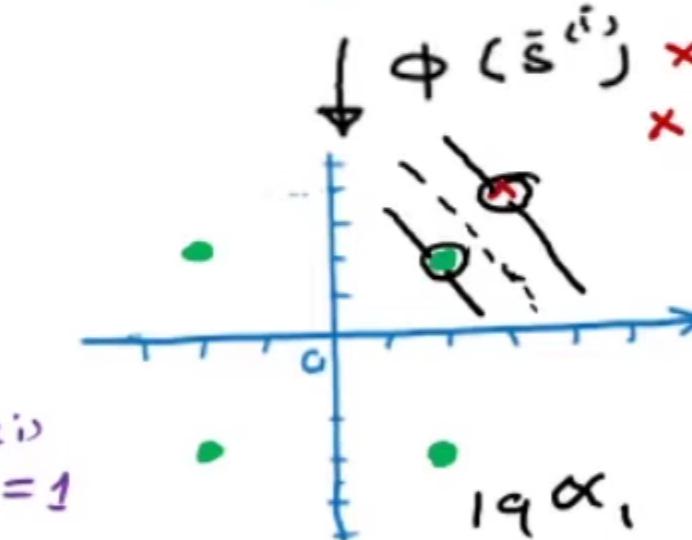
$$1000 \times 1000 \times 3 = 3M$$

$$\begin{aligned} \text{Sample-1} &\rightarrow \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \text{ Sample-2} \rightarrow \begin{pmatrix} 2 \\ -2 \end{pmatrix}, \text{ Sample-3} \rightarrow \begin{pmatrix} -2 \\ 2 \end{pmatrix} \} \rightarrow Y^{(i)} = 1 \\ \text{Sample-4} &\rightarrow \begin{pmatrix} -2 \\ -2 \end{pmatrix} \end{aligned}$$

$$\boxed{Y = w_1 x_1 + w_2 x_2 + b = -0.66x_1 - 0.66x_2 + 3.67}$$



$$\begin{aligned} Y^{(i)}(\bar{w} \cdot \bar{s}^{(i)} + b) &= 1 \\ \sum_{i=1}^2 \alpha_i Y^{(i)} (\bar{w} \cdot \bar{s}^{(i)}) &= 1 \end{aligned}$$



$$\begin{aligned} \alpha_1 \left(\begin{pmatrix} 3 \\ 3 \end{pmatrix} \right) + \alpha_2 \left(\begin{pmatrix} 3 \\ -3 \end{pmatrix} \right) &= 1 \\ \alpha_1 \left(\begin{pmatrix} 3 \\ 3 \end{pmatrix} \right) + \alpha_2 \left(\begin{pmatrix} -2 \\ 2 \end{pmatrix} \right) &= -1 \end{aligned}$$

$$19\alpha_1 + 13\alpha_2 = 1$$

$$13\alpha_1 + 9\alpha_2 = -1$$

$$\text{Solving} \rightarrow \alpha_1 = -8, \alpha_2 = 11.67$$

$$\begin{aligned} \tilde{w} &= \sum_{i=1}^2 \alpha_i \bar{s}^{(i)} Y^{(i)} \\ &= -8 \left(\begin{pmatrix} 3 \\ 3 \end{pmatrix} \right) + 11.67 \left(\begin{pmatrix} 3 \\ -3 \end{pmatrix} \right) + 11.67 \left(\begin{pmatrix} -2 \\ 2 \end{pmatrix} \right) = \begin{pmatrix} -0.66 \\ -0.66 \end{pmatrix} \end{aligned}$$

Making SVM capable of solving nonlinear classification task



Kernel Trick

- In general it is difficult to design data specific mapping function ϕ
- So we apply Kernel trick -

We write $k(x, z) = \underbrace{\phi(x)}_{\text{4x1}} \cdot \underbrace{\phi(z)}_{\text{4x1}} \quad \dots \quad (1)$

Now let $\phi(x) = \phi \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_{2 \times 1} = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_2 x_1 \\ x_2 x_2 \end{bmatrix}_{4 \times 1}$

$$\phi(z) = \phi \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}_{2 \times 1} = \begin{bmatrix} z_1 z_1 \\ z_1 z_2 \\ z_2 z_1 \\ z_2 z_2 \end{bmatrix}_{4 \times 1}$$

Then from (1) $\rightarrow k(x, z) = \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ x_2 x_1 \\ x_2^2 \end{bmatrix} \cdot \begin{bmatrix} z_1^2 \\ z_1 z_2 \\ z_2 z_1 \\ z_2^2 \end{bmatrix} = (x_1 z_1)^2 + x_1 x_2 z_1 z_2 + x_1 x_2 z_1 z_2 + (x_2 z_2)^2$
 $= (x_1 z_1 + x_2 z_2)^2$
 $= (x \cdot z)^2$

We call k is a quadratic kernel.

Kernel Trick



Take another example :

$$K(x, z) = \phi(x) \cdot \phi(z)$$

$$\text{Let } \phi(x) = \phi \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} x_1^3 \\ \sqrt{3} \cdot x_1 x_2^2 \\ \sqrt{3} \cdot x_1^2 x_2 \\ x_3^3 \end{bmatrix}$$

calculate

$$\phi(z) = \phi \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix} = \begin{bmatrix} z_1^3 \\ \sqrt{3} \cdot z_1 z_2^2 \\ \sqrt{3} \cdot z_2 z_1^2 \\ z_3^3 \end{bmatrix} \Rightarrow \phi(x) \cdot \phi(z) \\ \Rightarrow (x_1 z_1 + x_2 z_2)^3$$

Different types of Kernels:

1. Polynomial kernel $K(x, z) = (x \cdot z + 1)^d$, where d is the degree of polynomial.

2. Gaussian kernel, $K(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right)$

3. Gaussian Radial Basis function (RBF), $K(x, z) = \exp\left(-\gamma \|x - z\|^2\right)$

4. Laplace RBF kernel, $K(x, z) = \exp\left(-\frac{\|x - z\|}{\sigma}\right)$

5. Hyperbolic tangent kernel, $K(x, z) = \tanh(\beta_0 x \cdot z + \beta_1)$.

Kernel Trick



SVM= Optimal Classifier + Kernel Trick

$$\bar{s}^{(i)} \rightarrow x^{(i)} \rightarrow x; \quad \bar{s}^{(j)} \rightarrow x^{(j)} \rightarrow z$$

$$x \xrightarrow{\phi} \phi(x); \quad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in \mathbb{R}^{n=2} \xrightarrow{\text{O}(n)} \phi(x)$$

$$\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_2 x_1 \\ x_2 x_2 \end{bmatrix} \in \mathbb{R}^{h=4}$$

Considering all pairs of monomial terms.

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \in \mathbb{R}^{n=3} \rightarrow \phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix}$$

Because there are n^2 elements, you need $\text{O}(n^2)$ time to compute $\phi(x)$ or

$$\phi(x)^T \phi(z)$$

How kernel can reduce the complexity?

$$\text{Let, } K(x, z) = \underbrace{\phi(x)^T \phi(z)}_{\text{O}(n^2)} = \underbrace{(x^T z)^2}_{\text{O}(n)}$$

$$(x^T z)^2 = \left(\underbrace{\sum_{i=1}^n x_i z_i}_{x^T z} \right) \left(\underbrace{\sum_{j=1}^n x_j z_j}_{z^T x} \right) = \sum_{i=1}^n \sum_{j=1}^n x_i z_i x_j z_j$$

$$= \sum_{i=1}^n \sum_{j=1}^n (x_i \cdot x_j) (z_i \cdot z_j) \rightarrow \text{O}(n)$$

$$\in \mathbb{R}^{h=9}$$

Similarly $\phi(z) =$

$$\begin{bmatrix} z_1 z_1 \\ z_1 z_2 \\ z_1 z_3 \\ z_2 z_1 \\ z_2 z_2 \\ z_2 z_3 \\ z_3 z_1 \\ z_3 z_2 \\ z_3 z_3 \end{bmatrix} \in \mathbb{R}^{h=9}$$



Kernel Trick

If we choose $\phi(x)^T \cdot \phi(z) = K(x, z) = (x^T z + c)^d$; For $d = 2$

$$\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix} = \begin{bmatrix} n+d \\ d \end{bmatrix}$$

(Can have all
(n+d) features
of monomial
terms up to the
order of d. Roughly we can
have $(n+d)^d$ features. This is large -

$$\phi(x) = \begin{bmatrix} z_1 z_1 \\ z_1 z_2 \\ z_1 z_3 \\ z_2 z_1 \\ z_2 z_2 \\ z_2 z_3 \\ z_3 z_1 \\ z_3 z_2 \\ z_3 z_3 \end{bmatrix} ; \quad \phi(z) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix}$$

$$; \quad \phi(x) = \begin{bmatrix} \sqrt{2c} x_1 \\ \sqrt{2c} x_2 \\ \sqrt{2c} x_3 \\ c \end{bmatrix} ; \quad \phi(z) = \begin{bmatrix} z_1 z_1 \\ z_1 z_2 \\ z_1 z_3 \\ z_2 z_1 \\ z_2 z_2 \\ z_2 z_3 \\ z_3 z_1 \\ z_3 z_2 \\ z_3 z_3 \end{bmatrix}$$

But with our quadratic kernel, the computation is
still $O(n)$

Kernel Trick



How to make Kernels?

From Intuition if x and z are similar then

$$k(x, z) = \underbrace{\phi(x)^T \phi(z)}_{\text{is large}}$$

If x and z are dissimilar then $k(x, z)$ is small

$$K =$$



$$\left[\begin{array}{cccc} \phi(x_1)^T \phi(x_1) & \phi(x_1)^T \phi(x_2) & \dots & \phi(x_1)^T \phi(x_n) \\ \phi(x_2)^T \phi(x_1) & \phi(x_2)^T \phi(x_2) & \dots & \phi(x_2)^T \phi(x_n) \\ \vdots & \vdots & \ddots & \vdots \\ \phi(x_m)^T \phi(x_1) & \phi(x_m)^T \phi(x_2) & \dots & \phi(x_m)^T \phi(x_n) \end{array} \right]$$

Gram matrix \times

$$\left[\begin{array}{cccc} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_n) \\ k(x_2, x_1) & k(x_2, x_2) & \dots & k(x_2, x_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_m, x_1) & k(x_m, x_2) & \dots & k(x_m, x_n) \end{array} \right]$$

$k \rightarrow 0 \rightarrow \underline{\text{Symmetric}}$

SVM as linear machine

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j}^{m,m} y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

SVM as nonlinear machine

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j}^{m,m} y^{(i)} y^{(j)} \alpha_i \alpha_j \langle \phi(x^{(i)}) \cdot \phi(x^{(j)}) \rangle$$

SVM as nonlinear classifying with kernel trick

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j}^{m,m} y^{(i)} y^{(j)} \alpha_i \alpha_j$$

$$\underbrace{\phi(x)^T \phi(z)}$$

$$\boxed{k(x, z)}$$

How to design your own Kernel?



All Kernels more or less satisfy the small and large concept.

For example 1. Gaussian Kernel $\rightarrow K(x, z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$

If x & z are similar, then $K(x, z) \rightarrow 1$

If x & z are dissimilar, then $K(x, z) \rightarrow 0$

2. Radial basis Function kernel $\rightarrow K(x, z) = \exp(-\gamma \|x-z\|^2)$
Same conditions are obeyed.

All kernels must obey Mercer's theorem which states that :

For K to be a valid kernel function, there must exists ϕ such that
 $K(x, z) = \phi(x)^T \phi(z)$ and the corresponding kernel matrix K must be
symmetric and positive semidefinite i.e. $K \geq 0$, being constituted from a
finite set of m points and having dimension $m \times m$ with its $(i-j)$ entry is given by
 $K_{ij} = K(x^{(i)}, x^{(j)})$.

How to make your own kernel



Proof :

Let K is a valid kernel corresponding to some feature mapping ϕ .

Now consider some finite set of m points $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ and let a square matrix $K_{m \times m}$ be defined such that its (i, j) -entry is given by $K_{i,j} = K(x^{(i)}, x^{(j)})$. This matrix is kernel matrix.

Now if K is a valid kernel then $K_{i,j} = K(x^{(i)}, x^{(j)}) = \phi(x^{(i)})^T \phi(x^{(j)}) = \phi(x^{(j)})^T \phi(x^{(i)}) = K(x^{(j)}, x^{(i)})$

Also letting $\phi_k(x)$ to denote the k -th coordinate of the vector $\phi(x)$, then for any arbitrary vector z , we can write:

$$z^T K z = \sum_i \sum_j z_i K_{i,j} z_j \quad (\text{By the definition of matrix multiplication})$$

$$= \sum_i \sum_j z_i \underbrace{\phi(x^{(i)})^T \phi(x^{(j)})}_{\phi_k(x^{(i)}) \phi_k(x^{(j)})} z_j \quad \therefore a^T b = \sum_k a_k b_k$$

$$= \sum_i \sum_j z_i \sum_k \phi_k(x^{(i)}) \phi_k(x^{(j)}) z_j$$

$$= \sum_k \sum_i \sum_j z_i \phi_k(x^{(i)}) \cdot \phi_k(x^{(j)}) z_j = \sum_k \left(\sum_i z_i \phi_k(x^{(i)}) \right)^2 \geq 0$$

\Rightarrow Hence ' K ' is positive semi-definite i.e all its eigenvalues $\lambda_i \geq 0$



Generalised Kernel method

Kernel method is very general and can be applied to all the discriminative models you have learnt so far like:

- **Linear Regression**
- **Logistic Regression**
- **Perceptron learning**

Kernel based methods can also be applied to generative models like **PCA & LDA**



Properties of SVM

- Out of a large data set normally there are only a handful of support vectors and only support vectors are used to specify the separating hyper plane.
- Capable of handling large feature spaces:
 - computational complexity does not depend on the dimensionality of the feature space.
- Some flexibility of choosing kernel functions.
- Overfitting can be controlled by soft margin approach
- Has a nice mathematical foundation: The classification problem has been formulated as a constrained quadratic optimization problem, guaranteeing that the training parameters are tuned using global optimal values.



Limitations

- **It is sensitive to noise**
 - A relatively small number of mislabeled examples can dramatically decrease the performance
- **It only considers two classes**
 - how to do multi-class classification with SVM?
 - 1) with output m, learn m SVM's
 - **SVM 1 learns "Output==1" vs "Output != 1"**
 - **SVM 2 learns "Output==2" vs "Output != 2"**
 - :
 - **SVM m learns "Output==m" vs "Output != m"**
 - 2) To predict the output for a new input, just predict with each SVM and find out which one puts the prediction the furthest into the positive region.



SVM Applications

- SVM has been used successfully in many real-world problems:
 - Text (and hypertext) categorization.
 - Image classification.
 - Bioinformatics (Protein classification, Cancer classification, etc.).
 - Hand-written character recognition.



Comparison of SVM with other classifiers

S. No.	Name of the Classifier	Test error rate (%)
1.	K - Nearest Neighbours (well formulated)	3 - 5
2.	SVM (Gaussian Kernel / 4 degree polynomial)	0.5 -1.5
3.	NN (2 - layers, 300 Hidden units)	3 - 4.5
	6 - layers running on GPU	1
4.	CNN augmented with SVM	0.5