

Theory - SVM

Sou Sarathi
IPHD20

* Topics :

- ① SVMs
- ② Hard Margin SVM
- ③ Soft Margin SVM
- ④ Projected Gradient Descent [Incomplete]
- ⑤ Non Linear SVMs
- ⑥ Multi Class SVMs [Incomplete]

links :

- ① YouTube video
- ② YouTube - Statquest, multivariate
- ③ washington winter '17 Feisiga
- ④ cs.cmu Abadi & Georff.

Given the data, for supervised problem

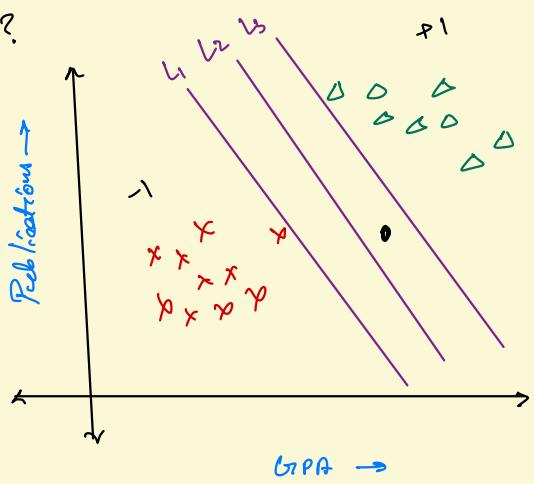
Sou Sarathli.

- ↳ Choose a model (Linear, Bilinear, Polynomial etc)
- ↳ Choose a suitable cost func. to find optimal parameters (least sq. error, logistic, sigmoidal etc)
- ↳ Choose a suitable iterative tech. to solve
(Grad. Descent, Stoch grad. descent)

Source: Riturik math SVM
YouTube

Support Vector Machines:

Why?



PHD Positions

+1 : Selected

-1 : not getting selected.

Question: which plane would you choose to separate the data? $\rightarrow L_2$ hyperplane. Why?
What's the prob. with L_3 ? It also separates the data.

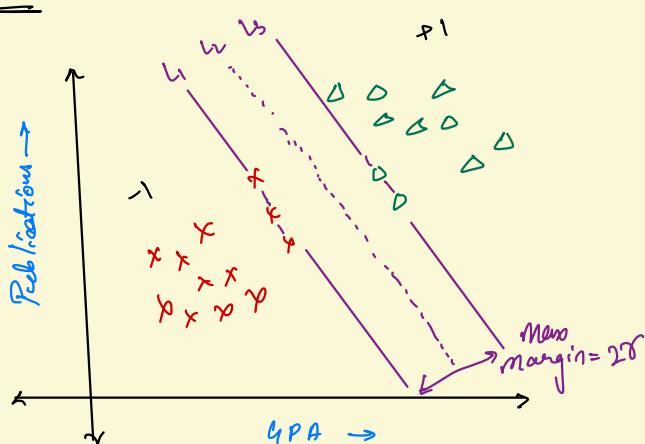
Suppose if we choose L_3 and we have a candidate in b/w $L_2 \times L_3$ (-) he can get rejected despite having good Prob. & C.R.A. Same is the case with L_1 .

we choose L_2 i.e we choose an hyperplane which is not dangerously close to both data to avoid misclassification.

In other words, we want to maximize the margin.

(Margin: is the shortest dist b/w data & the chosen hyperplane) Margin = 2r

Support Vectors:



we choose L_2 hyperplane & it is midway b/w the margin.
 L_2 = Decision boundary

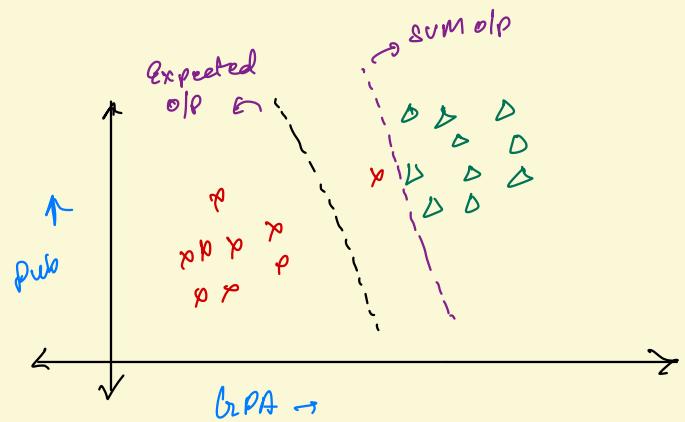
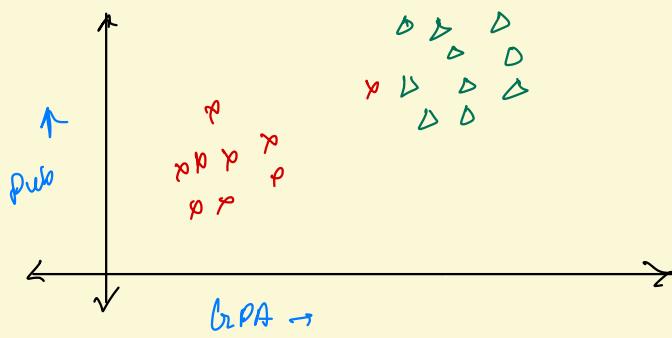
Support Vectors:

Some of the training data are called support vectors since they are the ones fully determining the margin.

Margin is supported by these support vectors. 1 data.

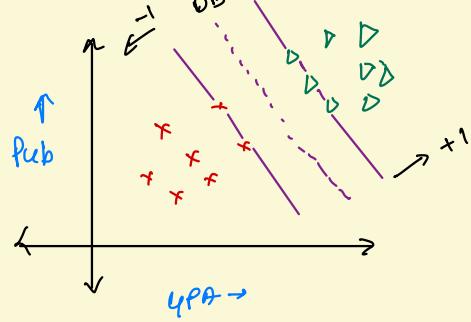
* Aim of SVM: is to maximize the margin.

One disadvantage of SVM: in case of an outlying data point



because SVM tries to maximize the margin

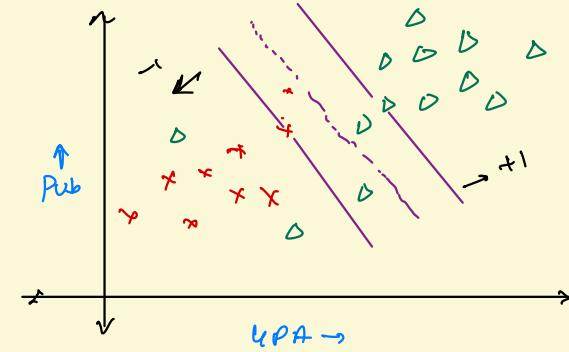
① Hard Margin



To right of DB: +1

To left of DB: -1

works for linearly separable



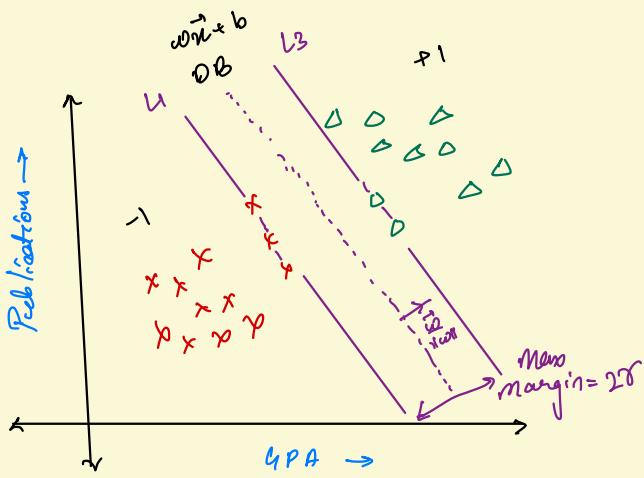
For data that are not linearly separable.

Classify data so that type I is on one side of Bound & type II is on the other side but with some misclassifications with a penalty for every mistake depending upon the mistake

SVM formulation:

- Margin If x on DB, $wx+b=0$.

to determine: γ



Let L_3 be the hyperplane

$$\vec{w} \cdot \vec{x} + b = 1$$

Given a x on DB : $\vec{w} \cdot \vec{x} + b$,

$$\vec{w} \cdot (\vec{x} + K \frac{\vec{w}}{\|\vec{w}\|}) + b = 1 \quad \text{for some } K.$$

$$\vec{w} \cdot \vec{x} + b + K \frac{\vec{w} \cdot \vec{w}}{\|\vec{w}\|} = 1$$

$$0 + K \frac{\|\vec{w}\|^2}{\|\vec{w}\|} = 1 \Rightarrow K = \frac{1}{\|\vec{w}\|}$$

$$K = \frac{1}{\|\vec{w}\|}$$

$$\text{So, } \delta = \frac{1}{\|\vec{w}\|}$$

\therefore Problem: Maximize 2δ = Maximize $\frac{2}{\|\vec{w}\|}$

$$= \text{minimizing } \|\vec{w}\| = \text{min } \|\vec{w}\|^2 = \text{min } \frac{\|\vec{w}\|^2}{2}$$

Constraints: ① If $y_i = 1$ we need $\vec{w} \cdot \vec{x}_i + b \geq 1$

② If $y_i = -1$ we need $\vec{w} \cdot \vec{x}_i + b \leq -1$

Hence we require, $y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1$ $\forall i$

Optimization problem. :

SVM

$$\min_{w,b} \frac{1}{2} w^T w \quad \text{subject to}$$

$$y_i (\vec{w}^T \vec{x}_i + b) \geq 1 \quad \forall i$$

\hookrightarrow (Convert as $1 - y_i (\vec{w}^T \vec{x}_i + b) \leq 0$)

Inequality constraint.

Step 1 $L(w, b, \alpha^*) = \frac{1}{2} w^T w + \sum_{i=1}^N \alpha^* (1 - y_i (\vec{w}^T \vec{x}_i + b))$, $\alpha^* \geq 0 \quad \forall i$

$$\alpha^*(\alpha^*) = \inf_{w,b} L(w, b, \alpha^*) \rightarrow \text{Dual prob.}$$

Step 2: To find infimum of L^* ,

$$\frac{\partial L}{\partial w} = 0, \quad \frac{\partial L}{\partial b} = 0$$

$$\frac{\partial L}{\partial w} = w - \sum_{i=1}^N d_i^o y_i^o x_i^o = 0 \quad \& \quad \frac{\partial L}{\partial b} = b \sum_{i=1}^N d_i^o = 0$$

$$w^* = \sum_{i=1}^N d_i^o y_i^o x_i^o \quad \text{and} \quad b \sum_{i=1}^N d_i^o = 0 \Rightarrow \sum_{i=1}^N d_i^o y_i^o = 0$$

Step 3:

Substitute in $L^*(x_i)$

$$L^*(d_i) = h(w^*, b^*, d_i)$$

$$h(w, b, d_i) = \frac{1}{2} w^T w + \sum_{j=1}^N d_j^o (1 - y_j^o (w^T x_j + b))$$

$$h(w, b, x_i) = \frac{1}{2} w^T w + \sum_{j=1}^N d_j^o - \sum_{j=1}^N d_j^o y_j^o w^T x_j + \sum_{j=1}^N d_j^o b$$

- $\frac{1}{2} w^{*T} w^* = \frac{1}{2} \left(\sum_{i=1}^N d_i^o y_i^o x_i^o \right)^T \left(\sum_{i=1}^N d_i^o y_i^o x_i^o \right) = \sum_{i=1}^N \sum_{j=1}^N d_i^o d_j^o y_i^o y_j^o x_i^T x_j$

- $\sum_{i=1}^N d_i^o y_i^o w^T x_i^o = w^T \sum_{i=1}^N d_i^o y_i^o x_i^o = \left(\sum_{i=1}^N d_i^o y_i^o x_i^o \right)^T \left(\sum_{j=1}^N d_j^o y_j^o x_j^o \right)$
 $= \sum_{i=1}^N \sum_{j=1}^N d_i^o d_j^o y_i^o y_j^o x_i^T x_j$

- $\sum_{i=1}^N b d_i^o = 0$

Hence $L^*(x_i) = \sum_{i=1}^N d_i^o - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N d_i^o d_j^o y_i^o y_j^o x_i^T x_j$

\therefore Our problem is reduced to:

$$L^*(\sum d_i) = \inf_{w, b} h(w, b, \sum d_i) = \sum_{i=1}^N d_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N d_i d_j y_i y_j x_i^T x_j$$

subject to $\sum_{i=1}^N d_i y_i = 0$

$$L^*(\sum d_i) = \sum_{i=1}^N d_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N d_i d_j y_i y_j x_i^T x_j$$

$y_i^o = \text{scalar}$

$x_i^o = \text{vector}$

SUMMARY ↴

Step 4:

Primal Prob.

$$\min_{w,b} \frac{1}{2} w^T w, \text{ subject to}$$

$$y_i(w^T x_i + b) \geq 1 \quad \forall i$$

Maximizing over w, b

Dual Prob.

$$\max_{\sum d_i = 1} \sum_{i=1}^m d_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N d_i d_j y_i y_j x_i^T x_j$$

$$\text{s.t. } \sum_{i=1}^m d_i = 1$$

Max over d_i

C dimension is red.

Recall $w^* = \sum_{i=1}^N x_i y_i x_i^*$ But $\tau = \frac{2}{\|w\|}$ τ is determined only by the support vectors.

But only vectors contributing to the margin are the support vectors.

It is necessary that $d_i = 0$ for the non support vectors.

$\therefore d_i = 0$ for non support vectors & dual problem simplifies.

$$b^* = y_i - w^{*T} x_i^*$$

$$\max_d \sum_{i=1}^m d_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N d_i d_j y_i y_j x_i^T x_j = \min_w \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N d_i d_j y_i y_j x_i^T x_j - \sum_{i=1}^N d_i$$

$$\begin{bmatrix} x_1^T x_1 & x_1^T x_2 \\ x_2^T x_1 & x_2^T x_2 \end{bmatrix}$$

$$Q = [Q_{ij}]_{n \times n}$$

$$[Q_{ij}]_{n \times n} = [y_i y_j^T]_{n \times n} \odot \begin{bmatrix} x_1^T x_1 & \dots & x_1^T x_n \\ x_2^T x_1 & \dots & x_2^T x_n \\ \vdots & & \vdots \\ x_n^T x_1 & \dots & x_n^T x_n \end{bmatrix}$$

$$y y^T = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \begin{bmatrix} y_1 & y_2 \end{bmatrix} = \begin{bmatrix} y_1 y_1 & y_1 y_2 \\ y_2 y_1 & y_2 y_2 \end{bmatrix}$$

Inner product b/w $x_i^T x_j$

$$\therefore y y^T \odot \begin{bmatrix} x_1^T x_1 & x_1^T x_2 \\ x_2^T x_1 & x_2^T x_2 \end{bmatrix} = \begin{bmatrix} y_1 y_1 x_1^T x_1 & y_1 y_2 x_1^T x_2 \\ y_2 y_1 x_2^T x_1 & y_2 y_2 x_2^T x_2 \end{bmatrix}$$

$$\begin{bmatrix} d_1 & d_2 \end{bmatrix} \begin{bmatrix} y_1 y_1 x_1^T x_1 & y_1 y_2 x_1^T x_2 \\ y_2 y_1 x_2^T x_1 & y_2 y_2 x_2^T x_2 \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \end{bmatrix} = \begin{bmatrix} x_1 y_1 y_1 x_1^T x_1 + d_1 y_2 y_1 x_1^T x_1 & x_1 y_1 y_2 x_1^T x_2 + d_2 y_2 y_1 x_1^T x_2 \\ x_2 y_1 y_1 x_2^T x_1 & x_2 y_1 y_2 x_2^T x_2 \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \end{bmatrix}$$

$$= \alpha_1^2 y_1 y_1 x_1^T x_1 + \alpha_1 \alpha_2 y_1 y_2 x_1^T x_2 + \alpha_2 \alpha_1 y_2 y_1 x_2^T x_1 + \alpha_2^2 y_2 y_2 x_2^T x_2 = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j$$

then $\sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j = \mathbf{d}^T Q \mathbf{d}$

- linear SVMs are convex optimization \Rightarrow strong duality
- the unique solution, i.e. $(\mathbf{w}^*, b^*, \{\alpha_i^*\})$, is a **saddle point**

$$\underbrace{\text{SVM1}}_{\text{prime problem}} \iff \underbrace{\max_{\{\alpha_i\}} L^*(\alpha_i)}_{\text{dual problem}} \text{ s.t. } \sum_{i=1}^N \alpha_i y_i = 0$$

- solving the **dual problem** yields $\{\alpha_i^*\}$
- **linear SVMs:** $\mathbf{w}^* = \sum_{i=1}^N \alpha_i^* y_i \mathbf{x}_i$ and $b^* = y_i - \mathbf{w}^{*T} \mathbf{x}_i$

\therefore Dual prob. can be eqn. written as:

$$\max_{\alpha} \mathbf{1}^T \alpha - \frac{1}{2} \alpha^T Q \alpha$$

subject to $y^T \alpha = 0$
 $\alpha \geq 0$

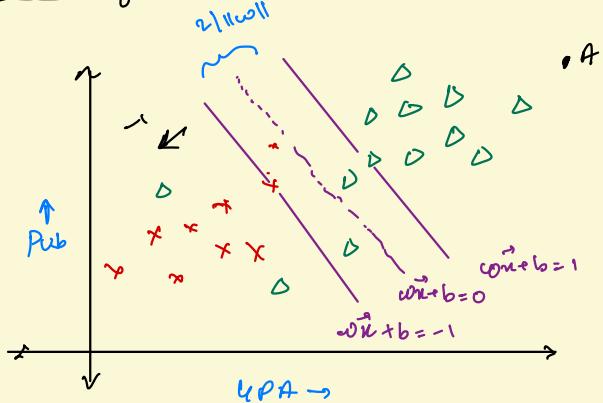
Insp. obs: For the optimization we just care about the inner product among the data.

$$Q = [Q_{ij}]_{n \times n}$$

$$[Q_{ij}]_{n \times n} = [y y^T]_{n \times n} \odot \begin{bmatrix} x_1^T x_1 & \dots & x_1^T x_n \\ x_2^T x_1 & \dots & x_2^T x_n \\ \vdots & \ddots & \vdots \\ x_n^T x_1 & \dots & x_n^T x_n \end{bmatrix}$$

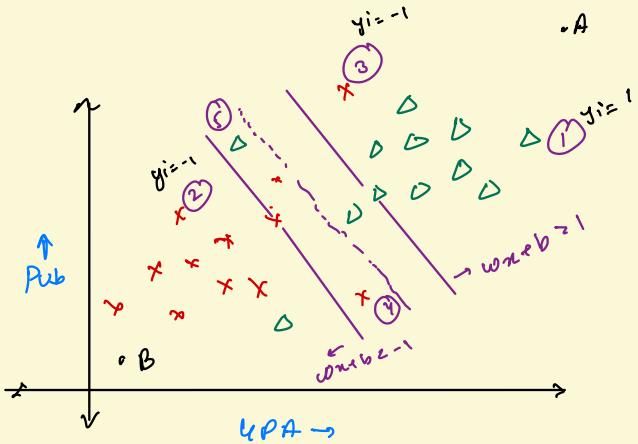
Involves only
inner product

② Soft Margin SVMs:



cannot use hard margin SVM to separate

$$\text{Hinge loss: } \max(0, 1 - y_i(w \cdot x_i + b))$$



Score @ pt A: $\gg 1$

Score @ pt B: $\ll 1$

\therefore for correct

classifications, $y_i^*(w \cdot x_i + b) \geq 0$

$\therefore \text{Score} = 0$.

Spirit of Soft Margin SVM:

Classify the data as well as possible but also willing to accept certain level of mistakes with penalties based on how big the mistake is.

$$\text{Score: } w \cdot x + b$$

If the obs. lies in b/w $w \cdot x + b = 0$ and $w \cdot x + b = 1$

$$0 < \text{Score} \leq 1 \dots$$

If obs. lies in b/w $w \cdot x + b = 0$ & $w \cdot x + b = -1$, $-1 < \text{Score} \leq 0$.

@ point A score is $\gg 1$

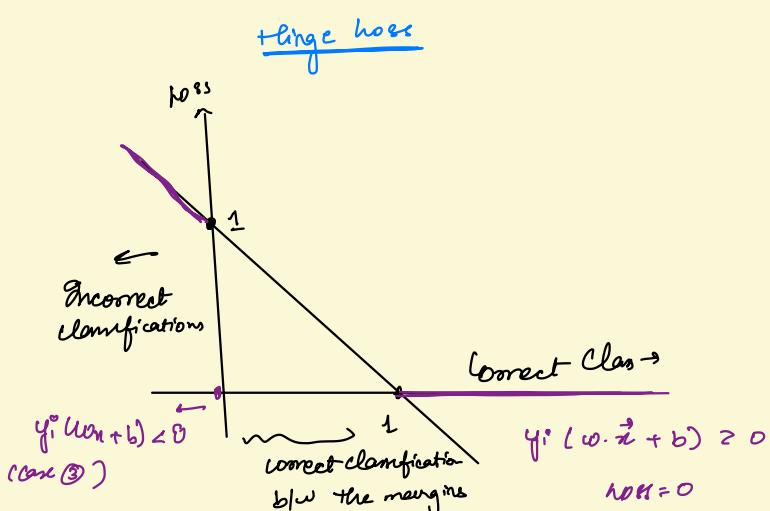
$$\textcircled{1} \max(0, 1 - 1(\geq 1)) = 0$$

$$\textcircled{2} \max(0, 1 + 1(\leq -1)) = 0$$

$$\textcircled{3} \max(0, 1 + 1(\geq 1)) = \geq 2$$

$$\textcircled{4} \max(0, 1 + 1(w \cdot x + b < -1)) = \text{something b/w 0 and 1}$$

$$\textcircled{5} \max(0, 1 - 1(w \cdot x + b < -1)) = \text{something b/w 1 & 2}$$



* soft margin SVM Mathematical formulation:

2 things: ① maximize the margin: minimize $\|w\|$

② minimize the hinge loss.

C = parameter that decides the tradeoff b/w max. margin or min. hinge loss.

C small \Rightarrow hinge loss given less importance, C large \Rightarrow focus on max. margin

$$\min_{w, b, \xi_i^0} \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i^0, \quad \xi_i^0 \xrightarrow{\text{hinge loss}} \max(0, 1 - y_i(w^T x_i + b))$$

Equivalently:

$$\begin{aligned} & \text{minimize}_{w, b, \xi_i^0} \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i^0 \\ & \xi_i^0 \geq 0 \end{aligned}$$

$$\xi_i^0 \geq 1 - y_i(w^T x_i + b) \quad \forall i \in \{1, \dots, N\}$$

$$\text{Derivation: } \left. \begin{array}{l} \text{Constraints:} \\ -\xi_i \leq 0 \\ 1 - y_i(w^T x_i + b) - \xi_i \leq 0 \end{array} \right\}$$

Step-1: Dual prob: for $d_j, \beta_j \geq 0$,

$$L(w, b, \xi_i, d_i, \beta_i) = \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N d_i \xi_i + \sum_{i=1}^N d_i - \sum_{i=1}^N \alpha_i y_i (w^T x_i + b) - \sum_{i=1}^N \beta_i \xi_i$$

$$\underline{\text{Step-2}} \quad \text{i)} \quad \frac{\partial L}{\partial w} = 0 \quad \Rightarrow \quad w - \sum_{i=1}^N \alpha_i y_i x_i = 0 \quad \Rightarrow \quad w^* = \sum_{i=1}^N \alpha_i y_i x_i$$

$$\text{ii)} \quad \frac{\partial L}{\partial \xi_j} = 0 \quad \forall j = 1, 2, \dots, N \quad \Rightarrow \quad C - \alpha_j - \beta_j = 0 \quad \Rightarrow \quad C = \alpha_j + \beta_j$$

$$\text{iii)} \quad \frac{\partial L}{\partial b} = 0 \quad \Rightarrow \quad - \sum_{i=1}^N \alpha_i y_i = 0 \quad \Rightarrow \quad \sum_{i=1}^N \alpha_i y_i = 0$$

$$b_i \geq 0 \Rightarrow c - \alpha_i \geq 0 \quad \forall i \in \{1, 2, \dots, N\}$$

$$C \geq \alpha_i \geq 0 \quad \forall i \in \{1, 2, \dots, N\}$$

Step 3 Substitute in Lagrangian

$$L(w^*, b^*, \xi_i, \alpha_i, \beta_i) = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j + \underbrace{\sum_{i=1}^N (c - \alpha_i - \beta_i)}_0 \xi_i + \sum_{i=1}^N \alpha_i - \sum_{i=1}^N \alpha_i y_i (w^T x_i + b)$$

$$L(w^*, b^*, \xi_i, \alpha_i, \beta_i) = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum_{i=1}^N \alpha_i (1 - y_i (w^T x_i + b))$$

$$\begin{aligned} \sum_{i=1}^N \alpha_i y_i w^T x_i &= w^T \sum_{i=1}^N \alpha_i y_i x_i = (\sum_{i=1}^N \alpha_i y_i x_i)^T (\sum_{j=1}^N \alpha_j y_j x_j) \\ &= \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j \quad \text{and} \quad \sum_{i=1}^N \alpha_i y_i b = 0 \end{aligned}$$

$$\therefore L(w^*, b^*, \xi_i, \alpha_i, \beta_i) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j$$

Dual Prob. : $L(w^*, b^*, \xi_i, \alpha_i, \beta_i) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j$ subject to

$$\sum_{i=1}^N \alpha_i y_i = 0 \quad \text{and} \quad 0 \leq \alpha_i \leq C$$

• Our Dual prob. can be eqn. written as:

$$\max_{\alpha} \alpha^T \alpha - \frac{1}{2} \alpha^T Q \alpha$$

$$\text{subject to } y^T \alpha = 0$$

$$0 \leq \alpha \leq C$$

Imp. obs: For the optimization we just care about the inner product among the data.

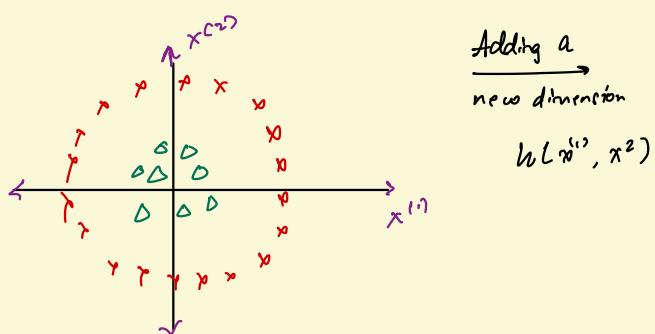
$$Q = [Q_{ij}]_{n \times n}$$

$$[Q_{ij}]_{n \times n} = [y y^T]_{n \times n} \odot \begin{bmatrix} x_1^T x_1 & \dots & x_1^T x_n \\ x_2^T x_1 & \dots & x_2^T x_n \\ \vdots & \ddots & \vdots \\ x_n^T x_1 & \dots & x_n^T x_n \end{bmatrix}$$

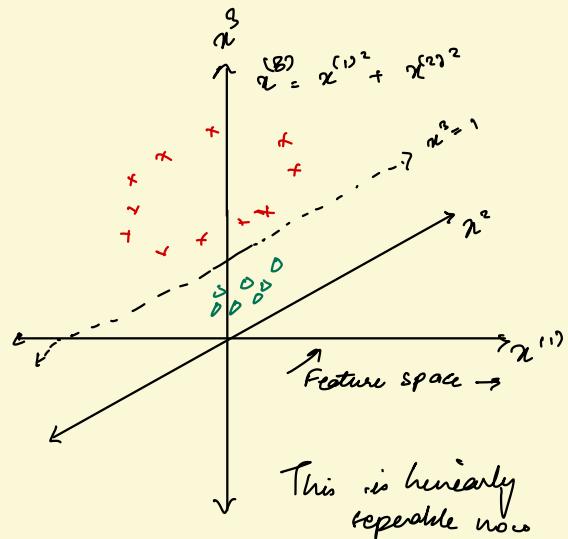
Involves only
inner product

* PROJECTED GRADIENT DESCENT (For the constraint $0 \leq \alpha \leq C$)

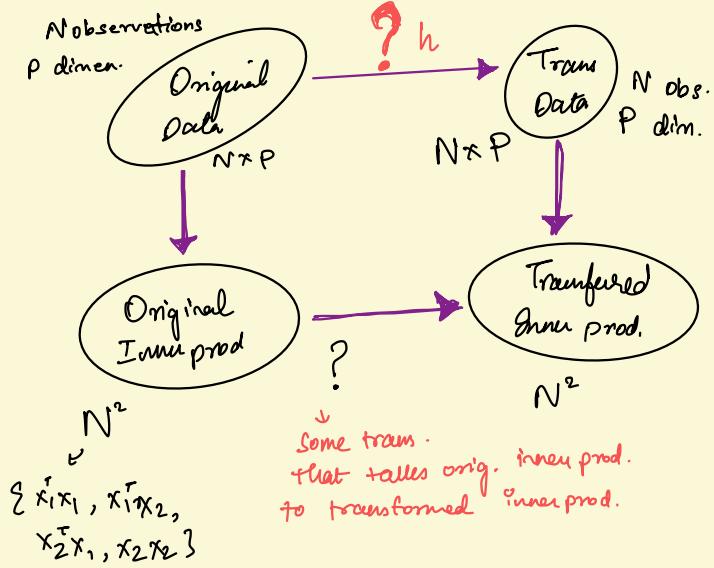
Non-linear SVMs:



We can't separate this by a hyperplane.



Idea: Project the data into a higher dim space so that it's linearly separable in the higher dimension



Recall: For the dual stage in SVM, we just cared about the inner products b/w the data. (for the Dual prob)

$$x_i \cdot x_j \rightarrow h(x_i) \cdot h(x_j)$$

Aim: We just need the inner products of the transformation. For example

for an input point x^0 ,

$$\begin{bmatrix} x_i^{(1)} \\ x_i^{(2)} \\ \vdots \\ x_i^{(P)} \end{bmatrix} = \begin{bmatrix} 1/2 \\ x_i^1 \\ x_i^2 \\ \vdots \\ x_i^P \\ x_i^1 x_i^1 \\ x_i^1 x_i^2 \\ \vdots \\ x_i^P x_i^P \\ x_i^1 x_i^P \end{bmatrix}$$

all possible
→ 2-way interactions.

Example of
projection in
7 dc

$N \times P$ matrix in memory is computationally very expensive than the $N \times P$. Hence we need not add new data, instead we can go for the dual problem calculate N^2 & carry on.

Aim: To need $h(x_i) \cdot h(x_j) =$

$$\begin{bmatrix} 1/2 \\ x_i^1 \\ x_i^2 \\ \vdots \\ x_i^P \\ x_i^1 x_i^1 \\ x_i^1 x_i^2 \\ \vdots \\ x_i^P x_i^P \\ x_i^1 x_i^P \end{bmatrix} \begin{bmatrix} 1/2 \\ x_j^1 \\ x_j^2 \\ \vdots \\ x_j^P \\ x_j^1 x_j^1 \\ x_j^1 x_j^2 \\ \vdots \\ x_j^P x_j^P \\ x_j^1 x_j^P \end{bmatrix} = \frac{1}{4} + x_i^1 x_j^1 + x_i^2 x_j^2 + (x_i^P x_j^P)^2 + (x_i^1 x_i^2 x_j^1 x_j^2) + (x_i^P x_j^P)^2 + (x_i^P x_i^1 x_j^1 x_j^P)$$

We want to get the above terms without explicitly computing the transformation.

Let Polynomial Kernel $\phi(x_i \cdot x_j) = (1 + x_i \cdot x_j)^2 = (\frac{1}{2} + x_i^1 x_j^1 + x_i^2 x_j^2)^2$

$$(\frac{1}{2} + x_i^1 x_j^1 + x_i^2 x_j^2)(\frac{1}{2} + x_i^1 x_j^1 + x_i^2 x_j^2) =$$

$$\frac{1}{4} + \frac{1}{2} x_i^1 x_j^1 + \frac{1}{2} x_i^2 x_j^2 + \frac{1}{2} x_i^1 x_j^1 + (x_i^1 x_j^1)^2 + (x_i^2 x_j^2 x_i^1 x_j^1) + \frac{1}{2} x_i^2 x_j^2 + (x_i^2 x_j^2 x_i^1 x_j^1) + (x_i^2 x_j^2)^2$$

$$= \frac{1}{4} + x_i^1 x_j^1 + x_i^2 x_j^2 + (x_i^1 x_j^1)^2 + (x_i^2 x_j^2)^2 + (x_i^2 x_j^2 x_i^1 x_j^1) + (x_i^2 x_j^2 x_i^1 x_j^1)$$

Polynomial Kernel

$$h: \mathbb{R}^m \rightarrow \mathbb{R}^n$$

hence we found

$$\phi(x_i, x_j).$$

$\phi: \text{Original IP} \rightarrow \text{Transformed IP.}$

$$\mathbb{R}^m + \mathbb{R}^m \rightarrow \mathbb{R}$$

* Polynomial Kernel: $\phi(x_i, x_j) = (x_i^T x_j + 1)^p$ $p = \text{dimension}$

↳ Overall picture:

① Introduce an additional dimension so that the non linear data is separable in the higher dimension. $h(x^{(1)}, x^{(2)}) - \text{obj could be any no. of dimensions.}$ [Not sure how to choose this dimension]

② Find a hyperplane in the higher dimension.
i.e. finding an optimal w, b in the higher dimension. How do we do it?

If it's linearly separable we can find using margin in SVM in the higher dim:

Primal Prob.

$$\min_{w, b} \frac{1}{2} w^T w, \text{ subject to}$$

$$y_i^p (w^T h(x_i) + b) \geq 1 \quad \forall i$$

minimizing over w, b

Dual Prob.

$$\max_{d \in \mathbb{R}^N} \sum_{i=1}^N d_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N d_i d_j y_i y_j h(x_i)^T h(x_j)$$

$$\text{s.t. } \sum_{i=1}^N d_i y_i = 0$$

max over d_i

(dimension is increased)

We use the dual problem. So here we just need to compute $h(x_i)^T h(x_j)$ which could be found by using Kernel function.

③ After finding w, b project it onto the current dimension
We have found the non-linear hyperplane.

(I'm not sure
how to do
this but
this is the
idea.)

RADIAL BASIS KERNEL:

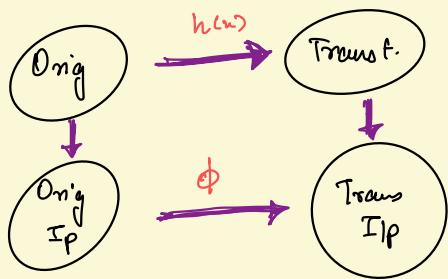
$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ \vdots \\ x_1^3 \\ x_2^2 \\ x_1^3 x_2^2 \\ \vdots \\ \infty \end{bmatrix}$$

$$\phi(x_i, x_j) = e^{-\frac{1}{2} \|x_i - x_j\|^2}$$

Wishes to consider all possible interactions b/w the variables.

It's Infinite dimensional.

How to take the inner product of this infinite dimensional vectors?



$$\begin{aligned} \phi(x_i, x_j) &= e^{-\frac{1}{2} \|x_i - x_j\|^2} \\ &= e^{-\frac{1}{2} ((x_i - x_j)^T (x_i - x_j))} \\ &= e^{-\frac{1}{2} (x_i^T x_i - 2x_i^T x_j + x_j^T x_j)} \\ &= e^{-\frac{1}{2} (x_i^T x_i + x_j^T x_j)} e^{x_i^T x_j} \\ &= C' e^{x_i^T x_j} \end{aligned}$$

$$\begin{aligned} \phi(x_i, x_j) &= C' e^{x_i^T x_j + 1} e^{-1} \\ &= C' e^{x_i^T x_j + 1} = C' \sum_{n=1}^{\infty} \left(1 + \frac{x_i^T x_j}{n}\right)^n = C' \sum_{n=1}^{\infty} \frac{\Phi_{\text{poly}}(x_i^T x_j)}{n!} \end{aligned}$$

This captures the no. of terms involved in the inner product

RBF Kernel: $\phi(x_i, x_j) = e^{-\gamma \|x_i - x_j\|^2}$

In summary, Non linear SVMs:

Feature space (higher dimensional space)

- First idea was: Learning a non-linear classifier using SVM
 - Calculate $h(x)$ for each training example – Find a linear SVM in the feature space.
 - Problems:
 - Feature space can be high dimensional or even have infinite dimensions, so storing $h(x)$ can be inefficient or even impossible

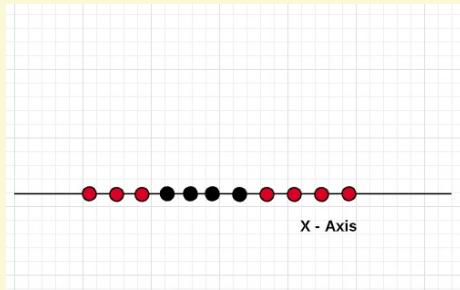
So we came up with the idea of “Kernels”:

- Kernels are similarity functions that return inner products between the images $h(x_i)$ of data points.
- Choosing kernel (ϕ) is equivalent to choosing $h(x)$ in the feature space .

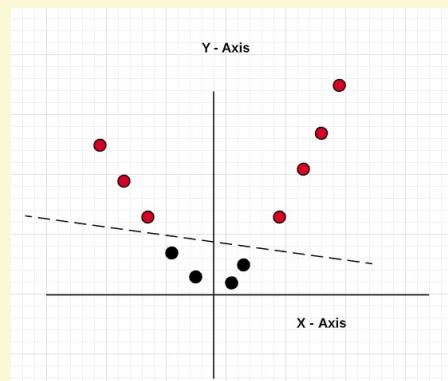
$$\begin{aligned} \text{- } \phi: x \times x \rightarrow \mathbb{R} \\ \phi(x_i, x_j) = h(x_i)^\top h(x_j) \quad (\text{inner product b/w } h(x_i) \text{ & } h(x_j)) \end{aligned}$$

The kernel trick

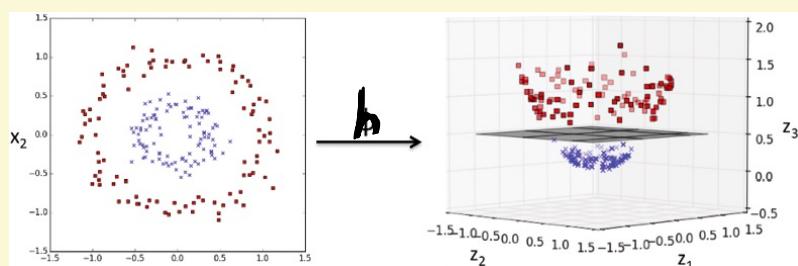
- No need to know what $h(x)$ is and what the feature space is.
- No need to explicitly map the data to the feature space. (memory constraint resolved)
- Define a kernel function ϕ and replace the dot product $\langle x, z \rangle$ with a kernel function $\phi(x, z)$ in both training and testing.



Linearly inseparable data in : one-dimension



Applying kernel method to represent data using 2-dimensions



linear

Maximize

$$L(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle \vec{x}_i, \vec{x}_j \rangle$$

Subject to $\alpha_i \geq 0$ and $\sum_i \alpha_i y_i = 0$

Non linear case

$$L(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \phi(\vec{x}_i, \vec{x}_j)$$

So problem is resolved without actually going to higher dimension, \therefore we needn't actually store $h(x)$. Just replace the orig. ip by ip in higher dim.

Equivalently

$$\max_{\alpha} \vec{1}^T \alpha - \frac{1}{2} \alpha^T Q \alpha$$

Sub. to $\vec{y}^T \alpha = 0$ (Hard SVM)
 $\alpha \geq 0$

$$Q = [Q_{ij}]_{n \times n}$$

$$[Q_{ij}]_{n \times n} = [\vec{y} \vec{y}^T]_{n \times n} \odot \begin{bmatrix} \vec{x}_1^T \vec{x}_1 & \dots & \vec{x}_1^T \vec{x}_n \\ \vdots & \ddots & \vdots \\ \vec{x}_n^T \vec{x}_1 & \dots & \vec{x}_n^T \vec{x}_n \end{bmatrix}$$

$$\text{Plane: } w^* = \sum_{i=1}^N d_i^* y_i \vec{x}_i$$

$$b^* = y_p - w^{*T} \vec{x}_p$$

$$\text{Plane: } w^* \vec{x} + b^*$$

$$= \sum_i \alpha_i y_i \langle \vec{x}_i, \vec{x} \rangle + b$$

Equivalently

$$\max_{\alpha} \vec{1}^T \alpha - \frac{1}{2} \alpha^T Q \alpha$$

Sub. to $\vec{y}^T \alpha = 0$
 $\alpha \geq 0$

$$\Phi = [\Phi_{ij}]_{n \times n}$$

$$[\Phi_{ij}]_{n \times n} = [\vec{y} \vec{y}^T]_{n \times n} \odot \begin{bmatrix} h(\vec{x}_1)^T h(\vec{x}_1) & \dots & h(\vec{x}_1)^T h(\vec{x}_n) \\ \vdots & \ddots & \vdots \\ h(\vec{x}_n)^T h(\vec{x}_1) & \dots & h(\vec{x}_n)^T h(\vec{x}_n) \end{bmatrix}$$

$$\text{Plane: } w^* = \sum_{i=1}^N d_i^* y_i h(\vec{x}_i)$$

$$b^* = y_p - w^{*T} \vec{x}_p$$

$$w^* \vec{x} + b^*$$

$$f(\vec{x}) = \sum_i \alpha_i y_i \phi(\vec{x}_i, \vec{x}) + b$$

$$\phi(\vec{x}_i, \vec{x}) = h(\vec{x}_i)^T h(\vec{x})$$

Summary so far

- Find the hyperplane that maximizes the margin.
- Introduce soft margin to deal with noisy data
- Implicitly map the data to a higher dimensional space to deal with non-linear problems.
- The kernel trick allows infinite number of features and efficient computation of the dot product in the feature space.
- The choice of the kernel function is important.

Common kernel functions

- Linear: $K(\vec{x}, \vec{z}) = \langle \vec{x}, \vec{z} \rangle$
- Polynomial: $K(\vec{x}, \vec{z}) = (\gamma \langle \vec{x}, \vec{z} \rangle + c)^d$
- Radial basis function (RBF): $K(\vec{x}, \vec{z}) = e^{-\gamma(||\vec{x}-\vec{z}||)^2}$

$$\vec{x} - \vec{z} = \begin{bmatrix} x_1 - z_1 \\ \vdots \\ x_N - z_N \end{bmatrix}$$

$$\|\vec{x} - \vec{z}\| = \sqrt{(x_1 - z_1)^2 + \dots + (x_N - z_N)^2}$$

Example

Linear Kernel: $h(x) = \alpha x$ for some α

$$h(x_1, x_2, \dots, x_N) = [\alpha x_1, \alpha x_2, \dots, \alpha x_N]$$

$$\begin{aligned} \phi(x, z) &= (h(x), h(z)) = (h(x))^T h(z) \\ &= [\alpha x_1, \dots, \alpha x_N] \begin{bmatrix} \alpha z_1 \\ \vdots \\ \alpha z_N \end{bmatrix} = \alpha^2 x_1 z_1 + \alpha^2 x_2 z_2 + \dots + \alpha^2 x_N z_N \end{aligned}$$

The Kernel Matrix (a.k.a. the Gram matrix)

$\phi(1,1)$	$\phi(1,2)$	$\phi(1,3)$...	$\phi(1,m)$
$\phi(2,1)$	$\phi(2,2)$	$\phi(2,3)$...	$\phi(2,m)$
...				
...				
$\phi(m,1)$	$\phi(m,2)$	$\phi(m,3)$...	$\phi(m,m)$

$\phi(i, j)$ means $\phi(x_i, x_j)$,

where x_i means the i-th training instance.

$$\phi(x_i, x_j) = (h(x_i))^T (h(x_j))$$

Just storing all possible innerproducts in the form of a matrix

Mercer's Theorem

- The kernel matrix is symmetric positive definite.
- Any symmetric, positive definite matrix can be regarded as a kernel matrix; that is, there exists a ϕ such that $\phi(x, z) = \langle h(x), h(z) \rangle$

SVM Learning Procedure (in a nutshell)

Given a training set as $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$

1. choose a kernel function $\Phi(\mathbf{x}_i, \mathbf{x}_j)$
2. build the matrices \mathbf{Q} , \mathbf{y} and $\mathbf{1}$ from \mathcal{D} and $\Phi(\mathbf{x}_i, \mathbf{x}_j)$
3. solve the quadratic programming problem to get:

$$\alpha^* = [\alpha_1^*, \dots, \alpha_N^*]^\top \text{ and } b^*$$

4. evaluate the learned model as:

$$y = \text{sign}\left(\sum_{i=1}^N \alpha_i^* y_i \Phi(\mathbf{x}_i, \mathbf{x}) + b^*\right)$$

↳ Multiclass SVMs

What if there are more than 2 classes for classification?

Assignment - 2

SVMs

$$l_i = \max(0, 1 - (y_i(\omega \cdot x_i + b)))$$

$$J = \frac{1}{2} \|\omega\|^2 + \frac{1}{n} \sum_{i=1}^n l_i = \frac{1}{2} \omega^\top \omega + \frac{1}{n} \sum_{i=1}^n l_i$$

where
 $y_i f(x_i) \geq 1$
 $l_i = 0$

$$\frac{\partial J}{\partial w_j} = w_j - \frac{1}{n} \left(\sum_{y_i f(x_i) < 1} y_i x_j \right)$$

$$\frac{\partial J}{\partial b} = -\frac{1}{n} \sum_{y_i f(x_i) < 1} y_i$$

$$y = \omega \cdot x + b$$

Plane:
 $y = \omega \cdot x + b$