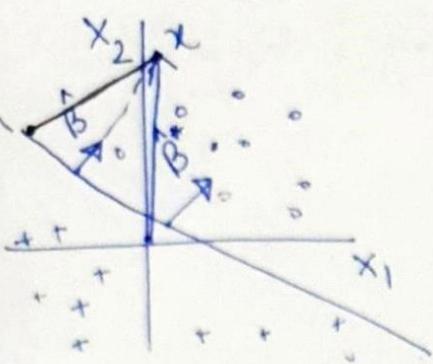


WEEK-4 - Lecture-20 - Perceptron Learning

Modelling a separating hyperplane



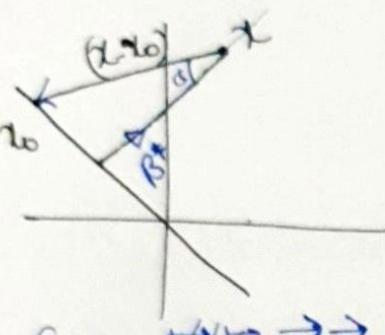
$$f(x) = \beta_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 = 0 = L \\ = \beta_0 + \beta^T x = 0$$

Properties:

$$1) \beta^T (\mathbf{x}_1 - \mathbf{x}_2) = 0 \Rightarrow \beta^* = \frac{\beta}{\|\beta\|}$$

$$2) \beta^T \mathbf{x}_0 = -\beta_0 \quad \forall \mathbf{x}_0 \in L \Rightarrow$$

$$3) \beta^T (\mathbf{x} - \mathbf{x}_0) = (\beta^T \mathbf{x} - \beta^T \mathbf{x}_0) \\ = \frac{1}{\|\beta\|} (\beta^T \mathbf{x} + \beta_0) = \frac{f(\mathbf{x})}{\|\beta\|}$$



$$\cos \theta = \frac{\vec{a} \cdot \vec{b}}{\|\vec{a}\| \|\vec{b}\|}$$

$$\frac{(\mathbf{x} - \mathbf{x}_0) \cdot \beta^*}{\|\beta^*\| \cdot \|\mathbf{x} - \mathbf{x}_0\|} = \cos \theta$$

$$f'(\mathbf{x}) = \beta$$

Distance of \mathbf{x} from hyperplane
Signed distance

$$= \frac{f(\mathbf{x})}{\|f'(\mathbf{x})\|}$$

Perceptron Learning Algorithm

Minimize distance of misclassified points to decision boundary.

$$y_i = 1, \text{ is misclassified} \Rightarrow \mathbf{x}_i^T \beta + \beta_0 < 0$$

$$y_i = 1, \text{ is misclassified} \Rightarrow \mathbf{x}_i^T \beta + \beta_0 > 0$$

$$y_i (\mathbf{x}_i^T \beta + \beta_0) < 0, \quad y_i = 1 \quad \text{Misclassified}$$

$$y_i (\mathbf{x}_i^T \beta + \beta_0) > 0, \quad y_i = 1 \quad \text{for sample } \mathbf{x}_i$$

For all samples

$$\textcircled{1} = - \sum y_i (\mathbf{x}_i^T \beta + \beta_0) \rightarrow \boxed{\text{Minimise}}$$

Assumption

$$\begin{cases} y_i = 1 \Rightarrow \\ \mathbf{x}_i^T \beta + \beta_0 > 0 \end{cases}$$

True
classification

[1]

Lecture-20

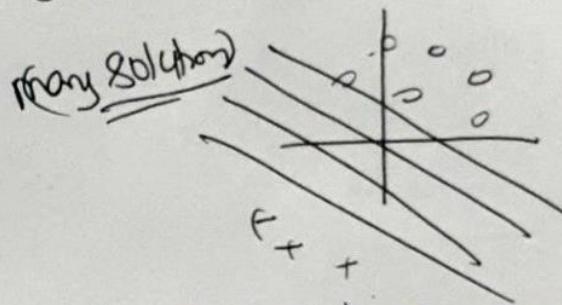
$\triangleright D(\beta, \beta_0) = - \sum_{i \in \text{M}} y_i (x_i^T \beta + \beta_0)$, (M - samples are misclassified)

(Minimize $M \rightarrow 0$) our motive

$$\frac{\partial D}{\partial \beta} = - \sum_i y_i (x_i), \quad \frac{\partial D}{\partial \beta_0} = \sum_i y_i$$

$$\begin{pmatrix} \beta \\ \beta_0 \end{pmatrix} = \begin{pmatrix} \beta \\ \beta_0 \end{pmatrix} - (\alpha) \left[\begin{pmatrix} y_i x_i \\ \sum y_i \end{pmatrix} \right] \quad \xrightarrow{\text{vector}} \quad \xrightarrow{\text{problems}}$$

- Linearly separable: then Converge to some solution.
- Can take along time if gap bet two classes are very less

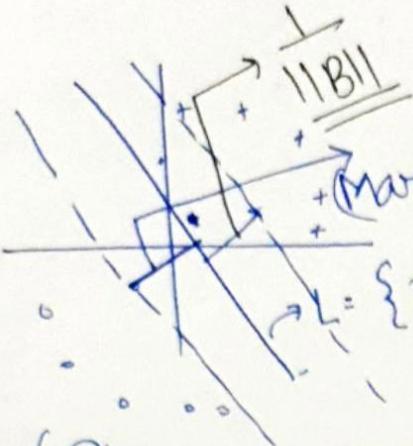


- Not linearly separable \rightarrow it enters in (loop) ∞

②

SVM FORMULATION

Lecture - 21



$$f(x) = \beta_0 + \beta^T x$$

(Margin)

$$\hat{L} = \{x : \hat{\beta}_0 + \hat{\beta}^T x = 0\}$$

Optimal Separating hyperplane
Maximize the distance
of closest point to
hyperplane.

(Thickness should be same for both hyperplanes)

$$\beta_0, \beta_i, \max M$$

$$(||\beta||=1) \quad \text{subject to } y_i(\mathbf{x}_i^T \beta + \beta_0) \geq M \text{ for all } i=1, \dots, N$$

β , should not be
larger.

We need to maximise
M, so that there is
a maximum separation
between classes.

distance of any x_i from hyperplane

$$\frac{f(x)}{\|f(x)\|} = \frac{(\mathbf{x}_i^T \beta + \beta_0)}{\|\beta\|} \rightarrow \text{signed distance}$$

multiply by y_i to make positive

Needn't worry about $\|\beta\|=1$

$$\frac{y_i(\mathbf{x}_i^T \beta + \beta_0)}{\|\beta\|} \geq M \Rightarrow y_i(\mathbf{x}_i^T \beta + \beta_0) \geq M \|\beta\|$$

~~$\|\beta\|$~~ Just
make sure

$$\|\beta\| = 1/M$$

only change in direction

$$\min \|\beta\| \\ \text{s.t. } y_i(\mathbf{x}_i^T \beta + \beta_0) \geq 1$$

$$\max M \\ \downarrow \\ \min \beta$$

[1]

SVM - interpretation & Analysis

$$S.t. y_i(x_i^T \beta + \beta_0) \geq 1 \quad \min_{\beta} \| \beta \| \Rightarrow \left\{ \min \frac{1}{2} \| \beta \|^2 \right\}$$

by our choice
so that it is easily differentiable

Lagrangian:

$$L_P = \frac{1}{2} \| \beta \|^2 - \sum_{i=1}^N \alpha_i [y_i(x_i^T \beta + \beta_0) - 1]$$

Setting derivatives to '0'.

$$\frac{\partial L_P}{\partial \beta} = 0 \Rightarrow \frac{1}{2} 2 \cdot \beta - \sum_{i=1}^N \alpha_i [y_i x_i] = 0 \quad \textcircled{1}$$

$$\beta = \sum_{i=1}^N \alpha_i x_i y_i$$

$$\frac{\partial L_P}{\partial \beta_0} = 0 - \sum_{i=1}^N \alpha_i y_i = 0 \quad \textcircled{2}$$

DQ91:

$$L_D := \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{k=1}^N \alpha_i \alpha_k y_i y_k x_i^T x_k \quad \textcircled{3}$$

Subject to $\alpha_i \geq 0, \forall i$

$$\alpha_i [y_i (x_i^T \beta + \beta_0) - 1] = 0 \quad \forall i \quad \textcircled{4}$$

$\alpha_i = 0$ For non-marginal points \rightarrow which are very far from hyperplane.

\downarrow will not contribute to calculating ' β '.

$\alpha_i \neq 0, x_i$ is on the margin. \rightarrow (support vectors)
 \hookrightarrow only points in the margin plane contribute to β .
 \downarrow points which lies on margin plane

\hookrightarrow Calculate β_0 from here.

$$f(x) = x^T \hat{\beta} + \hat{\beta}_0$$

$$= \sum \alpha_i y_i x_i^T x_i + \beta_0$$

SVM is stable

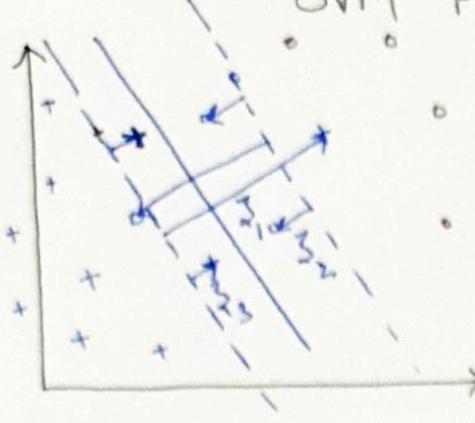
\downarrow less variance

points outside the margin will not affect the SVM. only

Support vectors will affect SVM. [2]

Lecture 23

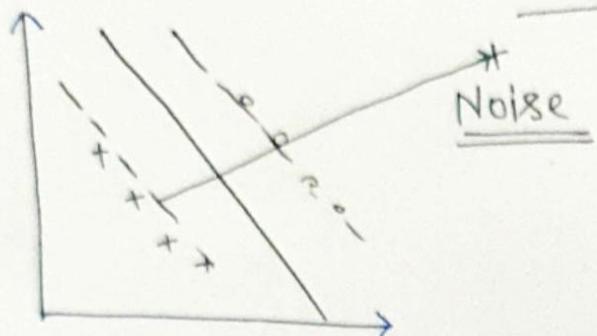
SVM FOR linear Non-separable data



1. Minimize these distances.

$\xi_1, \xi_2, \xi_3 \rightarrow$ misclassified samples

Ques. why not minimize the maximum distance rather than minimizing sum of distances.



Ans: classifier will try to minimize the noise, and hyperplane will get shifted towards Noise-like.

Formulation:

$$y_i(x_i^T \beta + \beta_0) \geq M(1 - \xi_i) \quad \begin{matrix} \rightarrow \text{slack} \\ \text{variable} \end{matrix}$$

1. Ideally we want to maximize M , and $\xi_i = 0$

2. but I want some sort of relaxation in maximizing M . why not $(M - \xi_i)$ Ans: Non-Convex optimization loss problem.

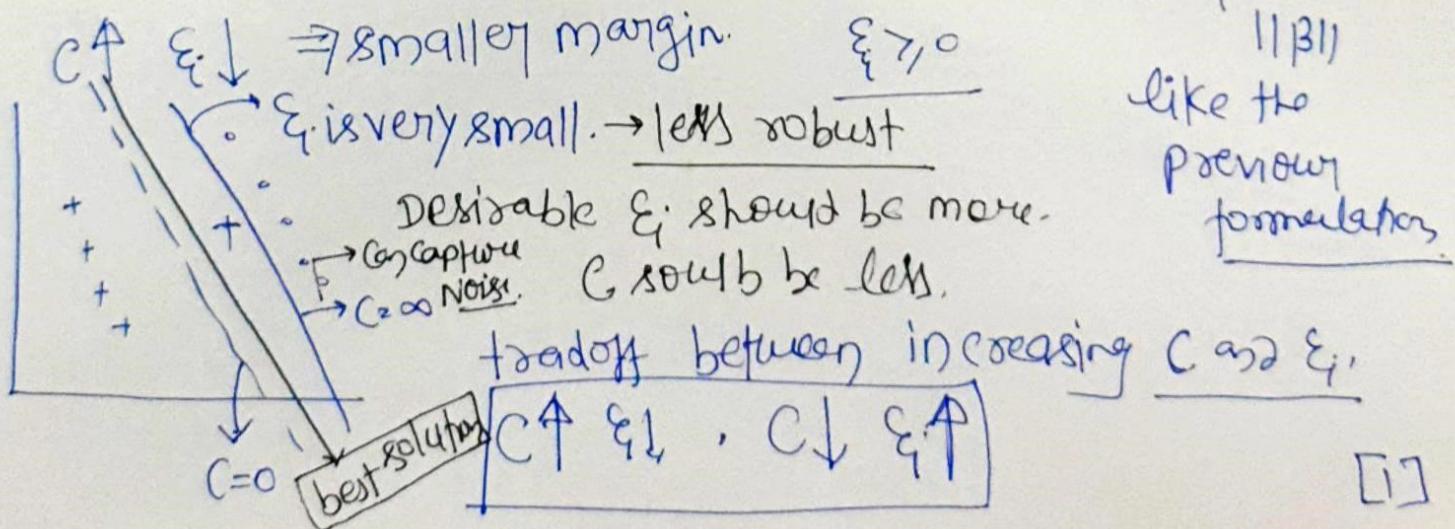
ξ_i : what fraction of margin (M)

$$\textcircled{1} \quad y_i \xi_i \geq 0, \quad \textcircled{2} \quad \sum_{i=1}^N \xi_i \leq \text{constant}$$

we don't want ξ_i to be very large.

$$\min \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^N \xi_i, \quad \text{Sub. to } y_i(x_i^T \beta + \beta_0) \geq (1 - \xi_i)$$

since $M = 1$



like the previous formulation

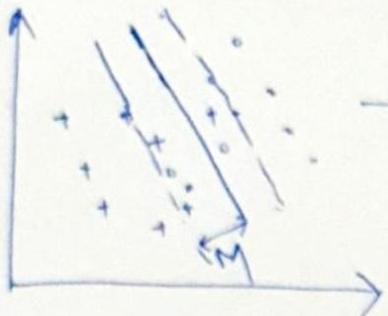
+ tradeoff between increasing C and ξ_i .

$$[C \uparrow \xi \downarrow, C \downarrow \xi \uparrow]$$

[1]

NOTE:

introducing ξ makes SVM more stable and less prone to noisy data. Suppose: we put $\xi=0$, then we will try to maximize M . as shown in figure.

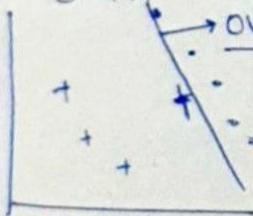


→ We should provide some slackness, because M is calculated via support vectors which lies on the boundary. It doesn't care about the points which lies within the margin. so, we need to provide some slackness in variable M .

by introducing ξ which is distance of points which are within the margin from margin plane. If ~~we make~~ ξ we need sum of ξ ~~to be constant~~ so it shouldn't overshoot.

If we make $C=\infty$, it means all ξ are very small, since classifier will try to minimize loss, in doing so, it has to make all $\xi=0$, \Rightarrow my margin will be very small. overfitting

with Capture Noise also as shown below



If $C=0$, ξ will be larger, more biased SVM classifier.

underfitting

Lecture-24 SVM KERNELS

$$f(x) = \mathbf{x}_i^T \beta + \beta_0 \Rightarrow \sum_{i=1}^N \alpha_i y_i \boxed{\mathbf{x}_i^T \mathbf{x}} + \beta_0$$

inner product.

is there any good way to compute inner product.

$x \rightarrow h(x)$ "Transformation on x "

$$L_D = \frac{1}{2} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle h(x_i), h(x_j) \rangle$$

similarly

$$f(x) = \sum_{i=1}^N \alpha_i y_i \boxed{\langle h(x), h(x_i) \rangle} + \beta_0$$

$$\langle h(x), h(x_i) \rangle = \underbrace{k(x, x_i)}_{\text{Kernel}} - \underbrace{\text{similarity measures}}_{\text{symmetric}}$$

\downarrow $x^T x \geq 0$

\downarrow (+ve) semi definite \downarrow semi definite

"popular choice of k "

$$\text{poly} : (1 + \langle x, x' \rangle)^d$$

$$\text{RBF} : \exp(-\gamma \|x - x'\|^2)$$

$$\text{ANN} : \tanh(\underbrace{k_1 \langle x, x' \rangle}_{\text{constant}} + \underbrace{k_2}_{\text{constant}})$$

POLY

$$(1 + \langle x_i, x \rangle)^d = (1 + x_1 x'_1 + x_2 x'_2)^2, d=2$$

$$X = \boxed{(x_1, x_2)} = \frac{1+2x_1 x'_1}{1} + \frac{2x_2 x'_2}{2} + \frac{(x_1 x'_1)^2}{3} + \frac{(x_2 x'_2)^2}{4} + \frac{2x_1 x'_1 x_2 x'_2}{6}$$

$$h_1(x) = 1, h_2(x) = \sqrt{2} x_1, h_3(x) = \sqrt{2} x_2$$

$$h_4(x) = x_1^2, h_5(x) = x_2^2, h_6(x) = \sqrt{2} x_1 x_2$$

(6-dimension space)

$$L_p = \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N (\alpha_i [y_i(x_i^T \beta + \beta_0) - (1 - \xi_i)]) - \sum_{i=1}^N \lambda_i \xi_i$$

Getting derivatives to 0.

$$\alpha_i, \lambda_i > 0$$

$$\beta = \sum \alpha_i y_i x_i \quad \text{①} \quad 0 = \sum_{i=1}^N \alpha_i y_i \quad \text{② we don't need}$$

$$\alpha_i = C - \lambda_i \quad \text{③}$$

$\sum \xi_i$ is constant

It's taken care by 'C' and optimization equations

Dual:

$$L_d = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j$$

Subject to

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

KKT Condition

$$C \xi_i = 0, \quad \alpha_i [y_i(x_i^T \beta + \beta_0) - (1 - \xi_i)] = 0, \quad \lambda_i \xi_i = 0 \quad \text{⑥}$$

$$y_i(x_i^T \beta + \beta_0) - (1 - \xi_i) \geq 0 \quad \text{① If } \alpha_i = 0 \Rightarrow x_i \text{ far away}$$

since $\xi_i = 0$ for faraway points

$$\text{if } y_i(x_i^T \beta + \beta_0) > 1 \Rightarrow \alpha_i = 0$$

$$\text{② if } \alpha_i < \infty, \xi_i = 0 \Rightarrow y_i(x_i^T \beta + \beta_0) = (1 - \xi_i)$$

$$\text{③ } y_i(x_i^T \beta + \beta_0) - (1 - \xi_i) < 1 \quad \text{Support vector}$$

Misclassified

$$\xi_i > 0, \Rightarrow \lambda_i = 0, \Rightarrow \alpha_i = 0$$

$$\text{if } \xi_i \neq 0$$

$$\alpha_i = C$$

in both cases

$\alpha_i \neq 0 \Rightarrow$ those x_i 's are support vectors.

Lecture-24 SVM kernel

Using kernels, we don't need to compute the basis function to transform vectors.

like in previous example, we don't need to compute $b(x)$.

C-SVM → we call it C-SVM - whatever we learnt so far.

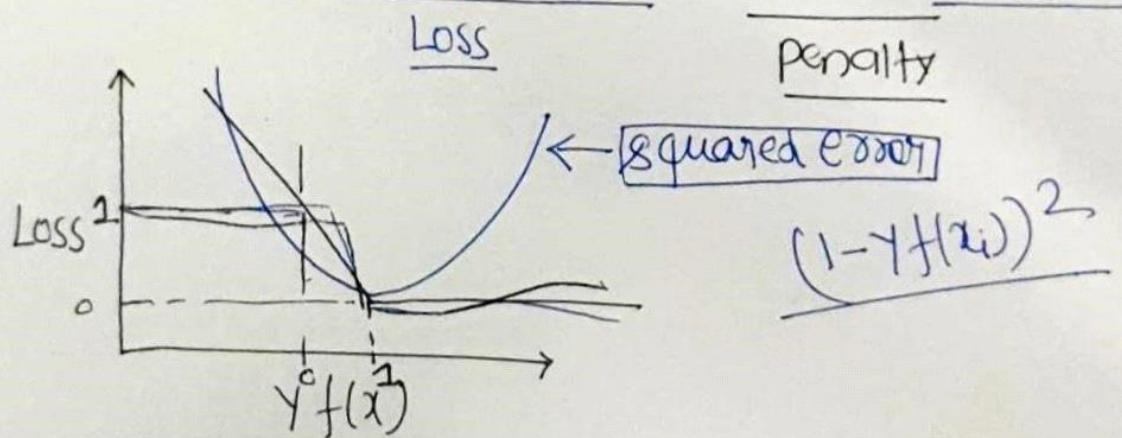
Different type of SVM's are possible where people use different optimisation constraints.

Lecture-25 - Hinge Loss Formulation

$$L_p = \frac{N}{\cancel{\sum_{i=1}^N}} \cdot \frac{1}{2} \|B\|^2 - \sum_{i=1}^N [Y_i (x_i^T \beta + \beta_0) - 1]$$

$$= \min_{\beta, \beta_0} \sum_{i=1}^N [1 - Y_i f(x_i)] + \frac{1}{2} \|\beta\|^2$$

Count only when it's positive.

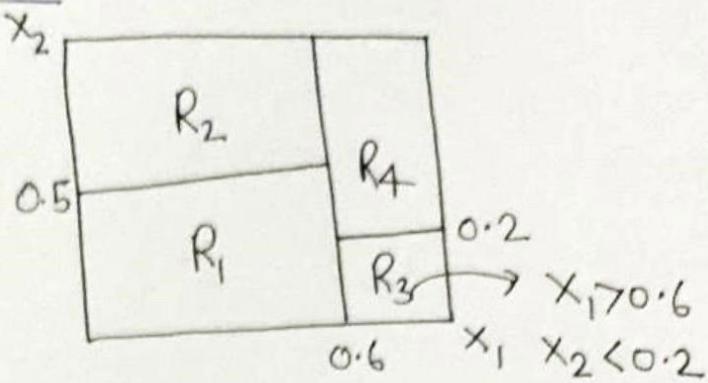
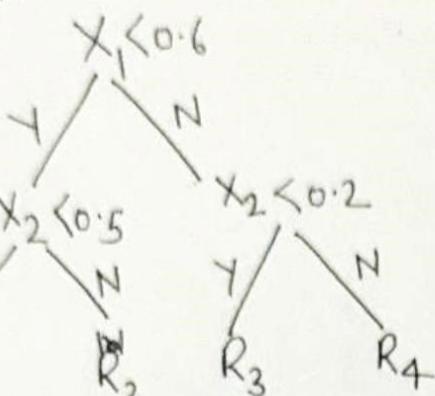


Lecture-33-1 Decision Trees - Introduction

Querying properties of the Data

- Partition I/P space into rectangles.

In theory regions could be of any shape.
we use high dimensional rectangular boxes.



- highly interpretable
- universal approximator
- Non-parametric

Regression Trees

- Fit a constant in each region.

$$\hat{f}(x) = \sum_{m=1}^M c_m I\{(x_1, x_2) \in R_m\}$$

$$(x_i, y_i); i=1, \dots, N, x_i \in \mathbb{R}^{d_x}, y_i \in \mathbb{R} \quad \Rightarrow \quad f(x) = \sum_{m=1}^M c_m I(x \in R_m) \\ x_i = \{x_{i1}, x_{i2}, \dots, x_{ip}\} \quad = c_1 I(x \in R_1) + c_2 I(x \in R_2) + c_3 I(x \in R_3) + c_4 I(x \in R_4)$$

- determine the 'M' regions

- Given regions, find Response. $c_m \rightarrow$ response for region 'm'.

Minimise \downarrow RSS:

$$= \sum_{j=1}^M \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2$$

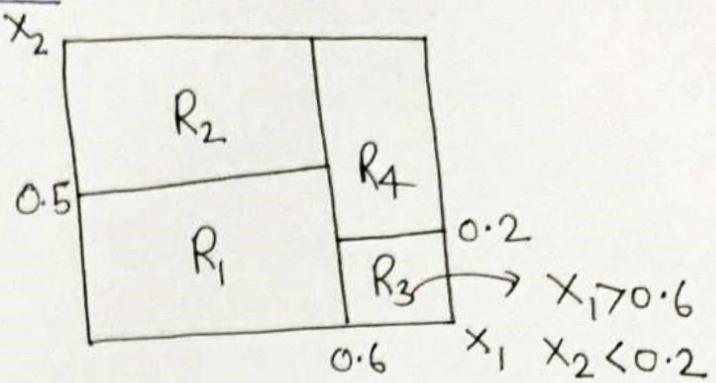
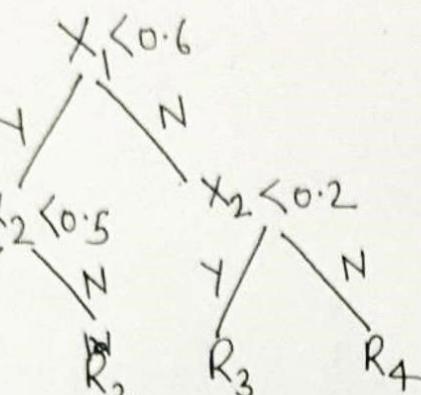
$$\hat{y}_{R_m} = \hat{c}_m = \frac{1}{|R_m|} \sum_{i \in R_m} f(x_i) \rightarrow \begin{matrix} \text{mean response for training} \\ \text{observations} \end{matrix}$$

Lecture-33-1 Decision Trees - Introduction

Querying properties of the Data

- Partition I/P space into rectangles.

In theory regions could be of any shape.
we use high dimensional rectangular boxes.



- highly interpretable
- universal approximator
- Non-parametric

Regression Trees

- Fit a constant in each region.

$$\hat{f}(x) = \sum_{m=1}^M c_m I\{(x_1, x_2) \in R_m\}$$

$$(x_i, y_i) : i=1, \dots, N, \quad x_i \in \mathbb{R}^p \quad \Rightarrow \quad f(x) = \sum_{m=1}^M c_m I(x \in R_m) \\ x_i = (x_{i1}, x_{i2}, \dots, x_{ip}) \\ = c_1 I(x \in R_1) + c_2 I(x \in R_2) + c_3 I(x \in R_3) + c_4 I(x \in R_4)$$

- Determine the 'M' regions

- Given regions, find Response. $c_m \rightarrow$ response for region m !

$$\text{Minimise } \underset{\substack{\downarrow \\ \text{RSS}}}{\sum_{j=1}^M} \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2$$

$$\hat{y}_{R_m} = \hat{c}_m = \underset{i \in R_m}{\sum} f(x_i) \rightarrow \text{mean response for training observations}$$

Space into J regions.

Approach (greedy)

1. All observations belong to a single region
2. Split feature space, i.e. split into two new branches.
3. At each step of tree building process, best split is made.
4. Consider a splitting variable j and cutpoint s , such the splitting the space into regions $R_1: \{x | x_j < s\}$, $R_2: \{x | x_j > s\}$ leads to greatest possible reduction in RSS.

$$R_1(j, s) = \{x | x_j \leq s\}, R_2(j, s) = \{x | x_j > s\}$$

Seek the value of J and S such that

$$\min_{J, S} \left[\sum_{i: x_i \in R_1(j, s)} (y_i - \bar{y}_{R_1})^2 + \sum_{i: x_i \in R_2(j, s)} (y_i - \bar{y}_{R_2})^2 \right]$$

\bar{y}_{R_1} Mean response

\bar{y}_{R_2} Mean response

We will do separately.

$$j \in \{1, 2, \dots, p\}$$

Fix J and find S

Strategy to find $S \rightarrow$ choose values from data points.
sort the points in ascending order and

choose ' s '-

$n \rightarrow$ Total points. J Total time: $n * p$ → For one split
 $p \in J$

$J \rightarrow$ Can be repeated while growing trees.

Lecture 35 - Stopping criteria & pruning

- * Early stopping - May lead into poor decision boundary.
- * Leaf of a tree is a region.

Tree pruning

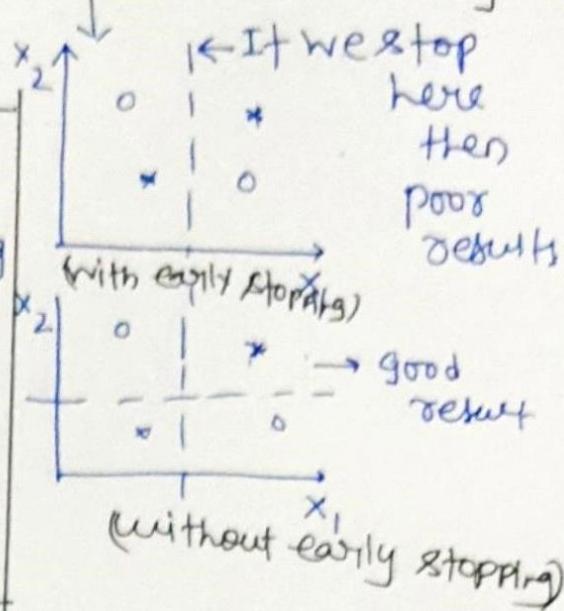
Above process will lead to overfitting

→ Tree will be very complex.

→ Variance higher

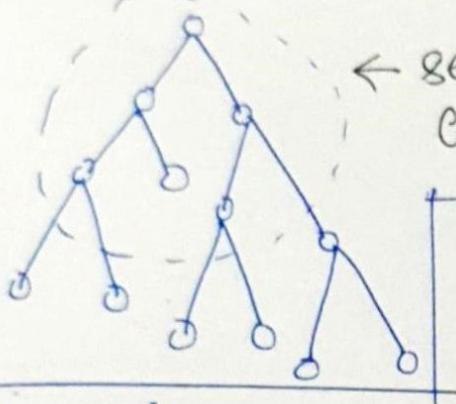
→ Simple Tree

less regions → lower variance.
better interpretation.



Strategy: Grow Tree T , and then prune it back in order to obtain subtree using (cross-validation)

← selecting best subtree will be very time consuming. → No. of subtrees are very large.



As $\alpha \uparrow \Rightarrow$ Number of Nodes pruned.

Cost-Complexity pruning: weakens link pruning

$$Q_\alpha(T) = \text{prediction error} + \alpha |T|$$

$$= \sum_{m=1}^{|T|} (y_i - \hat{y}_{R_m})^2 + \alpha |T| \text{ complexity penalty}$$

$|T| \rightarrow \underline{\text{Leaf Nodes}} = \text{No. of regions}$

$\alpha = 20,000$



$|T| = 3$

$|T| = 4$

$\alpha = 10,000$

$|T| = 4$

$\alpha = 22,000$

$|T| = 1$

$\alpha \uparrow \rightarrow$ (more penalty)

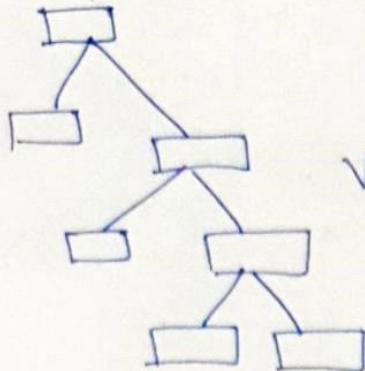
we want to penalize tree more having less no. of regions

① we already know that pruned tree will have more error.

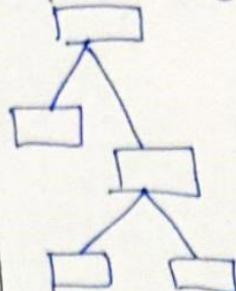
② we need to introduce one more term by which we can compare these trees. by introducing complexity penalty.

To compare

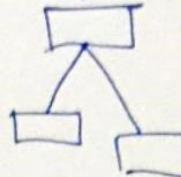
$SSR = 543.8$



$SSR = 5494.8$



$SSR = 19243$



$SSR = 28897.2$

vs

vs

vs

Let's take $\alpha = 10,000$,

Tree score = 40,543.8

Total score = $SSR + 10,000 * T$

Tree score = 35,994

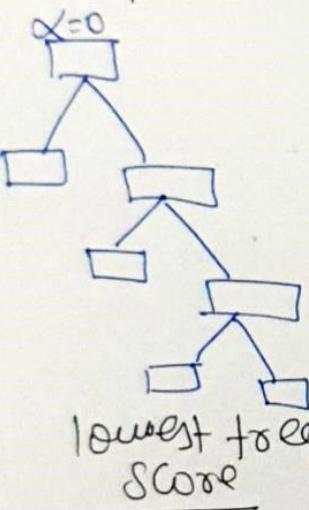
Tree score = 39,243.7

Tree score = 38,897.2

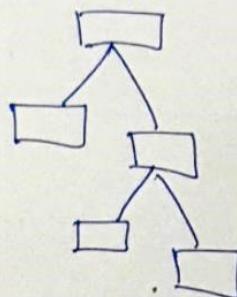
optimal subtree

using different α can give different result.
we need to find best value for alpha.

$\alpha = 0$, \Rightarrow Full sized tree will have lower Tree score.

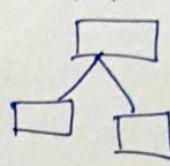


$\alpha = 10,000$



lower tree score

$\alpha = 15,000$



$\alpha = 22,000$

↓
gives
lower
tree score.

use testing data to get the optimal subtrees for different values of $|T|$.

Lecture 36

Decision Tree for classification - Loss function.

$$\hat{P}_{mk} = \frac{1}{|R_m|} \sum_{i \in R_m} I\{\hat{y}_i \neq k\}$$

Total Number of
 Correct Predictions
 Total Number of data points belonging
 to class k in region \underline{m}

Prediction

R_2	R_3
R_1	R_4

3 classes

$$b_{41}, b_{42}, b_{43}, b_{44}, b_{21}, b_{22}, b_{23}, b_{24}$$

$$b_{11}, b_{12}, b_{13}, b_{14}, b_{31}, b_{32}, b_{33}, b_{34}$$

if $x_i \in R_m$ then check which is greater.

$$[b_{41}, b_{42}, b_{43}, b_{44}]$$

$$k(m) = \arg \max_k \hat{P}_{mk} = \max_k [b_{41}, b_{42}, b_{43}, b_{44}]$$

$$k(4) =$$

↓ label for region 4

$$\text{Classification error} = \frac{1}{|R_m|} \sum_{i \in R_m} I(\hat{y}_i \neq k(m)) = (1 - \hat{P}_{mk(m)})$$

$$= (1 - \max_k \hat{P}_{mk})$$

$$\text{Gini Index} := \sum_{K=1}^K \hat{P}_{mk} (1 - \hat{P}_{mk})$$

Measure of variance across K -classes for a particular region.

$$\text{Cross Entropy} = - \sum_{K=1}^K \hat{P}_{mk} \log_{\hat{P}_{mk}}$$

estimated label

Information gain criteria

Lecture-37.

Decision Trees - Categorical attributes

q - unordered attributes:

Divide into 2 groups. Total No. of groups possible:

$$(q_1 + q_2 + q_3 + \dots + q_{q_{1/2}}) = 2^{q-1} - 1 \rightarrow \text{Not feasible to go and pick split points.}$$

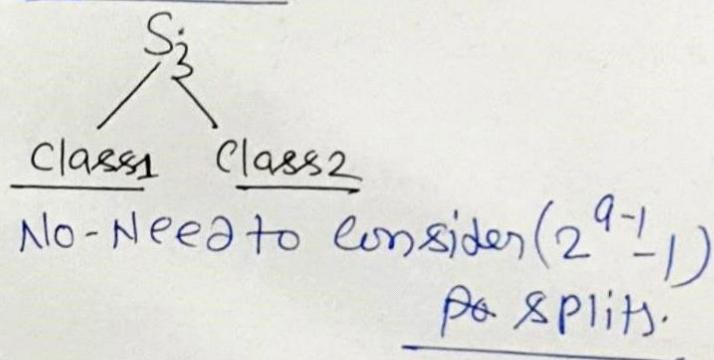
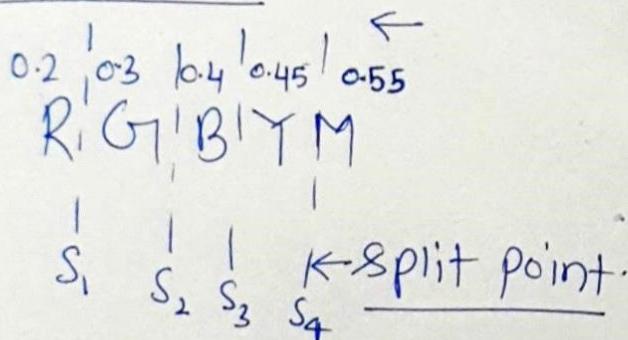
in case of regression: you have to look n points.
 $n \rightarrow \text{No. of data points.}$

0/1 - binary classification.

Suppose $q=5$ For color attribute: $q=0$ Red,

$P(\text{class1}/\text{Red})$, $P(\text{class1}/\text{green})$, $q=1$ green
 $P(\text{class1}/\text{blue})$, $P(\text{class1}/\text{yellow})$, $q=3$ blue
 $P(\text{class1}/\text{magenta}) \rightarrow$ Ascending order. and then make

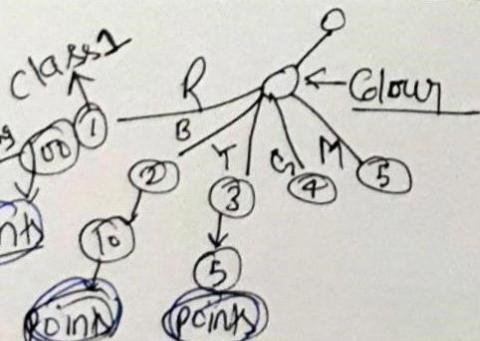
the split."



Multiclass

Heuristics:-

Multway-splits



Lecture-38.

Multway splits.

- Problem

→ Spawling - (spawning) (hard to interpret)

→ loss of interpretability

→ variance loss.

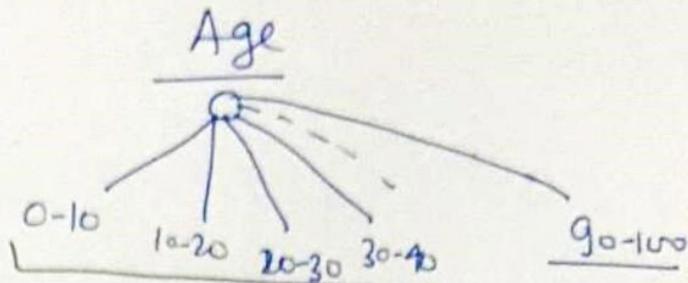
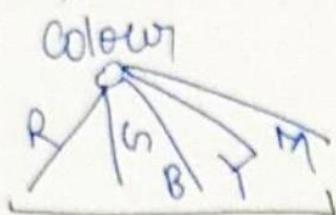
→ At each split, data points will be lower..

→ Tree will be become sparse.

→ Overfitting is the problem.

Lecture-38

- Favours attributes with more values.



ErrorColor > Error Age.

Less error in splitting

C4.5 - Algorithm - gain-Ratio

↓
what is the information gain
after splitting into 10 parts
rather than 6.

→ we can also use binary split to arrive at some
→ function which is predicted by multiway splits.

binary split is always better > multiway
split because
we can use recursive approach to make
multi-way split.

Multiway splits → avoids the choosing variable (S)
(split point)

Lecture-39.

Missing values, imputation & Surrogate result
imputation:

Mean, Regression, [predict missing points]

↳ Full information imputation.

Using all other feature to impute values.

Multiple imputation: - use distribution to predict
values.

Follows 3 steps.

1. Similar to single Imputation, missing values are imputed. However, the imputed values are drawn m times from a distribution rather than just once. → gives (m-datasets)

2. Analysis: At end of this step, there should be m - analyses.

3. Pooling: m results are consolidated into one result by calculating mean, variance, confidence intervals of the variable.

→ Single Imputation does not take care of uncertainty in imputed values whereas multiple imputation accounts for the uncertainty and range of values that the true value could have been.

→ Fitting it into a distribution and then sampling values using mean and variance of distribution

Imputation • Lecture-39: surrogate splits

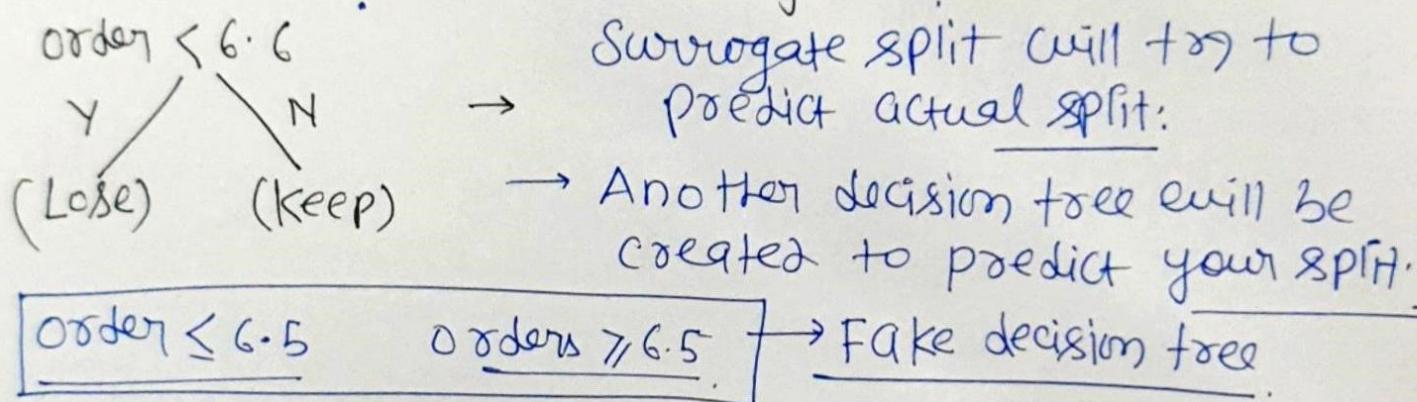
— New categorical feature value "missing".

<u>Age</u>	<u>Age</u>
20	20
30	30
--	"missing"
--	"missing"
40	40

} There are some reasons, why there is missing values in our data. If we impute it, we are going to add noise in the feature. So, it's best to just leave replace those values by "missing".

Surrogate splits

In decision tree, when a split is made, it depends on one variable. but what if that variable is missing. Suppose, we want to predict whether a customer will be kept or lost next year.



Fragment (specific to trees)

$$\underline{x_3 < 5 ?}$$

Lecture-40

Instability, smoothness & Repeated subtrees

Instability :- More variance. Prone to noises

Lack of smoothness :- Since decision tree has piecewise rectangular boxes. Clear cut boundaries.

Repeated subtrees :- Same set of split made at different points in decision tree.

$$\text{Loss } L_{\text{KK}} = \sum_{k=1}^K p_{mk} (\hat{p}_{mk}) = \sum_{k=1}^K \hat{p}_{mk} \hat{p}_{mk}$$

True estimated

: Advantages :-

- simple to understand and interpret.
- handle both numerical and categorical data.
- performs well with large datasets.

disadvantages :-

- Non-Robust: A small change in the training data can result in a large change in the tree and consequently final predictions.
- For data including categorical variables with different number of levels, Information-gain is biased in favour of attributes with more levels.
- Decision tree can learn very complex decision boundary. Overfitting.

Decision Tree tutorial

Dataset

	age	income	student	credit-rating	buys-computer	target
↓	high	yes		fair		
Youth	medium	No		excellent		
Middleaged	low					Total-data = 14

Senior

: impurity measure:

- cross entropy | gini-index

Multi-way split using cross entropy

Consider Attribute 'Age'

Youth → 3-No
2-Yes

Senior → 3-Yes
2-No

Middleaged → 4 Yes
0-No

$$\text{cross-entropy: } - \sum_{k=1}^K p_{mk} \log p_{mk}$$

$$K=2$$

Initially m=1

$$- (p_{11} \log p_{11} + p_{12} \log p_{12})$$

$$\text{For attribute age=Youth } p_{11} = \frac{2}{5}, \quad p_{12} = \frac{3}{5}$$

$$\text{cross-entropy}_{\text{Age}} = + (5/14) \left(-\frac{2}{5} \log \frac{2}{5} - \frac{3}{5} \log \frac{3}{5} \right) + \frac{4}{14} \left(-\frac{4}{5} \log \frac{4}{5} \right) + \frac{5}{14} \left(-\frac{3}{5} \log \frac{3}{5} \right)$$

$$= \underline{0.6935}$$

$$-\frac{2}{5} \log \frac{3}{5}$$

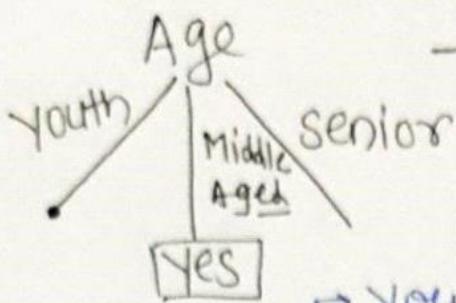
$$\text{cross-entropy}_{\text{credit_range}}(D) = \underline{0.8922}$$

$$\text{cross-entropy}_{\text{income}} = \underline{0.9111}$$

$$\text{cross-entropy}_{\text{student}} = \underline{0.7885}$$

less entropy →

Decision Tree

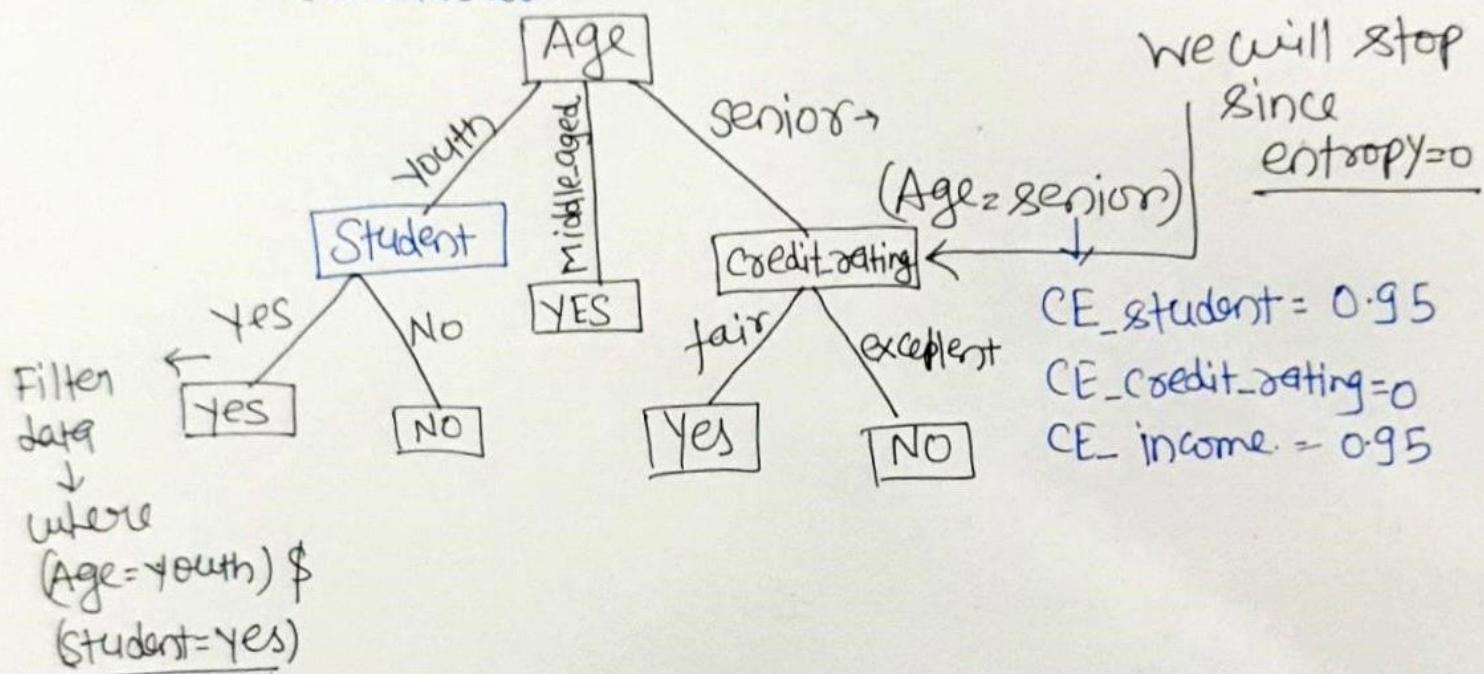


→ You will filter dataset where
age = youth.

cross_entropy ($\text{Age} = \text{youth}$) =
income

cross_entropy ($\text{Age} = \text{youth}$) = 0 → minimum value of entropy
student

cross_entropy ($\text{Age} = \text{youth}$) =
credit & Cox



Decision-Tree:

Binary split using Gini index

Consider 'Age'

Total data points = 14

(+ve) class-proportion: Age=Youth: 2/5

Age=Middle : 1

Age: Senior : 3/5

Ordering the probability.

Youth, senior, middle

$$\text{Gini-index} = \sum_{k=1}^K p_{mk} (1 - p_{mk})$$

Possible splits

{Youth}, {senior, middle} | {Youth, senior}, {middle}

⇒ We need to calculate impurity measure for both the above splits.

$$\begin{aligned} \text{Gini}_{\text{age} \in \{\text{Youth}\}} (\mathcal{D}) &= 2p(1-p) \\ &= \frac{5}{14} \left(2 \times \frac{2}{5} \times \frac{3}{5} \right) + \frac{9}{14} \left(2 \times \frac{7}{9} \times \frac{2}{9} \right) \\ &= 0.6508 \end{aligned}$$

$$\begin{aligned} \text{For '2' class problem } G_1 &= p(1-p) + (1-p)p \\ &= \boxed{2p(1-p)} \end{aligned}$$

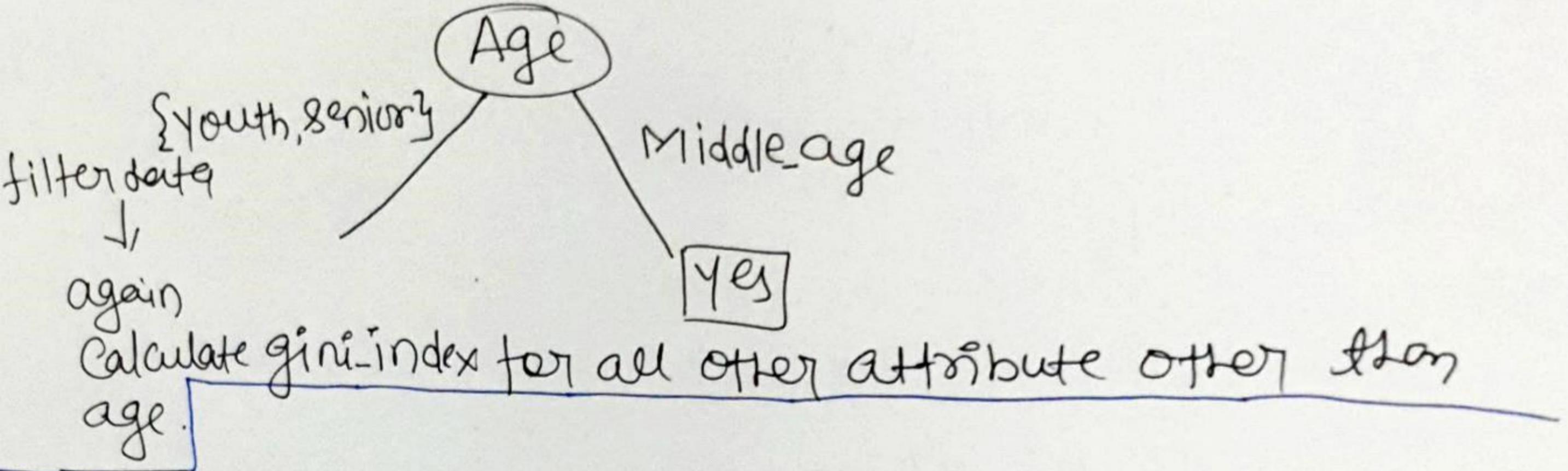
Gini-age = {Youth, senior} =

$$\text{Gini-Student} \{ \text{Yes} \} = \frac{7}{14} \left(2 \times \frac{3}{7} \times \frac{4}{7} \right) + \frac{7}{14} \left(2 \times \frac{2}{7} \times \frac{1}{7} \right)$$

Gini-income → Two splits =

$$\text{Gini-credit rating} \{ \text{fair} \} = \frac{8}{14} \left(2 \times \frac{6}{8} \times \frac{2}{8} \right) + \frac{6}{14} \left(2 \times \frac{3}{6} \times \frac{3}{6} \right)$$

Choose the one having less impure.: lesser value



statquest

Random Forest - Part 1

Building using, and Evaluating

- Decision trees are not flexing in order to classify new samples.

Building Random Forest

Step 1: Create a "bootstrapped" dataset.

Original dataset (Target)

Chest pain	Good blood	F3	F4	F5
No	No	No	125	Yes
Yes	Yes	Yes	180	No
Yes	Yes	No	210	No
Yes	No	Yes	167	Yes

Bootstrapped dataset

Chest pain	Good blood	F3	F4	F5
Yes	Yes	Yes	180	No
No	No	No	125	Yes
Yes	No	Yes	167	Yes
Yes	Yes	Yes	180	No

a) To create bootstrapped dataset, i.e of same size of original dataset, we just randomly select samples from the original dataset.

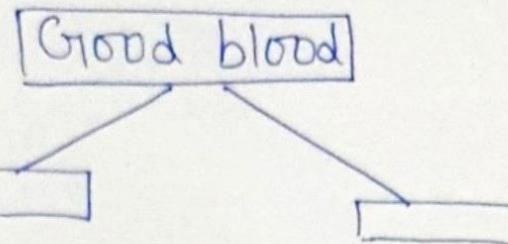
b) we are allowed to pick the same sample more than once.

Step 2: Create a decision tree using the bootstrapped dataset, but only use a random subset of variables (features) at each step.

In this case, we will randomly select two variables, "Chest" and "Good blood" ~~as well~~ as a candidate for root node.

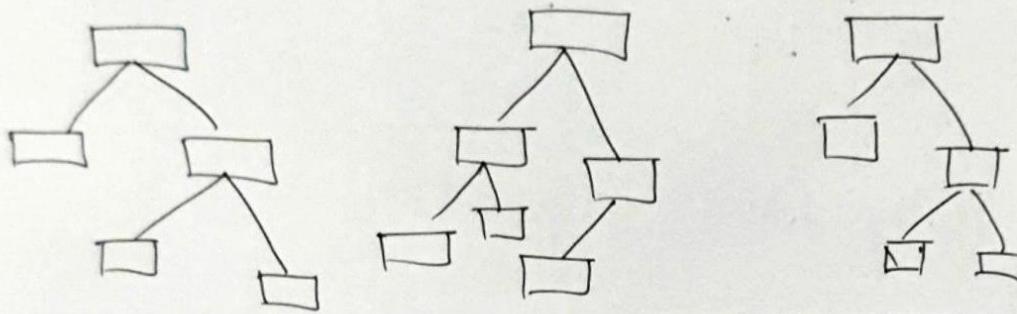
Assume, "Good blood" has lowest cross-entropy,

at this point
also, randomly
select two
variables from "3"
variables left.



and build the tree as usual.

Now, go back to step 1 and repeat: Make a new bootstrapped dataset and build a tree. Consider random subset of variables at each split.



This will result in wide variety of decision trees.

↓
makes Random forest more effective than individual decision trees.

How to predict?

Run the sample on all tree we have created and calculates voting of Yes vs No

Heart disease (Target)		⇒ class = Yes
Yes	No	
5	1	

Random Forests- Part 2

Missing data and sample clustering

Random forest Consider two types of missing data

1. Missing data in original dataset.
2. Missing data in New Sample.

→ We can make initial guess → refine our guess

↓
using mean, median or other imputation methods.

initial guess

Chest pain	Blood	Blocked	Weight	Target
Yes	Yes	No	167.5	No

How to refine our guesses

Step-1: Build Random forest...

Step-2: Run all data in all of distinct decision trees.

Find the samples which end up at same leaf Node as our data-sample which we want to make a guess.

$$\text{No. of Rows} = \text{No. of Columns} \\ = \text{Total Samples}$$

Build Proximity Matrix

	1	2	3	4
1		2	1	1
2	2		1	1
3	1	1		8
4	1	1	8	6

Suppose 3' and 4' end up in same leaf Node, we will add 1 in $P[3][4]$ and $P[4][3]$

→ we divide each value by $\frac{1}{4}$ Total No. of Trees.

1	0.2	0.1	0.1
2	0.2	0.1	0.1
3	0.1	0.1	0.8
4	0.1	0.1	0.8

Random Forest

Terminology:

Bootstrapping the data plus using the aggregate to make a decision is called "Bagging".

→ Typically $\frac{1}{3}^{\text{rd}}$ samples does not end up in bootstrapped dataset.

[33% of data] → called the (out-of bag dataset)
↓
doesn't appear in
bootstrapped dataset.

* Since out of bag datasets were not used in ~~class~~ building decision tree, we will use it to evaluate our model.

* The proportion of out of bag samples incorrectly classified by Random Forest is known as "out of bag error".

1. Build Random Forest

Change No. of variables

2. Estimate Accuracy of Random Forest

Usually, we start with $\sqrt{\text{No. of variables}}$.

Random Forest Part 2-2

Use proximity matrix to calculate missing values.

original dataset

Chest_pain	Good blood circ.	Blocked Arteries	weight	Heart_diseas
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	Yes	???	???	No

frequency(Yes) = $\frac{1}{3}$, f(No) = $\frac{2}{3}$

weighted_f(Yes) = $\frac{1}{3} \times (\text{The weight for Yes})^{\textcircled{1}} = 0.03$

The weight of Yes = $\frac{\text{Proximity of Yes}}{\text{All proximities}}$ value for sample 2

1	2	3	4
0.1	0.1	0.8	

$$= \left(\frac{0.1}{0.1+0.1+0.8} \right) = \frac{0.1}{1} = 0.1$$

weighted_f(No) = $\frac{2}{3} \times (\text{The weight for No})^{\textcircled{2}} = 0.6$

The weight for No = sample 1 and sample 3

$$\frac{0.8+0.1}{0.1+0.8+0.1} = \frac{0.9}{1} = 0.9$$

(No > Yes)

→ New guess = No

~~Lecture-23 SVM: for linearly non-separable data.~~ ③

Weightage-average of (weight) = $125 * \text{proximity weight}$

$$\text{Proximity weight}(125) = \frac{0.1}{0.1+0.1+0.8} = 0.1 + (180 * \text{proximity weight}) + (210 * \text{proximity weight})$$

$$\text{Proximity weight}(180) = \frac{0.1}{0.1+0.1+0.8} = 0.1$$

$$\text{Proximity_weight}(210) = \frac{0.8}{0.1} = 0.8$$

$$= 125 * (0.1) + 180 * (0.1) + 210 * (0.8) = \underline{\underline{(198.5)}}$$

After refining guesses, we run random forest again to refine guesses, until missing values converge. (6 to 7 times).