

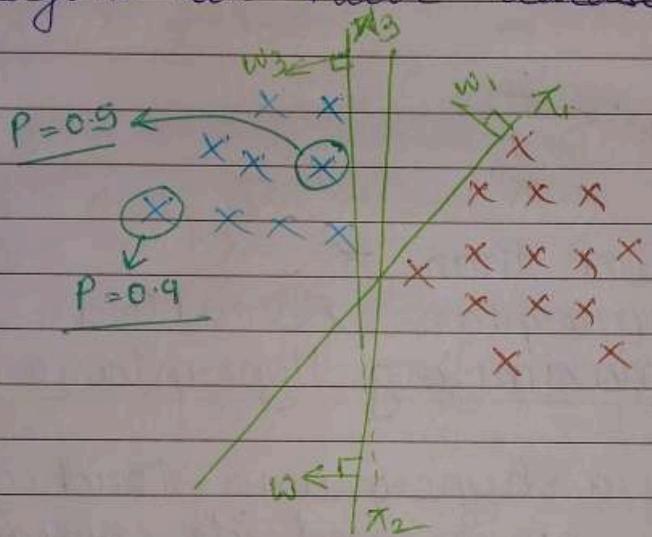
# SVM

36.1

SVM  $\rightarrow$  popular ML  $\rightarrow$  classification  
 $\rightarrow$  regression

## Geometric Intuition

Imagine we have dataset like this:

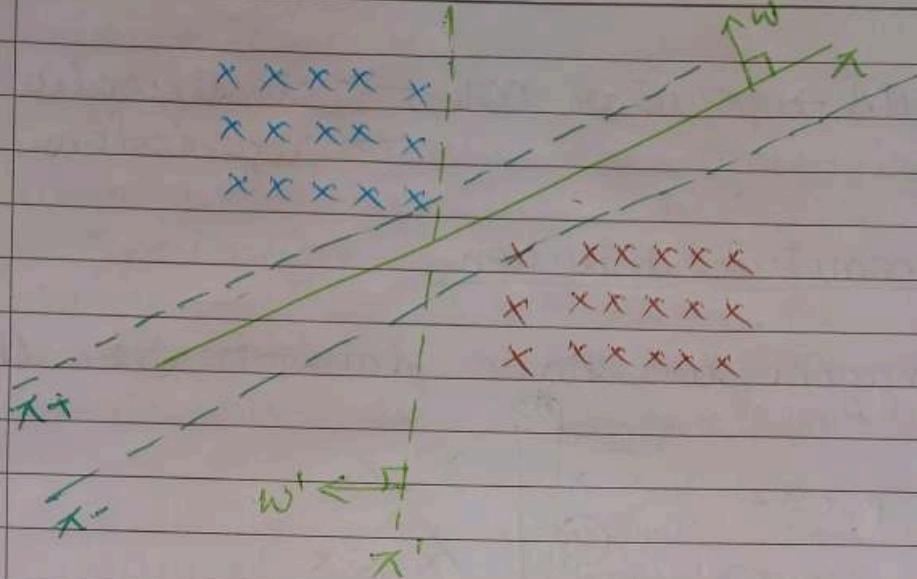


There are multiple hyperplanes separating the two classes of points

Consider the plane  $\pi_3$ . We see that the positive points are very close to that plane.

We have seen in Logistic Regression that if a point is close to hyperplane, then its probability of belonging to that class is 0.5.

We want a hyperplane which separates +ve points from -ve points as far/widely as possible. This is the core idea of SVM



$\pi$  is better than  $\pi'$ .  
Such a hyperplane is called  
Margin Maximizing Hyperplane.

Let  $\pi^+$  be a hyperplane touching a positive point and its parallel to  $\pi$   
Hence  $\pi^-$  touches a -ve point parallel to  $\pi$

Hence

$\pi^+$  and  $\pi^-$  are parallel to  $\pi$   
 $\pi^+$  and  $\pi^-$  are parallel to each other.

$\# \text{dist}(\pi^+, \pi^-) = \text{margin}$

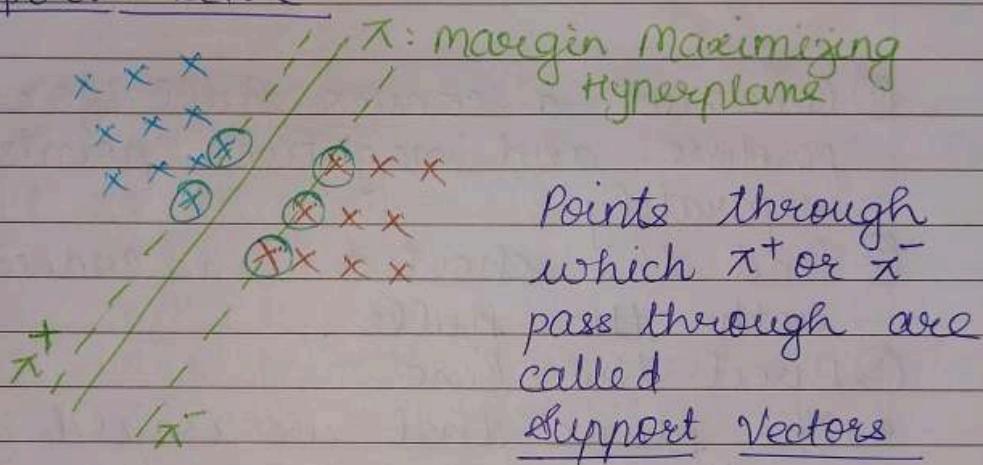
We want to maximize this.  
And wider this gap, the better it is

SVM attempt to find the hyperplane  
that maximizes the margin

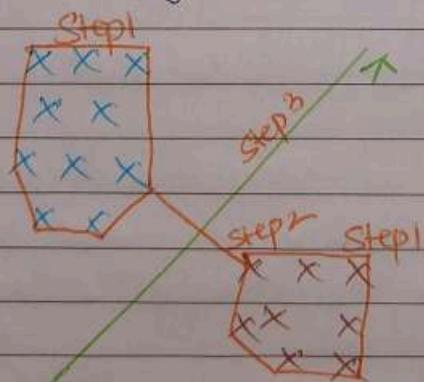
**\*\*\* -** If margin is high, the chance that points would be misclassified would reduce; means my generalization error would reduce (improve)

As Margin  $\uparrow$ ; Generaliz<sup>n</sup> Accuracy  $\uparrow$

### Support Vector



### Alternative geometric interpretation of SVM



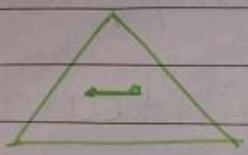
We have to draw a convex hull



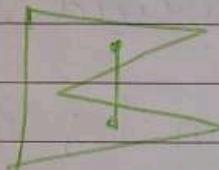
It's a convex polygon where all points are either in it or on the polygon

What is convex?

Line between two points in a polygon  
should lie within the polygon



Convex Polygon



Not Non convex Polygon

- ① Construct a convex hull for positive and negative points separately
- ② Find the shortest line connecting both the hulls.
- ③ Bisect the line
  - \* The plane that we would get is the marginal max plane

325

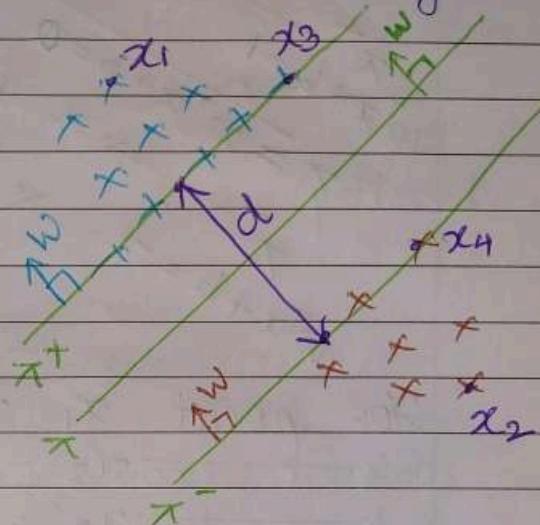
## 36.2 Mathematical Formulation of SVM

We have to find  $\pi$  with margin maximization

$$\pi : \text{margin-max}$$

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

$\pi, \pi^+, \pi^-$   
since  $w, w^+, w^-$  are parallel to each other,  $w$  will be  $\perp$  to all of them



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

$$\pi^+ : w^T x + b = 1 \quad ? \text{ we will see later}$$

$$\pi^- : w^T x + b = -1 \quad ? \text{ what it is}$$

We did not say  $w^T w = 1$ ; i.e. not a unit vec.

Margin :  $d = \frac{2}{\|w\|_2}$

We want to find

$$w^*, b^* = \underset{w, b}{\operatorname{argmax}} \left\{ \frac{2}{\|w\|_2} \right\}$$

such that

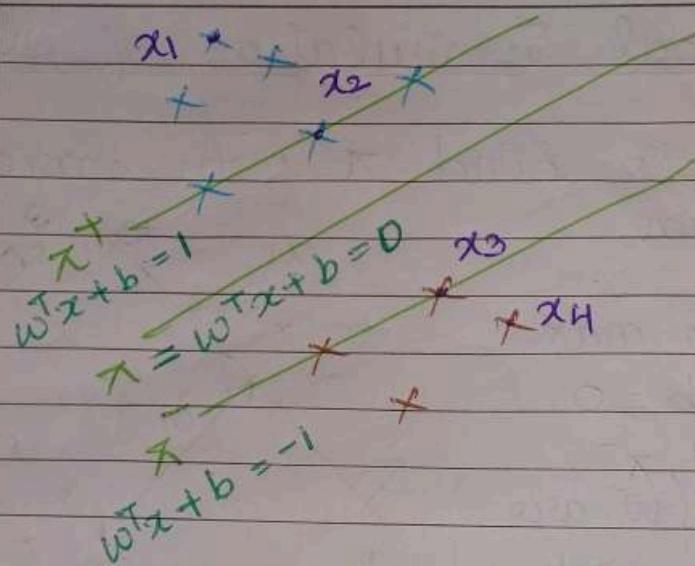
For point  $x_1$ ,  
 $(w^T x_1 + b) * y_1 \geq 1$

For point  $x_3$ ,  
 $y_3 (w^T x_3 + b) = 1$

For point  $x_2$   
 $y_2 (w^T x_2 + b) < -1$

For point  $x_4$   
 $y_4 (w^T x_4 + b) = +1$

$(-1) \quad (-1)$



For point  $x_1$ ;  $y_i(w^T x_i + b) > 1$

$\vdash x_2$ ;  $y_i(w^T x_i + b) = 1$

$\vdash x_3$ ;  $y_i(w^T x_i + b) = -1$

$\vdash x_4$ ;  $y_i(w^T x_i + b) \Leftarrow > 1$

$$\frac{(-1)}{(-1)} \frac{(-1)}{(-1)}$$

(less than -1) (Resultant product could be  $> 1$ )

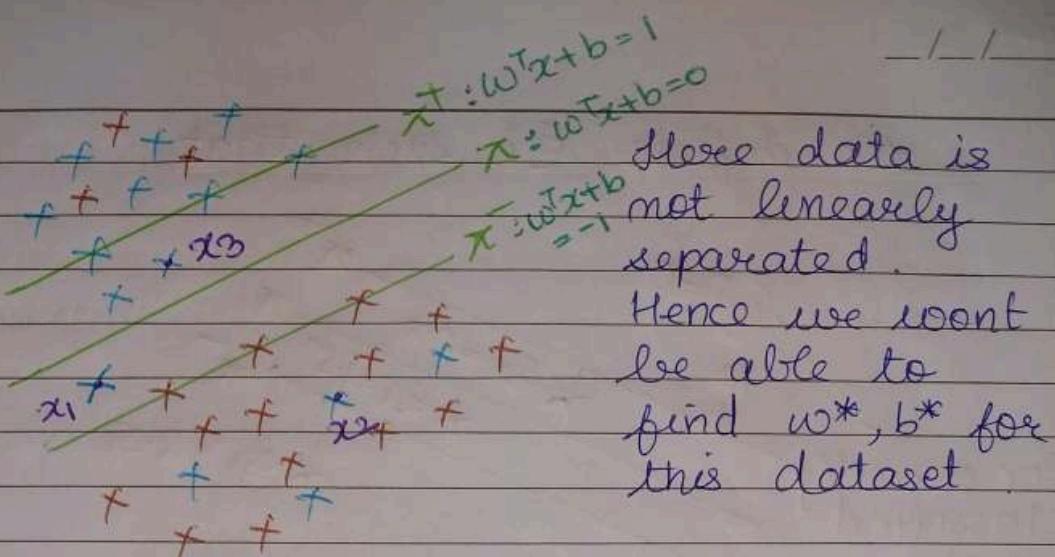
Our problem (optimiza<sup>n</sup>) is

$$\underset{\substack{\text{opt} \\ \text{prob.}}}{\text{Constraint}} (w^*, b^*) = \arg \max_{w, b} \frac{2}{\|w\|}$$

such that  $y_i(w^T x_i + b) \geq 1$  for all  $x_i$

The constraint is that all points should lie either in +ve or -ve region. Hence its not 1 constraint but 'n' constraints

Problem with this: Works well on data which is linearly separated.



Hence our optimization problem is called as Hard Margin SVM because it has a constraint of being linearly separable.

Can we modify the problem to suffice almost linearly separable problem?

Lets assume that point  $x_1$  is halfway between  $\pi^+$  and  $\pi^-$ .

$$\therefore y_1 (\mathbf{w}^T \mathbf{x}_1 + b) = -0.5 \dots \text{Its negative because } y_1 \text{ is positive class}$$

$$\therefore y_1 (\mathbf{w}^T \mathbf{x}_1 + b) = 1 - (1.5)$$

$\rightarrow$  Let this be called  $\xi_1$

For point  $x_2$ ;  $y_2 (\mathbf{w}^T \mathbf{x}_2 + b) = 1.5$   
assuming its halfway past  $\pi^-$

$$\therefore y_2 (\mathbf{w}^T \mathbf{x}_2 + b) = 1 - (2.5)$$

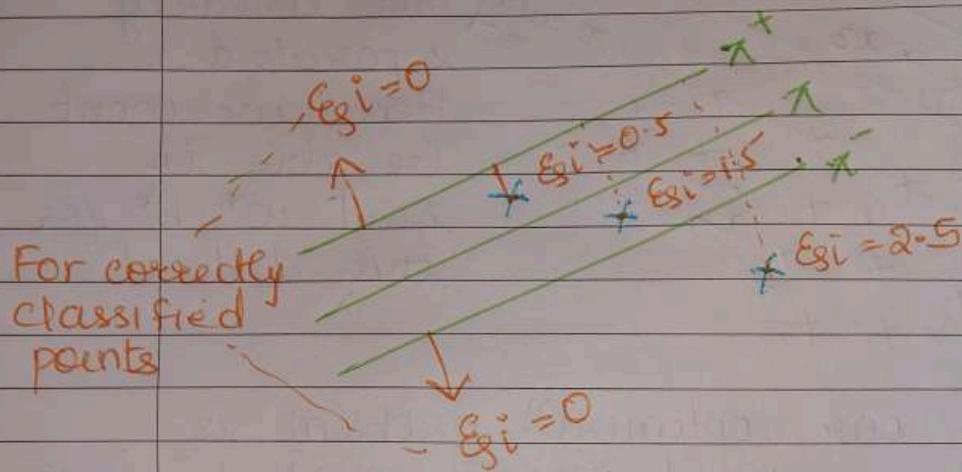
$\rightarrow$  Let this be  $\xi_2$

For point  $x_3$ ; its halfway between  $\pi$  and  $\pi^+$

$$\therefore y_3 (\mathbf{w}^T \mathbf{x}_3 + b) = 0.5 \\ = 1 - (0.5)$$

$\rightarrow$   $\xi_3$

We define  $\epsilon_i$  such that



$\epsilon_i = 0$  for all correctly classified point  
 $\epsilon_i > 0$  for misclassified points  
 depending on how far they  
 are from the plane of  
 the respective class.

For every  $x_i$ , we have  $\epsilon_i$   
 such that

$\epsilon_i = 0$  if  $y_i(w^T x_i + b) \geq 1$   
 i.e. correctly classified  
 points as per  $\pi^+$  and  $\pi^-$

$\epsilon_i > 0$  and is equal to some  
 units of distance away from  
 correct hyperplane ( $\pi^+$  or  $\pi^-$ )  
 in  $\epsilon$  incorrect direction.

It is not exactly the distance  
 because  $w$  is not a unit vector

Now, we know that

$$\arg \min (f(x)) = \arg \max \left( \frac{1}{f(x)} \right) \text{ where } f(x) \neq 0$$

$$\therefore \arg \max \left( \frac{2}{\|w\|_2} \right) = \arg \min \left( \frac{\|w\|_2}{2} \right)$$

Using  $\varepsilon_i$ , we can write optimization prob.

$$(w^*, b^*) = \arg \min_{w, b} \left[ \left( \frac{\|w\|}{2} \right) + C \cdot \frac{1}{n} \sum_{i=1}^n \varepsilon_i \right]$$

such that  $y_i(w^T x_i + b) \geq (1 - \varepsilon_i)$  for each  $x_i$   
 and  $\varepsilon_i \geq 0$

for all misclassified points, we can  
 write like this. As we have seen  
 for  $x_i$ :  $(y_i)(w^T x_i + b) = 1 - (1 - \varepsilon_i)$

For correctly classified points,  $\varepsilon_i = 0$

We want to minimize errors + misclassification  
 i.e.  $\frac{1}{n} \sum_{i=1}^n \varepsilon_i$

Now  $\frac{1}{n} \sum_{i=1}^n \varepsilon_i \Rightarrow$  Average distance of  
 misclassified points from correct  $\pi$ s

We want to minimize this distance,  
 hence we have added it to the term  
 for minimization

what is  $C$ ?

$C$  is the hyperparameter

$$(\omega^*, b^*) = \underset{w,b}{\operatorname{arg\,min}} \frac{\|w\|}{2} + C \cdot \frac{1}{n} \sum_{i=1}^n \xi_i$$

Regulariza<sup>n</sup> Loss

margin      average dist of misclassified points

constraint such that  $y_i(\omega^T x_i + b) \geq 1 - \xi_i \quad \forall i$   
and  $\xi_i \geq 0$

This formula<sup>\*</sup> is called soft margin SVM  
Hence we have

$$\underset{w}{\operatorname{Min}} \quad C \cdot \text{Loss} + \text{Regulariza}^n$$

As  $C \uparrow$ ; the tendency to make mistakes on DTrain  $\downarrow$   
 $\Rightarrow$  overfit  $\Rightarrow$  High Variance

~~comment on video~~  
If  $C$  is high, the loss term would dominate the regularization term. In such case, even though the optimization problem tries to minimize the sum of the loss and the regularization term, as the regularization term is negligible when compared to the loss term, the optimization becomes only minimizing the loss term, as much as possible, thereby there are chances for ensuring no misclassification, leading to overfitting.

As  $C \downarrow$ ; its underfit  $\Rightarrow$  High-Bias

$C$  behaves exactly opposite to  $\lambda$  in Logistic Regression

36.3 Why we take values +1 and -1 for Support Vector Machines.

Because, we are not saying  $\|w\|_2 = 1$   
Its any vector apart from unit vector.

Our task is to maximize the margin

Assume, I take

$$\begin{aligned} \pi^+ : w^T x + b &= k \\ \pi^- : w^T x + b &= -k \end{aligned} \quad \left. \begin{array}{l} k > 0 \\ \end{array} \right\}$$

Why are we taking  $+k$  and  $-k$  & not  $k_1$  and  $-k_2$ ? Because we want both  $\pi^+$  and  $\pi^-$  to be equally far away i.e.  $\text{dist}(\pi^+, \pi) \cong \text{dist}(\pi^-, \pi)$

$\therefore \text{Márgen} = \frac{2k}{\|w\|}$  We know that if

$$\text{if } k=8; \text{ márgen} = \frac{8}{\|w\|}$$

It would change as per  $k$ .

so  $k$  doesn't carry much significance

$\therefore$  We are taking +1 and -1 as way of convinience

From algebraic point,

$$\text{Let's say: } \pi^+ : w^T x + b = k$$

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

$$\Rightarrow (w')x + (b') = 1$$

Hence we can divide by  $k$ , because  $w$  is not a unit vector.

Reason to take +1 and -1: To simplify the math.

36.4

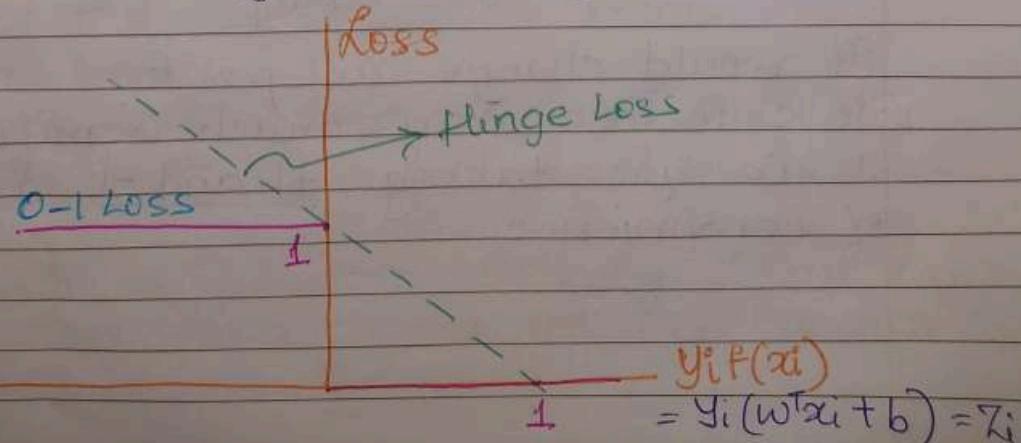
Loss Function (Hinge Loss) based Interpretation

When we study Logistic Regression,

Logistic Regression = Logistic + Regularization loss  
equivalent to 0-1 Loss

Linear Reg: Linear Loss + Regularization

SVM: Hinge Loss + Regularization

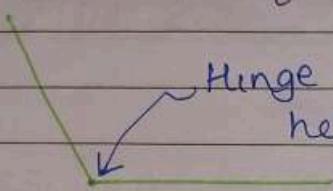


$$\text{let } (\omega^T x_i + b) y_i = z_i$$

$z_i > 0$ :  $x_i$  is correctly classified  $\rightarrow 0$  loss

$z_i < 0$ :  $x_i$  is incorrectly classified

Hinge loss is a straight line from  $-\infty$  to 1; after that it becomes 0.



Hinge Loss is not differentiable here as its not smooth curve

But there are some optimization

$\therefore$  Hinge Loss :  $z_i \geq 1$ ; hinge loss = 0  
 $z_i < 1$ ;  $\rightarrow = \underline{1-z_i}$

One way to simplify the def  
 Hinge Loss :  $\max(0, 1-z_i)$

Case 1:  $z_i \geq 1$ ;

$1-z_i$  is negative value.

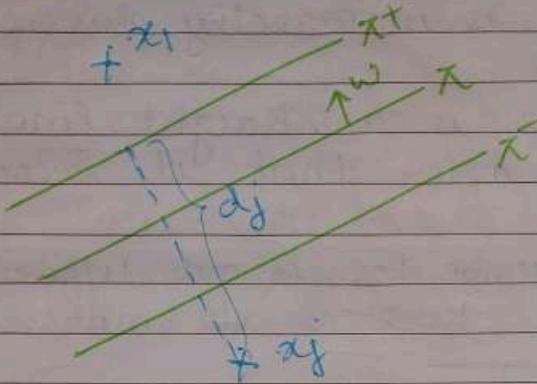
$$\therefore \max(0, 1-z_i) = 0$$

Case 2:  $z_i < 1$

$$\therefore (1-z_i) > 0$$

$$\therefore \max(0, 1-z_i) = (1-z_i)$$

### Geometric Intuition



$x_i$  is correctly classified :  $\epsilon_{gi} = 0$

$x_j$  is incorrectly

Let the distance of  $x_j$  from  $x^+ = d$

$$d_j = 1 - y_j (\underline{w^T x_j + b})$$

value would be -ve as

$\rightarrow x_j$  is on opp side of  
w direction

$d_j$  is distance of misclassified positive pt

$$d_j = 1 - \frac{y_j (w^T x_j + b)}{\downarrow +ve \quad \downarrow -ve} = z_j (1 - z_j)$$

$\therefore d_j$  would be +ve

We know

$$\epsilon_{gj} = \text{dist to } x_j \text{ from } x^+ = d_j = 1 - z_j$$

$$\therefore \underline{\epsilon_{gj} = 1 - z_j}$$

$\therefore \epsilon_{gj} = 0 \rightarrow \text{correctly classified}$

$\therefore \epsilon_{gj} = 1 - z_j \rightarrow \text{misclassified}$

$$\therefore \max(0, 1 - z_i) = \varepsilon_i$$

Soft SVM formula:

$$\textcircled{I} \quad \underset{\omega, b}{\text{Min}} \left[ \frac{\|\omega\|}{2} + C \sum_{i=1}^n \varepsilon_i \right] \text{ such that } (1 - y_i(\omega^T x_i + b)) \geq \varepsilon_i \quad \varepsilon_i \geq 0$$

Loss Min formula:

$$\textcircled{II} \quad \underset{\omega, b}{\text{Min}} \sum_{i=1}^n \max(0, (1 - y_i(\omega^T x_i + b))) + \lambda \|\omega\|^2$$

Let's prove that I and II are same

$$1. \frac{\|\omega\|}{2} \text{ is same as } \|\omega\|^2$$

Since  $\|\omega\| > 0$ , minimizing  $\|\omega\|$  is similar to minimizing  $\|\omega\|^2$

We can put constant to  $\|\omega\| \Rightarrow \frac{\|\omega\|}{2}$  or remove it.

2. We have hyperparameter  $C$  &  $\lambda$
- $C \rightarrow$  multiplied to Loss
- $\lambda \rightarrow$  Regularizer

<u>Soft SVM</u>	<u>Loss Min</u>	<u>Behaviour</u>
As $C \uparrow$	As $\lambda \downarrow$	Overfit
As $C \downarrow$	As $\lambda \uparrow$	Underfit

Typically we call it  $C \rightarrow$  multiplied with loss func<sup>n</sup> &  $\lambda \rightarrow -$  with regularizer. And the behaviour is inverted.

36.5

### Dual form of SVM formulation

$$\underset{w, b}{\text{Min}} \quad \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

such that

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

&  $\xi_i \geq 0$

This formulation is soft Margin SVM  
Also called as Primal formula of SVM

This Primal is equivalent to other formula<sup>m</sup>:

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

such that  $0 \geq \alpha_i \geq 0$

and  $\sum_{i=1}^n \alpha_i y_i = 0$

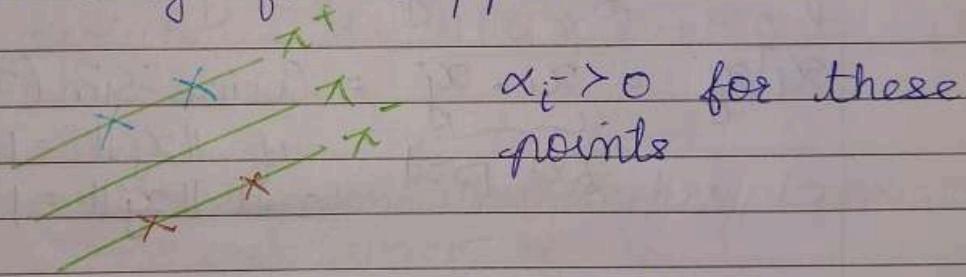
→ called Dual form of SVM

- For every  $x_i$ , we have corresponding  $\alpha_i$
- $x_i^T x_j$  only occur in product form.
- $x_i$  occur only in form of  $\alpha_i^T x_j$
- We want that given query point, we want to compute  $(w^T x_q + b)$ .  
If its positive/negative determines the class label?

$$\therefore f(x_q) = \sum_{i=1}^n \alpha_i y_i x_i^T x_q + b$$

$\underbrace{\qquad\qquad}_{w^T x_q}$

- $\alpha_i > 0$  only for support vectors



$\alpha_i = 0$  for non support vectors

$$f(x) = \sum_{i=1}^n \alpha_i y_i x_i^T x_q + b$$

$\underbrace{\qquad\qquad\qquad}_{}$

Support Vec:  $\alpha_i > 0$

Non-Supp:  $\alpha_i = 0$ ; whole term becomes 0

∴ To compute  $f(x)$ ; only thing that matters are Support Vector.  
Hence the name Support Vector Machine.

The Non support Vector dont really change the class

Hence, the formulation

$$\max_{\alpha_i} \left\{ \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j \right\}$$

such that  $\alpha_i \geq 0 \rightarrow \alpha_i = 0$  for Non SV  
 and  $\sum_{i=1}^m \alpha_i y_i = 0$  for S.V.

We also said that  $x_i$  occur in form of  $x_i^T x_j$ . Why is that imp?

$$x_i^T x_j = \underbrace{x_i \cdot x_j}_{\text{Dot prod}} = \text{Cosine-Sim}(x_i, x_j)$$

if  $\|x_j\| = 1$   
 and  $\|x_j\| = 1$

We get the similarity between  $x_i$  and  $x_j$  only in Dual form  
 Often times the similarity is written as  $K(x_i, x_j)$   
 where  $K$  = Kernel function

$K(x_i, x_j)$  = Similarity function  
 = between  $x_i$  and  $x_j$

36.6 Kernel Trick

Dual formulation of SVM

$$\underset{\alpha}{\text{Max}} \left[ \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \boxed{x_i^T x_j} \right]$$

$\rightarrow \text{sim}(x_i, x_j)$

such that  $\sum_{i=1}^m \alpha_i y_i = 0$  Also called  
 $K(x_i, x_j)$   
Kernel function

and  $\alpha_i \geq 0$

The most imp idea in SVM is the Kernel Trick.

If we use  $x_i^T x_j \rightarrow$  linear SVM

If we replace it by  $K(x_i, x_j) \rightarrow$  Kernel SVM

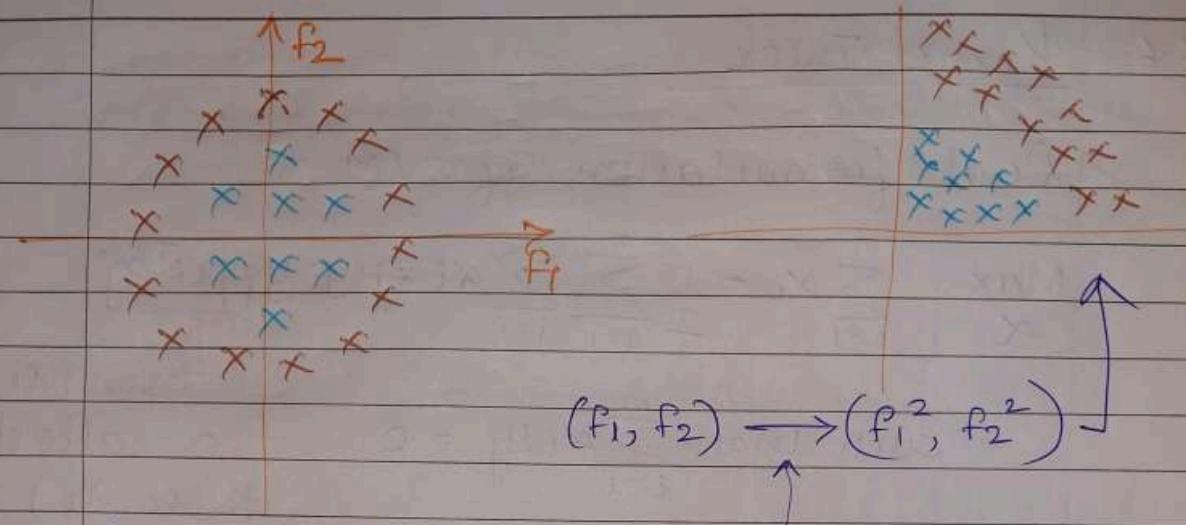
Linear SVM :- We try to find Margin Max hyperplane in space of  $x_i$

Logistic Regression : Minimize logistic loss in space of  $x_i$

World changing idea in area of SVM is Kernel function

Using Linear SVM, we find a hyperplane

But it cannot solve all kinds of data points



In this dataset

Linear SVM : fail

Logistic Regression : fail

Logistic Regression + Feature Transformation = Succeed

Kernel SVM : will find right = succeed.  
kernel

There is a difference between  
feature transformation & kernelization.

Hence,

Kernelization  $\Rightarrow$  Enables SVM handle  
Non linearly separable  
dataset .

36.7

POLYNOMIAL KERNEL

In logistic regression,  
we transformed  
the features into

$$(f_1, f_2) \longrightarrow (f_1^2, f_2^2)$$

feature  
Transforma<sup>n</sup>. ↓

Logistic Reg



How does kerneliza<sup>m</sup> handles this?

Given two datapoints  $x_1$  and  $x_2$ ,  
the polynomial kernel is defined as

$$K(x_1, x_2) = (x_1^T x_2 + c)^d \Rightarrow c \text{ and } d \text{ are constants}$$

$$\text{Eg: } K(x_1, x_2) = (1 + x_1^T x_2)^2 \Rightarrow c=1 \text{ and } d=2$$

Consider  $x_1 = \langle x_{11}, x_{12} \rangle$  and  
 $x_2 = \langle x_{21}, x_{22} \rangle \Rightarrow d=2$

$$K(x_1, x_2) = (1 + x_{11} \cdot x_{21} + x_{12} \cdot x_{22})^2$$

$$\begin{aligned}
 &= 1 + x_{11}^2 x_{21}^2 + x_{12}^2 x_{22}^2 \\
 &\quad + 2 \cdot x_{11} x_{21} x_{12} x_{22} \\
 &\quad + 2 \cdot x_{11} \cdot x_{21} \\
 &\quad + 2 \cdot x_{12} \cdot x_{22}
 \end{aligned}$$

We can split this into two vectors

$$\Rightarrow d' = 6$$

Let

$$\begin{aligned} \mathbf{x}_1' &= [1, x_{11}^2, x_{12}^2, \sqrt{2}x_{11}x_{12}, \sqrt{2}x_{11}, \sqrt{2}x_{12}] \\ \&\mathbf{x}_2' = [1, x_{21}^2, x_{22}^2, \sqrt{2}x_{21}x_{22}, \sqrt{2}x_{21}, \sqrt{2}x_{22}] \end{aligned}$$

$$\therefore K(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{x}_1')^\top (\mathbf{x}_2')$$

So, given features datapoints

$$\mathbf{x}_1 \quad \mathbf{x}_2$$

$$\Downarrow \text{feature transformation}$$

$$\mathbf{x}_1' \quad \mathbf{x}_2'$$

and then we can perform  $(\mathbf{x}_1')^\top \mathbf{x}_2'$

What Kerneliza<sup>n</sup> does internally is

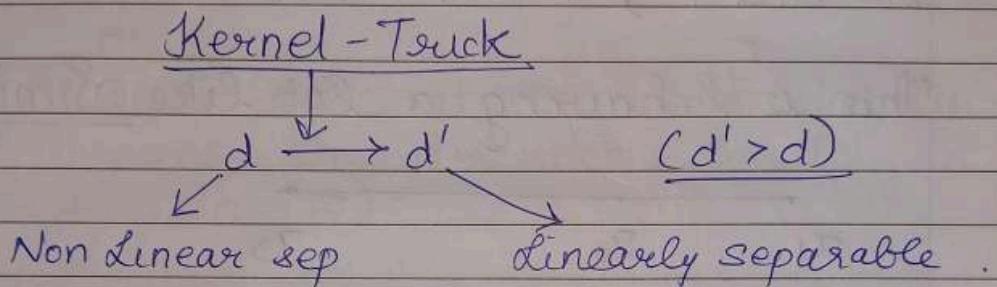
$d$  Feature  $\times$   $d'$   
Transforma<sup>n</sup>  
(internally)

Here, in our eq.  $d = 2$  and  $d' = 6$   
such that  $d' > d$

So, we went from  $d = 2$  to  $d = 6$  and  
in 6-D data, we have a quadratic  
term equivalent to  $(F_1, F_2) \Rightarrow (F_1^2, F_2^2)$

There is a theorem called Mercer's Theorem:

It says that Kernel trick takes d-dim data which is not linearly separable and transforms it into  $d'$ -dim which is linearly separable.



The explicit feature transformation in Logistic Regression is replaced by Kernelization trick ('implicit feature X-forma')

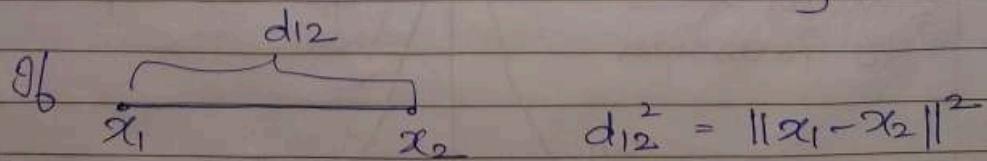
Challenge: What is the right Kernel.

### 36.8 RBF (Radial Basis func<sup>n</sup>) Kernel

SVM: Most popular/general-purpose RBF

Given two datapoints  $x_1$  and  $x_2$ ;

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

$\sigma$  

$$d_{12} \quad d_{12}^2 = \|x_1 - x_2\|^2$$

$\sigma$  = Hyperparameter

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

where  $d_{12} = \|x_1 - x_2\|^2$

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

This is behaving a lot like Similarity

$$\cdot \quad \cdot \quad \cdot \quad \cdot$$

$x_1 \quad x_2 \quad x_3$

①  $K(x_1, x_2) > K(x_1, x_3)$

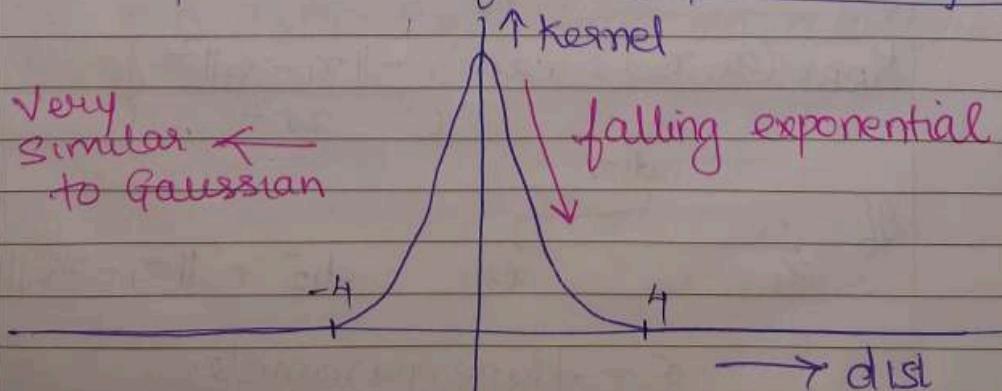
② Impact of  $\sigma$  (for a constant  $d_{12}$ )

$$\underline{\sigma = 1}; \quad K(x_1, x_2) = \exp\left(-\frac{d_{12}^2}{2}\right)$$

case 1: If  $d_{12} = 0$

$$\therefore K(x_1, x_2) = \exp(0) = 1$$

As  $d_{12} \uparrow$  in positive/negative direction  
the  $K$  value falls exponentially



Kernel value is equivalent to similarity  
As the distance reduces, the similarity value increases.

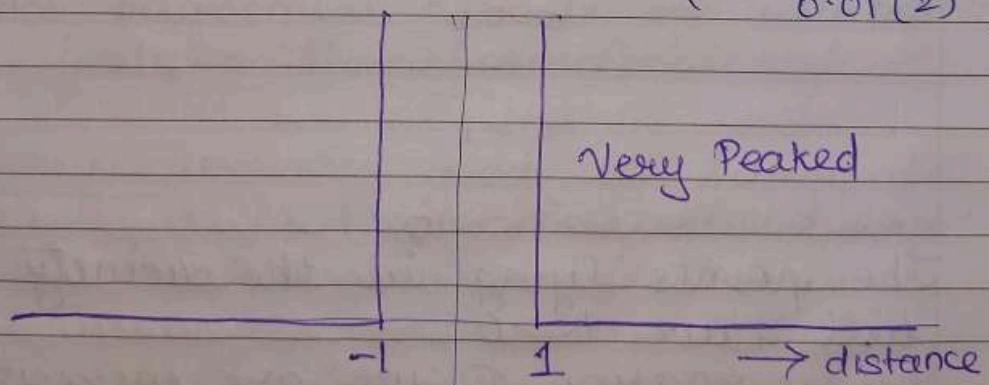
When  $d_{12} = 0.5$ , we get high Kernel Value.  
When  $d_{12}$  is high, the similarity reduces down to 0.

If  $K(x_1, x_2) = 0$ ; it means similarity = 0,  
It means the distance should be high ( $d_{12} = \infty$ )

Intuitively ; kernel  $\rightarrow$  similarity  
distance  $\rightarrow$  Dissimilarity.

$$\sigma = 0.1$$

$$\therefore \sigma^2 = 0.01 \quad \therefore K(x_1, x_2) = \exp\left(-\frac{\text{dist}(x_1, x_2)}{0.01(2)}\right)$$

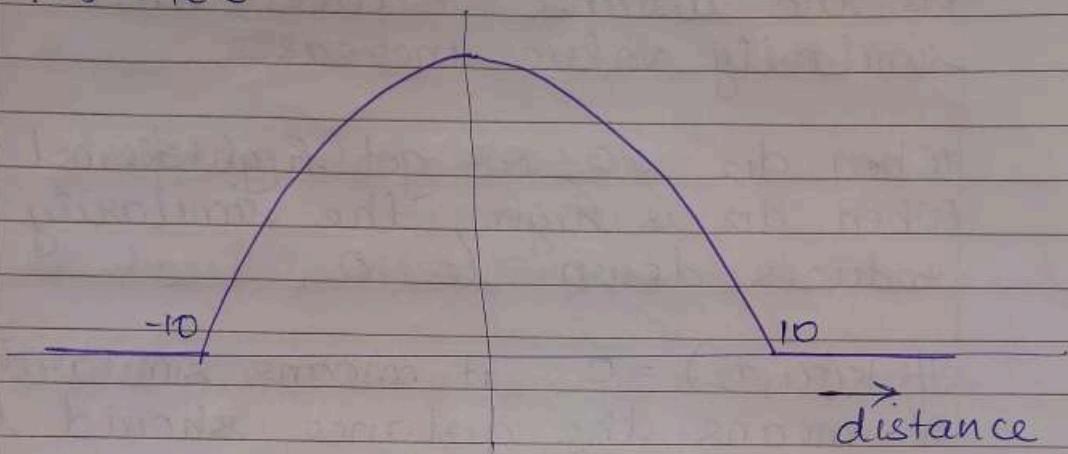


when  $d > 1$ ;  $K = 0$

346

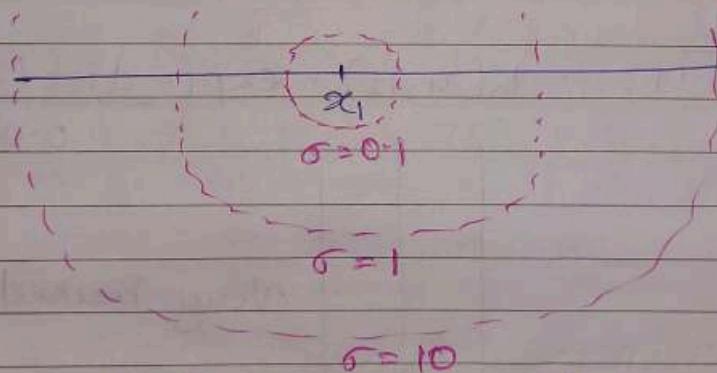
-/-/-

$$\begin{aligned}\sigma &= 10 \\ \therefore \sigma^2 &= 100\end{aligned}$$



When  $d_{12} > 10$ ;  $K = 0$

When  $\sigma \uparrow$ , the vicinity of points having  $K > 0$  increases



The points lying in the vicinity will have  $K > 0$ .

By increasing  $\sigma$ , we are increasing vicinity

Similarity and distance are opposite to each other

If we observe, this is similar to K-NN

There as well, we take a large value of  $k$  where we want to increase the distance of neighbourhood.

The  $k$  reduces if the neighbourhood is closer.

$\sigma = 0.1 \Rightarrow$  Very small  $k$  in K-NN

$\sigma = 10 \Rightarrow$  Large  $k$  in K-NN

### K-NN and RBF-Kernel

$\sigma \uparrow$  in RBF =  $k \uparrow$  in K-NN

- Interview Q :
- ① How is SVM related to K-NN?
  - ② How is SVM related to Logistic Reg?
  - ③ How polynomial kernels are related to feature transforms

K-NN are not good for large dataset since we had to store all points for Locality Sensitive Hashing.

RBF-SVM : We just need to have the support vectors and corresponding  $\alpha_i$

K-NN : Lot of time : Time Compl :  $O(nd)$

Hence RBF-SVM is nice approximation to K-NN

348

-/-/-

Hence, if we don't know the "Best" Kernel to use



Simply use RBF-SVM

Soft.  $\leftarrow C$        $\sigma \rightarrow$  RBF  
margin optimization

We have to do research on both  $C$  &  $\sigma$   
i.e. we have to do Grid Search OR  
Random Search on both  $C$  &  $\sigma$ .

36.9

## Domain Specific Kernels

Kernel trick is similar to the Feature Transformation.

And we know that feature Xforma<sup>m</sup> is domain specific and its the Art. There are lot of domain specific kernels in ML.

→ String Kernels → Text classification

→ Graph →

If we don't have domain specific kernels, we can use RBF

If we want to classify, we can check for any available domain specific Kernels.

Eg:

Amazon food reviews  $\rightarrow$  String Classification

Features : hard  $\rightarrow$  Partially replaced by finding the right Kernel  
Transform Part in ML

If we don't find any, RBF can be used for Backlog.

### 36.10 Train & Run Time Complexity of SVM

We can train SVM using S.G.D but there are better algorithms  $\rightarrow$  Sequential Minimal Optimization (SMO)

LibSVM : Best libraries of training SVM.

Training Time Complexity :  $\sim O(n^2)$  for Kernel SVM

There was a more optimized algo published around 2007 :  $O(nd^2)$  if  $d < n$

If  $n$  is large  $\rightarrow O(n^2)$  is very very large & hence most data scientists do not use SVM when  $n$  is large.

$n$  is Large  $\rightarrow$  Internet applica<sup>m</sup>  
(Billions of data)

### Run Time Complexity

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

$\alpha_i = 0$ ; for non support vectors.

Hence the run time complexity depends on the number of support vectors.

If #S.V = k; then

RunTime Complexity :  $O(kd)$

The number of S.V. could be from 1 to n. i.e  $1 \leq k \leq m$

Imagine if k is high.

Let's say  $n = 100k$  and  $k = 1000$ .

$\therefore$  RunTime Compl :  $O(1000d) \rightarrow$  high

Comparing to Logistic Regression

RunTime Compl :  $O(d)$

Now  $d \ll k$

Hence Logistic Regression has much smaller RunTime Complexity than SVM

And we don't have any means to control the Support Vectors

### 36.11. nu-SVM

We had studied till now was C-SVM

C-SVM  $\rightarrow$  Original Formula<sup>n</sup>;  $C \geq 0$



Hyperparameter

There is an alternative formula<sup>n</sup> of SVM  
nu-SVM: Hyperparameter  $\Rightarrow$  nu.

$$\therefore 0 \leq \text{nu} \leq 1$$

and  $\text{nu} \geq \text{fraction of errors}$

$\text{nu} \leq \text{fraction of support Vectors}$ .

Support if I want to build SVM  
 on DTrain and

① If I want less than 10% error  $\Rightarrow \text{nu} = 0.1$

②  $\frac{1}{n} \rightarrow 1\% \Rightarrow \text{nu} = 0.01$

When we say

$\text{nu} = 0.01 \Rightarrow \% \text{ of errors} = 1\% (\text{Max})$

$\Rightarrow \# \text{S.V.S} \geq 1\% \text{ of } n$

Hence, in a way, it's very useful since  
 it's controlling the number of errors

But for Runtime Complexity, we want  
 fewer SVs; unfortunately nu-SVM  
 doesn't control that.

But, mu-SVM is one way to get  
 → no of Support Vectors  
 → no. of Misclassified points in D<sub>Train</sub>

### 36.12 SVM - Regression

In classifica<sup>n</sup>;  $y_i \in \{+1, -1\}$

When SVM Regression is used,  
 Math. formula<sup>n</sup>:

Linear  
form  
of  
SVR

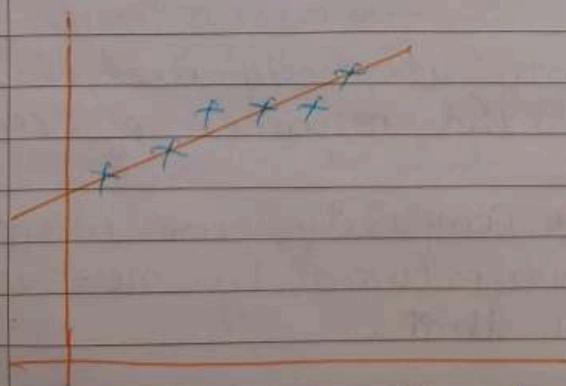
$$\text{Min}_{w} \frac{1}{2} \|w\|^2$$

such that  $y_i - (w^T x_i + b) \leq \epsilon$   
 and  $(w^T x_i + b) - y_i \leq \epsilon \quad \forall i$   
 for  $\epsilon \geq 0$  Hyperparameter

We can Kernelize the formula<sup>n</sup>  
 to fit any shape.

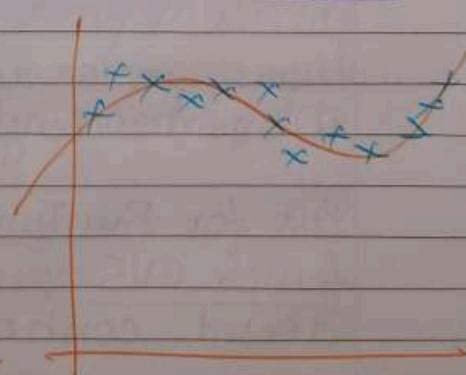
Because

LR - SVR



Will be able to fit only  
Planes

kernel SVR

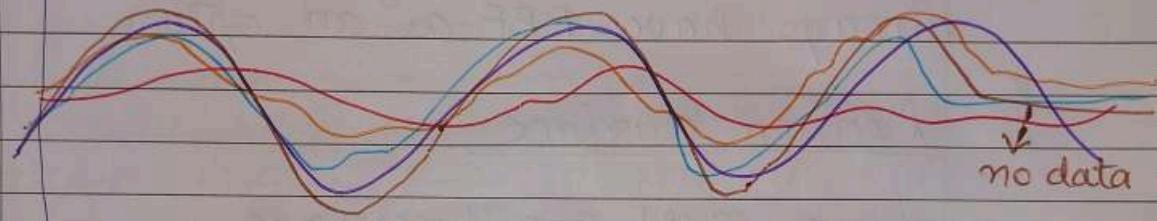


would fit all  
shapes  
 $\epsilon$ : hyperparam

If  $\epsilon$  is  $\downarrow$ ; errors are low on DTrain  
; overfitting  $\uparrow$

If  $\epsilon$  is  $\uparrow$ ; errors are high on DTrain  
; underfit  $\uparrow$

Eq:



True Curve

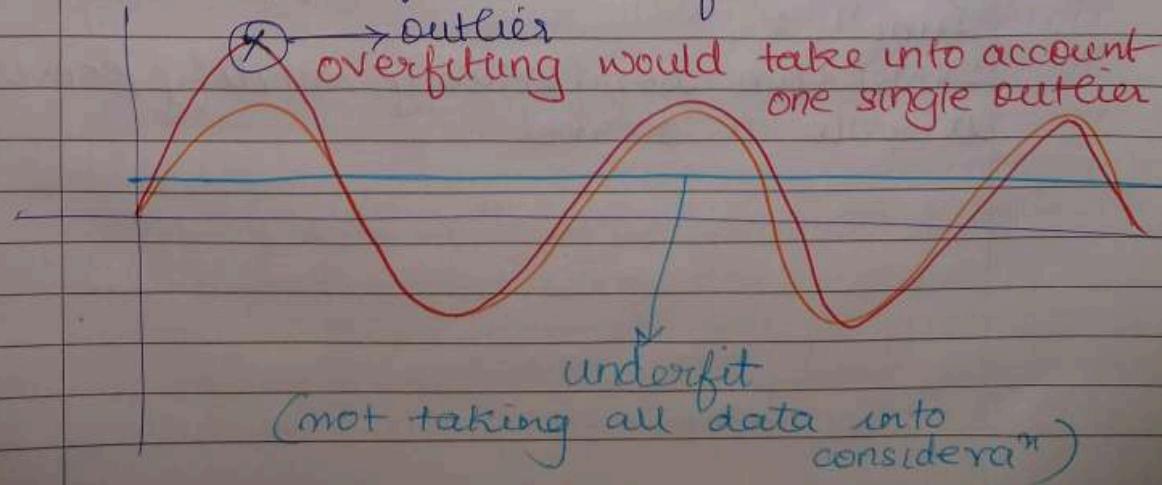
$\epsilon = 0.01 \rightarrow$  overfit

$\epsilon = 0.1$

$\epsilon = 1.0$

$\epsilon = 2.0 \rightarrow$  underfit

In SVM, overfit & underfit means



36.13 Cases : (SVM)1 Feature Engineering and feature Transformations

Replaced by Kernel design,  
It's all about finding the right Kernel.

If we don't have right kernel, we always have RBF as an option

2 Decision Surface .

Linear SVM :  $\rightarrow$  Hyperplane

Kernel SVM : In the original space, i.e. the d-dimensional space in which the points lie ; the decision surface could be a non linear surface.

Using Kernel trick, we can convert  $d \rightarrow d'$  dimensions ; by which the decision surface changes from Non linear  $\rightarrow$  linear .

(3) If we are given Similarity /distance func<sup>n</sup>.

SVMs can comfortably modify /convert it into a kernel .

(4) Big challenge with SVM is about Kernel Interpretability and feature importance

→ Now way to get F.I. directly from SVM

→ we can get using forward feature selection

→ FI and interpretability can be thought of exactly like Logistic Regression

(5) Outliers

→ have very little impact on model

→ Note that if we using RBF with small  $\sigma$ ; it behaves similar to K-NN with small  $K$ .

Since K-NN with small  $K$  can be impacted with outliers.

SVM with  $\sigma$  small can also be impacted

→ In general, impact is very small as only Support Vectors matter

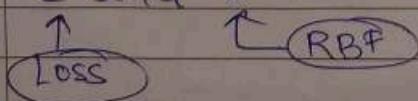
(6) Bias Variance Tradeoff

As  $C \uparrow$ ; we would Overfit  $\Rightarrow$  High Variance  
 $C \downarrow$ ;  $\rightarrow$  underfit  $\Rightarrow$  High Bias

We have Grid Search / Random Search techniques to find optimal  $C$ .

In RBF-SVM we have two Hyperparameters

$C$  and  $\sigma$



(7) What about Large d

→ Very Good for SVM

(8) Best Cases for SVM

① If we have the right Kernel

(9) Worst Case

① 'n' is large.

→ Train Time is high

→ Not used in Internet Based App

② 'k' is large

→ Means #SVs is large

→ Cannot have low latency

→ Time to compute  $f(x_q)$  is large

Typically when 'n' is high, i.e. millions / billions, people end up using Logistic Regression with feature transforms.

BF-1 Decision Trees

KNN → Instance Based Method  $\Rightarrow$  We are given data points and we work on that  
 N.B. → Probabilistic Method  
 Log Regression }  
 Linear Reg. } inspired by ideas of Geometry, hyperplane  
 SVM }  $\rightarrow$  Kernel Based.

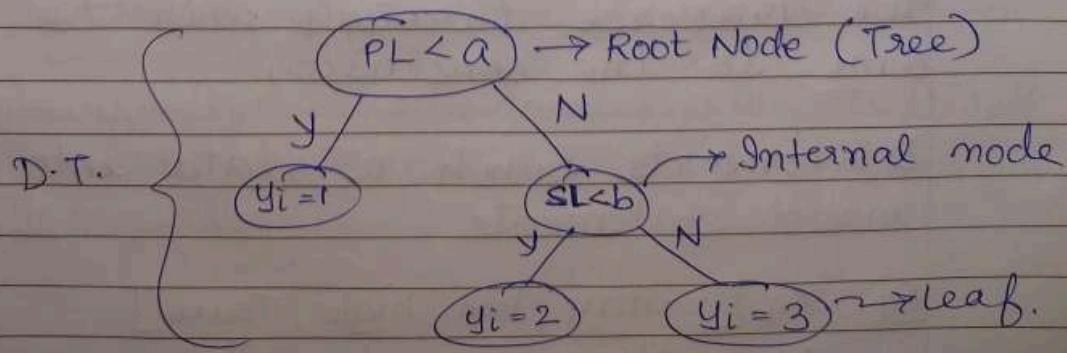
Thus, we have seen 3 types of models.  
 4th Model: Decision Trees.  
 ↳ nested if..else condition

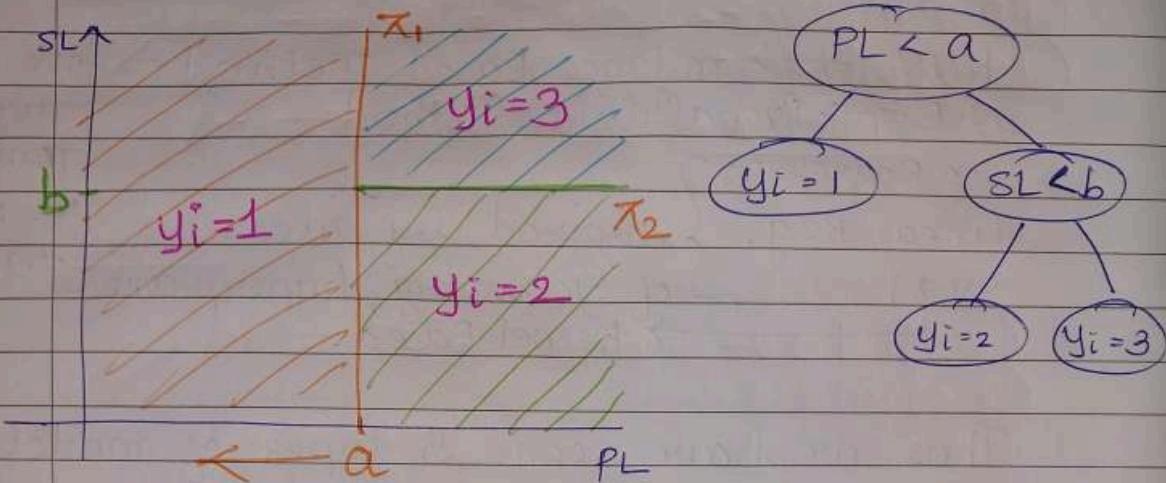
Ex: EDA: Iris dataset  
 where  $y_i = \{1, 2, 3\} \Rightarrow \{\text{setosa, Virginica, Versicolor}\}$   
 and feature =  $\{\text{SL, SW, PL, PW}\}$

We formulated as:

if  $(PL < 5)$       } simple condition  
 then  $y_i = 1$

If my model can be written as nested if else conditions, then we can draw the decision Tree



Geometric Intuition

There would be a  $\pi$  for every decision

V-Imp.: All hyperplanes are axes parallel. in a D-T.  
(parallel to either X/Y axis)

D-T: Set of axis parallel hyperplanes which divide the space into hyperslabs

37-2 Sample D.T.

Eq: Play Tennis

All features are categorical and  $y_i \in \text{Binary}$

We traverse through path to come to the conclusion

D.T. converts bunch of data into insert statements

Eq:  $x_9[\text{sunny, hot, high, True}]$

$\therefore y_9 = \text{No}$

### 37.3 Building D.T. (Entropy)

#### Entropy

→ concept of information Theory, Physics

Suppose we are given r.v. Y which can take 'K' value

Entropy:

$$H(Y) = - \sum_{i=1}^K P(Y_i) \log_b (P(Y_i))$$

where  $b = 2$  or

$$b = e = 2.718$$

Considering the Tennis example.

$$H(Y) = - \frac{9}{14} \log \left( \frac{9}{14} \right) - \frac{5}{14} \log \left( \frac{5}{14} \right) = 0.94$$

$$\uparrow \\ P(Y+)$$

$$\uparrow \\ P(Y-)$$

= No of +ve pts  
by total pts

- % of +ve pts

= % of -ve points

#### Properties

Assume that Y can take two values ( $Y_+, Y_-$ )

Two categories only.  
Two classes.

Case 1:

$$\begin{array}{l} \text{D} \swarrow \\ y+ \rightarrow 99\% \end{array} \quad \left. \begin{array}{l} H(y) = -0.99 \log(0.99) \\ -0.01 \log(0.01) \end{array} \right\} = 0.0801$$

$$\begin{array}{l} \text{D} \searrow \\ y- \rightarrow 1\% \end{array}$$

Case 2:

$$\begin{array}{l} \text{D} \swarrow \\ y+ \rightarrow 50\% \end{array} \quad \left. \begin{array}{l} H(y) = -0.5 \log(0.5) \\ -0.5 \log(0.5) \end{array} \right\} = 1 \\ \begin{array}{l} \text{D} \searrow \\ y- \rightarrow 50\% \end{array} \end{math>$$

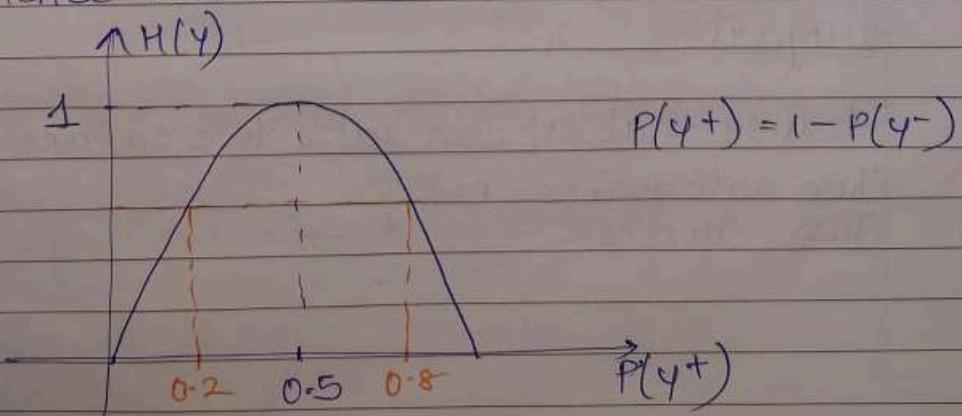
Case 3:

$$\begin{array}{l} \text{D} \swarrow \\ y+ \rightarrow 0\% \end{array} \quad \left. \begin{array}{l} H(y) = 0 \\ y- \rightarrow 100\% \end{array} \right\}$$

from this we understand

→ If the values of class are equiprobable  
i.e. case 2; then we get Max Entropy.→ If one class dominates over other,  
then entropy is 0

Hence Plot looks like:



361

Curve is Symmetric  
 $\therefore H(0.2) = H(0.8)$

$$H(0.2) = -0.2 \log(0.2) - 0.8 \log(0.8)$$
$$H(0.8) = -0.8 \log(0.8) - 0.2 \log(0.2)$$

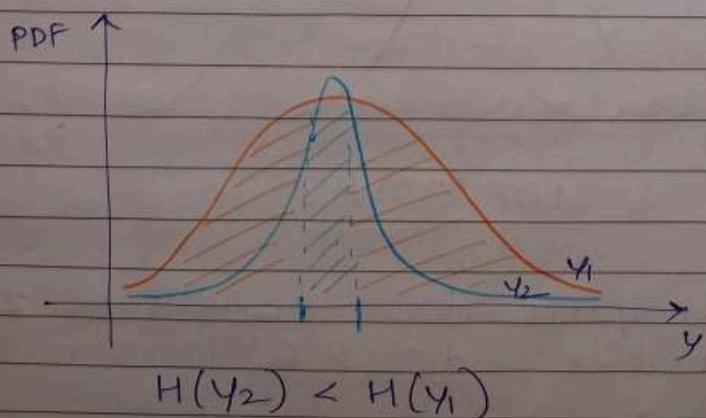
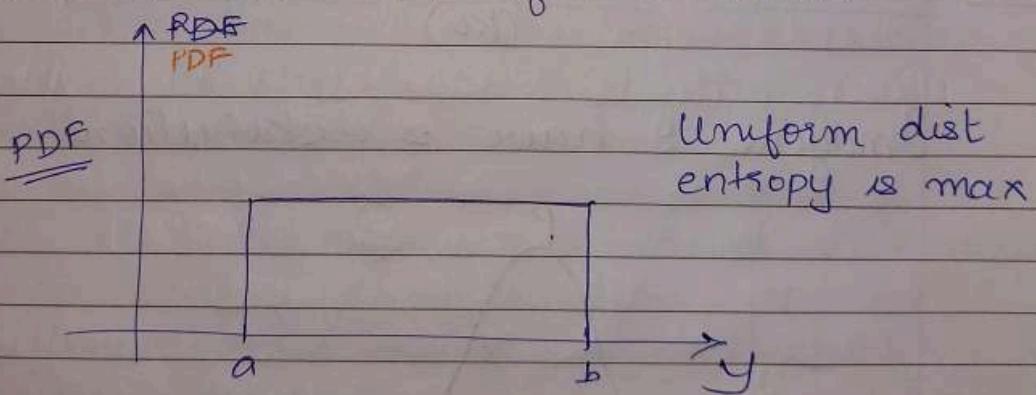
For multi-class r.v.

$$Y \rightarrow y_1, y_2, \dots, y_K$$

if its equiprobable  $\rightarrow$  entropy is max

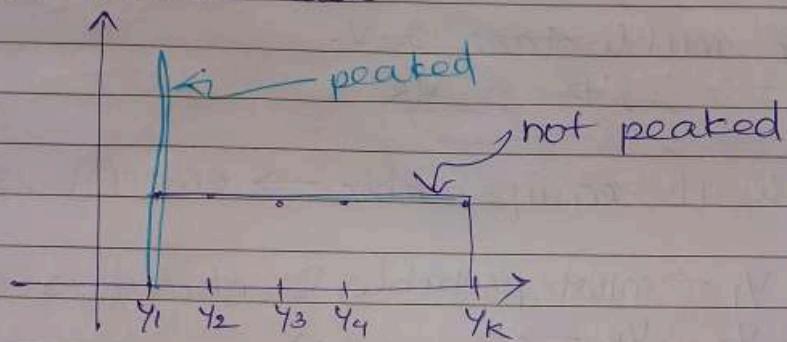
if  $y_1 \rightarrow$  most probable. } entropy is min  
 $y_2, y_3, \dots, y_K \rightarrow 0$  }

Let's draw PDF of the same



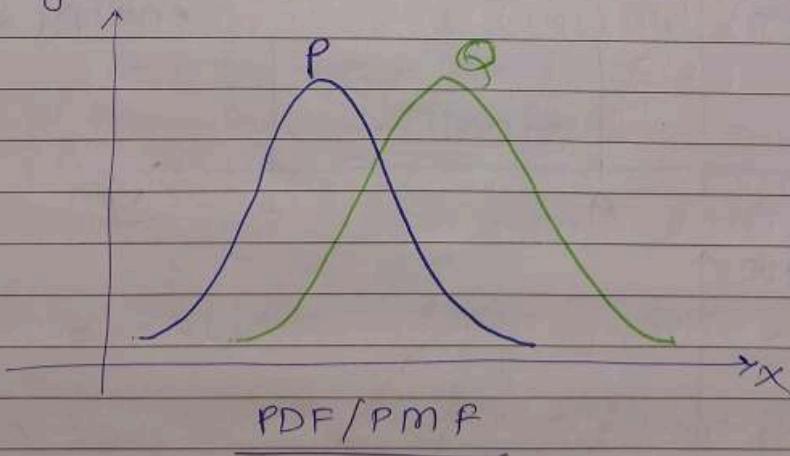
The more peaked the distribution is; less we would see more points in  region. We would see less points in spread region 

For discrete case



### Kullback - Leibler divergence (KL)

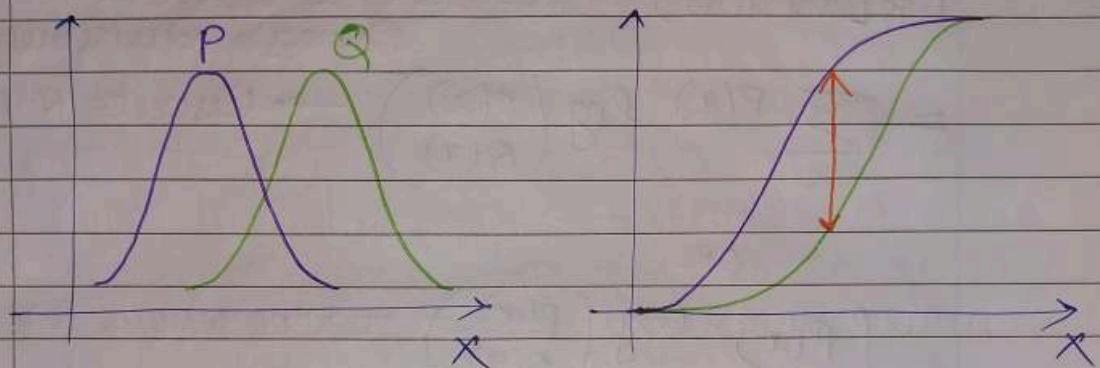
Imagine I have 2 distributions



Can I compute distance between 2 distributions : dist(P, Q)

### One Idea: K-S statistics

We have already studied this.  
It gives the supremum or the maximum gap between the CDF of two distributions.



$$\text{K-S statistics} = \sup (|P'(x) - Q'(x)|)$$

The more closer (similar) two func<sup>n</sup>s/distributions are, the more overlapping would be their CDF and the gap will be smaller.

In all ML algo we studied till now, we require the function to be differentiable.

We cannot differentiate supremum

$$\sup (|P'(x) - Q'(x)|)$$

↓       $\underbrace{\qquad\qquad\qquad}_{\text{differentiable}}$   
 Not       $\underbrace{\qquad\qquad\qquad}_{\text{not}}$

The moment we cannot differential K-S stat, we cannot use it as Loss func<sup>n</sup> in Optimiza<sup>n</sup> Based Method.

Hence, I would need a distance measure which is differentiable

### Info-Theoretic Measure

$D_{KL}(P \parallel Q) \rightarrow$  KL distance between two Probab distributions

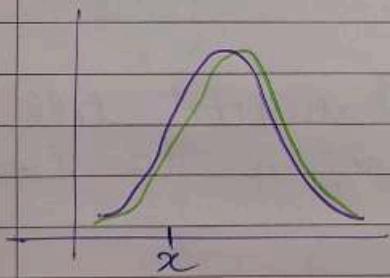
$$= \sum_x P(x) \log \left( \frac{P(x)}{Q(x)} \right) \rightarrow \text{Discrete R.V}$$

OR

$$\int P(x) \log \left( \frac{P(x)}{Q(x)} \right) \rightarrow \text{continuous R.V}$$

#### Case 1:

Distributions  $P$  and  $Q$  are very similar



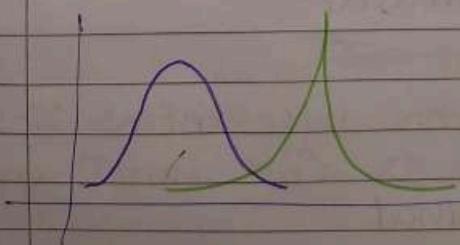
So,  $P(x)$  would be very similar to  $Q(x)$

$$\therefore \frac{P(x)}{Q(x)} \approx 1$$

$\therefore \log$  would be 0

Hence  $D_{KL}$  would be very small.

#### Case 2: $P$ and $Q$ are very different



$D_{KL}$  is very Large

Why there is log?

KL divergence is also called Relative Entropy

$$D_{KL}(P||Q) = \sum_x P(x) \log \left( \frac{P(x)}{Q(x)} \right)$$

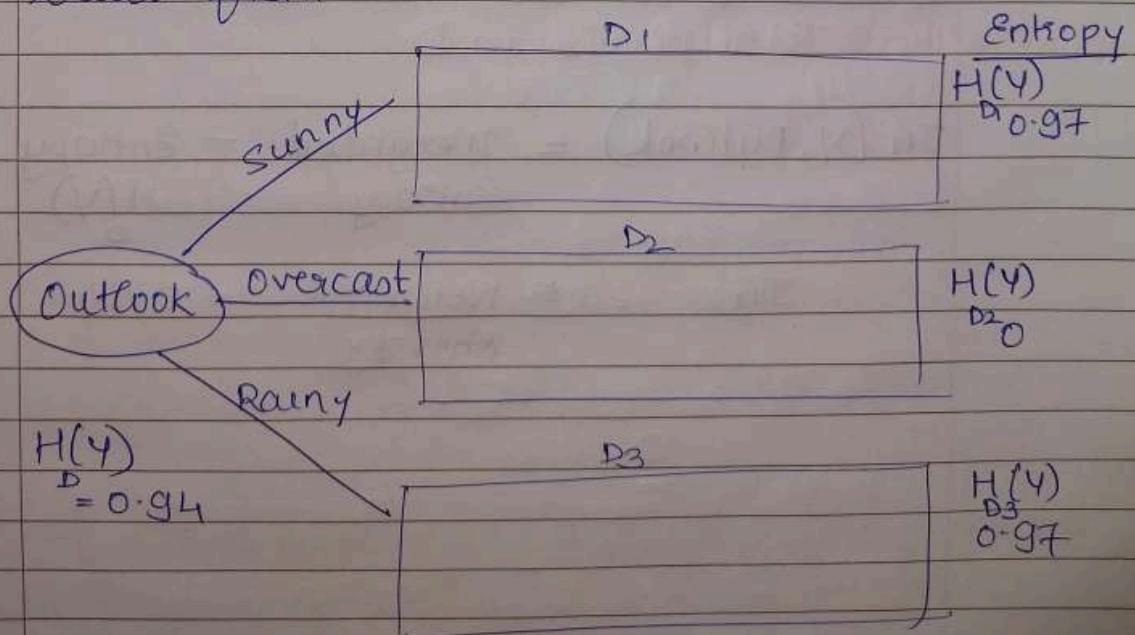
$$= \sum_x P(x) \log(P(x)) - \sum_x P(x) \log(Q(x))$$

This term looks very similar to Entropy.  
Hence we have the log term

Now the whole DKL term is differentiable

### 37.5 Information Gain

Consider the Tennis dataset and DT.  
built from that



- ① We calculated the Entropy of entire dataset with respect to "Outlook".  $H(Y) = 0.94$
- ② Then we calculated split the dataset D into 3 datasets as per the class of Outlook
- ③ We calculated the entropy of each sub dataset.
- ④ Calculate the average entropy.  
(also called Weighted Entropy)

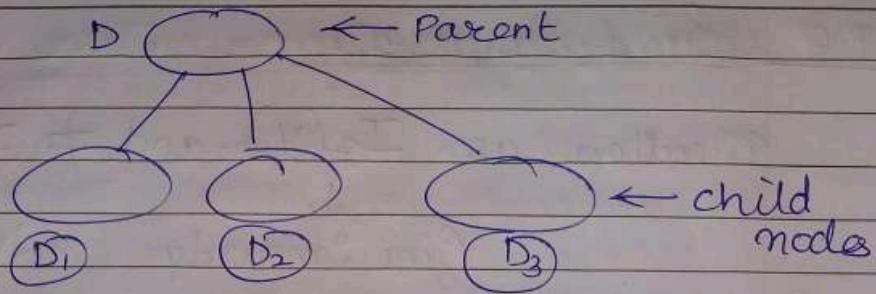
$$\text{Weighted Entropy} = \left( \frac{5}{14} * 0.97 \right) + \left( \frac{4}{14} * 0 \right) + \left( \frac{5}{14} * 0.97 \right)$$

$\downarrow$  Giving the weight as per the number of points       $\downarrow$  Entropy of  $D_1$

$5 \rightarrow$  No. of points in  $D_1$   
 $14 \rightarrow$  Total points in D

$$IG(Y, \text{Outlook}) = \text{Weighted Entropy} - H(Y)$$

$$IG = \text{Weighted Entropy} - 0.94$$



$$\left( -H_D(Y) \right) + \left( \frac{|D_1|}{|D|} * H_{D_1}(Y) + \frac{|D_2|}{|D|} * H_{D_2}(Y) + \frac{|D_3|}{|D|} H_{D_3}(Y) \right)$$

$\leftarrow$  Entropy of Y       $\leftarrow$  Weighted Entropy  $\rightarrow$

$\therefore IG(Y)$

$\therefore$  To sum it up.

If r.v.  $Y \xrightarrow[\text{broken into}]{\text{variable}} Y_1, Y_2, \dots, Y_k$ .

$$\therefore \underline{IG(Y, \text{Variable})} = \left[ \sum_{i=1}^k \frac{|D_i|}{|D|} * H_{D_i}(Y) \right] - H_D(Y)$$

## 37.6 Gini Impurity

Written as  $I_G(Y)$  and  $\neq IG(Y)$

$\downarrow$        $\downarrow$

Gini Impurity

Info. Gain

For  $Y \rightarrow Y_1, Y_2, \dots, Y_k$

$$I_G(Y) = 1 - \sum_{i=1}^k (P(Y_i))^2$$

Assume that a r.v.  $Y$  takes two classes  $Y+$  and  $Y-$

Case 1:

$$P(Y+) = 0.5$$

$$P(Y-) = 0.5$$

$$I_G(Y) = 1 - (0.25 + 0.25)$$

$$= 0.5$$

$$H(Y) = 1$$

Case 2:

$$P(Y+) = 1$$

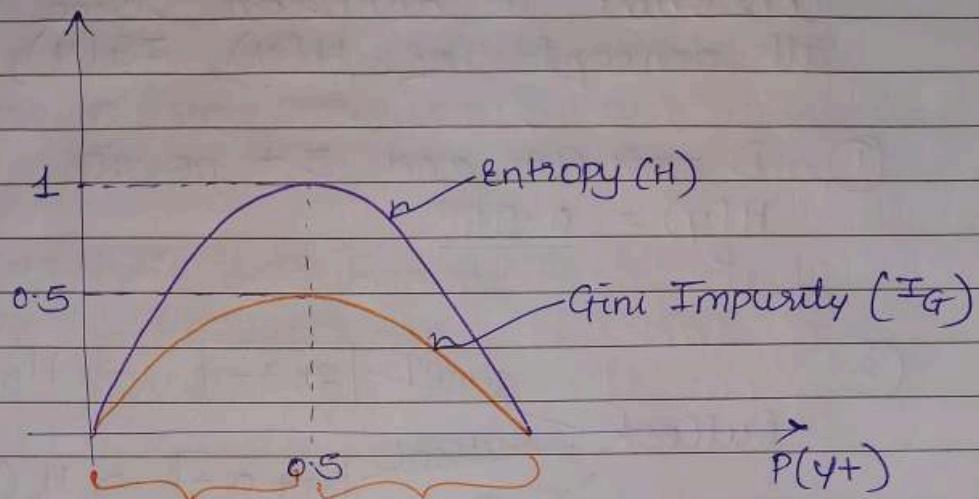
$$P(Y-) = 0$$

$$I_G(Y) = 1 - (1 + 0)$$

$$= 0$$

$$H(Y) = 0$$

Assume we have a 2-category case  
 $(Y+ \text{ and } Y-) ; P(Y+) = 1 - P(Y-)$



As  $P(Y+)$  increases,  
 Both  $H(Y)$  and  
 $I_G(Y)$  increases at  
 a different rate

As  $P(Y+)$  decreases,  
 Both  $H(Y)$  and  $I_G(Y)$   
 decreases

Q: Why are we studying Entropy?

Consider the formula of Entropy

$$H(Y) = -P(Y+) \log(P(Y+)) - P(Y-) \log(P(Y-))$$

In this case, we are computing  $\log$   
 which is more time consuming.

$$I_G(Y) = 1 - (P(Y+)^2 + P(Y-)^2)$$

We are not computing  $\log$ .  
 Hence  $I_G(Y)$  is more computationally  
 efficient than  $H(Y)$ .

370

11

37-7

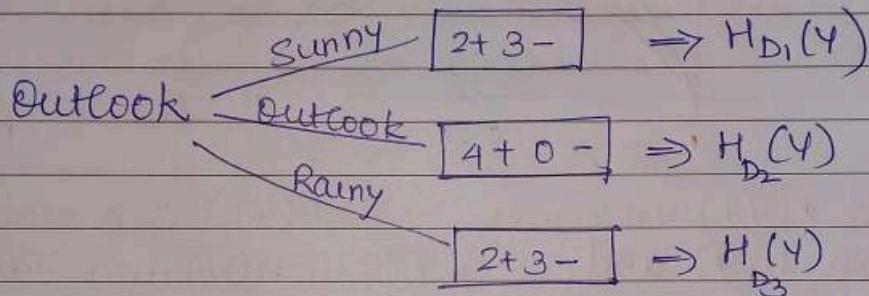
## Building a Decision Tree

Construct a decision tree using all concepts of  $H(Y)$ ,  $IG(Y)$ ,  $IG_1(Y)$

①  $D \rightarrow 9+$  and  $5+$  points

$$H_D(Y) = \underline{0.94}$$

②



$$\text{Weighted Entropy} = 0.69$$

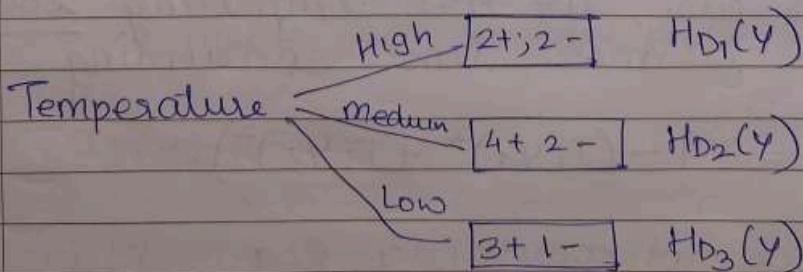
Difference of Entropy at this level with entropy at above level

$$= 0.94 - 0.69$$

$$= \underline{0.25} = IG(Y)$$

③

Breaking  $D$  w.r.t Temperature



And we can find  $IG$  at well

④ Similar process I can do for humidity and windy as well

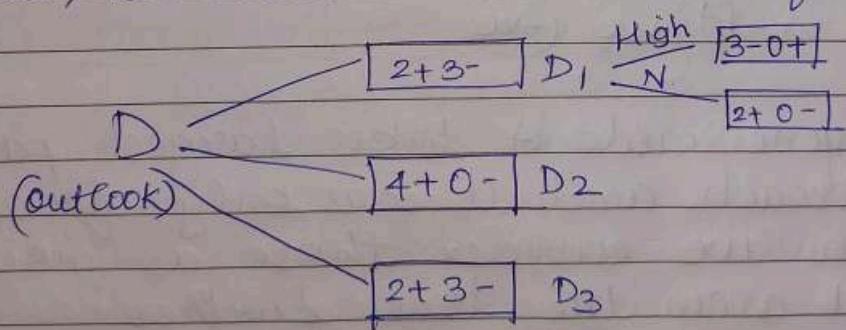
Hence, we performed following steps:

→ choosing Root Node and we computed details.

$$\left. \begin{array}{l} IG(4, \text{Outlook}) = 0.25 \\ IG(4, \text{Temperature}) = \\ IG(4, \text{Humidity}) = \\ IG(4, \text{Windy}) = \end{array} \right\} \begin{array}{l} \text{Pick feature} \\ \text{with Max IG} \end{array}$$

Basically  $IG = \text{entropy at Parent Level} - \text{Entropy at child Level}$

Let's say "outlook" has max entropy.  
Now, we break dataset into following:



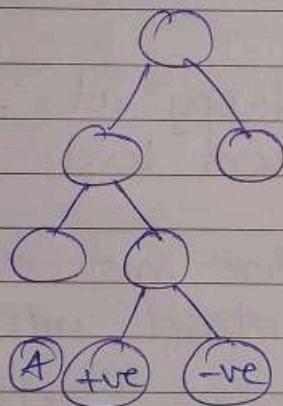
Consider  $D_2 \rightarrow$  It has only +ve points.  
Such a node is called PURE Node.  
because we stop the computation  
and we assign  $y_i = +ve$

Hence, we recursively break each node using IG as criteria.

We start with root node and till the time we get a pure node, we continue the traversal.

OR

If we can't grow the tree because of lack of points, then we won't grow the node further.



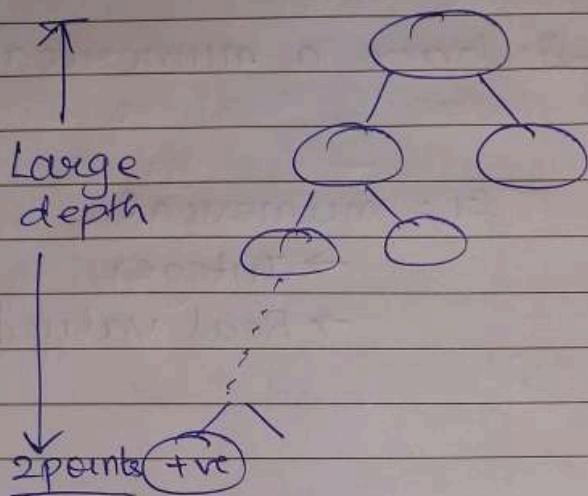
Imagine node A takes largest path to reach and it has only few points, then we consider those as noise and don't grow the tree further.

Hence, we stop computing (growing tree)

if:

- ① Pure node
- ② Few points at a node
- ③ If we are too deep

Imagine  $n = 10k$ . and we construct a tree such that



Hence as depth of tree increases, the chances of overfitting also increases

If depth is small  $\Rightarrow$  underfit.

In D.T. :- hyperparameter  $\Rightarrow$  depth of tree

37.8

### Splitting Numerical features

We have learnt to

Construct a :- Splitting a  $\rightarrow$  IG  
 D.T  $\underbrace{\text{Node}}_{\text{Imp step}}$

We have seen that we use/prefer Gini Impurity over Entropy because:-

- $\rightarrow$  It is computationally efficient
- $\rightarrow$  It has almost same behaviour as Entropy

Till now, we have seen discrete R.V or categorical feature

Let's assume I have a numerical feature

$f_1$	$y$	
2.2	1	$f_1$ : numerical
2.6	1	→ Integer
3.5	0	→ Real valued
3.8	0	
4.6	1	
5.3	0	

Here the split based on categorical feature is not intuitively. It is like, every point has a distinct value  $\Rightarrow$  each point would be a category.

Steps we take:

- 1 Sort the numerical features in ascending order
- 2 Then build the decision D.T condition as per the values

Eq:  $f_1 < 2.2$

$f_1 < 2.6$

$f_1 < 3.5$

$f_1 < 3.8$

$f_1 < 4.6$

$f_1 < 5.5$

consider

$f_1 < 3.8$

Y

N

There are only  
two values  
(Yes/No)

Here, I would have atmost  $n$ -possible variab<sup>ns</sup>. like here we had 6 variab<sup>ns</sup>. For numerical features; we can have  $f_i < T_i$  for each feature

### Numerical features:

$$f_1 < T_1 \quad \begin{cases} D_1 \\ D_2 \end{cases}$$

$$f_2 < T_2 \quad \begin{cases} D_1 \\ D_2 \end{cases}$$

$$f_3 < T_3 \quad \begin{cases} D_1 \\ D_2 \end{cases}$$

:

$$f_n < T_n \quad \begin{cases} D_1 \\ D_2 \end{cases}$$

We evaluate all of them. We check all possible splitting criteria. And we pick the feature with Max IG.

→ Super Time Consuming

→ Non Trivial Computation

There are hacks/computational tools around it to minimize the complexity

376

— / —

37.9

## Feature Standardization

We have learnt that

KNN  
Logistic Regression } Feature Std.  
SVM

We used this

because we were  
computing some  
form of distance

	$f_j$	
$x_i$	$x_{ij}$	
	$\mu_j \sigma_j$	

$$x'_{ij} = \frac{x_{ij} - \mu_j}{\sigma_j}$$

In case of Decision Trees ;

It is not a distance based method.

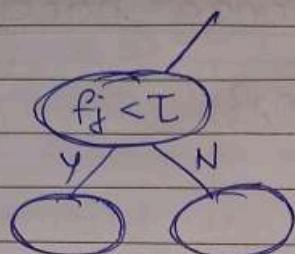
Suppose if I have a real valued feature

	$f_j$	

→ Sort the column

We would check for each of values  
as potential threshold

Hence, as part of D.T we would traverse the path.



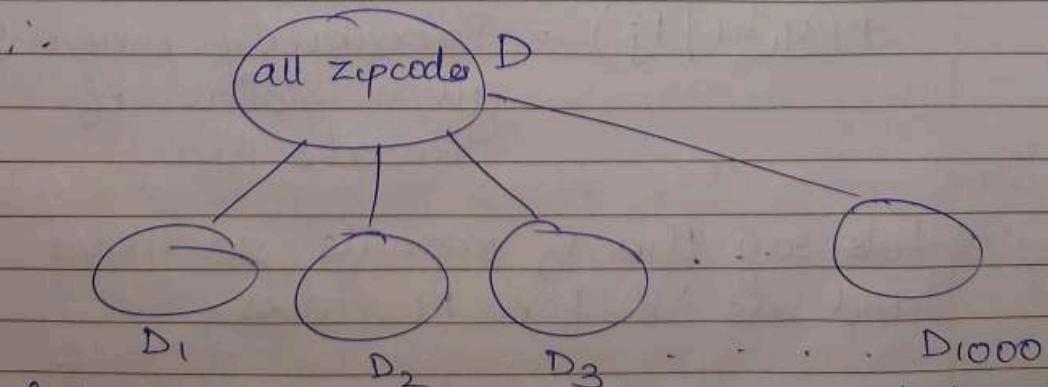
Hence, In D.T we do not care about actual values of features and hence, we don't need to perform Feature Standardizati<sup>n</sup>.

### 37.10 Building a D.T : Categorical features with many possible values

Ex: Pincode / Zipcode

We have 1000s of categories.

We cannot really compare two zipcodes: Ex: 94087 and 94089  
So, every zipcode is a category.



lot of categories and every dataset would be extremely small.

This happens a lot in ML.

Hack: Feature Engg hack.

Pincode	$y_i = \{0, 1\}$
$P_1$	
$P_2$	
$P_3$	
$P_1$	
$P_2$	
:	
:	
-	

Hack would be to convert the categorical feature  $\rightarrow$  numerical feature.

① Assign a value  $y_i$  to a zipcode.

② Calculate Probability as

$P(y_i=1 | P_j)$  = Probability when  $y_i=1$   
given value of  
Pincode  $P_j$

Let's say the  $P_j$  repeats 20 times  
and its '1' for 19 times

∴

$$P(Y_i=1 | P_j) = \frac{\# Y_i=1 \text{ and } P_j}{\# P_j}$$

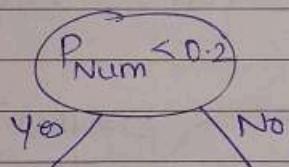
$$P(Y_i=1 | P_j) = \frac{19}{20}$$

So, we would convert

$P_j(Pincode)$  → Numerical feature  
(eg: Probability).

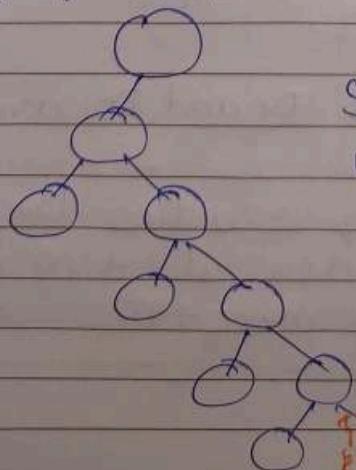
Let's say 2 Pincode have same  
Numerical feature = 0.2.

It would be easier to form a D.T.

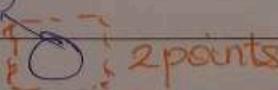


### 37.11 Overfitting and Underfitting

As depth ↑, the possibility of having very few points at leaf node ↑.



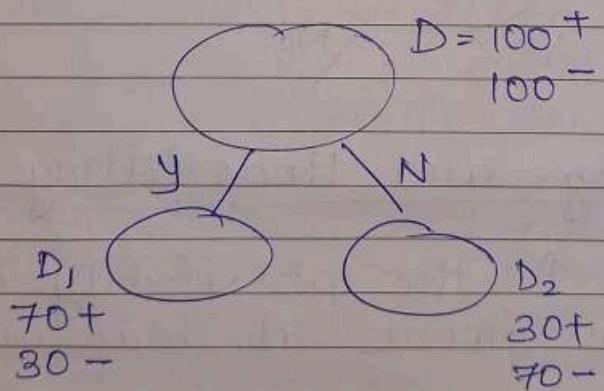
If the 2 points are noisy,  
it would be overfitting

 2 points → noisy

As depth  $\uparrow$ ; the interpretability of model  $\downarrow$

Since there are too many if..else conditions, the interpretability decreases, it becomes hard to understand

At the same time, lets say  
depth = 1



Both of the modes  $D_1$  and  $D_2$  are not Pure  
pulse

Hence if point  $x_0$  reaches  $D_1$ , we would consider it +ve as majority of points are +ve.

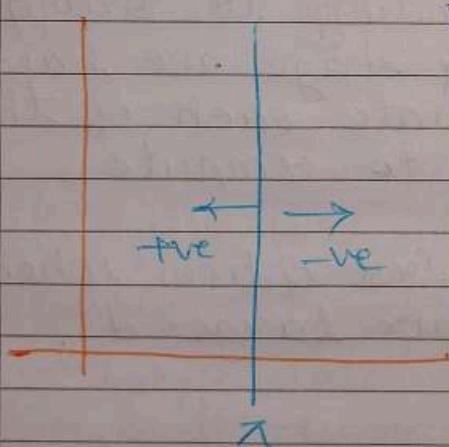
Similarly, if it reaches  $D_2$ , it would be classified as -ve.

Hence the tree with depth=1 is called as Decision Stump.

So, we'll So, when we don't have pure mode, the majority points in leaf nodes are considered.

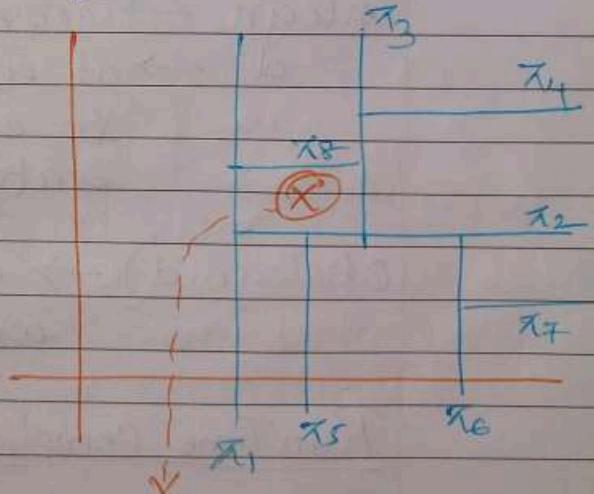
Hence we want to determine right depth  $\rightarrow$  cross Validation

Too shallow depth  $\Rightarrow$  underfit



Very few  
Hypercubes &  
Hyperspaces

Too deep depth = overfit



It has only one point  
which may be a noise

Too many Hypercubes &  
Hyperspaces

## 37.12 Train & RunTime Complexity

### Train Time Complexity:

$$\sim O(n \log n * d)$$

where  $n = \# \text{ points in } D_{\text{Train}}$   
 $d = \text{dimensionality}$

We may think that the numerical features may take long time to compute as there is Threshold associated with it.

But there are certain algo methods to take care of the same.

$n \log n \rightarrow$  corresponding to sorting.  
 $d \rightarrow$  at every stage we have to evaluate each of the feature to compute.

$O(n \log n d) \rightarrow$  not best option when we have large ' $d$ '.

### RunTime Complexity:

We have to save DT<sub>ree</sub>. We convert DT<sub>ree</sub> into nested if... else condition which is very much space efficient.

Space Complexity: Num of Internal nodes  
 $+ \frac{n}{r} \rightarrow$  Leaf nodes

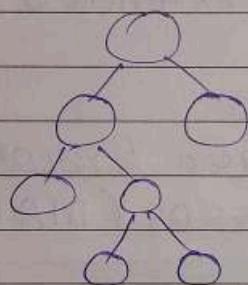
Space Complexity :  $O(\text{nodes})$

Typically we don't train decision trees of more than depth 5; Atmost 10 on very massive datasets.

Hence Run Time Space Complexity : Reasonable.

Run-Time Time Complexity :

While computing  $x_q \rightarrow y_q$ .



Max time taken is the max num of traversals.  
 $= K$   
 $= \text{max depth of any Leaf mode}$

Hence Run Time Complexity is very reasonable

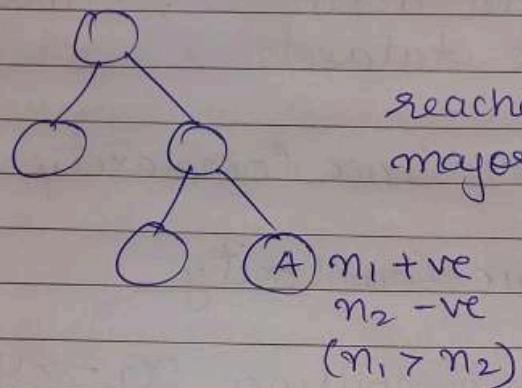
Hence D.T. is super useful during  
 $\rightarrow$  Large data  
 $\rightarrow$  small dimensionality  
 $\rightarrow$  Low latency.

We would see variants of the same like  
 $\rightarrow$  Random Forest  
 $\rightarrow$  GBDT

which are very popular in internet applications

### 37.13 Regression Using D.T.

We have seen D.T. for classification



When the control reaches A, it checks the majority class, and assigns  $y_q$  to  $x_q$ .

It can be easily extended to Regression

We use Mean Square Error (MSE) or Median Absolute Error (MAE)

MSE(D)

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad D \Rightarrow$$

MSE at this stage  
is like Entropy in classification.

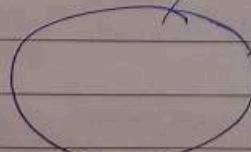
$y_i \in \mathbb{R}$

Compute

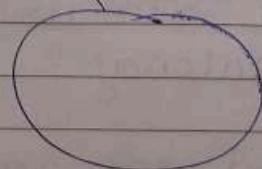
$$\text{Avg} = \text{Mean/Median}(y_i)$$

$$\hat{y}_i = \text{Mean}_D(y_i)$$

Compute  
MSE(D<sub>1</sub>)



D<sub>2</sub>  
compute  
MSE(D<sub>2</sub>)



The feature which gives least MSE would be considered

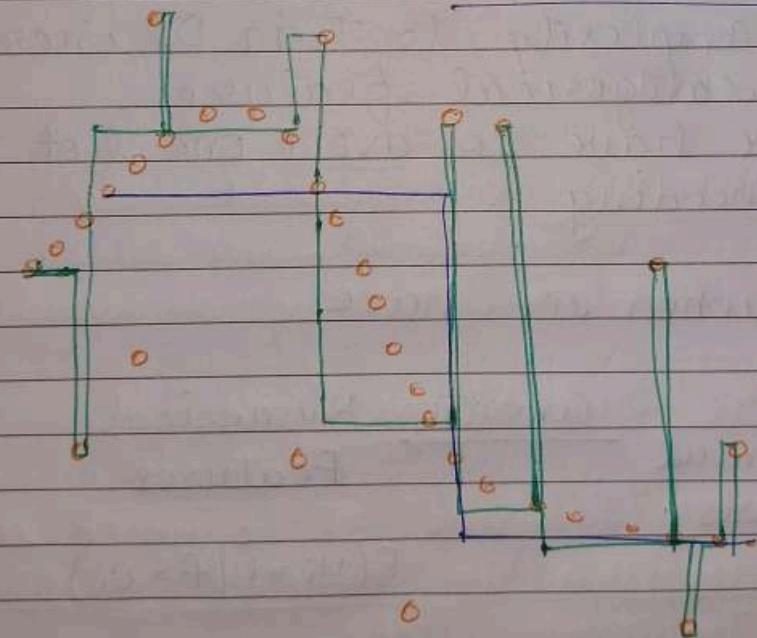
Re. which  
Classification  $\Rightarrow$  Pick a feature  $\Rightarrow$  Reduce Entropy  
 $f_i$  (Min Entropy = 0)

Regression  $\Rightarrow$  Reduce MSE  $\Rightarrow$  Pick feature.

And

$$\text{MSE}_D(y_i) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{MAE} = \text{MAD} = \text{Median}(|y_i - \hat{y}_i|)$$



o  $\rightarrow$  data

—  $\rightarrow$  max depth = 2  $\rightarrow$  underfit

—  $\rightarrow$  max depth = 5  $\rightarrow$  overfit

386

37-14 Cases

① Imbalanced data

- We HAVE TO Balance data
- Upsampling /downsampling
- Imbal. data impacts the entropy calculations /MSE calcula<sup>n</sup>.

② Large dimensionality

- We have to evaluate each of features at each node
- Time complexity to Train DT increases
- For categorical feature
  - we have to avoid one hot encoding

Hence when we have.

Lot of .       $\xrightarrow{\text{Converts}}$  Numerical  
categorical      Features  
features . . . . .  
 $P(y_i = 1 | f = c_i)$

③ Similarity Matrix

- D.T. wont work
- D.T. need features explicitly
- otherwise we cannot find entropy

#### ④ Multiclass Classification

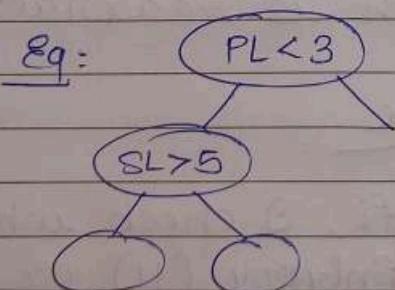
- need not do One Vs Rest (OVR) for D.T.
- Entropy, IG, IG can be computed when we have more than 1 category
- It is straight forward.
- D.T. naturally can be extended

#### ⑤ Decision Surface

- non linear decision surfaces
- axes parallel hyper cuboid
- Each of the internal node has a plane corresponding to it.

#### ⑥ Feature Interaction

- Possible in decision Trees.
- As we have two features interacting with each other.



My condition would be:

$$\underline{(PL < 3)} \text{ and } \underline{(SL > 5)}$$

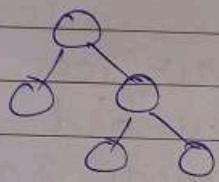
↳ Two features ↑  
interact<sup>n</sup>.

- Feature interact<sup>n</sup> come as an added inbuilt benefit for D.T.
- In Logistics Regression, we didn't have explicit feature interact<sup>n</sup>.

→ In LR, we had feature transformation  
~~Eq:~~  $(f_i * f_j)$  OR ;  $(f_i \text{ AND } f_j)$

### ⑦ Outliers

→ As depth  $\uparrow$ ; outliers would impact



will make Tree unstable

→ [2 pts]

### ⑧ Interpretability

→ D-T are very interpretable.

→ very much readable

→ when depth is reasonable

### ⑨ Feature Importance

→ For every feature  $f_i$ , I check what is reduction in Entropy ( $H$ ) or Gini Impurity ( $I_G$ ).

→ I can add up (Normalized sum) all the reductions in  $H$  due to  $f_i$

→ SKLearn →  $f_i$  has F.I. associated with it