# Support Vector Machines

◇ __Geometric-Intuition__



X +ve pts
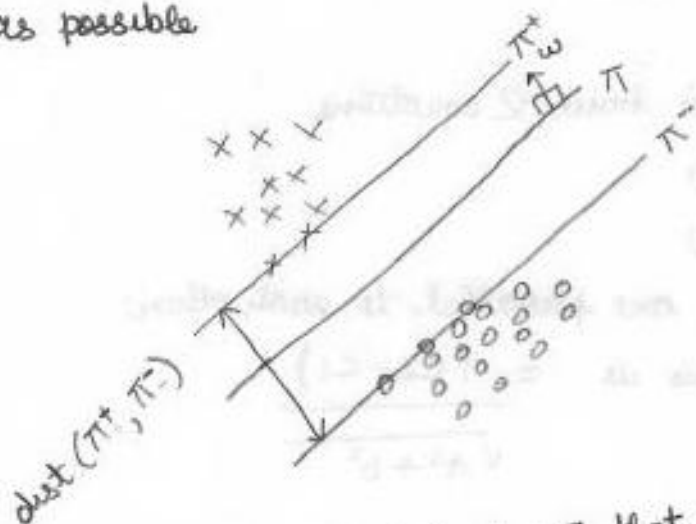O -ve pts

many $\Pi$s that seperates +ve pts. and -ve pts.
A point that are very close to hyperplane would be low.

▽ Key idea of Support Vector Machine
  $\Pi$ seperates the hyperplanes from -ve and +ve pts as widely as possible



$\Pi$ : margin maximizing hyperplane

$\Pi^+$ and $\Pi^-$ are both parallel to $\Pi$.

SVM : Try to find a $\Pi$ that maximizes the margin = dist $(\Pi^+, \Pi^-)$

If my margin is high chance of misclassification decreases.
Points through which $\Pi^+$ and $\Pi^-$ goes through are called __support vectors.__
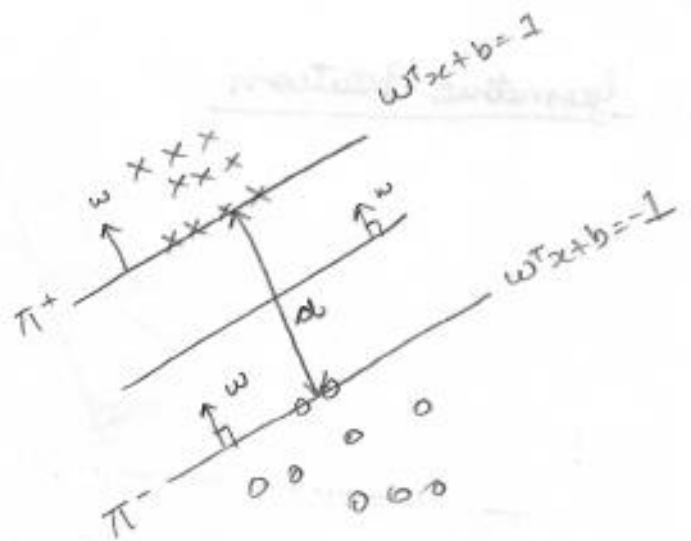
## ✧ Mathematical derivation

$\pi :-$ margin-maximization

$\pi : w^T x + b = 0$

$\pi^+ : w^T x + b = 1$

$\pi^- : w^T x + b = -1$

margin: $d = \dfrac{2}{\|w\|}$

Objective :- $(w^*, b^*) = \underset{w, b}{argmax} \dfrac{2}{\|w\|}$

s.t $(w^*, b^*) = \underset{w, b}{argmax} \dfrac{2}{\|w\|} = margin.$

Proof :

for example, if we have 2 equations

$ax + by + c1 = 0$

$ax + by + c2 = 0$

here both lines are parallel to each other

So distance between lines is $= \dfrac{|c2 - c1|}{\sqrt{a^2 + b^2}}$

now we are using 2 hyperplanes equations
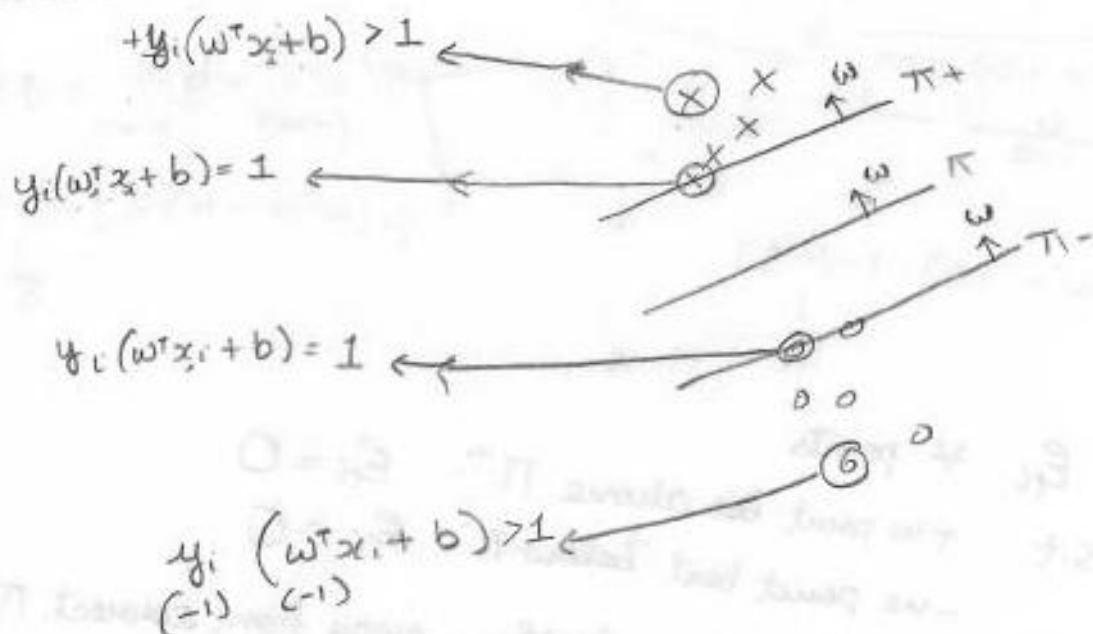
$w^T x + b = 1$

$w^T x + b = -1$

$w^T x + b - 1 = 0$

$w^T x + b + 1 = 0$

Comparing eq of plane to eq. of lines

distance between 2 planes $= \dfrac{|\text{difference of vector components}|}{||w||}$

diff of vector components $= |w^T x + b - 1 - w^T x - b - 1|$

$$= |-2|$$

$$\boxed{\text{distance between 2 planes} = \dfrac{2}{||w||}}$$

$+y_i(w^T x_i + b) > 1$

$y_i(w^T x_i + b) = 1$

$y_i(w^T x_i + b) = 1$

$\underset{(-1)}{y_i} \underset{(-1)}{(w^T x_i + b)} > 1$

$(w^*, b^*) = \underset{w, b}{\text{argmax}} \dfrac{2}{||w||} = \text{margin}$

$\left.\begin{array}{c}\end{array}\right\}$ Constr optimization problem of SVM

s.t $y_i(w^T x_i + b) \geq 1$ for all $x_i$

works when data is linearly separable

$\underset{+ve}{(w^T x_i + b)} \underset{-ve}{(-y_i)} < -1 \longleftarrow$

$(w^T x_i + b)(y_i) < -1$

$\underset{-ve}{\phantom{(}} \underset{+ve}{\phantom{)}}$

No errors allowed
This is known as hard margin SVM
1) no point between the margins.
2) no -ve point on +ve side of hyperplane and vise versa as these are not two separable by hard margin SVM.

$y_i(w^T x_i + b) = 1.5$

$y_i(w^T x_i + b) = 0.5 = 1 - 0.5$

$y_i(w^T x_i + b) = 0.5$
↓
-ve
↓
+ve

$y_i(w^T x_i + b) = 1 - (1.5)$
|
$\xi_i$

$y_i(w^T x_i + b) = -1.5 \rightarrow$ (assume)

$y_i(w^T x_i + b) = 1 - (2.5)$
|
$\xi_i$

$\xi_i \ \forall$ points

s.t    +ve point lies above $\Pi^+$   $\xi_i = 0$
       -ve point lies below $\Pi^-$   $\xi_i = 0$

$\Rightarrow \xi_i \uparrow$, pt. is farther away from correct $\Pi$ in incorrect direction.

$\xi_i = 0$   if $y_i(w^T x_i + b) \geq 1$

$\xi_i > 0$   else  it is equal to the some units of dist away from correct hyperplane in incorrect direction

Objective:-

$(w^*, b^*) = \underset{w, b}{\arg \max} \dfrac{2}{\|w\|} = \underset{w, b}{\arg \min} \dfrac{\|w\|}{2}$

$\xi_i$'s

$$(w^*, b^*) = \underset{w,b}{\text{argmin}} \ \frac{\|w\|}{2} + C \cdot \left( \frac{1}{n} \sum_{i=1}^{n} \xi_i \right)$$

→ average distance of misclassified pts from $\pi$'s /loss

$$\Downarrow$$

margin
(regularization)   hyperparameter

s.t $y_i (w^T x_i + b) \geq 1 - \xi_i \quad \forall i$ ] Correctly classified pt $\xi_i = 0$

$$\xi_i \geq 0$$

minimize errors / misclassifications

$$\|$$

$$\min \ \Sigma \ \xi_i$$

$$\cancel{s.t \ y_i (w^T x_i + b) \geq 1 - \xi_i \ \forall i \ \xi_i \geq 0}$$

$C \uparrow$ : giving more importance to make mistakes (reduces) on training data which leads to overfitting (variance)
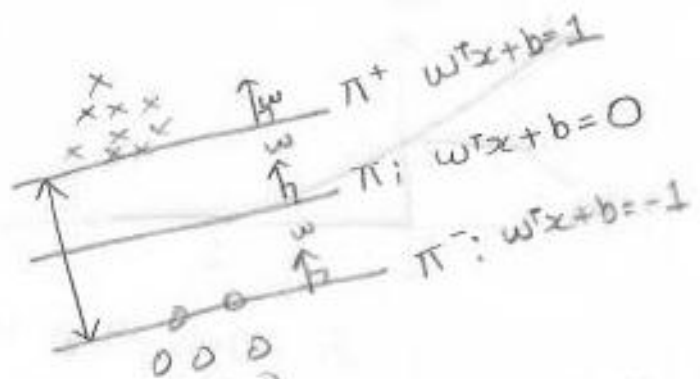
$C \downarrow$ high bias or underfit model.

$$\boxed{(w^*, b^*) = \underset{w,b}{\text{argmin}} \ \frac{\|w\|}{2} + C \cdot \frac{1}{n} \sum_{i=1}^{n} \xi_i}$$ Soft margin SVM

---

Q  SVM:  why we take values +1 and -1 for SVM vector planes for hard margin. SVM?

margin $\div \frac{2}{\|w\|}$

$$w^*, b^* = \underset{w,b}{\text{arg max}} \ \frac{2}{\|w\|}$$



$\pi^+ \ w^T x + b = 1$
$\pi : \ w^T x + b = 0$
$\pi^- : \ w^T x + b = -1$

$\{ \ \|w\| \neq 1$ (any vector) need not to be unit vector.

① $\Pi^+ :- \quad w^T x + b = K$

$\Pi^- :- \quad w^T x + b = -K$ $\qquad K > 0$

We are taking $+k$ and $-k$ as we want 2 hyperplanes $\Pi^+$ and $\Pi^-$ to be equally far away from $\Pi$.

$$\text{margin} = \frac{2K}{||w||}$$

$$\underset{w, b}{\text{argmax}} \; \frac{2}{||w||} = \underset{w, b}{\text{argmax}} \; \frac{2K}{||w||} = \frac{8}{||w||}$$

② $\Pi^+ : \quad w^T x + b = K$ $\qquad w \perp \Pi$

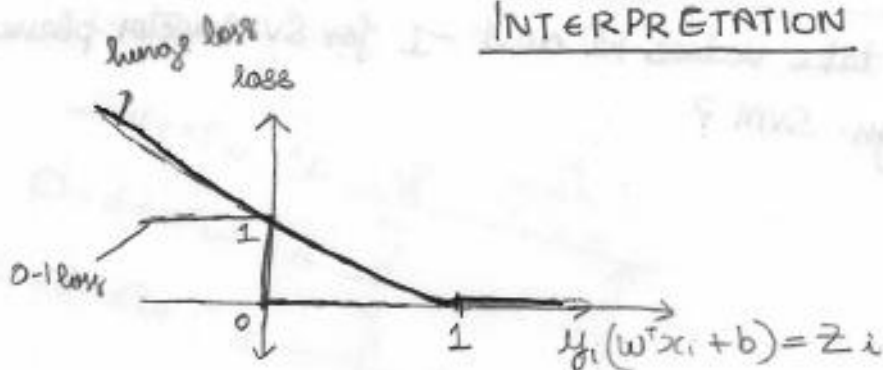$$\left(\frac{w}{K}\right)^T x + \frac{b}{K} = 1$$

$$(w')^T x + b' = 1$$

Reason to take $+1$ and $-1$

is to simple the math.

## Loss Function (Hinge Loss) Based

### Interpretation



$$\begin{cases} z > 0 : x_i \text{ is correctly classified} \\ z_i < 0 : x_i \text{ is incorrectly classified} \end{cases}$$

hinge loss is not differenciable as it is not continous.
approximated 0-1 loss by hinge loss

hinge loss :- $\begin{cases} z_i \geq 1 ; & \text{hinge loss} = 0 \\ z_i < 1 ; & \text{hinge loss} = 1 \end{cases}$
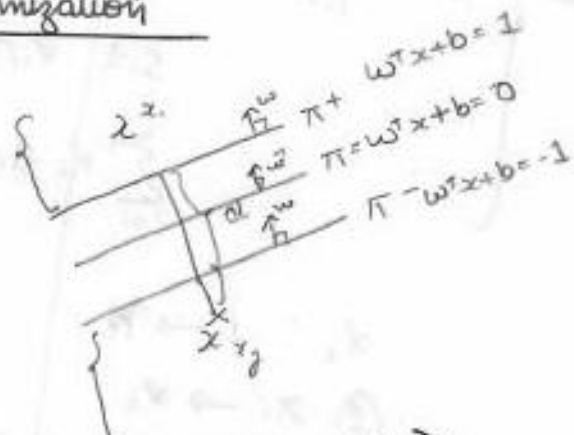
$\quad\quad \longrightarrow \max(0, 1 - z_i)$

Case I: $z_i \geq 1$; $1 - z_i$ is -ve value $\Rightarrow \max(0, 1 - z_i)$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad = 0$

Case II: $z_i < 1$, $1 - z_i$ is +ve value $\Rightarrow \max(0, 1 - z_i)$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad = 1 - z_i$

✧ <u>Geometric Formulation & loss minimization</u>

$\xi_i = 0 \quad \leftarrow x_i$



$d = 1 - y_i(w^T x_i + b)$

<u>soft SVM</u>

$\quad\quad\quad \longrightarrow$ loss (hyperparameter)

$\underset{w, b}{\min} \dfrac{\|w\|}{2} + c \sum_{i=1}^{n} \xi_i$

s.t $(1 - y_i(w^T x_i + b)) \geq \xi_i \; \forall i$

$\quad\quad\quad\quad\quad \xi_i \geq 0$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \longrightarrow$ reg.

<u>loss min</u> $\quad \underset{w, b}{\min} \sum_{i=1}^{n} \max(0, 1 - y_i(w^T x_i + b)) + \lambda \|w\|^2$

$\|w\| \geq 0 \Rightarrow \min \|w\|$ is same as $\min \|w\|^2$

Note → If we multiply hyperparameter with loss function it will
cause overfitting $c\uparrow$ = overfit
$\quad\quad \lambda\uparrow \Rightarrow$ underfit

3) Hard And Soft Margin SVMs
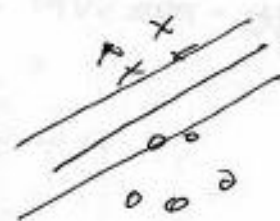
a) Hard Margin SVM

1.) Key idea of SVM maximise the margin

For the supports vectors, we know:

$$y_i * (wx + b) = 1$$

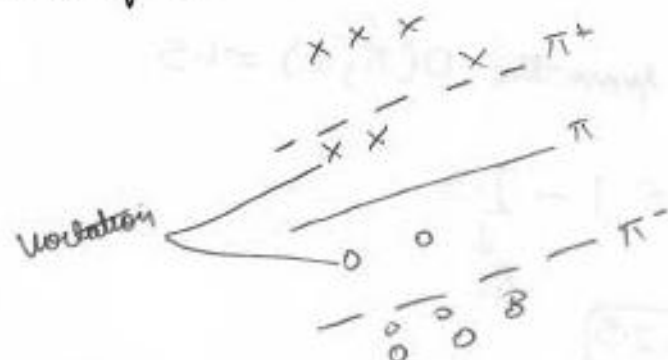for non support vectors in the above image

$$y_i * (wx + b) > 1$$

So optimization problem becomes:

$$\max(w) \left\{ \frac{2}{||w||} \right\} \text{ such that } y^* (wx_i + b) \geq 1$$

As we know the hard margin SVMs are optimal svm for linearly seprable data where +ve pts are above $\pi^+$ and -ve pt are below $\pi^-$. There is no points in the margin area or we can say no points violating the margin.
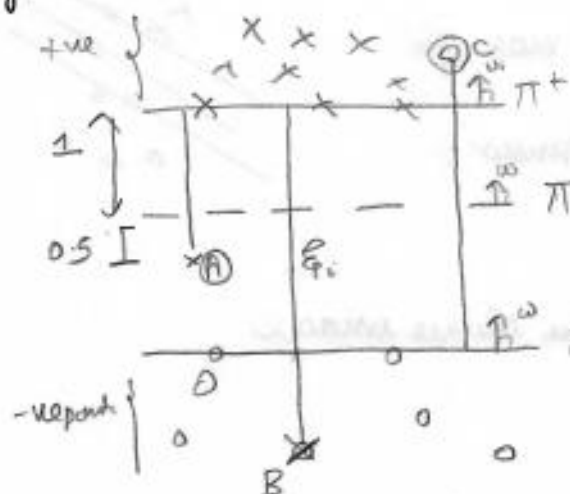
◇ what if pts lie between the margin?



② The solution for this SVM is soft margin SVM?
Soft Margin SVM.

11) **Soft marguin SVM**

When the points are almost linearly seperable we consider soft - max SVM as. the best idea to implement.



$$\max(\omega)\left\{\frac{2}{||\omega||}\right\} \text{ such that}$$

$$y_i(\omega x_i + b) \geq 1$$

**For point A :** The distance from the $\pi$, $D\{\pi, A\} = -0.5$

we can write this distance as $D\{\pi, A\} = 1 - (1.5)$

lets call this $1.5$ as zeta $\{i\}$
Symbolized as $\xi_i = 1.5$

**For point B :** The distance from the, $D(\pi, B) = -1.5$

$$\xi_i = 2.5$$

$$D(\pi, B) = 1 - \underset{\underset{\xi_i}{\downarrow}}{2.5}$$

**For point C :** The distance from the $D(\pi, C) = -1.5$

$$D(\pi, C) = 1 - \underset{\underset{\xi_i}{\downarrow}}{2.5}$$

$$\boxed{\xi_i = 2.5}$$

**Note** for the SVM support vectors and the correctly classified pts we have $\xi_i = 0$ and for misclassified pts $\xi_i > 0$.

So the total error in our case is : (Soft margin)

$$C \sum_{i=1}^{m} \xi_i$$

As we know $\max(f(x)) = \min(-f(x))$

$$\max(f(x)) = \frac{1}{\min f(x)}$$

So we can write $\max(w) \left\{ \frac{2}{||w||} \right\}$ as $\min(w) \frac{1}{2} ||w||$

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

Subject to $y_i (w^T x_i + b) \geq 1 - \xi_i \quad \forall \; \xi_i \geq 0$

C: hyperparameter : It is the hyperparameter which tunes how much error rate we have to optimize. As C value $\uparrow$ the importance of loss term increases.
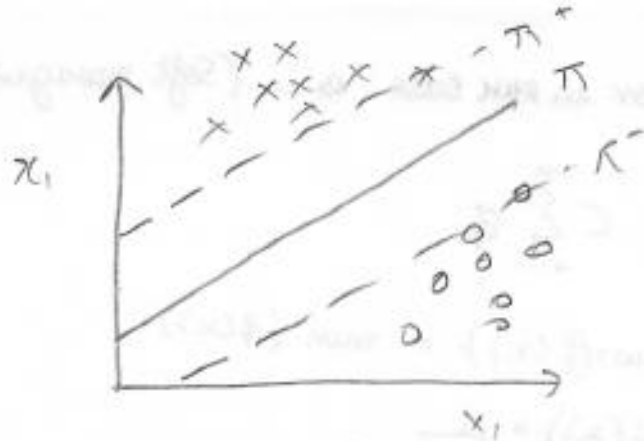
$C \uparrow$ overfitting $\uparrow$
$C \downarrow$ underfitting.

↦ <u>Primal and Dual Forms of SVMs</u>

Hard and Soft Margin is called constrained optimization problem

To get Dual form of the optimization problem we use Lagrange's multipliers

i.) <u>Hard Margin SVM Primal form and Dual Form:</u>

$x_1$

optimization problem for hard SVM

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

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

This is primal form of SVM

$\Downarrow$ convert

Dual form we use (Lagrange multiplier $\alpha$)

rewriting the primal form with Lagrange multiplier

$$\underset{(w,b)}{argmin} \quad \frac{1}{2} \|w\|^2 - \alpha_n \left\{ \text{Constraint eq}^n \right\}$$

where cons eq$^n$

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

$$\underset{(w,b)}{\overset{arg}{min}} \left\{ \frac{1}{2} w^T w - \sum_{i=1}^{n} \alpha_n \ (y_i (w^T x_i + b) - 1) \right\}$$

use min wrt $(w,b)$ & max$^n$ wrt $\alpha_n \geq 0$

$$L \ (w,b,\alpha) = \left\{ \frac{1}{2} w^T w - \sum_{i=1}^{n} \alpha_n \ (y_i (w^T x_i + b) - 1) \right\}$$

Please note here we are minimizing the $L$ wrt to $(w,b)$ but we always maximize $L$ w.r.t $\alpha$

_Gradient definations_ → derivative of the fn.

After finding the gradient w.r.t to $w$ we get:

$$\nabla_w L = w - \sum_{n=1}^{N} \alpha_n y_n x_n = 0 \quad \text{———①}$$

$$\frac{\partial L}{\partial b} = -\sum_{n=1}^{N} \alpha_n y_n = 0 \quad \text{———②}$$

} using lagrange

$$w = \sum_{n=1}^{N} \alpha_n y_n x_n \quad \& \quad \sum_{n=1}^{N} \alpha_n y_n = 0$$

In (1),(2) → partial derivative w.r.t $w$ and partial derivative w.r.t $b$.

Putting $w$ in $\alpha(w,b,\alpha)$

$$\ell(w,b,\alpha) = \frac{1}{2} w^T w - \sum_{n=1}^{N} \alpha_n (y_n(w^T x_i + b) - 1)$$

$$\Rightarrow \sum_{n=1}^{N} \alpha_n - \frac{1}{2} \sum_{n=1}^{N} \alpha_n (y_n(w^T x_A))$$

$$\sum_{n=1}^{N} \alpha_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} \alpha_n \alpha_m y_n y_m x_n^T x_m$$

$$\Rightarrow \boxed{\sum_{n=1}^{N} \alpha_n - \frac{1}{2} \sum_{n=1}^{N} \alpha_n y_n x_n w^T}$$

$$\sum_{n=1}^{N} \alpha_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} \alpha_n \alpha_m y_n y_m x_n^T x_m$$

$$\text{maximize} \quad \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j y_i y_j x_i^T x_j$$

$$\text{subject to} \quad \sum_{i=1}^{n} \alpha_i y_i = 0, \quad \alpha_i \geq 0 \; \forall i \quad \text{Dual form of hard margin}$$

KKT condition    We have for $n = 1, \ldots, N$

$$\alpha_n \left( y_n \left( w^T x_n + b_0 \right) - 1 \right) = 0$$

$\underbrace{\hspace{3cm}}$

Interior points                              $\rightarrow$ SVs

$(w^T x_n + b) > 1$              $w^T x_n + b = 1$

$\Downarrow$                                         $\downarrow$

$\alpha = 0$                                      $\alpha_n > 0$

so    $w^* = \sum\limits_{n \in SV}^{N} \alpha_n y_n x_n$

# Solve for b:-

using $y_n (w^T x_n + b) = 1$    for SVs

Put $w^* = \sum\limits_{n \in SV}^{N} \alpha_n y_n x_n$   and   get 'b'

$$y_n \left( \sum\limits_{n \in SV}^{N} \alpha_n y_n \, x_n \cdot x_n + b \right) = 1$$

$$\sum\limits_{n \in SV}^{N} \alpha_n (y_n)^2 \, x_n^T x_n + b = 1$$

$$\boxed{b = \left( \frac{1}{y_n} \right) - \sum\limits_{n \in SV}^{N} \alpha_n (y_n) * x_n^T x_n}$$

Now we have optimal $w$ & $b$    but we donot have $\alpha$ value

$\Diamond$ How to find $\alpha$ values?

We compute $\alpha$ using the concept of quadratic programming which takes minimization problems.

$$\min_{\alpha} \frac{1}{2} \sum \sum \alpha_n \alpha_m y_n y_m x_n^T x_m - \sum \alpha_y$$

Solution of quadratic programming gives

$$\min_{\alpha} \frac{1}{2} \alpha^T Q \alpha - 1^T \alpha \quad \text{subject to} \quad y^T \alpha = 0 ; \ \alpha \geq 0$$

$\alpha$ is a vector where each $\alpha_i$ corresponds to each $x_i$

---

✦ Soft - More Margin SVM Primal and Dual form.



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

subject to $y_i (w^T x_i + b) \geq 1 - \xi_i \quad \forall i \ \xi_i \geq 0$

2 lagrange multiplier $\alpha$ and $\beta$.

$$L(w, b, \xi, \alpha, \beta) = \frac{1}{2} w^T w + C \sum_{n=1}^{N} \xi_n - \sum_{n=1}^{N} \alpha_n (y_n (w^T x_n + b) - 1 + \xi_n)$$

$$- \sum_{n=1}^{n} \beta_n \xi_n$$

minimize w.r.t $w, b$ & $\xi$    maximize w.r.t $\alpha_n \geq 0$ & $\beta_n \geq 0$

gradient w.r.t 'w' and partial derivative wrt 'b' and $\xi_i$

$$\nabla_w \alpha = w - \sum_{n=1}^{n} \alpha_n y_n x_n = 0 \quad \text{—(1)}$$

$$\frac{\partial L}{\partial b} = - \sum_{n=1}^{n} \alpha_n y_n = 0 \quad \text{—(11)}$$

$$\frac{\partial L}{\partial \xi_n} = (C - \alpha_n - \beta_n) = 0 \quad \text{—(III)}$$

Plug in equation into $L\{w, b, \alpha, \beta\}$
we got same dual form as we computed for hard margin SVM.

$$\max_\alpha L(\alpha) = \max_\alpha \left\{ \cdot \alpha_n - \frac{1}{2} \sum_{n=1} \sum_{n=1} y_n y_m \alpha_n \alpha_m x_n^T x_m \right.$$
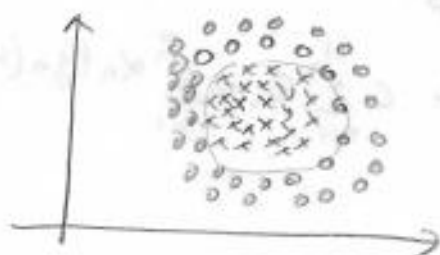
$$s.t \quad 0 \le \alpha \le C \quad , \quad \sum \alpha_n y_n = 0$$

$$\# \quad w = \sum_{n=1} \alpha_n y_n x_n$$

$$\text{Subject to} \quad \sum_{i=1}^n \alpha_i y_i = 0 \quad , \quad C \ge \alpha_i \ge 0 \quad \forall_i$$

Hard margin SVM to separate the linearly separable dataset and soft margin SVM to separate the almost linearly separable dataset.

## Kernels TRICK



Kernel → It takes data points in X space in d dimensional and transform the points in 2 space in d' dimension.

$$\max_{\alpha_i} \sum \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \underbrace{(x_i^T x_j)}_{} \rightarrow \text{similarity function}$$
$$\rightarrow K(x_i, x_j)$$

$$st \quad \sum_{i=1}^n \alpha_i y_i = 0 \quad ; \quad \alpha_i \ge 0$$

$$f(x_a) = \sum_{i=1}^n \alpha_i y_i \underbrace{K(x_i, x_q)}_{} + b$$

Kernelization ÷ SVM handle non-linear separate datasets. -9 -

✡₁) **Polynomial Kernels**

$$k(x_1, x_2) = (x_1^T x_2 + c)^d$$

eg $k(x_1, x_2) = (1 + x_1^T x_2)^2$

⇕

quadratic Kernel

$$x_1 = \langle x_{11}, x_{12} \rangle$$
$$x_2 = \langle x_{21}, x_{22} \rangle$$

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

$$= \left[ 1, x_{11}^2, x_{12}^2, \sqrt{2} x_{11}, \sqrt{2} x_{12}, \sqrt{2} x_{11} x_{12} \right] : x_1'$$
$$\left[ 1, x_2^2, x_{22}', \sqrt{2} x_{21}, \sqrt{2} x_{22}, \sqrt{2} x_{21} x_{22} \right] : x_2'$$

$$= (x_1')^T (x_2')$$

Kernelization doing internally is feature transformation

✡ **Radial Basis Functions (RBF)**

SVM : most popular / general purpose : RBF

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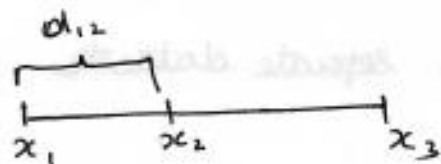
↳ hyperplane parameter

$$d_{12} = \|x_1 - x_2\|_2$$

$$K(x_1, x_2) = \exp\left( \dfrac{-d_{12}^2}{2\sigma^2} \right)$$

$d\uparrow, d^2\uparrow, e^{d^2}\uparrow$

$\dfrac{1}{e^{d^2}}\downarrow$

1.) $d_{12}\uparrow$ ; $K(x_1, x_2) \downarrow$

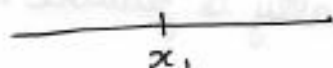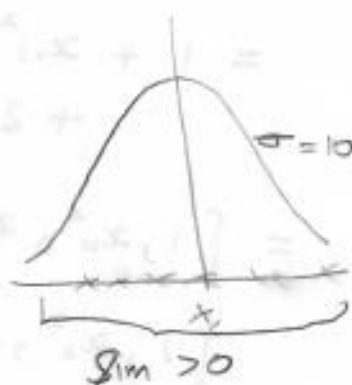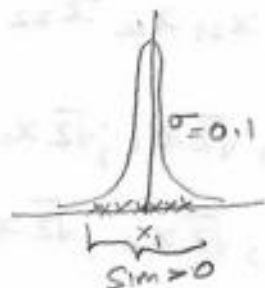$$K(x_1, x_2) > K(x_1, x_3)$$
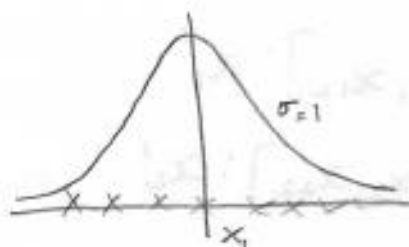
(2) $\sigma$

RBF ~ guassian kernel $(\mu, \sigma^2)$

$$\sigma = 1 \text{ to } \sigma = 0.1$$

$$\sigma = 0.1$$

$$\sigma^2 = 0.01$$

dist $> 1 ; K = 0$



RBF approximation is similar to KNN.

If you don't know which kernal to use you use RBF kernels

⭐ <u>Domain specific Kernels</u> → These are specialized kernels for specified tasks

for eg graph kernel → specific kernel → lot of specilized graph kernel

$Ex \ 1.8$

# Train & Run Time Complexity of SVM.

Train → SGD
  ↳ Specialized algo ( ) → Sequential minimal optimization
                                    (SMO)

Training Time ~ $O(n^2)$ for kernel SVM

{ more opt  ~(2007) ÷ $O(nd^2)$ if $d < n$.
  algo

{ if you have large data $n$ is large → $O(n^2)$ ↑↑
              ↑ Typically do not use SVM when $n$ is large.
                                                ↓
                                            internet
                                            applications

$O(Kd)$

---

## nu - SVM

alternative formulation of SVM

            $0 \leq nu \leq 1$
hyperparameter — nu

        nu ≥ fraction of errors
      nu ≤ fraction of SV's
lower bound.

  nu = 0.01 ⇒ %age of errors = 1%
              # SVs ≥ 1% of n points

SVM.  Dtrain
  10% errors
  nu = 0.1
  1% error
  nu = 0.01

## SVM Regressions

### SVM Classification : SVC $\quad y_i \in \{+1, -1\}$

Mathematical formulation:

br form of
SVR

$$\min_{w, b} \frac{1}{2} \|w\|^2 \quad \longrightarrow \text{regularization}$$

$$\text{S.t} \quad y_i - (w^T x_i + b) \leq \mathcal{E} \quad \longrightarrow \text{hyperparameter}$$

$$(w^T x_i + b) - y_i \leq \mathcal{E}$$

$$\mathcal{E} \geq 0$$

$$\boxed{f(x_i) = w^T x_i + b}$$
$$= \hat{y}_i$$

error $\boxed{\begin{aligned} &= y_i - \hat{y}_i \leq \mathcal{E} \\ &\hat{y}_i - y_i \leq \mathcal{E} \end{aligned}}$

If not kernelized ⟵



If kernelized



Kernel SVM

$\mathcal{E} \downarrow \Rightarrow$ errors are low on training data

$$\Rightarrow \text{overfitting} \uparrow$$

$\mathcal{E} \uparrow \Rightarrow$ errors on $D_{train} \uparrow \Rightarrow$ underfit $\uparrow$.

RBF SVR $\rightarrow$ KNN —reg