# Support Vector Machines
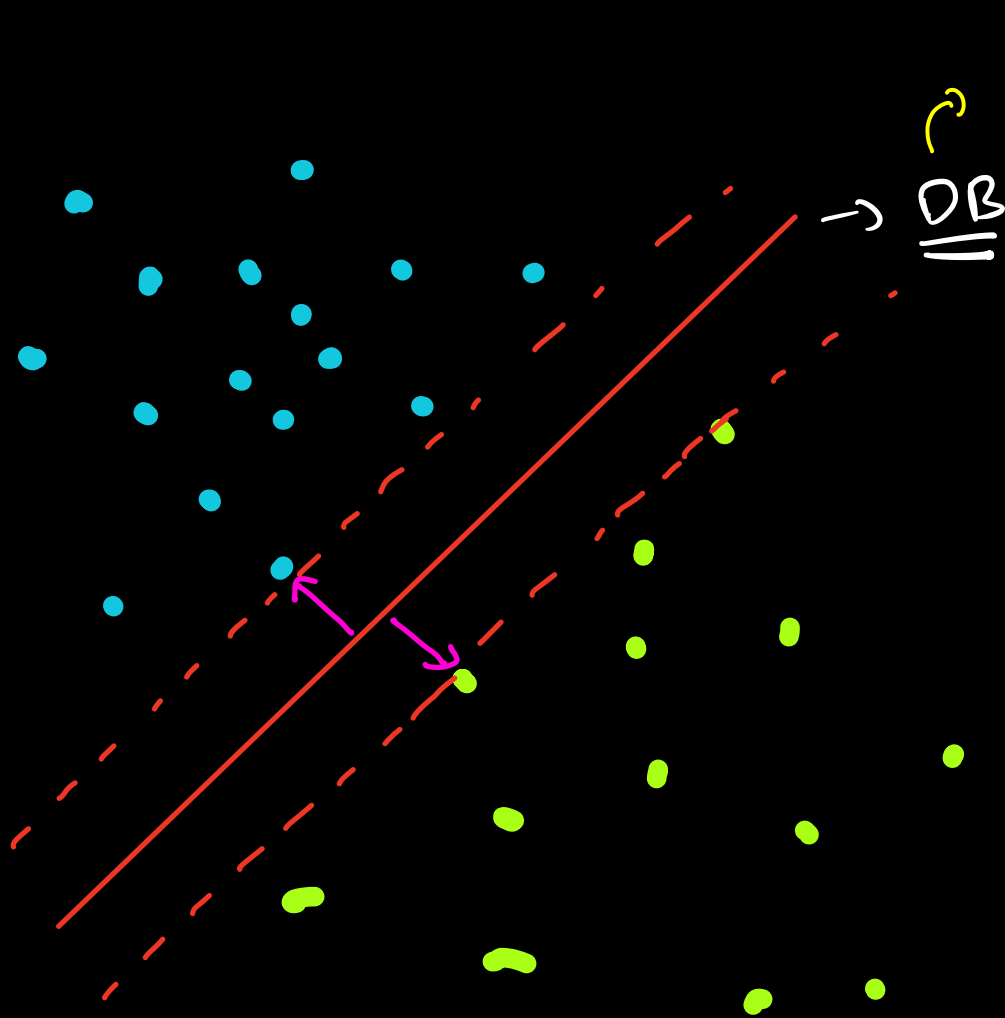## [ ML - 1 ]

→ Recap SVM-1

→ Dual form

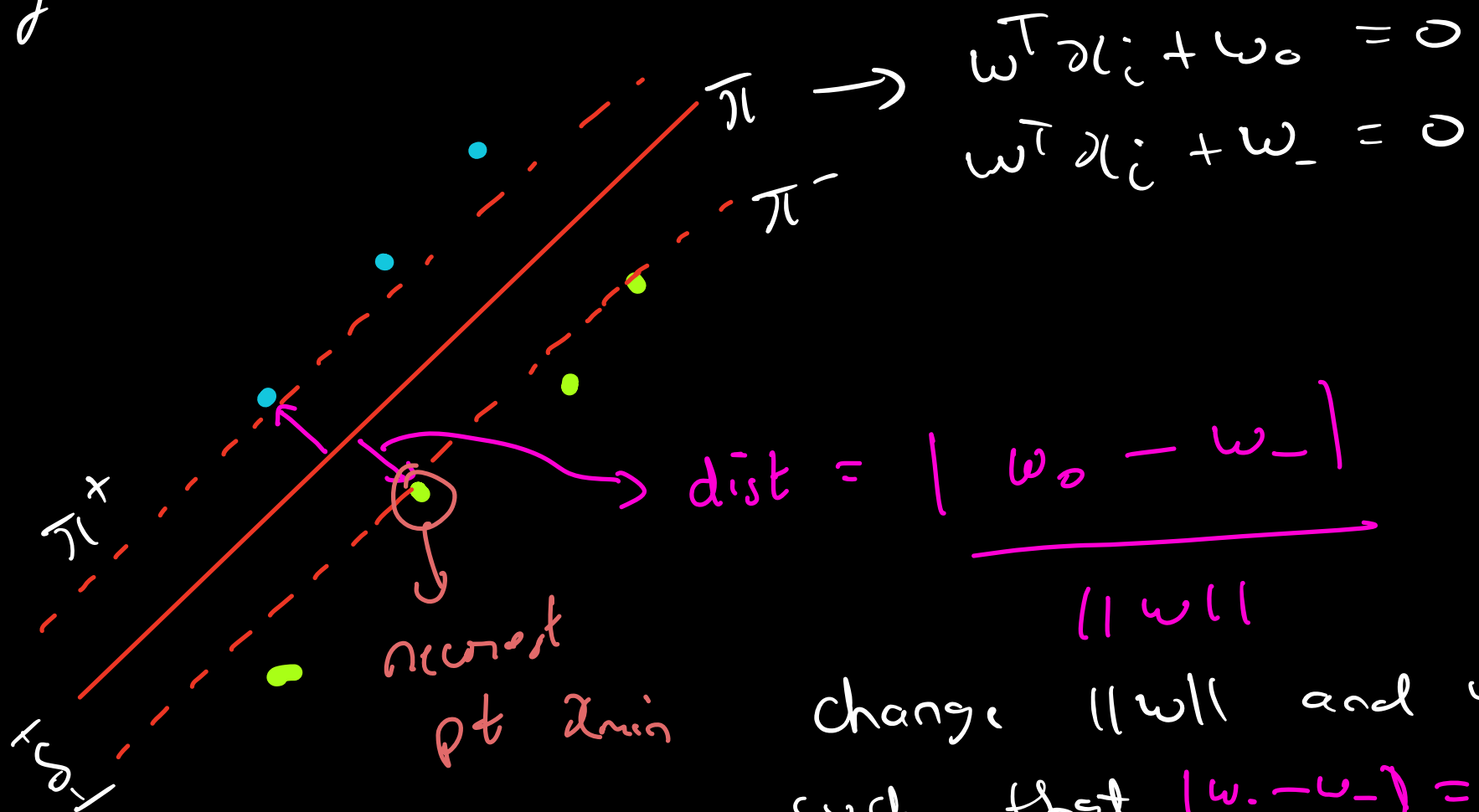→ Kernel Trick

# Recap



$$w^T x_i + w_0 = 0$$

→ DB

( linear )

( perfectly separable

DB should be
as far as possible
from either of the
nearest point

Maximise   margin
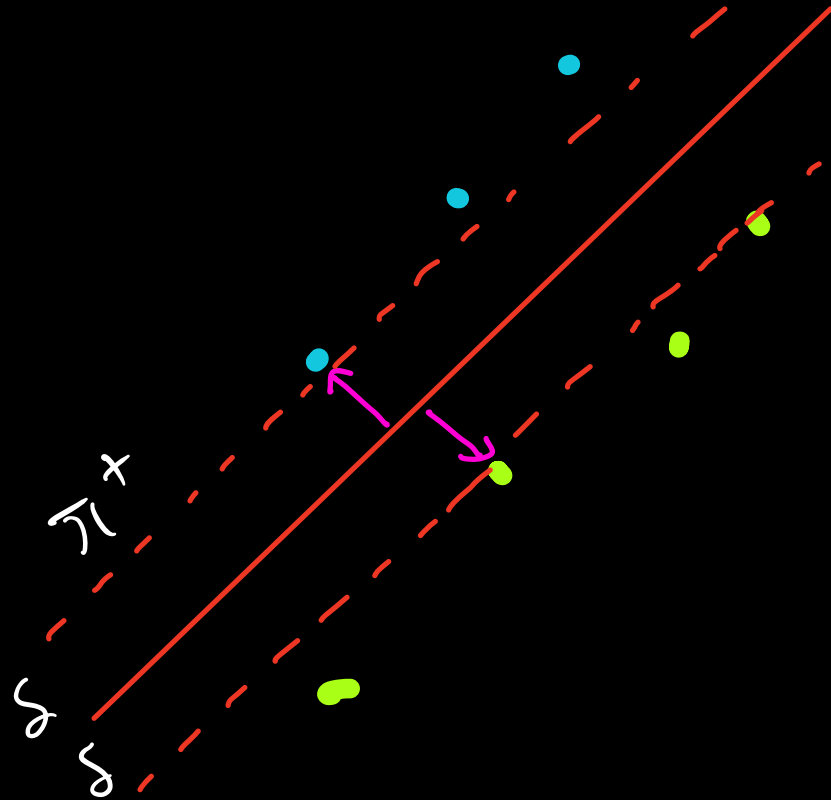
→ We take nearest points because they deserve to have a say. Further pts, may have outliers.

$\pi \longrightarrow$   $w^T x_i + w_0 = 0$

$\pi^-$   $w^T x_i + w_- = 0$

$\pi^+$

$\pi_+$

$\pi_-$

nearest pt $x_{min}$

$dist = \dfrac{|w_0 - w_-|}{||w||}$

change $||w||$ and $w_0$ such that $|w_0 - w_-| = 1$

Change s.t. min dist,

i.e $\delta = 1$:

$w^T x + (w_0 + 1) = 0$

$w^T x + w_0 = 0$

$w^T x + (w_0 - 1) = 0$

For simplicity

$\pi^+$

$\delta$

$\delta$

Objective:

$$\max_{\vec{w}, w_0} : \frac{2}{||w||}$$

# Hard Margin SVM
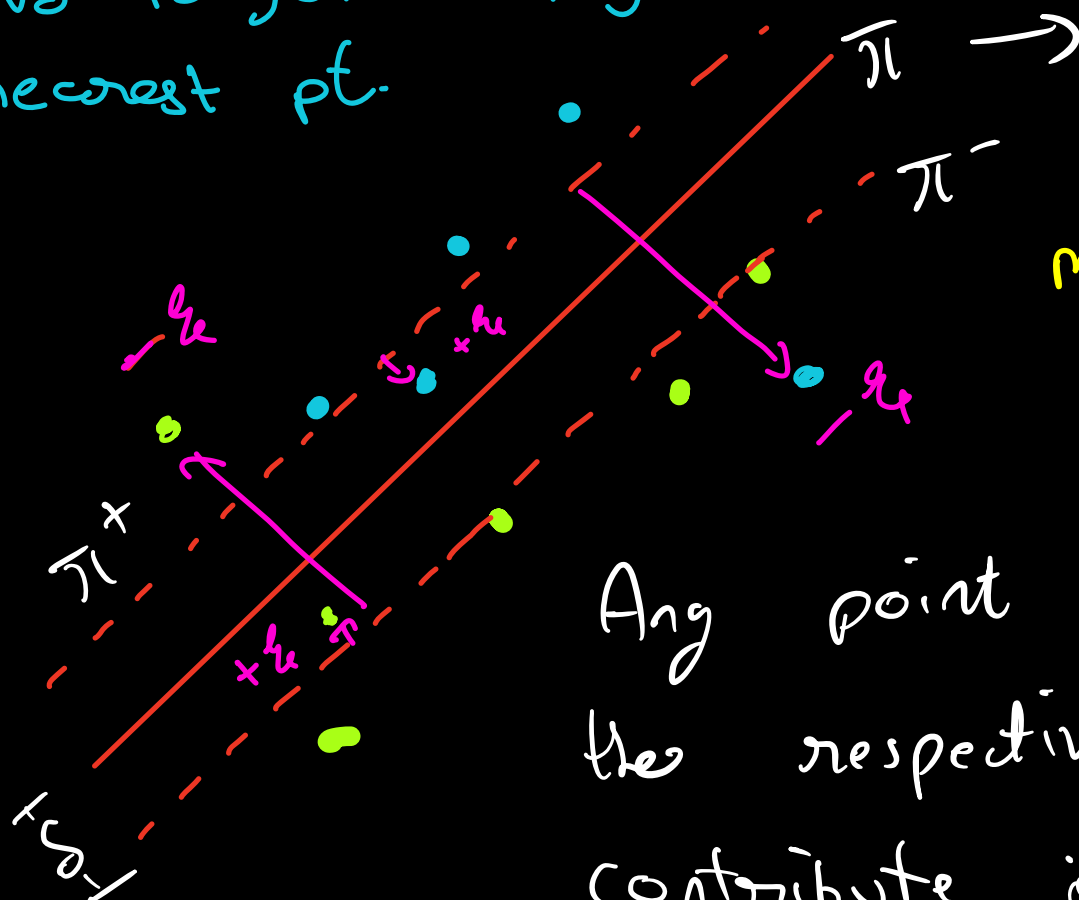
$$\max_{\vec{w}, w_0} \frac{2}{||w||} \quad ; \text{ s.t. } (w^T x_i + w_0) y_i \geq 1$$

$$\forall i: 1 \to n$$

But, this does not work if there is any mix / overlap of classes.

$$= \min_{\vec{w}, w_0} \frac{||w||}{2}$$

## Soft Margin SVM

No longer taking nearest pt.

$\pi \longrightarrow$ Here we care about dist from margin lines, not DB.

$\pi^-$

$\pi^+$

Any point which is outside the respective margin line will contribute in computing the DB.

$\xi > 0$  correctly classified support vectors

$\xi < 0$  incorrectly classified support vectors

$$\min_{w, w_0, \xi_i} \frac{||w||}{2} + C \sum_{i=1}^{n} \xi_i$$

only support vectors

max margin ⟷ min crossing of margin

balance

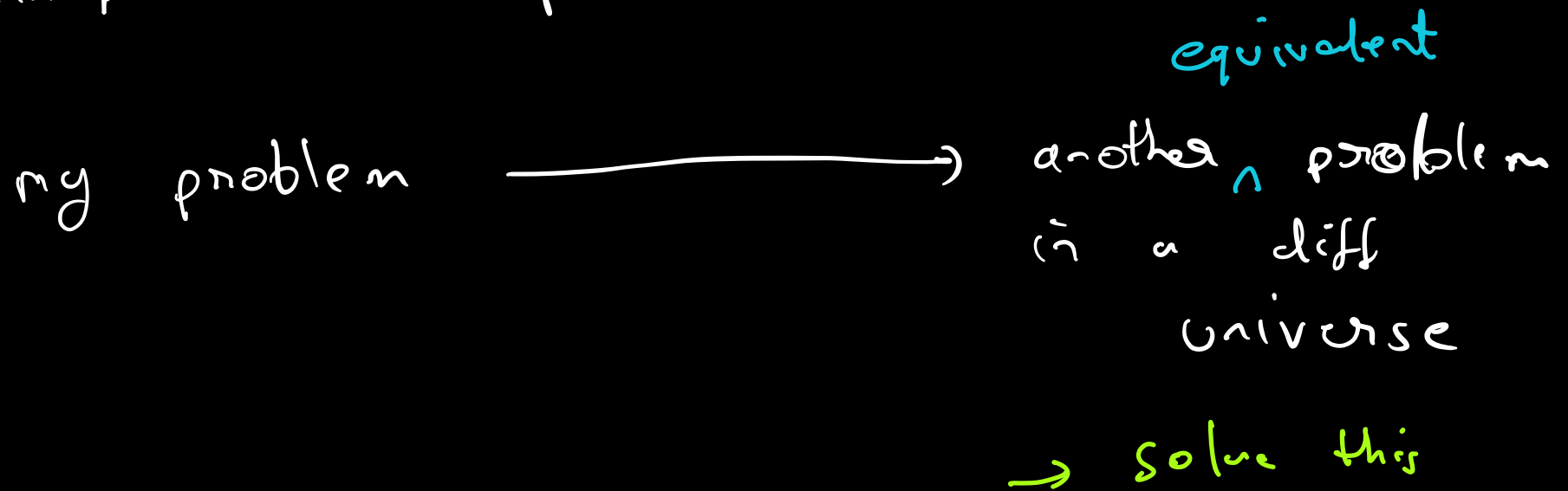$$st: \quad (w^T x_i + w_0) y_i \geq 1 - \xi_i$$

$$\xi_i > 0$$

$$\forall i: 1 \longrightarrow n$$

Solving it in this form is possible, however
converting it to a dual form would make
it possible to introduce non-linearity

## Dual Form of Objective Func$^n$

Principal in optimisation

my problem $\longrightarrow$ another $\underset{\wedge}{}$ problem
                                    equivalent
                                    in a diff
                                        universe

$\rightarrow$ solve this

$$\max_{\alpha_i} \quad \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j \, y_i \, y_j \, x_i^T x_j$$

$$\text{s.t} \quad 0 \le \alpha_i \le C \qquad \Biggr] \quad \forall i : 1 \to n$$

$$\sum_{i=1}^{\hat{}} \alpha_i y_i = 0 \qquad \Biggr]$$

Here,

prediction eqⁿ : → $f(x_q) = \sum_{i=1}^{\hat{}} \alpha_i y_i \, x_i^T x_q$  ⟵ query point

$\underbrace{\qquad\qquad\qquad}_{\text{support vectors}}$

Observe that $\vec{x}$ never appears alone anymore!

it's always $x_a^T x_b$

$\alpha_i = 0 \quad \rightarrow$ non support vect

$\alpha_i > 0 \quad \rightarrow$ support vect

## Full Picture

Man Maximisation

$\downarrow$

Soft margins $\longrightarrow \xi_i$

$\downarrow$

Optimisation

$\downarrow$

Dual $\rightarrow \alpha_i$

# Kernel SVM

$$\max_{\alpha_i} \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j \, y_i \, y_j \, \underbrace{x_i^T x_j}$$

$\downarrow$

Dot product

Let us re-look this!
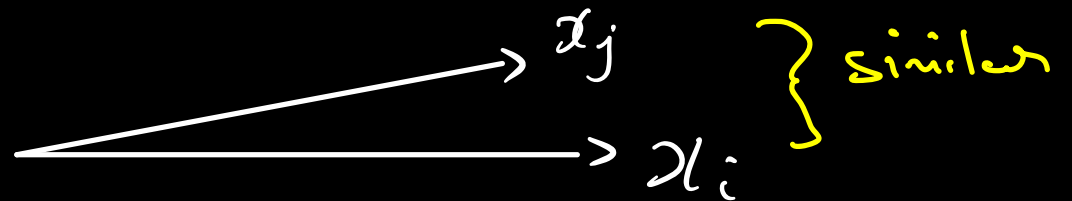
"Dot product" can be interpreted as similarity

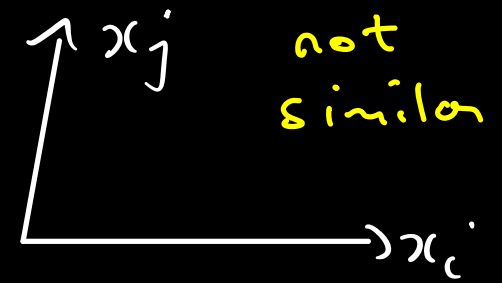Mathemetitians realised that they can change this to any notion of similarity

$\rightarrow$ cosine

$\rightarrow$ custom ..

$\rightarrow$ etc ..



} similar

Now we can write:

$$\sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j \, y_i \, y_j \, \underbrace{K(x_i, x_j)}_{\text{Kernel Function}}$$

$x_j$ not similar

$x_i$

$\rightarrow$ predict func$^n$ : $\sum_{i=1}^{\hat{} } \alpha_i \, y_i \, K(x_i, x_q)$

The beauty here is that the kernel function helps provide a convinient place to introduce non-linearity.

Examples :
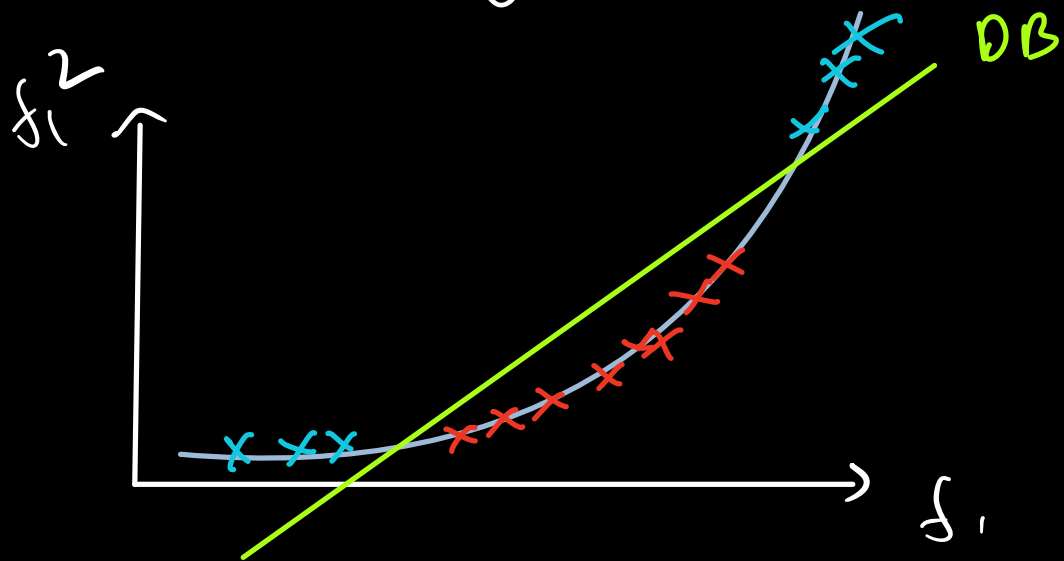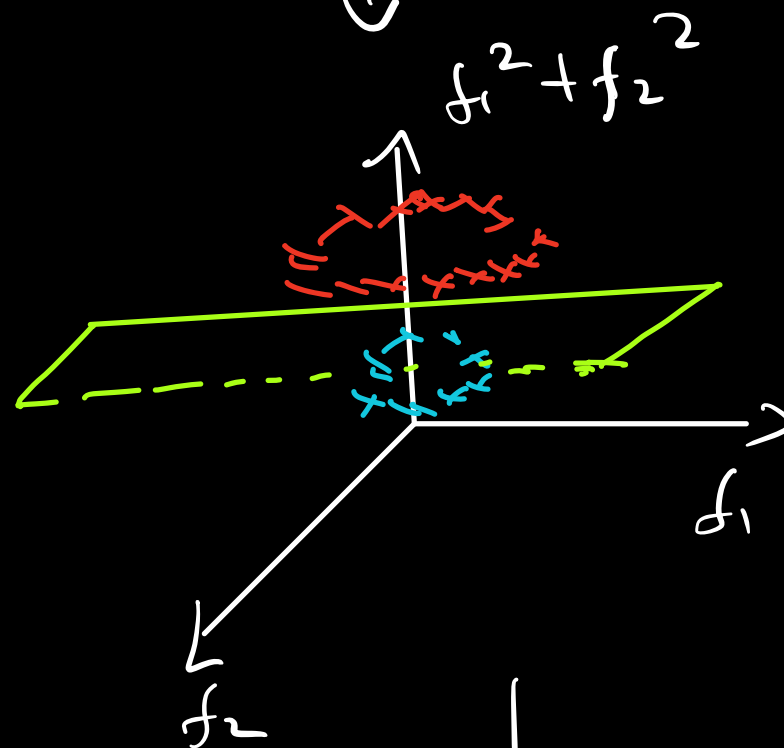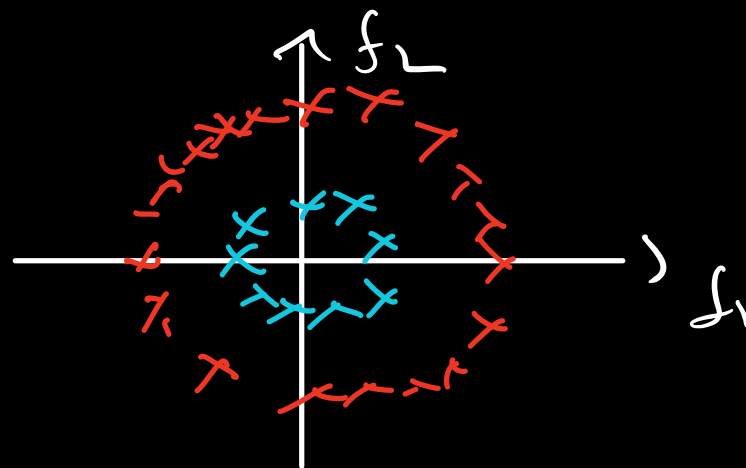
Low - Dim

$f_1$

not    seperable

High - Dim

$f_1^2$

DB

$f_1$

$f_2$

$f_1$

$f_1^2 + f_2^2$

$f_1$

$f_2$
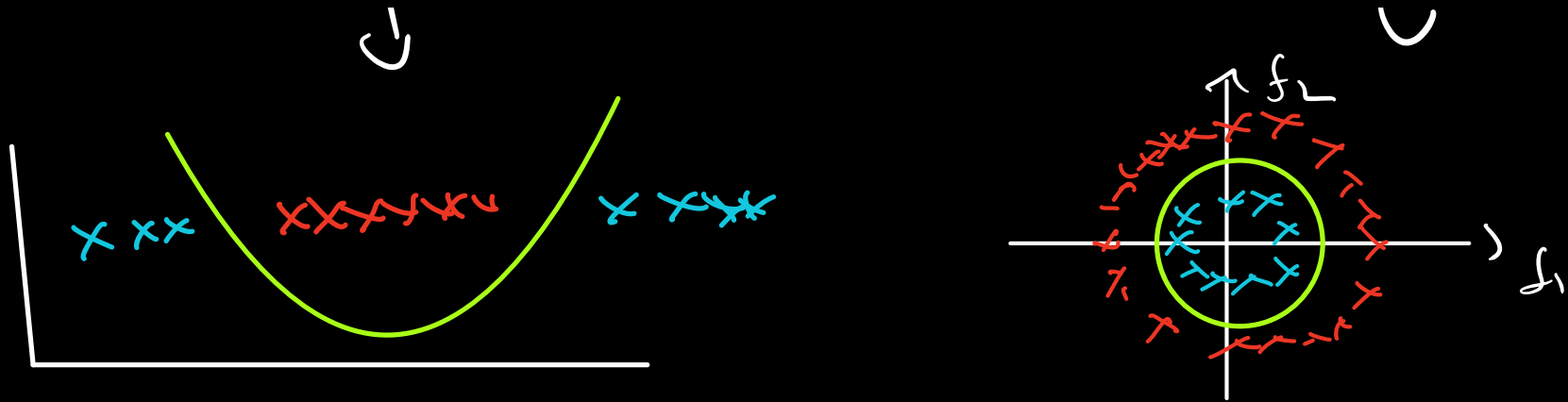
So kernels may be used to create more features (non-linear) where SVM can now try a linear decision boundary in a high dimensional space. When we come back to original space we get a non-linear DB

How does kernel accomplish this?

## Polynomial Kernel

$$K(x_i, x_2) = (x_1^t x_2 + 1)^m$$

$m =$ degree of polynomial

Eg: $m = 2$ (quadratic)

$$\vec{x_1} = \begin{bmatrix} x_{11} \\ x_{12} \end{bmatrix}, \quad \vec{x_2} = \begin{bmatrix} x_{21} \\ x_{22} \end{bmatrix}$$

$$\therefore K(x_1, x_2) = \left( \begin{bmatrix} x_{11} & x_{12} \end{bmatrix} \begin{bmatrix} x_{21} \\ x_{22} \end{bmatrix} + 1 \right)^2$$

$$= (1 + x_{11} x_{21} + x_{12} x_{22})^2$$

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

$$= 1 + x_{11}^2 x_{21}^2 + x_{12}^2 x_{22}^2 + 2x_{11}x_{21} + 2 x_{12}x_{22}$$

$$+ 2 x_{11}x_{21}x_{12}x_{22} \qquad — Ⓐ$$

Now consider:

$$x_1' = [1, x_{11}^2, x_{12}^2, \sqrt{2}x_{11}, \sqrt{2}x_{12}, \sqrt{2}x_{11}x_{12}]$$

$$x_2' = [1, x_{21}^2, x_{22}^2, \sqrt{2}x_{21}, \sqrt{2}x_{22}, \sqrt{2}x_{21}x_{22}]$$

$$\underbrace{\qquad}_{1} \quad \underbrace{\qquad}_{x_i^2} \quad \underbrace{\qquad}_{x_i} \quad \underbrace{\qquad}_{x_i x_j}$$

bias    non-linear    linear    non-linear

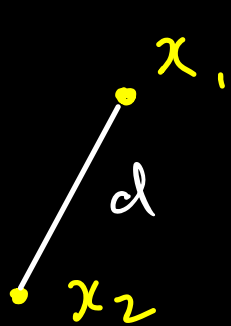$$\text{Magic} \rightarrow \vec{x_1'} \cdot \vec{x_2'} = Ⓐ \quad [2d \rightarrow 6d]$$

So using this small $K(w_1, w_2)^m$ we are able to get all powers of original features and all multiplicative combinations
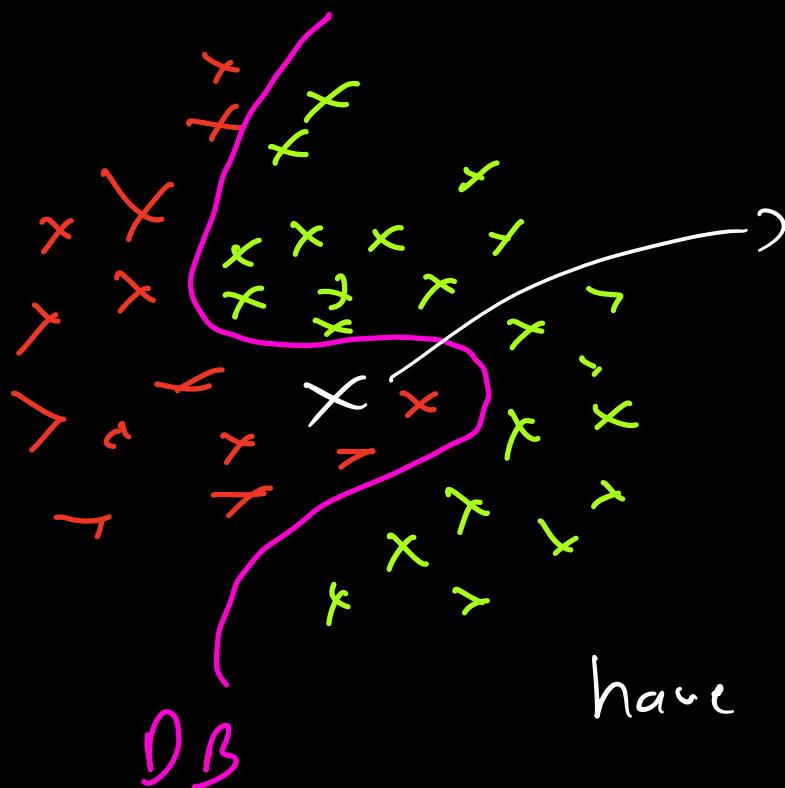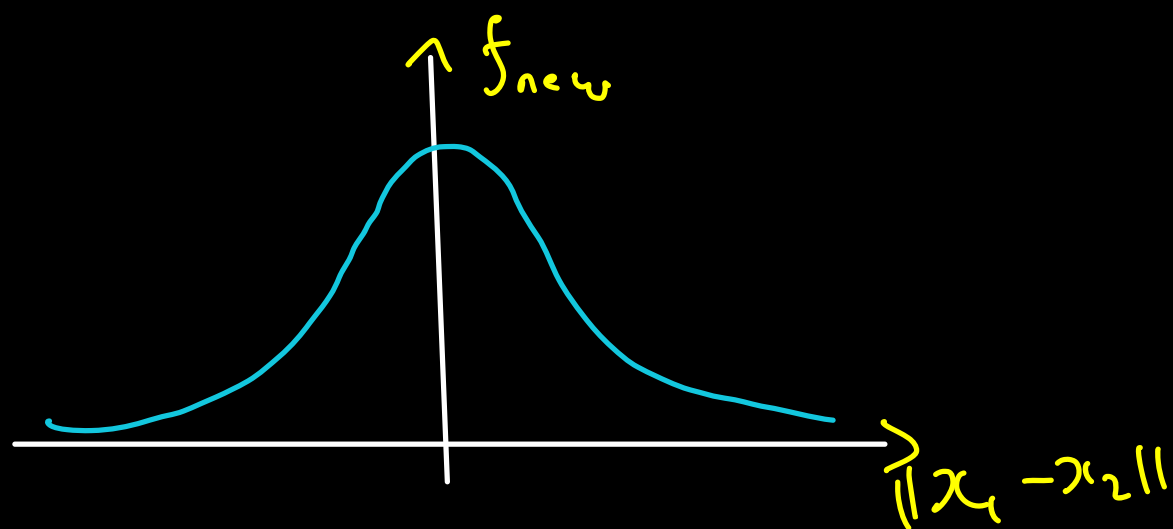
RBF Kernel $\rightarrow$ Most Popular

(Radial Basis Func$^n$ / Gaussian) $\rightarrow$ euclidean dist

$$K_{RBF}(x_1, x_2) = \exp\left(\frac{-\|x_1 - x_2\|^2}{2\sigma^2}\right)$$

$\hookrightarrow$ H.P

$$K_{RBF}(x_1, x_2) = e^{\frac{-d^2}{2\sigma^2}}$$

$f_{new}$

$\|x_1 - x_2\|$

DB

query point

→ dist b/w nearest points?)

→ close points have high $f_{new}$.

→ KNN behaviour

KNN $\rightarrow$  K $\uparrow$   underfit ;  K$\downarrow$  } overfit

RBF $\rightarrow$  $\sigma \uparrow$   underfit ;  $\sigma \downarrow$

RBF $\sim$ Polynomial $(m \rightarrow \infty)$

How?

m=1  m=2  m=3

m=7

$\rightarrow$

$m \rightarrow \infty$
$\rightarrow$ RBF
can make
such DB ✓