

Module-4

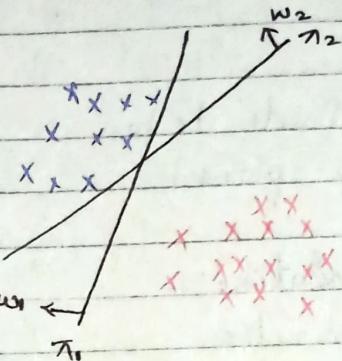
CH - SVM

→ Geometric intuition:-

SVM is a popular ML algorithm, it can do both classification and regressions.

There are many hyperplanes that separate the +ve & -ve points.

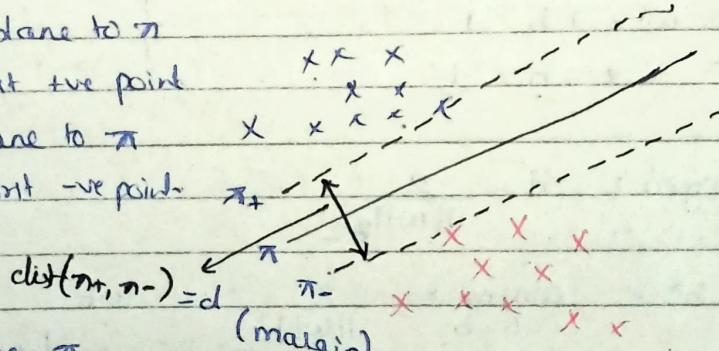
so, which of one is to be preferred.



→ The key idea of SVM: try to find the hyperplane that separates the +ve & -ve points as wide as possible. In above plot π_2 is more preferable because π_1 is very close to points. Such a plane which separates the points with maximum margin is called margin maximizing hyperplane.

π_+ is a parallel plane to π which touches the first +ve point

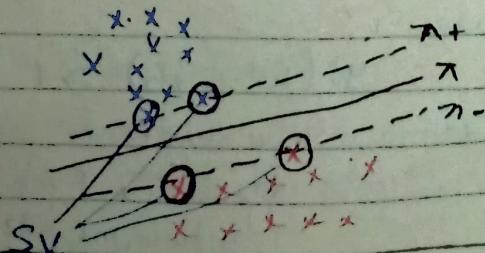
π_- is a parallel plane to π which touches the first -ve point



We want a hyperplane π s.t. $\text{dist}(\pi+, \pi-)$ is maximum.

SVM: Try to find the hyperplane that maximizes the $\text{dist}(\pi+, \pi-)$ because as margin \rightarrow , there is less chance of misclassification.

SV (support vectors):- The points through which π_+ & π_- pass are called support vector.



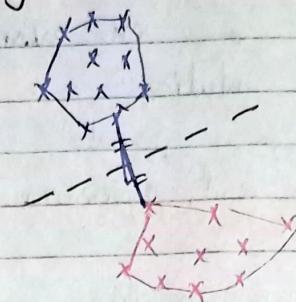
→ Alternative geometric interpretation of SVM:-

① Build a convex hull for the point -ve

② find the shortest line consisting these hulls

③ Bicut the line

④ The plane that bisects the line is margin maximizing plane.



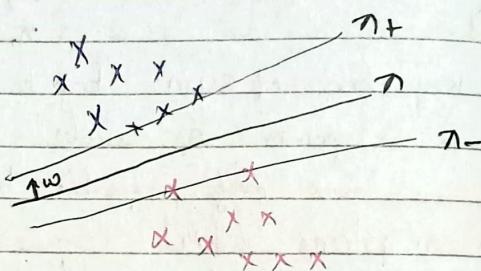
* Mathematical foundation:-

π :- margin maximization

$$\pi = \mathbf{w}^T \mathbf{x} + b = 0$$

As π^+ & π^- are normal to

\mathbf{w} to π the \mathbf{w} will also be normal to π^+ & π^-



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

$$\text{if } \pi^+: \mathbf{w}^T \mathbf{x} + b = 1$$

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

} \mathbf{w} need not be unit vector.

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

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

s.t. all the positive points lie above π^+ and all negative points lie below π^-

Let x_i be a positive point. Then distance of this point from π will be ≥ 1 because $\mathbf{w}^T \mathbf{x}_i + b = \pi^+$. $\pi^+: \mathbf{w}^T \mathbf{x} + b = 1$ and π^+ is boundary of the positive region. Hence any point in the region \mathbf{w} will have

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

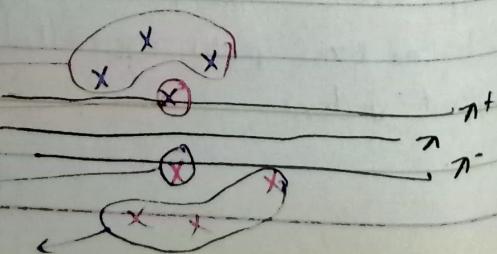
$$[x_i \leftarrow y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1]$$

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1$$

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) = 0$$

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) < 1$$

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) > 1$$



Hence for SV: $y_i(w^T x_i + b) = 1$, for all other points: $y_i(w^T x_i + b) > 1$

$$w^*, b^* = \underset{w, b}{\operatorname{argmax}} \frac{2}{\|w\|_2} \text{ s.t. } y_i(w^T x_i + b) \geq 1 \quad \forall x_i$$

so, we have constraint optimization problem.

But there is one problem. we have considered only linearly separable data. what if our data is non-linearly separable or almost linearly separable in that case our optimization eq will fail.

Hence the above equation is called Hard margin SVM

$$\boxed{w^*, b^* = \underset{w, b}{\operatorname{argmax}} \frac{2}{\|w\|_1} \text{ s.t. } y_i(w^T x_i + b) \geq 1 \quad \forall x_i},$$

Hard margin SVM.

so, can we change our equation slightly so that our model also works for almost linearly separable data.

$$y_i(w^T x_i + b) = 0.5 \xrightarrow{\text{approx}} \begin{array}{cccc} x & x & x & x \\ & x & x & \end{array} \xrightarrow{\text{mt}} \text{mt}$$

$$y_i(w^T x_i + b) = -0.5 \text{ (approx)} \xrightarrow{\text{mt}} \begin{array}{cccc} x & & & \\ x & x & x & x \\ x & x & x & \end{array} \xrightarrow{\text{mt}} \text{mt}$$

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

ξ_i (slack i)

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

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

Since we add a new variable ξ_i & the points s.t. if

the point lies in the region $\xi_i = 0$

-ve " " " -ve " $\xi_i = 0$

~~too~~ point ~~too~~ another

ξ_i for all other point is +ve.

As $\xi_i \uparrow$ the point is further away from its correct hyperplane in incorrect direction.

$$\forall x_i \rightarrow \xi_i = 0 \text{ if } y_i(w^T x_i + b) \geq 1$$

$\xi_i \geq 0$ & if is equal to some threshold and away from correct hyperplane in incorrect direction.

we also know that $\arg \max f(u) = \arg \min \frac{1}{f(u)}$

$$w^*, b^* = \arg \min_{w, b} \frac{\|w\|_2^2}{2} + C \cdot \frac{1}{n} \sum_{i=1}^n \xi_i;$$

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

we want to minimize error or misclassification and ξ_i is > 0 for incorrectly classified point.

C is hyperparameter.

$$w^*, b^* = \arg \min_{w, b} \left(\frac{\|w\|_2^2}{2} + C \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \right)$$

hyperparameter
 \rightarrow distance of misclassified point
 \rightarrow average distance of misclassified point

\downarrow margin
 \rightarrow s.t. $y_i(w^T x_i + b) \geq 1 - \xi_i, \forall i$ and $\xi_i \geq 0$
 \rightarrow Loss: This has to be minimized. (Hinge loss)

\rightarrow This can be think as regularization

As $C \uparrow$ more weightage is given to loss term, hence for even a small error/loss the term will become more. Hence there is a tendency to make mistake decrease hence we get overfit.

If $C \downarrow$ the more importance is given $\frac{\|w\|_2^2}{2}$, then we have underfit.

The above formulation of sum is called soft-margin SVM.

$$w^*, b^* = \arg \min_{w, b} \frac{\|w\|_2^2}{2} + C \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \quad \text{s.t. } y_i(w^T x_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0$$

soft-margin SVM.

* Why $+1$ and -1 for π^+ , π^-

$$\pi: w^T x_i + b \leq \pi^+: w^T x_i + b = (+1) \quad \pi^-: w^T x_i + b = (-1)$$

why?

* w need not be unit vector.

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

① Let's assume instead of π_1 we look K , $K > 0$

$$\pi^+ : w^T x_i + b = K$$

$$\pi^- : w^T x_i + b = -K$$

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

$$\underset{w, b}{\operatorname{argmax}} \frac{2}{\|w\|} = \underset{w, b}{\operatorname{argmax}} \frac{2K}{\|w\|} \quad \text{because } (2 \times 2K \text{ are constant, hence no effect)}$$

∴ we can take any value of K , But for simplicity we take $K=1$

* loss-minimization interpretation :-

while studying LR, we learnt that

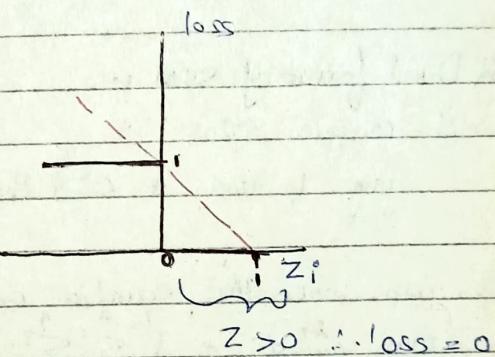
As C-loss is not continuous we took an approximation and that approximation

is logistic loss.

Logistic reg :- logistic loss + reg

Linear Reg :- linear loss + reg

SVM :- hinge loss + reg



z_i in SVM is $y_i(w^T x_i + b)$

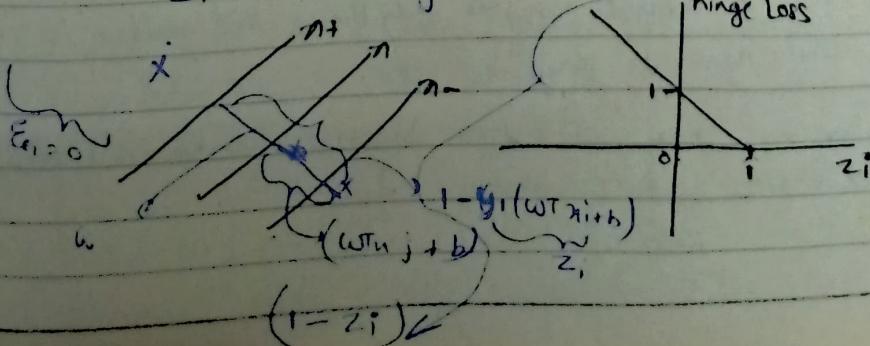
If $y_i(w^T x_i + b) > 0$ that means point is correctly classified. hence loss = 0

when $z_i < 0$:- x_i is incorrectly classified hence in C-1 loss we give it a loss of 1.

And for SVM we use hinge loss as approximation
So, how does our hinge loss behave.

If $z_i \geq 1$, hinge loss = 0

$z_i < 1$, hinge loss = $1 - z_i$



Hence hinge loss is same as ϵ_i

• less min. sum formulation is

$$\min_{w,b} \sum_{i=1}^n \max(0, 1 - y_i(w^T x_i + b)) + \frac{\gamma}{2} \|w\|^2$$

↑ underfit
↓ overfit
Regularization

is same as soft-margin sum, behave opposite (↑ overfit, ↓ underfit)

$$\min_{w,b} \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \epsilon_i \quad \text{s.t. } (1 - y_i(w^T x_i + b)) \geq \epsilon_i, \epsilon_i \geq 0$$

* Dual form of SVM:

soft-margin SVM:

$$\min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \epsilon_i \quad \text{s.t. } y_i(w^T x_i + b) \geq 1 - \epsilon_i, \epsilon_i \geq 0$$

we can write this equation as

$$\max_{x_i} \sum_{i=1}^n d_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n d_i d_j y_i y_j x_i^T x_j \quad \text{s.t. } d_i \geq 0, \sum_{i=1}^n d_i y_i = 0$$

solving both the equation is equivalent.

→ This formulation is called Dual form.

(i) x_i we have d_i

(ii) x_i only occur in form of $x_i^T x_j$

$$(iii) \hat{y}_q = f(x_q) = \sum_{i=1}^n d_i y_i x_i^T x_q + b$$

(iv) $d_i > 0$ only for SV, $d_i = 0$ for non support vectors.

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

, if $d_i = 0$ whole thing is zero.

Hence to compute x_q the only thing that matters is SV and non-SV plays no role.

That's why the name is SVM because only SV play role in classification!

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

s.t. $d_i > 0$ $\begin{cases} d_i = 0 \text{ for non-SV} \\ d_i > 0 \text{ for SV} \end{cases}$

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

$$x_i^T x_j = x_i \cdot x_j = \text{cosine similarity b/w } x_i \text{ & } x_j \quad \left\{ \text{if } \|x_i\| = 1, \|x_j\| = 1 \right\}$$

we can replace $x_i^T x_j$ with any similarity not just cosine-similarity.
That makes the sum super popular because we can plug in
any similarity measure in dual form equation.

after this we use $K(x_i, x_j)$ it is called Kernel func.

$$\max_{d_i} \sum_{i=1}^n d_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n d_i d_j y_i y_j K(x_i, x_j)$$

Kernel func

$$f(x_2) = \sum_{i=1}^n d_i y_i K(x_i, x_2) + b$$

* Kernel trick :-

o Dual form:-

$$\max_{d_i} \sum_{i=1}^n d_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n d_i d_j y_i y_j x_i^T x_j \quad \text{s.t. } \sum_{i=1}^n d_i y_i = 0; d_i > 0$$

Replacing $x_i^T x_j$ with $K(x_i, x_j)$ is
called Kernel trick.

can be replaced by any similarity
func. One popular class of
func is Kernel.

→ The most important idea in SVM is Kernel-trick.

→ if we not replace $x_i^T x_j$ with any other Kernel then it is called
Linear sum. otherwise it is called Kernel SVM.

⇒

Linear-sum:- find margin-max hyperplane in x_i 's space
Logis-reg:- min logistic loss in x_i 's space.

Both of them approx give similar results. But the world
changing idea in SVM become the Kernel Trick.

Generalization:-

In linear ^{SVM} space we try to separate with hyperplane. But
if data is not linearly separable, linear SVM will fail.
In Logistic regression also we need feature transform to succeed.
But Kernel sum will succeed if we identify the
correct Kernel.

Kernel trick do a similar task as that of feature engineering
in logistic regression.

Hence Kernel-SVM can separate non-linear ~~data~~ also.

Note * There is a diff b/w feature engineering & Kernel-trick

* Polynomial Kernel:-

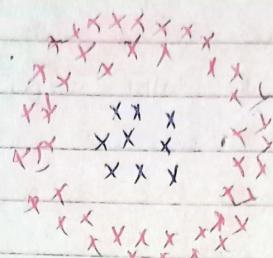
In logistic Reg we had seen that we can apply feature transformation

$$f_1 = f_1^2$$

$$f_2 = f_2^2$$

to make the data linearly separable

Now, let's look how Kernelization works on this.



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

Ex:- Quadratic Kernel.

$$\begin{aligned} K(x_1, x_2) &= (1 + x_1^T x_2)^2 && \text{let } x_1 : \langle x_{11}, x_{12} \rangle \\ &= (1 + x_{11}x_{21} + x_{12}x_{22})^2 && x_2 : \langle x_{21}, x_{22} \rangle \\ &= 1 + \cancel{x_{11}^2} \cancel{x_{21}^2} + \cancel{x_{12}^2} \cancel{x_{22}^2} + \underline{2x_{11}x_{21}} + \underline{2x_{12}x_{22}} + \underline{2x_{11}x_{22}} \cancel{x_{12}x_{21}} \\ &\rightarrow [1, x_{11}^2, x_{12}^2, \sqrt{2}x_{11}, \sqrt{2}x_{12}, \sqrt{2}x_{11}x_{12}] \rightarrow x'_1 \\ &\rightarrow [1, x_{21}^2, x_{22}^2, \sqrt{2}x_{21}, \sqrt{2}x_{22}, \sqrt{2}x_{21}x_{22}] \rightarrow x'_2 \end{aligned}$$

$$\rightarrow (x'_1)^T (x'_2)$$

$$\begin{matrix} x_1 & x_2 & = x_1^T x_2 \\ \downarrow & \downarrow & \downarrow \\ x'_1 & x'_2 & = x'_1^T x'_2 \end{matrix}$$

Kernelization :- To a internal feature transformation.
we don't do external feature transformation.

Kernel trick :- take a d dim data and change it to d' data internally using feature transform.

$$\begin{matrix} d & \rightarrow & d' \\ d' & \rightarrow & d \end{matrix}$$

But we need right kernel.

* Radial Basis Function (RBF) Kernel -

These are the most imp & general kernel in SVM:-

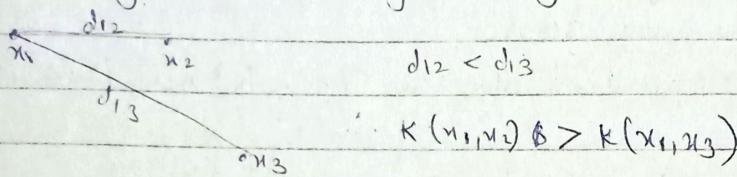
$$K_{RBF}(x_1, x_2) = \exp\left(\frac{-\|x_1 - x_2\|^2}{2\sigma^2}\right)$$

$x_1 \rightarrow \curvearrowright \quad x_2 \rightarrow \curvearrowright$

$\|x_1 - x_2\|^2 \Rightarrow$ square of distn
b/w x_1, x_2

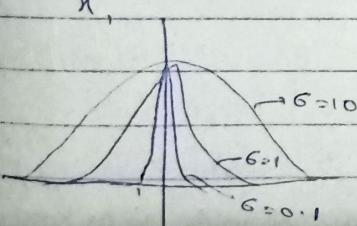
σ = hyper parameters.

as distance increase $K_{RBF}(x_1, x_2) \downarrow$ of $K(x_1, x_2) = \frac{1}{e^{(||x_1 - x_2||^2)/\sigma^2}}$
This is behaving like some sort of similarity.



That means if point are very far, $K(x_1, x_2)$ approaches 0 i.e. similarity becomes 0

If σ increases, more points here are similar to here if σ is very large then all the points will be similar and hence overfit.
If σ is very less then all the points will be non-similar hence underfit.



σ can be treated as K-NN.

K-NN is similar to RBF Kernel - SVM.
Hence RBF is said general Kernel SVM because it is very much similar to KNN.

If we remember, time complexity was only drawback of KNN, but RBF Kernel do a similar task with very less time complexity.

So RBF SUM are nice approximation to KNN.

→ If you don't know what Kernel to use, just simply use RBF.

but remember, if we have a \rightarrow hyperparameters if we use RBF, i.e. σ & C from soft margin SVM equation,

* Domain specific Kernels:-

There are lot of Domain specific Kernels:- like
String kernel for text classifiers.

Genome kernel for genome prediction problem.

Graph Kernel for Graph based problems.

* Time complexity:-

Train:- \rightarrow using SGD

\rightarrow specialized algo like (sequential minimal optim) (smo)

$\approx O(n^2)$ for Kernel SVMs

when n is very large $O(n^2)$ is very time consuming

$$\text{Test/Runtime} := f(x_2) = \sum_{i=1}^{n-1} y_i k(x_i, x_2) + b$$

$\rightarrow O$ for non-SV

$\cdots O(kd)$
 \downarrow # of SV dimensionality.
 $1 \leq k \leq n$

* nu-SVM:- standard SVM formulation with C as hyperparameter

Here we only said $C \geq 0$ nothing else. There is alternative formulation of SVM, called nu-SVM

nu lies b/w 0 & 1 $0 \leq \text{nu} \leq 1$

• nu is an upperbound of fraction of error

nu is an lowerbound of Support vectors.

Ex:- we don't want $> 10\%$ error. we simply set nu=0.1

nu max fraction of error that can happen.

If $\text{nu} = 0.01 \Rightarrow \% \text{ of error} \leq 1\%$.

of SV $\geq 1\% \text{ of } n \text{ data points}$.

* Support Vector Regression (SVR)

We can do SVR with slight changes.

$$\min_{w,b} \frac{1}{2} \|w\|^2 \quad \text{s.t. } y_i - (w^T x_i + b) \leq \epsilon \\ (w^T x_i + b) - y_i \leq \epsilon \\ \epsilon > 0$$

This is the linear formulation of SVR, we can also have kernel formulation.

ϵ is hyper parameter.

If $\epsilon \downarrow$ then errors are low on training data : overfitting ↑

$\epsilon \uparrow$ underfitting may increase.

* Cases:-

1. Feature engg & FT :- Kernel design.

2. Decision surface :- Hyperplane for linear sum, non-linear surface ~~for RBF~~ using Kernel

3. Similarity & Distance func :- convert to kernel & then plug it in SVM.

4. Interpretability & Feature imp :- No implicit way, hence forward feature selection is necessary.
Kernel case.

In linear SVM it is easy to get features.

5. Outliers :- very little impact.

But if we are using RBF with small σ then it may impact.

6. Bias-variance tradeoff :- As $c \downarrow$ overfit, $c \uparrow$ underfit.

7. Large d :- It's good for SVM. They work very well in higher dimension.

8. Bad case :- If we have a right kernel

9. Worst case :- ① n is large \rightarrow train time is high,
② K is large $\# \{K = \text{no. of SV}\}$

CH Decision Trees

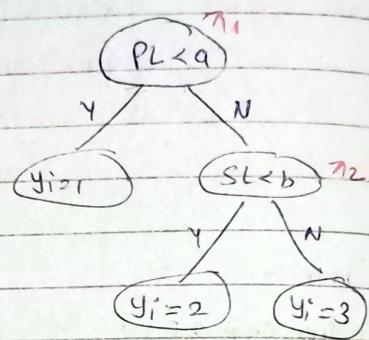
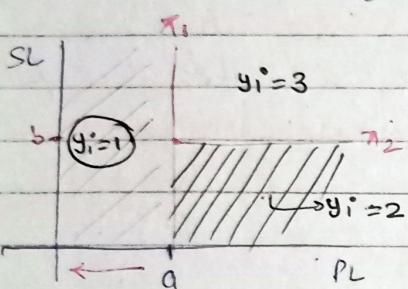
* Geometric Intuition:- Since we have learned algorithms like

KNN, NB, Log reg, Lr, reg, SVM
 (instance based) ↓
 (probabilistic) ↓
 (geometric, because they try to find hyperplane)

Decision Tree is ~~some~~ what like if the condition.

DT is nested If else classifier.

Let Example DT of this dataset

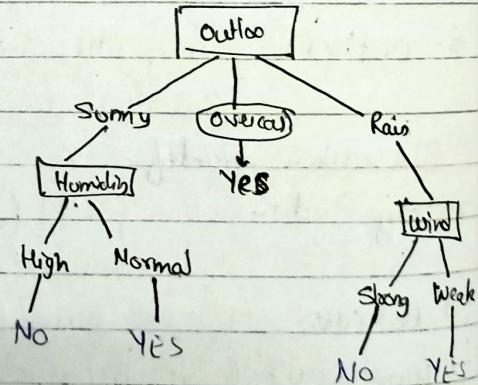


• all the hyperplanes are axis parallel.

∴ DT is set of axis parallel hyperplane.

Example:-

	outlook	Temp	Humidity	Windy	Play
S	Sunny	Hot	High	False	No
S	H	H	H	True	No
O	Overcast	H	H	False	Yes
R	Rainy	Mild	H	F	Yes
R	Cool	Normal	N	F	Yes
R	C	N	T	F	No
O	C	N	T	F	Yes
S	M	H	F	F	No
S	C	N	F	F	Yes
R	M	N	F	F	Yes
S	M	N	T	T	Yes
O	M	H	T	F	Yes
O	H	N	F	T	No
R	M	H	T	T	No



x_i

$y_i \in \{0, 1\}$

* Building a Decision Tree: In order to build a DT, we need to have understanding of some of the concepts

i. Entropy: Let we have random variable y which can take K values y_1, y_2, \dots, y_K . Then entropy of y is

$$H(y) = - \sum_{i=1}^K p(y_i) \log_b(f(y_i)) \quad \begin{cases} b=2 \text{ or} \\ b=e=2.718 \end{cases}$$

$$p(y_i) = P(Y=y_i)$$

In our previous example, suppose $y =$, $p(y_i) = \{y_1, y_2\}$

$$p(y=y_1) = \frac{9}{14} \quad p(y=y_2) = \frac{5}{14}$$

$$H(y) = - \left[\frac{9}{14} \log \left(\frac{9}{14} \right) \right] - \left[\frac{5}{14} \log \left(\frac{5}{14} \right) \right] \rightarrow \text{value of } -\infty \text{ points}$$

$$\frac{\text{# of true points}}{\text{Total # points}} H(y) = 0.94$$

\approx of true points

Properties of Entropy: - If $y \in \{0, 1\}$ (2 class R.V)

Case 1: If 99% of points are y_1 and only 1% of points are y_2 in given dataset.

$$H(y) = -0.99 \log(0.99) - 0.01 \log(0.01) = 0.0801$$

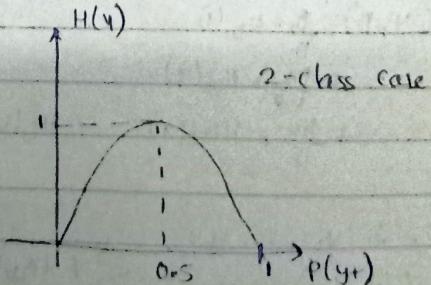
$$\text{Case 2: } D \begin{cases} y_1 \rightarrow 50\% \\ y_2 \rightarrow 50\% \end{cases} \quad H(y) = -0.5 \log 0.5 - 0.5 \log 0.5 = 1$$

$$\text{Case 3: } D \begin{cases} y_1 \rightarrow 0\% \\ y_2 \rightarrow 100\% \end{cases} \quad H(y) = -0 \log(0) - 1 \log(1) = 0$$

Now, lets look at maximum multiplicity case.

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

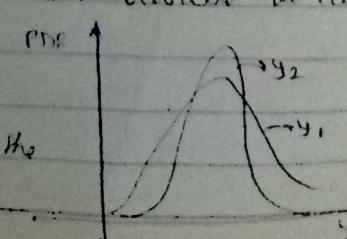
If all these equally probable
Then entropy is maximum.



By seeing the PDF of a random variable we can guess the entropy.

$$H(y_2) < H(y_1)$$

more the peaked PDF less will be the entropy



② Information gain :- Now took back to our example dataset.

Original dataset:

Outlook	Temp	Hum	Wind	Play
S	H	H	F	No
S	H	H	T	No
S	M	H	F	No
S	C	N	F	Yes
S	C	N	T	Yes

Information gain calculations:

- For 'sunny':
Subset 1 (Sunny):

Outlook	Temp	Hum	Wind	Play
S	H	H	F	No
S	H	H	T	No
S	M	H	F	No
S	C	N	F	Yes

 $E = 0$
- For 'overcast':
Subset 2 (Overcast):

Outlook	Temp	Hum	Wind	Play
O	H	H	F	Yes
O	C	N	T	Yes
O	M	H	T	Yes
O	H	N	F	Yes

 $E = 0.97$
- For 'rainy':
Subset 3 (Rainy):

Outlook	Temp	Hum	Wind	Play
R	M	H	F	Yes
R	C	N	F	Yes
R	C	N	T	No
R	M	N	F	Yes
R	M	H	T	No

 $E = 0.97$

Average information gain:

$$\text{Average} = \frac{0.97 + 0 + 0.97}{3} = 0.97$$

Information gain: $IG(Y, \text{outlook}) = \frac{(5 \times 0.97) + 4 \times 0 + 5 \times 0.97}{14} = 0.94$

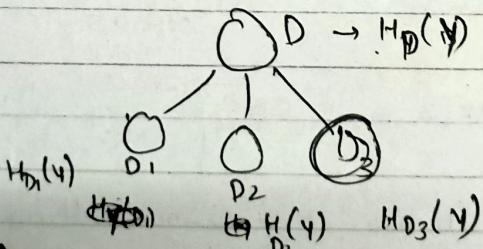
Annotations for weighted entropy calculation:

- Number of points in that division
- Total points
- Entropy of that division
- Weighted entropy
- Entropy of whole Data set

$$IG(Y, \text{outlook}) = 2 \times \frac{5}{14} \times 0.97 - 0.94$$

~~IG(Y, D)~~

$$IG(Y, D) = \frac{D_1}{D} \cdot H_{D_1}(Y) + \frac{D_2}{D} \cdot H_{D_2}(Y) + \frac{D_3}{D} \cdot H_{D_3}(Y) - H_D(Y)$$



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

$$IG(Y, \text{var}) \Rightarrow \sum_{k=1}^K \frac{|D_k|}{|D|} H_{D_k}(Y) - H_D(Y)$$

② Gini Impurity:- It is very similar to entropy

$$I_G(y) = 1 - \sum_{i=1}^x (P(y_i))^2$$

$$(x = y \in \{y+, y-\})$$

(case 1:- $P(y+) = 0.5$ $P(y-) = 0.5$)

$$I_G(y) = 1 - (0.5)^2 + (0.5)^2 \Rightarrow 0.5$$

$$H(y) = 1$$

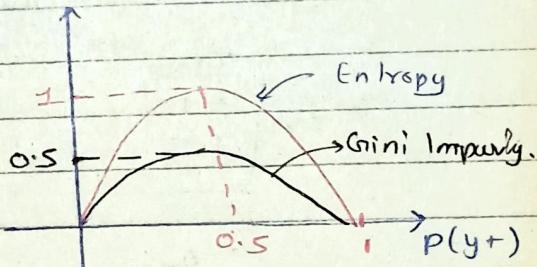
(case 2:- $P(y+) = 1$ $P(y-) = 0$)

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

$$H(y) = 0$$

As we can see that both entropy and gini impurity are very much similar only diff is max value of entropy is 1 and gini val = 0.5

So, why we use gini impurity.



$I_G(y)$ $1 - \{ P(y^+)^2 + P(y^-)^2 \}$	$H(y)$ $-P(y+) \log P(y^+) - P(y-) \log (P(y^-))$
---	--

Entropy need log calculation which is more time taking than square calculation. Hence $I_G(y)$ is more computationally efficient than entropy.

→ Construct a DT :- In our dataset of play or not we have

$$D \begin{cases} \text{Sunny} \rightarrow 9 \\ \text{Overcast} \rightarrow 4 \\ \text{Rainy} \rightarrow 5 \end{cases} \quad H_D(y) = 0.94$$

Now, we have 4 features to choose using which we want to break the Data:-

① Outlook

- Sunny $\begin{cases} \text{Yes: } 2 \\ \text{No: } 3 \end{cases} \rightarrow H_{D1}(y) = 0.97$
- Overcast $\begin{cases} \text{Yes: } 4 \\ \text{No: } 0 \end{cases} \rightarrow H_{D2}(y) = 0$
- Rainy $\begin{cases} \text{Yes: } 3 \\ \text{No: } 2 \end{cases} \rightarrow H_{D3}(y) = 0.97$

$$I_G \Rightarrow \text{weighted Entropy} - H_{D1}(y) = 0.25$$

② Temperature

- Hot $\begin{cases} \text{Yes: } 2 \\ \text{No: } 2 \end{cases} \rightarrow H_{D4}(y) = 0.1 \quad I_G \Rightarrow$
- Mild $\begin{cases} \text{Yes: } 4 \\ \text{No: } 2 \end{cases} \rightarrow H_{D5}(y) =$
- Cool $\begin{cases} \text{Yes: } 3 \\ \text{No: } 1 \end{cases} \rightarrow H_{D6}(y) =$

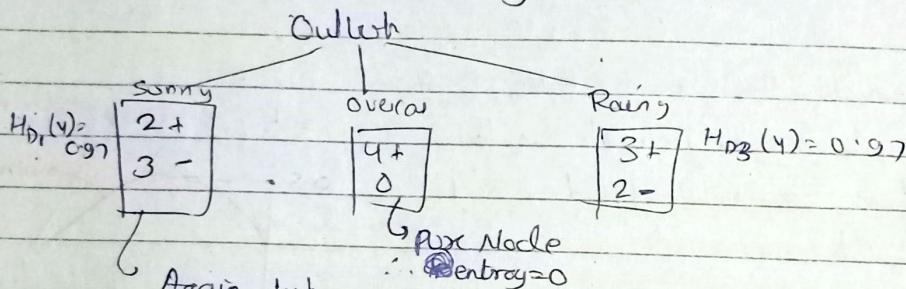
Similarly we can do for windy, humidity:-

We have I_G for each feature.

$I_G(Y, \text{Outlook})$
 $I_G(Y, \text{Temp})$
 $I_G(Y, \text{Windy})$
 $I_G(Y, \text{Humidity})$

, max, pick the feature with most information gain and break the dataset using that feature

Suppose outlook gives most information gain, here we break the dataset using outlook.



Again look at remaining 3 features and break this using features which give most information gain.

At any point we have 3 cases

- (1) pure-node \Rightarrow stop growing the node
- (2) cannot grow the tree anymore because of lack of points.
we stop growing when
 - (i) got a pure node
 - (ii) Node has very few points
 - (iii) if we are too deep.

As depth of tree \uparrow chance of overfitting \uparrow

If depth is \downarrow " underfit \uparrow

In DT, depth is the hyperparameter.

* splitting numerical features:-

Till now we have only seen categorical feature now, let's

A see how to deal with numerical features.

2.2	1
2.6	1
3.5	0
3.8	0
4.2	1
5.3	0

- (1) sort the numerical feature in ascending order.
- (2) we check for each data point as splitting criteria

$$\begin{aligned} f < 2.2 &\rightarrow \text{y}_0 \\ f < 2.6 &\rightarrow \text{y}_0 \\ f < 3.5 &\rightarrow \text{y}_0 \\ f < 3.8 &\rightarrow \text{y}_0 \\ f < 4.6 &\rightarrow \text{y}_0 \\ f < 5.3 &\rightarrow \text{y}_0 \end{aligned}$$

} max Again Recursively do the same thing on smaller dataset.

* Feature standardisation:-

Till now in each Alg of Logistic Regr., SVM, etc) we have done feature standardisation. But in the case of Decision Tree we are not computing distance i.e. it is not a distance based algorithm. Here we don't need feature standardization because we don't care about scales here.

* Categorical features with many categories:-

Suppose we have a feature pincode, it is a numerical feature but really it is a categorical feature. It may have 1000's of categories. There may be some categories which may contain very few points.

We can apply feature engineering:-

pincode/zipcode \rightarrow categorical

\downarrow
convert to numerical feature.

So how do we convert it to numerical feature.

$$P(y_i=1 | \text{pincode } j)$$

i.e. probability of $y_i=1$ given my pincode j

Pincode	Yield
P ₁	Y ₁
P ₂	Y ₂
⋮	⋮
P _n	Y _n

$\rightarrow H(y_i=1 | P_i)$ Now, we are converting every pincode j into numerical feature which

is equivalent to probability of $y_i=1$ given the pincode j . Now, converting pincode to probability is very beneficial. we can treat them as normal numerical feature.

* Overfitting & underfitting :-

As depth ↑ :- possibility of having very few points @ a leaf node increases and hence they may be noisy points and we are attempting to correctly classify them and hence overfitting happens.
As depth increases interpretability decreases.

(i) if depth is very low then we are underfitting.

The correct depth can be predicted using cross validation.

* Time & space complexity

Training :- Time : $O(n(\log n)d)$ At every stage we need to evaluate each feature IG.
 $n = \# \text{ points}$
 $d = \text{dim}$
Sorting of numerical features

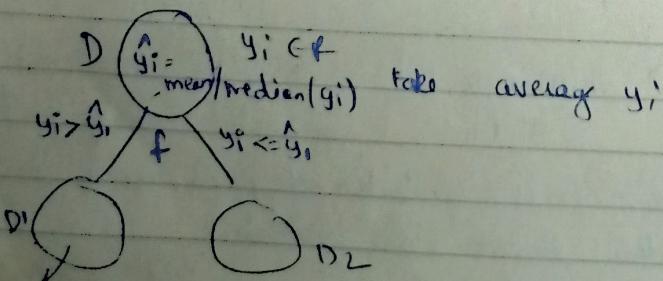
Testing: Time :- $O(D)$ → Depth of tree.

Space :- $O(\text{no. of nodes in tree})$ ~~nodes~~

∴ DT are good for, large data, dim is small, low latency

* Regression using DT:-

We have seen DT as a classifier. We can extend DT as Regression. We don't use IG here, instead of that we use Mean square error or MAD to break the data.



We use ~~use~~ check the feature with lowered MSE or MAD

* Cases:

- ① Imbalanced data :- balance it using upsampling or down sampling because Imbalance data impact entropy calculation
- ② Large dim :- @ each node we have to evaluate each feature
Time complexity increases
avoid one-hot encoding for categorical feature because suppose feature have 10 categories then 1 feature will be converted to 10 features :-
~~If feature have lot of categories it's useful to convert them to numerical feature of probability.~~
- ③ Similarity matrix :- DT cannot handle similarity matrix; ~~This~~
They need data points
- ④ Multiclass class :- DT can handle multi-class. If there are more than one type of points at a node choose the majority one.
- ⑤ decision surface :- Axis parallel types cuboids.
- ⑥ feature interaction :- logical feature interaction is inbuilt in DT like if outlook=sunny & humidity then feature interaction plays Yes
- ⑦ Outlier :- if depth is more ; outlier impact the DT.
- ⑧ Interpretability :- ~~DT~~ are very much interpretable.
- ⑨ feature importance :- We can do feature importance by calculating the reduction in entropy at each node - Sum up all the reduction due to a feature , if sum is large then it is more important.

CH- Ensemble models

→ Ensembles = In english ensemble means group of things. In ML ensemble means multiple model & combined together to build a powerful model.

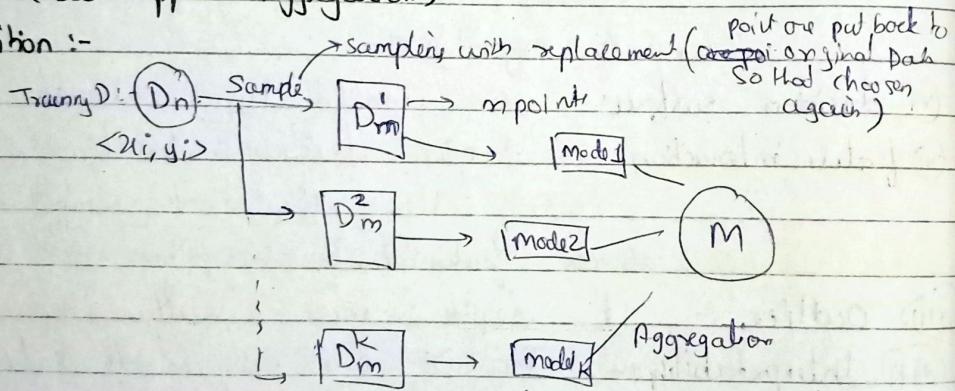
There are 4 types of Ensembles

1. Bagging (Bootstrapped Aggregation)
2. Boosting
3. Stacking
4. Cascading

Key concept = more different the models are, the better you can combine them.

① Bagging (Bootstrapped Aggregation)

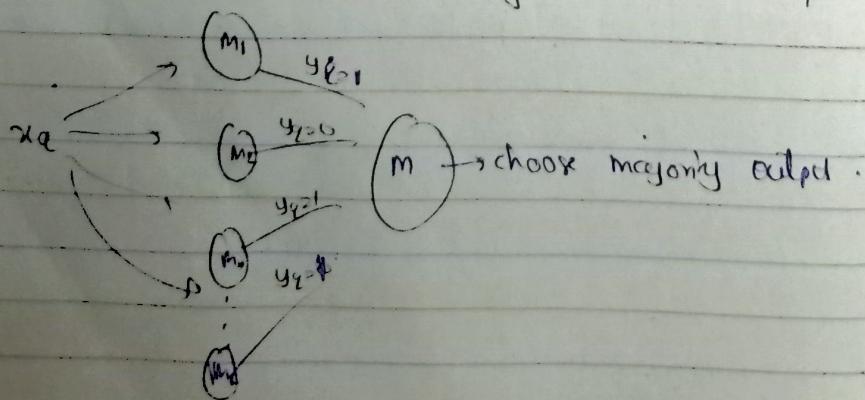
* Intuition :-



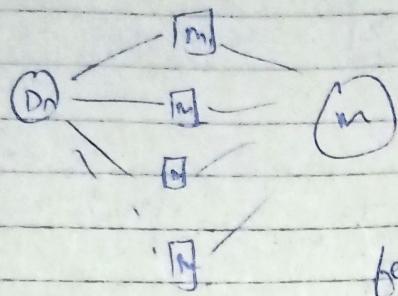
Each model m_i is built using D_m^i of size m ($m \leq n$)
⇒ each model trained on different subset of data.

A typical way of Aggregation is in classification :- Majority vote
in Regression := Mean/modes.

Ex:-



- * Bagging can reduce variance in model: Let see how variance is \rightarrow how your model change as the dataset changes.



Each of the models are not seeing the whole data, they are based only on a small subset of data. If the data set is changed then only few of them the model will change and hence there is less chance that our final output after aggregation gets changed. Hence even though if there is change in dataset final output may not change and hence variance decreases.

- * Bagging can reduce the variance in model without impacting the Bias.

Now suppose each of our M_i is low bias high variance model.

** Bagging (M_i 's) \Rightarrow low bias, reduced variance model.

i.e. Bagging says take a bunch of low bias & high variance model & combine them using Bagging and you get a low bias and reduced variance model.

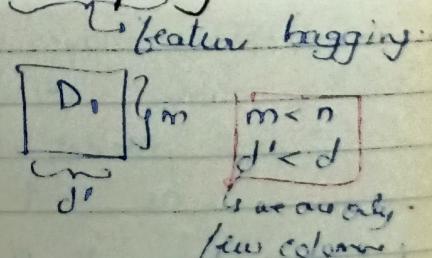
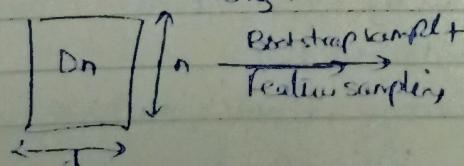
Ex:- of a low bias, high variance model is DT.

\rightarrow Random Forest:- it is the most popular Bagging model.

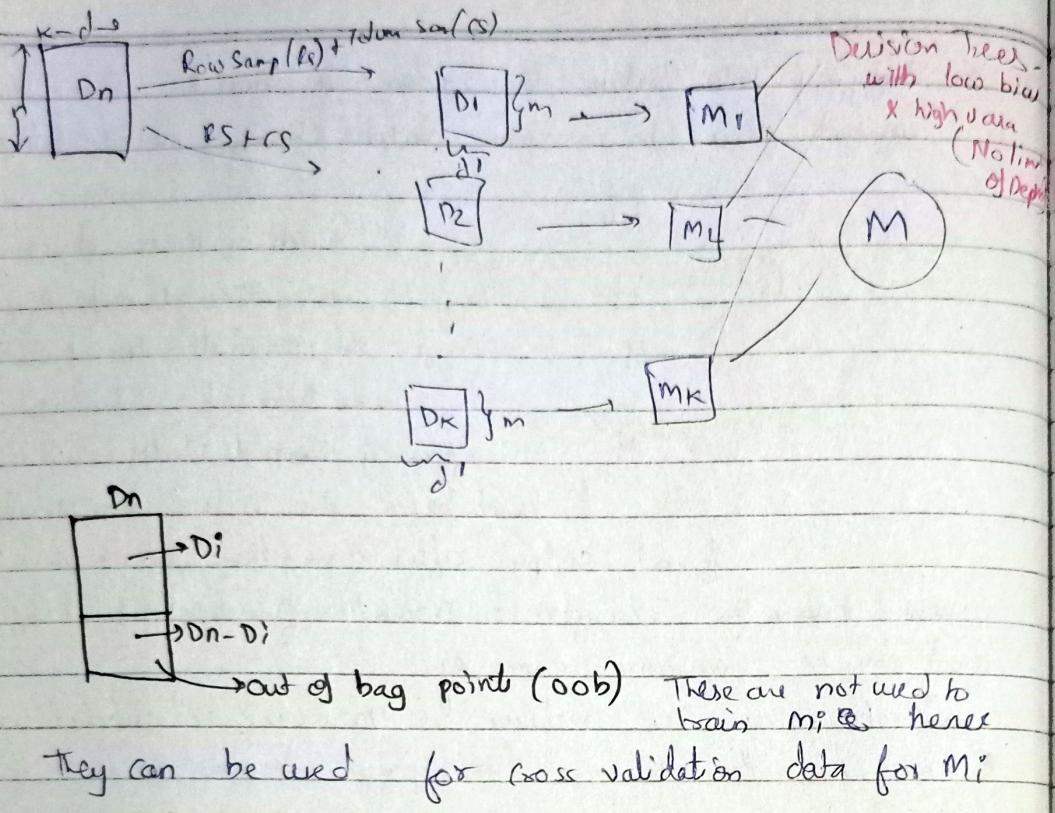
RF take DT as base model and apply bagging and also apply a unique tech called column sampling.

RF = DT + Bagging + col. sampling

Total Data



The idea of RF is, to sample both rows & cols.



Ex:- Suppose ~~for~~ for Model 1, we chose D_1 as D_1 is subset of D_n then there must be some point which are not part of D_n . Hence that points can be used as cross validation for Model 1. Same apply for all Model i.

RF: $\cdot D_i + \text{Row sampling} + \text{Col sampling} + \text{Aggregation}$

\downarrow
fully grown
Decision trees

\downarrow
Bagging

\downarrow
Feature sampler
(Majority/mean)

1. No limit on height
to make the model high variance.

→ Bias-variance tradeoff:-

RF : have very low bias because base model (M_i 's) are low bias.

Model = agg (M_1, M_2, \dots, M_k)
as $k \uparrow$ variance \downarrow

more no. of models used in agg lesser will be the variance.

One more hyperparameter is dimension of sampler.

$$D_n \rightarrow \left[D_i \right]_{i=1}^m \quad d' < d \quad m^d \leq n$$

If these are small ~~and~~ dataset in each model will be more different hence variance ↓

Time complexity :-

Train :- $O(n(\log n) d \cdot k)$ $n \rightarrow$ no. of points

But we can parallelize the training of base model's to reduce the time complexity.

When we have large amount of data with reasonable features. we can train more base model parallelly by creating batches and parallelly train models in batch:-

Test :- Time: $O(\text{depth} \cdot k)$

Space: $O(DT \cdot k)$

* Extremely Random Trees:-

There is a variation of RF
In RF:- row sampling, colⁿ sampling, aggregation

ExtremeRF:- don't we all the values in a numerical feature to for decision making. As we have learned in DT that if a feature is numerical then we sort it and check for each value in that feature if it can give a better threshold. Instead of that extreme tree only uses random sample of possible values to determine threshold.

Extreme tree:- colⁿ sampling + s.s + agg + randomization
when selecting threshold in numerical feature Extreme tree tends to reduce variance better than RF but bias may increase.

(2) Boosting:- (base model) Randomization

In Bagging :- high-variance, low-bias + (CS, RS) → aggregation

Boosting : low-variance, high-bias + additively combine.

$\underbrace{\quad\quad\quad}_{\text{reduce bias while keeping}} \underbrace{\quad\quad\quad}_{\text{variance low.}}$

Core-idea :- how boosting reduces bias

$$\text{Train} = \{x_i, y_i\}_{i=1}^n$$



Step 0-(a) Build a model M_0 on whole of D_{Train} .

M_0 need to have high bias low variance

$(M_0) \leftarrow \{x_i, y_i\}$ for example a DT with small depth,

holu) As we have high bias model we encounter large training error.

(b) calculate error $_i = y_i - \hat{y}_i$ {Regression}

∴ At point in train data x_i, y_i we have error $_i$
 $\{x_i, y_i, \text{error}_i\}_{i=1}^n$

Step 1: - train model M_1 with take $\{x_i, \text{error}_i\}$ as its training data.

$$(M_1) \leftarrow \{x_i, \text{error}_i\}_{i=1}^n$$

$F_1(n)$ = model at end of step 1.

$$f_1(x) = \alpha_0 h_0(n) + \alpha_1 h_1(x)$$

Step 2: - $\{x_i, \text{error}_i\}$

$$\downarrow y_i - f_1(x_i)$$

$$F_2(x) = \alpha_0 h_0(n) + \alpha_1 h_1(n) + \alpha_2 h_2(x)$$

$$\text{Step } k: F_k(n) = \sum_{i=0}^k \alpha_i h_i(x)$$

$\underbrace{\quad\quad\quad}_{\text{additive weights model}}$

Each model is trained to fit the residual error at end of the previous stage.

$$F_k(x) = \sum_{i=1}^n d_i h_i(x)$$

ends up having a low residual error.

Train error ~~decreases~~ reduces because at each step we try to fit the residual error & hence bias reduces.

There are many boosting algorithm some of them are:

- ① Gradient Boosted DT (GBDT)
- ② AdaBoosting (Adaptive boosting)

* Residuals, Loss-func. & Gradients:-

$$F_k(x) = \sum_{i=1}^n d_i h_i(x)$$

residuals:- ~~error~~ residual at end of stage K is:

$$\text{err}_i = y_i - F_k(x_i)$$

↳ residual.

$$\text{loss-func.} : L(y_i, F_k(x)) = (y_i - F_k(x_i))^2$$

$$\frac{\partial L}{\partial F_k(x)} = -2 * (y_i - F_k(x_i))$$

↳ residual.

$$\frac{\partial L}{\partial F_k(x)} = -2 * \text{residual}$$

$$\frac{\partial L}{\partial F_k(x)} = 2 * \text{residual.}$$

negative deviate/gradient

* negative gradient \approx residual.
of loss func. w.r.t. $F_k(x)$.

↳ pseudo residual

Now instead of using residual we can train of pseudo-residual.

The advantage pseudo-residual helps us to choose any loss func.
as long as the loss func. is differentiable.

for example RF cannot be used to min hinge loss.
 But gradient boosting can be used to minimize any loss
 That's a big advantage.
 So algorithm now looks like this
 At any stage i : we have already built a model $F_{i-1}(x)$

$$(m_i) \leftarrow (x_i, e_{\text{ori}})$$

↳ pseudo-residual
 $\frac{\partial L}{\partial F_{i-1}(x)}$

Gradient Boosting:-

Algorithm:

Input: training set of x_i, y_i, γ_i , a differentiable loss funcn $L(y, F(x))$, number of iteration M

Algo: # of base models (m_1, m_2, \dots, m_M)

1. Initialize model with a constant value:

$$F_0(x) = \underset{\gamma}{\operatorname{arg\min}} \sum_{i=1}^n L(y_i, \gamma) \quad || \text{find a constant } \gamma \text{ s.t error on training data is minimized}$$

Learn all model

2. for $m = 1$ to M :

1. Compute so-called pseudo-residuals:

$$\gamma_{im} = - \left[\frac{\frac{\partial L(y_i, F(m_i))}{\partial F(x_i)}}{\partial F(x_i)} \right] \quad \text{for } i=1 \text{ to } n$$

$F(x) = F_{m-1}(x)$

2. Fit a base learner (e.g tree) $h_m(x)$ to pseudo-residuals

i.e. train it using $\{x_i, \gamma_{im}\}_{i=1}^n$.

3. Compute multiplica γ_m by solving 1-Dim optimiza prob.

$$\gamma_m = \underset{\gamma}{\operatorname{arg\min}} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i))$$

4. Update the model:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x)$$

5. Output $F_m(x)$

Ex ↗ Logistic loss
 ↗ Squared loss
 ↗ Hinge loss

So, what a base learner can be, the most popular learner is Decision Tree. we get Gradient Boosting Decision Tree. Decision Trees in GB DT have very small depth.

* Regularization by shrinkage:-

so the formulation of Gradient Boosting

$$F_m(x) = h_0(u) + \sum_{i=1}^m \gamma_i h_i(x)$$

$m = \# \text{ base-models}$

as $\# \text{ base-models} \uparrow$, variance \uparrow & has overfit.
But there is a concept called shrinkage.

An important part of gradient boosting methods is regularized by shrinkage which consists in modifying the update rule as follows:

$$f_n(x) = F_{m-1}(x) + v \cdot \gamma_m h_m(x) \quad 0 \leq v \leq 1$$

where v is called learning rate.

Empirically it has been found that small learning rates (such as $v < 0.1$) yield dramatic improvement in model's generalization ability over gradient boosting without shrinking ($v \geq 1$). However, it comes at the price of increasing computational time both during training & querying: lower learning rate req. more iterations.

So we have 2 hyper parameters M & v .

* Time complexity:-

Train : Time: $O(n \log n d \cdot M)$

But remember RF was trivially parallelizable but GB DT is not easy to parallelize. Because this is a serial algo.

Predict : Time: $O(\text{depth} \cdot m)$

↳ small in GB because in GB DT we take DT with very small depth.

Space: $O(DT \cdot m + \gamma_m)$

* XGBoost + Randomization:-

XGBoost : ~~GBDT~~ is an implementation of GBDT, in its implementation, in addition to normal GBDT, it also includes Random Sampling & Column sampling.

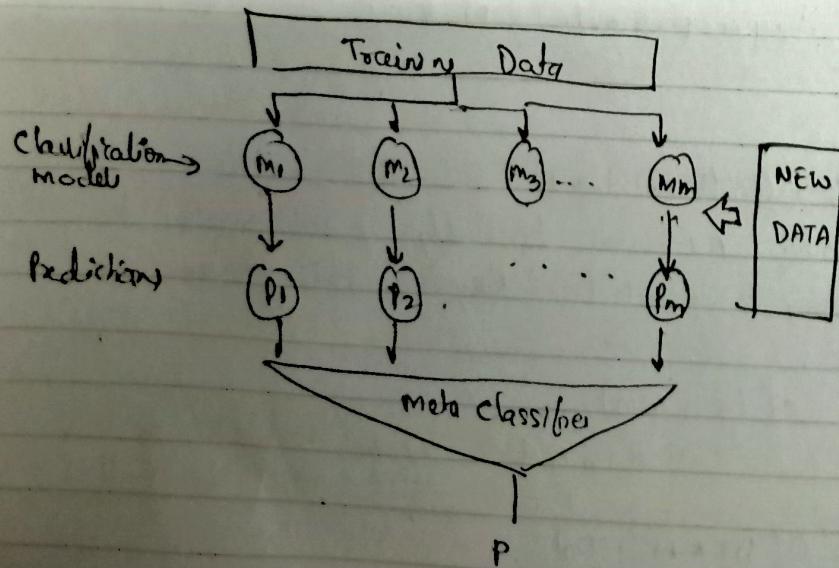
XGBoost implementation is faster than sklearn implementation of GBDT.

It is developed by PhD students of university of Washington.

* AdaBoost :- it is similar to Gradient Boost. It is mostly used in image detection. In case of GBDT we used pseudo-examples. Here all the misclassified points are assigned more weight.

(3) stacking classifier:-

Stacking is an ensemble learning technique to combine multiple classification models via a meta-classifier. The individual classification models are trained based on the complete training set, the meta-classifier is fitted based on the o/p - meta feature - of the individual classification models in ensemble. The meta classifier can either be trained on the predicted class labels or probabilities from the ensemble.



This is same as RF but the diff. is, in RF each base-model is DT but in stacking each of the model can may be different more different. The models are more good the working is. And in RF we have sampled the data but here we create various types of model on whole training data.

Algo: Input: $D = \{x_i, y_i\}_{i=1}^m$ ($x_i \in \mathbb{R}^n$)

Output: An ensemble classifier H .

Step 1: Learn the first level classifier.

for $t \leftarrow 1$ to T do

 Learn a base classifier h_t based on D

Step 2: Construct new data sets from D

 for $i \leftarrow 1$ to m do

 Construct a new dataset contain $\{x'_i, y'_i\}$

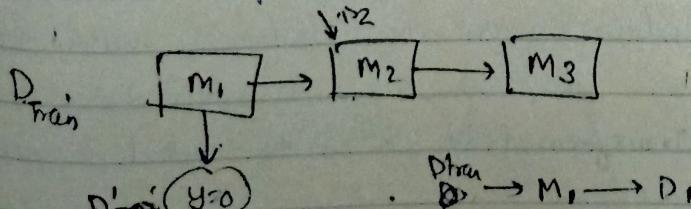
$x'_i = \{h_1(x_i), h_2(x_i), \dots, h_T(x_i)\}$

Step 3: Learn a second level classifier.

 Learn a new classifier h' based on newly constructed data set.

return $H(x) = h'(h_1(x), h_2(x), h_3(x), \dots, h_T(x))$

1) Cascading models:- They typically used when the cost of making a mistake is high i.e. we want a very very accurate prediction, we want to be 100% sure.



$D_{train} - D_1 \rightarrow M_2 \rightarrow D_2$

$D_{train} - D_1 - D_2 \rightarrow M_3 \rightarrow D_3$

We are reducing the training data, we are keeping the data which are correctly classified in previous model. It is mostly used in ~~critical~~ critical situations like card fraud detection in banks and medical diagnosis.