

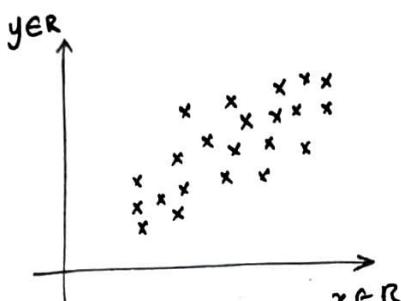
INTRODUCTION TO MACHINE LEARNING:

Consider a data, $\mathcal{D} = \{(x_i, y_i) : x_i, y_i \in \mathbb{R}^2, i=1(1)n\}$

	x	y
1	x_1	y_1
2	x_2	y_2
:	x_3	y_3
:	:	:
n	x_n	y_n

$x_i, y_i \in \mathbb{R}^2$

The datapoints can be plotted in a real plane.



GOAL: To build a predictive model to predict y based on x

(ii) to create a function $f(x)$ that accurately predicts y .

(ii) $y_i = f(x_i)$, $\forall i = 1(1)n$ \forall for all

↳ finding a y_i for every value of x_i

However, the equality in the above equation never holds true because in real life it's extremely difficult to capture the exact effect.

x = Area (sqft living)

y = Price

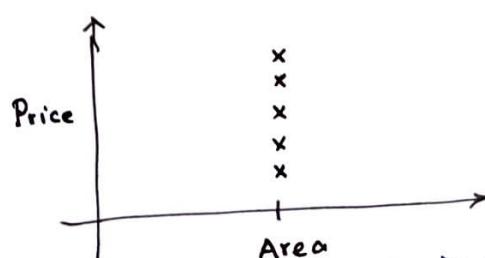
$$y = f(x)$$

Whatever equation or function you use to predict the price, you won't be able to predict with only the area.

$$y = a + bx$$

$$y = a + bx + cx^2$$

$$y = a + bx + cx^2 + dx^3$$



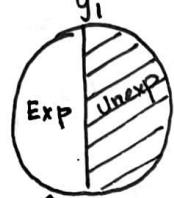
For a given area, the house can be underconstruction, ready to move in, semifurnished, furnished, 1st floor, 10th floor which that these affect the homeprice.

$y = f(x_1, x_2, \dots, x_6)$, if we do so we may be able to reach closer to y but the equality will still not hold.

→ Whatever may be the case, exact prediction is not possible if given the inputs and will end up making some error every time we make a prediction.

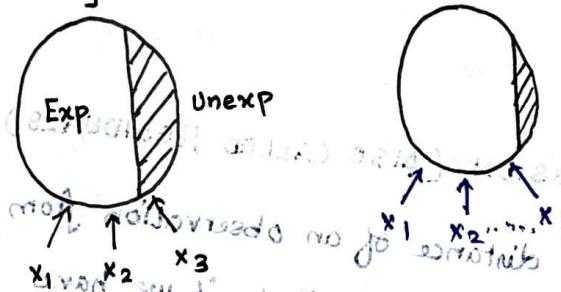
MODIFIED EQUATION:

$$y_i = f(x_i) + \epsilon_i, i=1(1)n \text{ and } \epsilon_i \text{ are random noise.}$$



→ There may be no reasons that we couldn't account for (quick sale, waiting for price shoot up, suicide).

whereas ϵ_i is called random noise or error.



This equation holds true

$$y_i = f(x_i) + \epsilon_i$$

whereas

$$y_i \neq f(x_i), \forall i=1(1)n$$

ERROR:

→ Error or random error is the unexplained part of y .

→ If ϵ_i is +ve then $f(x_i)$ has underestimated y_i [$y_i - f(x_i) < 0$]

→ If ϵ_i is -ve, then $f(x_i)$ has overestimated y_i . [$y_i - f(x_i) > 0$]

→ Since errors are inevitable (it's unavoidable), we always focus on designing statistical frameworks that solves the problem of prediction by minimizing the errors.

Design a model that's minimizing the errors.

$$M(\beta) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_1 - \dots - \beta_p x_p)^2$$

→ Minimizing $M(\beta)$

LINEAR REGRESSION MODEL:

→ Considering the same data $D = \{(x_i, y_i) : i=1(1)n, x_i, y_i \in \mathbb{R}^2\}$

Consider a linear model:

$$y_i = b_0 + b_1 x_i + \epsilon_i \quad \leftarrow \text{consider called a linear model because the dependent variable } y_i \text{ is linear in the parameters}$$

Is the model: $y_i = b_0 + b_1 x_i^2 + \epsilon_i$ a linear model?

It's a linear model because it is linear in parameters.

Linear \rightarrow flat surface in p dimensional space.

$\rightarrow f$ and

In higher dimensions, linear regression is a plane.

Higher the shape dimensions, then linear regression will take different shapes.

ERRORS IN LINEAR REGRESSION (ALSO CALLED RESIDUALS)

Errors are the vertical distance of an observation from the linear regression line which means that if we have

two variables x and y ,

Error = (Actual y - predicted y)

Principles no need, you do not predicted value for x to predict y .

PARAMETERS IN SIMPLE LINEAR REGRESSION MODEL:

↑ prior simple-one predictor.

$$y_i = b_0 + b_1 x_i + \epsilon_i, i=1(1)n$$

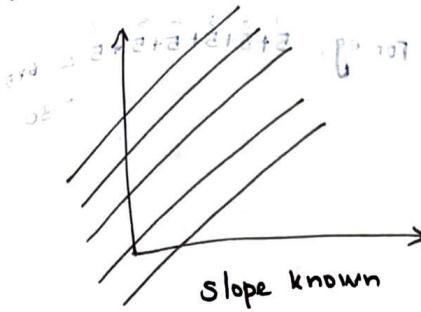
$y = b_0 + b_1 x$	X
$y = b_0 + b_1 x$	X

Incorrect

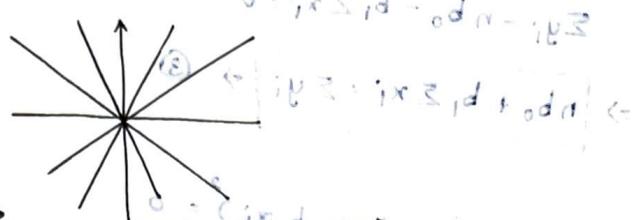
b_0 = y-intercept

b_1 = slope of the line

A line depends on these two numbers.



either fixed no. $b_0 + b_1 x$



$$y = b_0 + b_1 x \quad (1)$$

$$y = b_0 + b_1 x \quad (2)$$

Estimate the parameters.

Estimation of Parameters

(Optimization objective in linear regression) Estimate the parameters such that the sum of squared error is minimum. (ii) $\sum_{i=1}^n \epsilon_i^2$ is minimum.

choose that particular line with min. SSE.

$$\epsilon_i = ? , \quad \epsilon_i = y_i - \hat{y}_i$$

$$\epsilon_i = y_i - (b_0 + b_1 x_i) = y_i - b_0 - b_1 x_i$$

$\Rightarrow \epsilon_i$ is a function that depends on b_0 and b_1 .

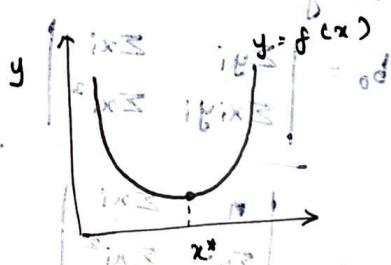
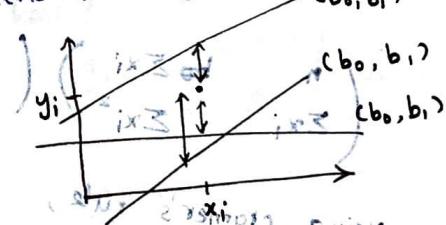
$$\epsilon_i = f(b_0, b_1) = y_i - (b_0 + b_1 x_i) \quad (\text{where } x_i \text{ and } y_i \text{ are fixed})$$

$$\sum_{i=1}^n \epsilon_i^2 = \sum (y_i - b_0 - b_1 x_i)^2$$

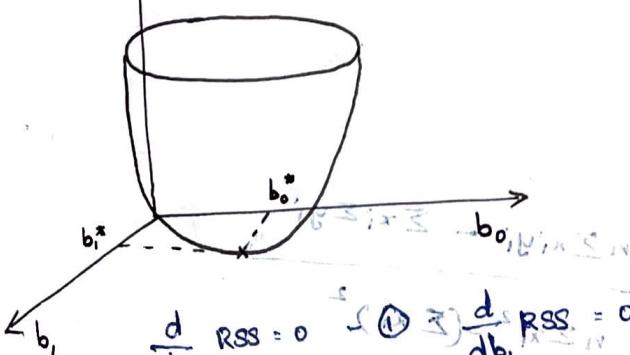
We want to minimize $\sum \epsilon_i^2$ by choosing

b_0, b_1

$$\sum \epsilon_i^2 = \text{RSS} \quad (\text{Residual sum of squares})$$



$$\frac{d}{dx} f(x) = 0 \quad \text{and solve for } x.$$



$$\frac{d}{db_0} \text{RSS} = 0 \quad (1) \quad \frac{d}{db_1} \text{RSS} = 0 \quad (2)$$

$$(1) \Rightarrow \frac{d}{db_0} \sum (y_i - b_0 - b_1 x_i)^2 = 0 \Rightarrow \sum \frac{d}{db_0} (y_i - b_0 - b_1 x_i)^2 = 0$$

$$\Rightarrow \sum 2 (y_i - b_0 - b_1 x_i) (-1) = 0$$

$$\Rightarrow -2 \sum \epsilon_i = 0$$

÷ by -2 on both sides

$$\sum y_i - nb_0 - b_1 \sum x_i = 0$$
$$\Rightarrow [nb_0 + b_1 \sum x_i : \sum y_i] \rightarrow ③$$

Adding constant n times
is same as nb_0 ,
For eg, $5+5+5+5+5+5 = 6 \times 5 = 30$

$$(2) \frac{d}{db_1} \sum (y_i - b_0 - b_1 x_i)^2 = 0$$

$$\sum \frac{d}{db_1} (y_i - b_0 - b_1 x_i)^2 = 0$$

$$\sum 2(y_i - b_0 - b_1 x_i)(-x_i) = 0$$

$$\Rightarrow -2 \sum (x_i y_i - b_0 x_i - b_1 x_i^2) = 0$$

$$\Rightarrow \sum x_i y_i - b_0 \sum x_i - b_1 \sum x_i^2 = 0$$

$$\Rightarrow b_0 \sum x_i + b_1 \sum x_i^2 = \sum x_i y_i - ④$$

$$nb_0 + b_1 \sum x_i = \sum y_i - ③$$

$$b_0 \sum x_i + b_1 \sum x_i^2 = \sum x_i y_i - ④$$

$$\begin{pmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix} \begin{pmatrix} b_0 \\ b_1 \end{pmatrix} = \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix}$$

Using cramer's rule,

$$b_0 = \frac{\begin{vmatrix} \sum y_i & \sum x_i \\ \sum x_i y_i & \sum x_i^2 \end{vmatrix}}{\begin{vmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{vmatrix}} = \frac{\sum y_i \sum x_i^2 - \sum x_i \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$b_1 = \frac{\begin{vmatrix} n & \sum y_i \\ \sum x_i & \sum x_i y_i \end{vmatrix}}{\begin{vmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{vmatrix}}$$

$$= \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$= \frac{n(\sum x_i \sum y_i - \sum x_i \sum y_i)}{n \sum x_i^2 - (\sum x_i)^2}$$

$$= 0 = 30$$

The linear regression, every point in the line

is considered to be the mean or average

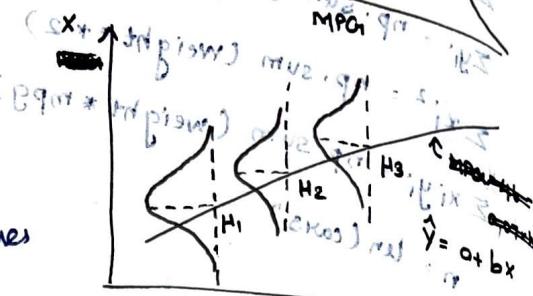
such that MPG follows a normal distribution.

$$y = a + bx$$

Slope

With one unit increase in x , y increases by b units.

$$\hat{MPG} = 46.325 - 0.007 \times WT$$



Interpreting slope

With 1 unit increase in wt, the MPG increases by -0.007 miles.

"", the MPG decreases by $+0.007$ miles.

With 1 lb increase in wt, the MPG decreases by $\underline{\underline{-0.007}}$ miles.

With 1000 lb increase in wt, the MPG decreases by $\underline{\underline{-0.007}}$ miles.

↳ Most interpretable.

$y = a + bx \rightarrow$ intercept $y = a$ when $x = 0$.

Interpreting y -intercept

For an equation, $y = a + bx$,

a can be interpreted as the

value y takes when $x = 0$,

For eg consider the popular

regression model in economics

$\rightarrow a$ is the expenditure incurred when income = 0 . This

Shows that even there's no earning there may be a fixed expenditure.

However intercept in this way isn't possible for many examples. As in our case when we write $\hat{MPG} = 46.325 - 0.007 \times WT$.

Since the weight of a car cannot be 0 pound, the intercept cannot be interpreted in the same way as before.

DECISION TREE:

→ KNN does not use any parameters.

→ Linear regression uses parameters $\rightarrow y_i = b_0 + b_1 x_{i1} + \dots$

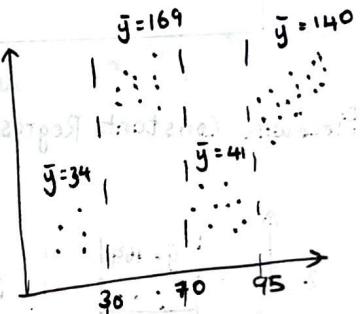
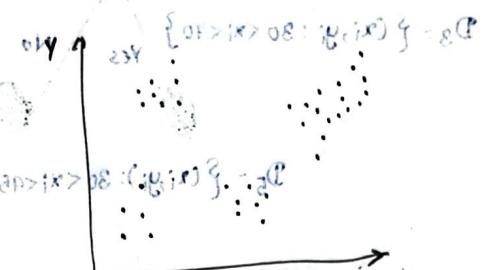
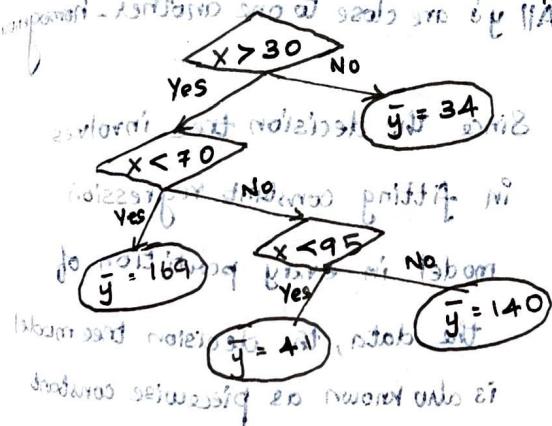
CART → Classification and Regression

Seeing graph we can't say x and y are related. The issue is that x and y are not linearly related.

But the relation can be expressed as

$y_i = b_0 + b_1 x_i + b_2 x_i^2 + \dots$. I.e., if x_i increases by b_2 units for unit change in x_i^2 , the value of y_i increases by b_2 units.

keeping all others constant \rightarrow not interpretable.



Decision tree can be identified as cluster.

DT: Identifies how data can be approximated in a region, the predicted

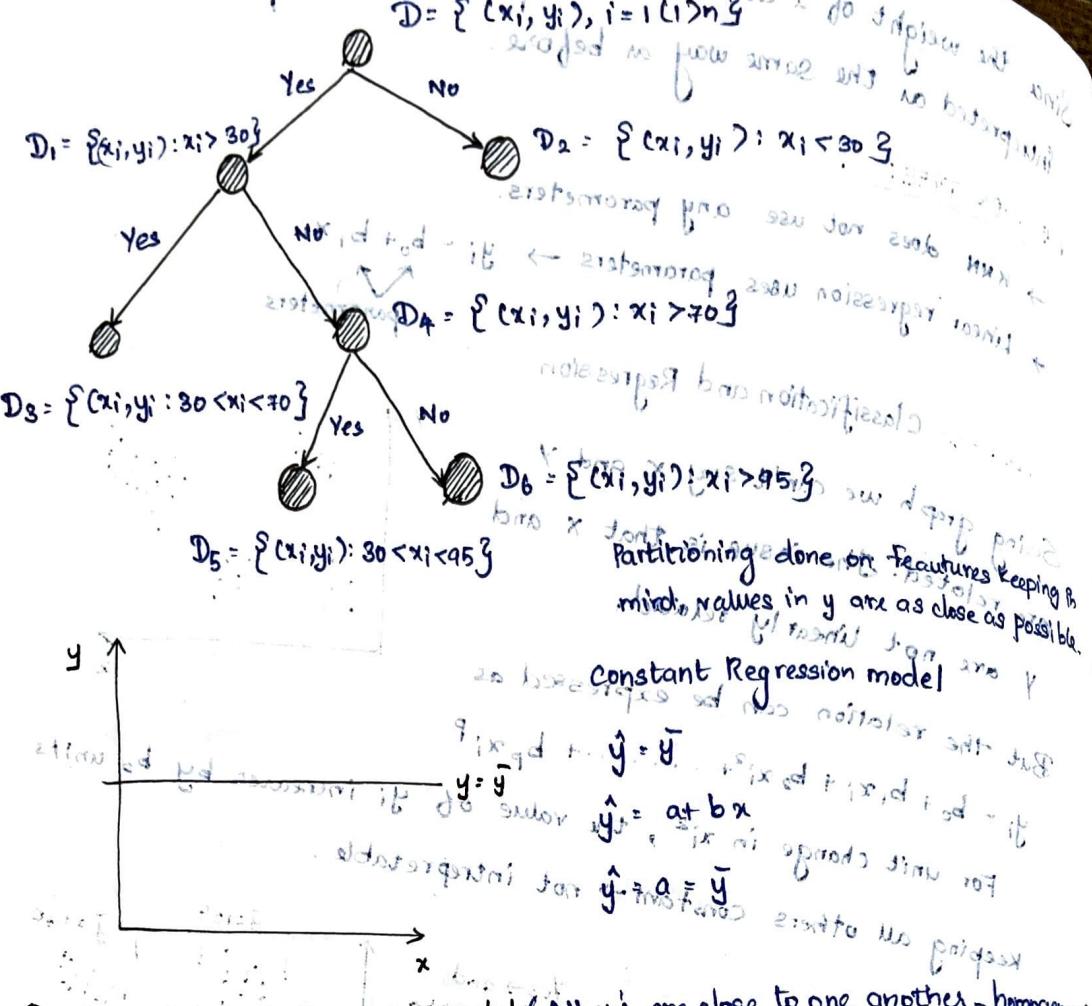
Assumes that variability in a region is approximated to the average values

of y in that region and this is called a node.

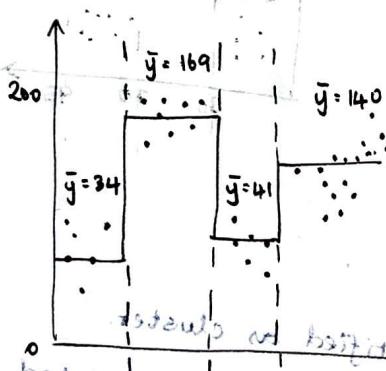
Each decision is represented as a subset of

Every node in a decision tree is associated with a subset of data.

Each node with attributes and conditions is also called a leaf node.



Piecewise Constant Regression Model (All y 's are close to one another - homogeneous)



Since the decision tree involves in fitting constant regression model in every partition of the data, the decision tree model is also known as piecewise constant regression model.

→ We usually ignore sign in continuous distributions since $p(x=30)$

→ Splitting can be done only based on the features (x) and homogeneity of target variable (y).

- X.X.X. NOTES:
- (1) Every node of a decision tree contains a subset of data.
 - (2) The root node of a decision tree contains the entire data.
 - (3) A CART model can only have two arms, for every decision nodes.



It's right arm always corresponds to value of True and the left arm corresponds to a value of False.

A decision tree partitions the data into a number of nodes based on the values of y (i.e. the homogeneity is observed according to the values of y).

The partitioning of the data happens in the feature space only.

A decision tree begins with one root node and may have several leaf nodes.

A decision tree with one root node and 2 leaf nodes is called a "decision stump".

Decision Tree
makes prediction by predicting mean

Measure of node Impurity for Regression Tree:

looks like variance.

$$D_{node} = \sum_{y \in \text{node}} (y - \bar{y})^2$$

looks like variance.

$$D_{node} = \sum_{y \in \text{node}}$$

giving the SSE of each interval.

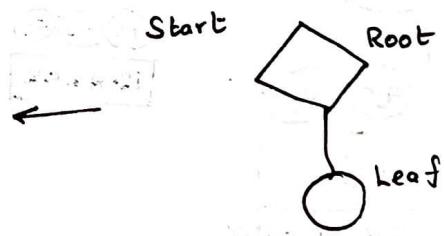
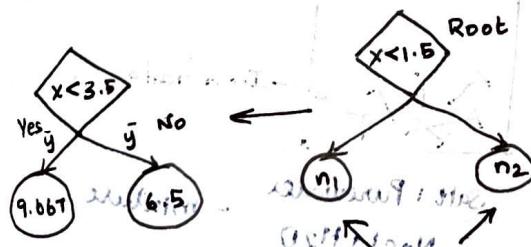
minimum overall deviance

If the value of y is same for different values of x in intervals, then $y = \bar{y}$, then $D_{node} = 0$.

minimum overall deviance

The deviance of each node must be minimum and overall deviance must also be minimum.

A decision tree is a combination of decision stumps.



Lesser the value of D_T , more the homogenous is the data split.

Decision Tree algorithm needs to be stopped (or) otherwise is pruning letting decision tree grow and cutting unwanted leaves.

All conditions are bounded (ANDED).

1st 2nd 3rd
Linear Regression, KNN

Most interpretable: Decision Tree, Linear Regression, KNN

In KNN we are cluster data based on value of X and in decision tree we cluster based on value of Y (homogeneity).

→ Different models used different approaches.

Feature selection can't be usually interpretable using decision trees.

→ Not OK → using decision tree or random forests for feature selection and using those features in linear regression. (NEVER USE) Europe

CLASSIFICATION TREE

ISLR Book

↳ 20-30 days

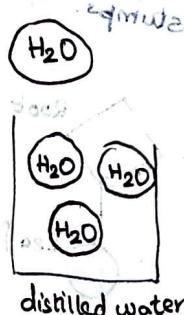
Recursive partitioning → A decision (i) CART algorithm partitions the data in a recursive fashion and therefore it's also known as a recursive partitioning algorithm.

Course era → University of Washington Machine Learning

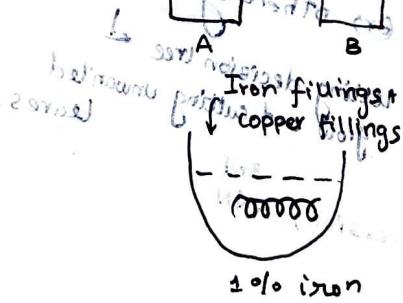
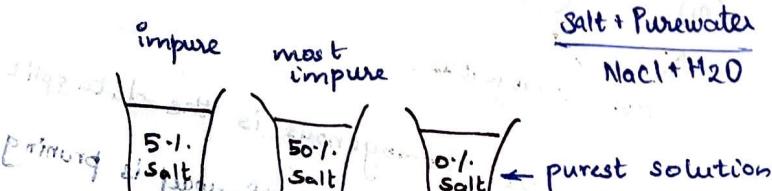
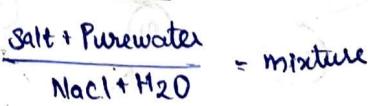
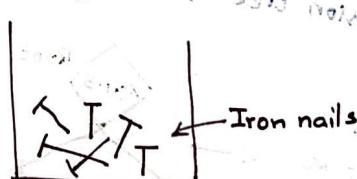
Assignments in GitHub or Apply for compensation.

1. Mixture → impure substance

2. Compounds/Elements → pure substance.



IRON BAR



When y is categorical in nature, homogeneity is calculated based on which feature the split is made.

↳ out of 2 features, 1 is chosen for splitting. (A) 7.9

The probability of tumour being malignant given the age of person between 28 and 50 and smokes cigarettes is $4/5$. (A) 7.9

MEASURES OF SELECTING THE BEST SPLIT (FOR CLASSIFICATION TREE).

Example data:

CID	Savings	Assets	Credit Risk
1.	M	H	G1
2.	L	L	B
3.	M	M	B
4.	L	M	G1
5.	H	M	G1
6.	L	L	G1
7.	M	M	G1
8.			

To predict credit risk based on Savings and Assets.

Regression	Classification
Homogeneity measured	Purity of node measured
-1 - (Using Total Deviance)	Dr

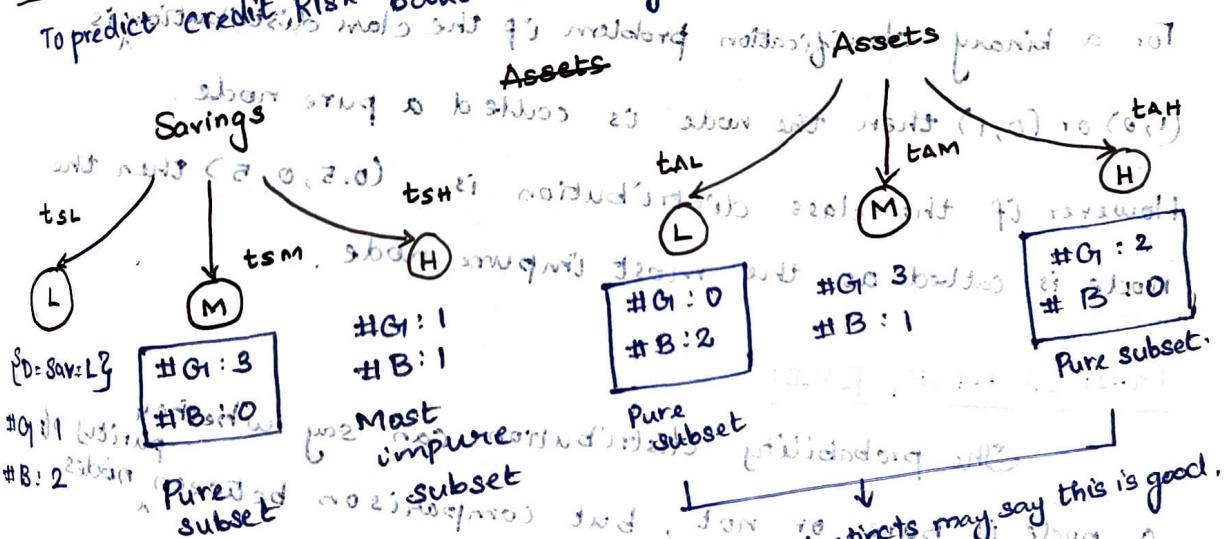
Individual labels are

can be used to partition the data.

Labels are

can be used to partition the data.

Labels are



Labels are

Labels are

NOTATION:

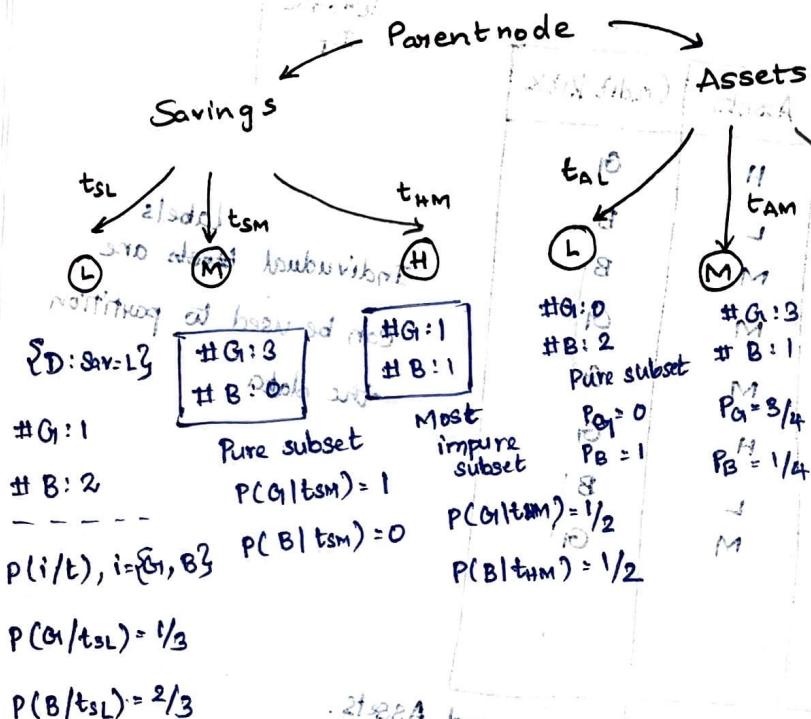
$p(i/t)$ = proportion of observations belonging to the i^{th} class at node t .
 $(e.g.) p(0/t) = \text{proportion of observations belonging to class } 0 \text{ at node } t$

note: t is time

For a binary classification problem, where $y_i \in \{0, 1\}, \forall i$

$$p(0/t) + p(1/t) = 1.$$

two classes $\rightarrow G_1$ and B



$$P(G_1|t_{SL}) = 1/3$$

$$P(B|t_{SL}) = 2/3$$

for a given node

For a binary classification problem if the class distribution is

$(1,0)$ or $(0,1)$ then the node is called a pure node.

However if the class distribution is $(0.5, 0.5)$ then the

node is called as the most impure node.

GINI SPLITTING RULE:

The probability distribution can say

whether purity of a node is present or not, but comparison between nodes

Although the conditional probability distribution of the classes help us to understand the purity of a node, but it's difficult to use such distributions to compare

the purity of multiple nodes.
 For the purpose of comparison, we need one single number that can help us to understand if one node is purer compared to another.

Genie Splitting Rule helps us to come up with such numbers that are useful for comparing the purity of multiple nodes.

$$Gini(t) = 1 - \sum_i p_i^2 (i/t) - (e^t) - (e^t) + 1 = (impure)$$

$$Gini(t) = 1 - p^2 (0/t) - (e^t) - (e^t) + 1 = (pure)$$

Note: In binary classification, a pure node has a genie = 0.

1. For a binary classification problem, if $t \in \{0, 1\}$ then
 e.g. if a node t has a class distribution of $(1, 0)$ then

$$Gini(t) = 1 - 1^2 - 0^2 = 0$$

The same holds to if class distribution of t is $(0, 1)$.

$$Gini(t) = 1 - 0^2 - 1^2 = 1 - 0.5 = 0.5$$

2. For a binary classification problem an impure node bias $Gini = 0.5$

If a node t has a class distribution of $(1/2, 1/2)$ then

$$Gini(t) = 1 - (1/2)^2 - (1/2)^2 = 1/2 = 0.5$$

3. For a multiclass classification problem, a pure node has

$Gini$ equal to 0. (i) $Gini = 0$, (ii) $Gini = 1$, (iii) $Gini = 2$, (iv) $Gini = 3$, (v) $Gini = 4$, (vi) $Gini = 5$, (vii) $Gini = 6$, (viii) $Gini = 7$, (ix) $Gini = 8$, (x) $Gini = 9$, (xi) $Gini = 10$, (xii) $Gini = 11$, (xiii) $Gini = 12$, (xiv) $Gini = 13$, (xv) $Gini = 14$, (xvi) $Gini = 15$, (xvii) $Gini = 16$, (xviii) $Gini = 17$, (xix) $Gini = 18$, (xx) $Gini = 19$, (xxi) $Gini = 20$, (xxii) $Gini = 21$, (xxiii) $Gini = 22$, (xxiv) $Gini = 23$, (xxv) $Gini = 24$, (xxvi) $Gini = 25$, (xxvii) $Gini = 26$, (xxviii) $Gini = 27$, (xxix) $Gini = 28$, (xxx) $Gini = 29$, (xxxi) $Gini = 30$, (xxii) $Gini = 31$, (xxiii) $Gini = 32$, (xxiv) $Gini = 33$, (xxv) $Gini = 34$, (xxvi) $Gini = 35$, (xxvii) $Gini = 36$, (xxviii) $Gini = 37$, (xxix) $Gini = 38$, (xxx) $Gini = 39$, (xxxi) $Gini = 40$, (xxii) $Gini = 41$, (xxiii) $Gini = 42$, (xxiv) $Gini = 43$, (xxv) $Gini = 44$, (xxvi) $Gini = 45$, (xxvii) $Gini = 46$, (xxviii) $Gini = 47$, (xxix) $Gini = 48$, (xxx) $Gini = 49$, (xxxi) $Gini = 50$, (xxii) $Gini = 51$, (xxiii) $Gini = 52$, (xxiv) $Gini = 53$, (xxv) $Gini = 54$, (xxvi) $Gini = 55$, (xxvii) $Gini = 56$, (xxviii) $Gini = 57$, (xxix) $Gini = 58$, (xxx) $Gini = 59$, (xxxi) $Gini = 60$, (xxii) $Gini = 61$, (xxiii) $Gini = 62$, (xxiv) $Gini = 63$, (xxv) $Gini = 64$, (xxvi) $Gini = 65$, (xxvii) $Gini = 66$, (xxviii) $Gini = 67$, (xxix) $Gini = 68$, (xxx) $Gini = 69$, (xxxi) $Gini = 70$, (xxii) $Gini = 71$, (xxiii) $Gini = 72$, (xxiv) $Gini = 73$, (xxv) $Gini = 74$, (xxvi) $Gini = 75$, (xxvii) $Gini = 76$, (xxviii) $Gini = 77$, (xxix) $Gini = 78$, (xxx) $Gini = 79$, (xxxi) $Gini = 80$, (xxii) $Gini = 81$, (xxiii) $Gini = 82$, (xxiv) $Gini = 83$, (xxv) $Gini = 84$, (xxvi) $Gini = 85$, (xxvii) $Gini = 86$, (xxviii) $Gini = 87$, (xxix) $Gini = 88$, (xxx) $Gini = 89$, (xxxi) $Gini = 90$, (xxii) $Gini = 91$, (xxiii) $Gini = 92$, (xxiv) $Gini = 93$, (xxv) $Gini = 94$, (xxvi) $Gini = 95$, (xxvii) $Gini = 96$, (xxviii) $Gini = 97$, (xxix) $Gini = 98$, (xxx) $Gini = 99$, (xxxi) $Gini = 100$.

E.g.: If a node t has 4 classes with a class distribution of $(1/4, 1/4, 1/4, 1/4)$ then $Gini(t) = 1 - 0^2 - 1/4^2 - 1/4^2 - 1/4^2 = 0$.

We see that the value of $Gini$ decreases with increase in purity of a node and increases with an increase in impurity of a node. Therefore, $Gini$ can be thought as the measure of "badness of a split". (i) higher the value of $Gini$ worse is the split.

Homework: Compute the $Gini$ of each node arising due to splits by savings and assets. Express the numbers in decimals.

By savings and assets: Express the numbers in decimals.

Sibarit Dujitum Jo Wiway

Comparing purity of 2 parent nodes
of basegmas purg at sbar and 2i berasberg at 26 9/3/14

$$\text{Gini}(t_{SL}) = 1 - \sum p^2(i)$$

$$= 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 1 - \frac{1}{9} - \frac{4}{9} = \frac{9-5}{9} = \frac{4}{9} = 0.44$$

$$\text{Gini}(t_{SM}) = 1 - (1)^2 - (0)^2 = 0$$

$$\text{Gini}(t_{SH}) = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 1/2 = 0.5$$

$$\text{Gini}(t_{AL}) = 1 - (0)^2 - (1)^2 = 0$$

$$\text{Gini}(t_{AM}) = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = 1 - \frac{9}{16} - \frac{1}{16} = \frac{16-10}{16} = \frac{6}{16} = \frac{3}{8}$$

$$\text{Gini}(t_{AH}) = 1 - (1)^2 - (0)^2 = 0.$$

~~LOGISTIC REGRESSION~~ instead of node splits for now

→ Individual gene values tell us about the purity of the terminal nodes.

→ Savings & Assets are the parent node

To know which of these parent nodes gives us a purer splits, we make use of "Weighted Gini"

$$\text{Weighted Gini} = \sum \text{Prob} \times \text{Gini}$$

$$= P(\text{Savings} = L) \times \text{Gini}(\text{Savings} = L) + P(\text{Savings} = M) \times \text{Gini}(\text{Savings} = M) + P(\text{Savings} = H) \times \text{Gini}(\text{Savings} = H)$$

$$= \left(\frac{3}{8} \times 0.44\right) + \left(\frac{3}{8} \times 0\right) + \left(\frac{2}{8} \times 0.5\right) = 0.1875$$

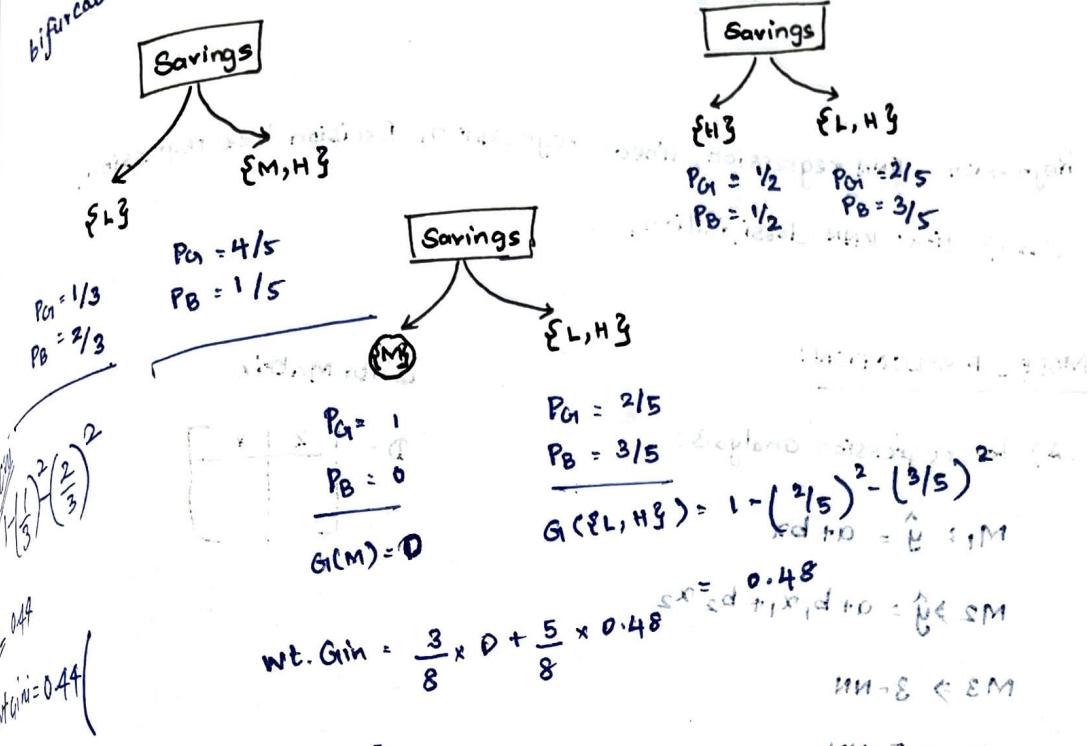
$$\text{Weighted Gini} = \sum \text{Prob} \times \text{Gini}$$

$$= \left(\frac{7}{8} \times 0\right) + \left(\frac{1}{8} \times \frac{3}{8}\right) + \left(\frac{1}{8} \times 0\right)$$

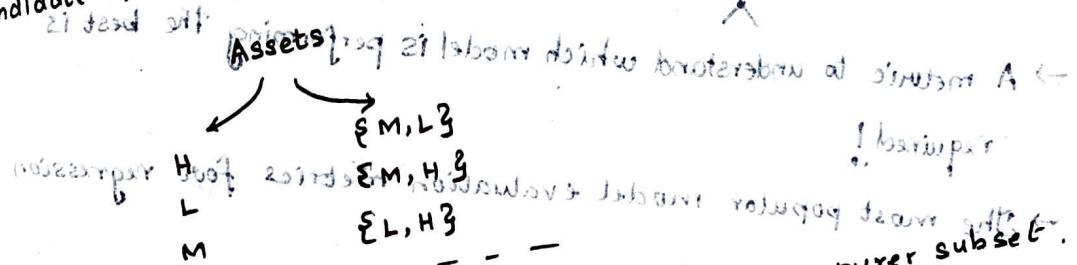
$$= \frac{12}{64} = 0.1875$$

Since $\text{wt. Gini}(\text{Assets}) < \text{wt. Gini}(\text{Savings})$, we confirm that splitting on the variables "Assets" gives a purer split.

CREATING A DECISION STUMP USING BINARY SPLIT ONLY (CART MODEL):
 Cart doesn't allow us to do multiple splits ($splits > 2$).
 To solve the same problem using cart, we need to bifurcate the data as follows (Candidate splits on the variable savings)



Candidate splits on the Variable Assets



Homework: Find out which candidate split gives a purer resulting subset.

$$\begin{aligned}
 & \text{Assets Node: } G(M) = 0 \\
 & \text{S Node: } G(L, H) = 0.375 \\
 & \text{L, M Node: } G(L, M) = 0.5, G(H) = 0.5 \\
 & \text{Final Gini: } 0.5 \times 0.5 + 0.5 \times 0.5 = 0.5
 \end{aligned}$$

$$\begin{aligned}
 & \text{Assets } \$L3, \$M3, \$3 \\
 & Wt. Gini = (2/8 \times 0 + 6/8 \times 0.277) = 0.437 \\
 & \quad (\< 0.437) \text{ justify splitting} \\
 & \quad \text{ob} = \frac{4}{8} \times [G(M) + \frac{4}{8} \times G(L)] \\
 & \quad = 0.207 \\
 & \quad \text{at } 0.207 \text{ split point} \\
 & \text{Assets: } \$H3, \$L3, \$M3 \\
 & Wt. Gini = (2/8 \times 0) + (6/8 \times 0.5) = 0.375 \\
 & \quad (\> 0.375) \text{ no split}
 \end{aligned}$$

Regression: KNN regression, linear regression, Decision Tree regression.

Classification: KNN classification,

MODEL EVALUATION:

(A) For regression analysis:

$$M_1: \hat{y} = a + bx$$

$$M_2: \hat{y} = a + b_1x_1 + b_2x_2$$

M3 \Rightarrow 3-NN

M4 \Rightarrow 5-NN

M5 \Rightarrow



Data Matrix

$$D = \begin{bmatrix} x & y \\ \vdots & \vdots \\ 0 & (M_1)_0 \end{bmatrix}$$

\rightarrow A metric to understand which model is performing the best is required!

\rightarrow The most popular model evaluation metrics for regression problems are:

$$1. \text{ MAE (Mean Absolute Error)} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| = \frac{1}{n} \sum_{i=1}^n |e_i|$$

Most interpretable one

It says what is the average deviation of the predicted values from the actual values.

$$2. \text{ MSE (Mean Squared Error)} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{n} \sum_{i=1}^n e_i^2$$

$$3. \text{ RMSE (Root Mean Squared Error)} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} = \sqrt{\text{MSE}}$$

$$4. \text{ MAPE (Mean absolute Percentage Errors)} = \frac{1}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right| \times 100\% = \frac{1}{n} \sum_{i=1}^n \left| \frac{e_i}{y_i} \right| \times 100\%$$

MAPE helps us to express the error in percentage form. $\rightarrow \left(\frac{\text{error}}{\text{actual}} \right) \times 100$

Since all these are error metric the lesser the value, the better it is.

Therefore, among all the models that are fitted on a data, the model with the least value of RMSE should be considered and is expected to give better results.

(B) For Classification models

We compute something called as Confusion Matrix

A confusion matrix is a bivariate table between the actual values and the predicted values. For example, for a binary classification problem with 2 classes "Positive" and "Negative", a confusion matrix would look like:

		Predicted	
		+ve	-ve
Actual	+ve	#TP	#FN
	-ve	#FP	#TN

x	y	g
+	+	TP
+	-	FN
-	-	TN
-	+	FP
5	+	

TP + # FN = Total no. of actual positive (+ve) cases

FP + # TN = Total no. of actual negative (-ve) cases

TP + # FP = Total no. of predicted positive (+ve) cases

FN + # TN = Total no. of predicted negative (-ve) cases.

TP \rightarrow Truly classified into +ve class
 FN \rightarrow False classified into -ve class
 TN \rightarrow Truly classified into -ve class
 FP \rightarrow False classified into +ve class

POPULAR EVALUATION METRICS FOR CLASSIFICATION PROBLEMS:

1. Overall error rate: $\frac{\# \text{mistakes}}{\# \text{obs}}$

2. Overall accuracy = $1 - \text{Overall error} = \frac{\# \text{correct predictions}}{\# \text{obs}}$

Overall accuracy = $\frac{\# \text{TP} + \# \text{TN}}{n}$

001 x (75773) \leftarrow most, spams in years 2013-2014
3. Sensitivity (or Recall) = True Positive Rate (TPR) \rightarrow 29/31

It means

Actual fraudulent transactions ≥ 150

Out of 150 fraudulent transactions, 100 were detected right.

		Prediction	
		True	False
Actual	True	100	50
	False	150	9700

$\frac{100}{150} \rightarrow$ Even though the transaction alarm goes off as a false alarm, it is not fraudulent, model predicts it as fraudulent.

$$\text{Sensitivity} = \frac{\text{# correct predictions}}{\text{# actual positives}}$$

$$= \frac{\text{# TP}}{\text{# TP} + \text{# FN}}$$

\rightarrow Sensitivity measures the proportion of positive cases that are correctly predicted by the model.

$$\rightarrow \text{In our example, Sensitivity} = \frac{100}{100+50} = \frac{100}{150} = 0.667$$

which means 66.67% of the positive cases are identified correctly by the model.

4th Specificity or True negative rate (TNR)

Specificity measures the proportion of negative cases that are correctly predicted by the model.

$$\text{In our example, Specificity} = \frac{9700}{9700+150} = 0.984$$

That is, 98.4% of the negative cases are identified correctly by the model.

5. Precision: Helps to understand the number of precise predictions.

$$\text{Precision} = \frac{\text{# no. of correct (true) prediction}}{\text{Total no. of (true) prediction}} = \frac{\# TP}{\# TP + \# FP}$$

Higher precision, lower no. of false alarms

$$\text{In our example, precision} = \frac{100}{100+150} = \frac{100}{250} = 2/5 = 0.4$$

which means 40% of the positive predictions made by the model are correct predictions.

→ Note that, higher the value of the precision the less will be the no. of false alarms and vice versa.

6. F1 score (most popular)?

F1 Score = H.M (Precision, Recall)

Property of arithmetic mean: A. Mean of 10 and 100 = $\frac{110}{2} = 55$.
Mean towards higher value.

It says either precision or recall is high.

Property of harmonic mean: H. mean of 10 and 100

says Both precision and recall is high.

higher value of harmonic mean says higher the both values.

Since precision and recall are 2 rate measures

it's appropriate to use harmonic mean to compute the average of these two rates.

Moreover, harmonic mean ensures that, a higher value of F1 score would mean both precision and recall are high.

In general we don't compute accuracy in training data. Because it gives biased results as model has already seen the training data.

The moment you get the data, divide the data randomly into 2 parts (training and test), even before data exploration.

First, secure your test data,

$$\text{avg speed} = \frac{\text{Total Dist}}{\text{Total Time}}$$

$$\begin{aligned} &= \frac{2d}{\frac{d}{60} + \frac{d}{40}} \\ &= \frac{2}{\frac{1}{60} + \frac{1}{40}} \\ &= \text{HM}(60, 40) \end{aligned}$$

$$\text{avg speed} = \frac{60+40}{2}$$

Data exploration and data cleaning

VALIDATION

TRAIN TEST SPLIT: (TRAIN, VALIDATION, TEST SPLIT)

do not mix training and test data

→ We randomly split the data into 2 parts. We call one

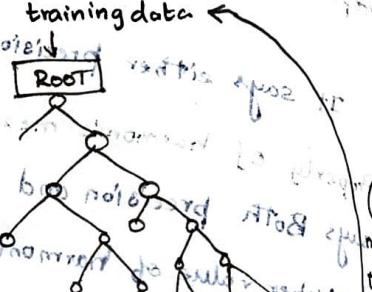
part as the training data and the other one as test data.

→ The training data is used to train a model and the test data is used to evaluate the out-of-sample performance of the model.

CAUTION: The test data should never be exposed to the user or to the model until at the final stage of model validation.

→ Somewhere we take whispers from data and make our model being exposed to test data even though we're training our model

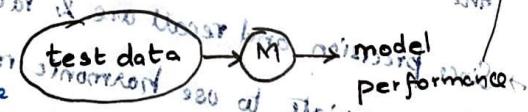
on the test data.



training data → ROOT → d=1 → d=3 → leaf
d=2 → d=5 → leaf

→ Here we pass test data into the model and evaluate model's performance. Then based on model performance we decide the

model performance we decide the model to grow more.



test data → M → model performance

→ Then we observe that when test data is passed model performance is ↑. We come to intuition making decision tree grow results are good. But once again when we allowed the

model performance comes down (↓) randomly a tree to grow, the model performance goes down (↓) randomly.

→ But this is not preferred because of \bullet .

IDEAL FRAMEWORK:

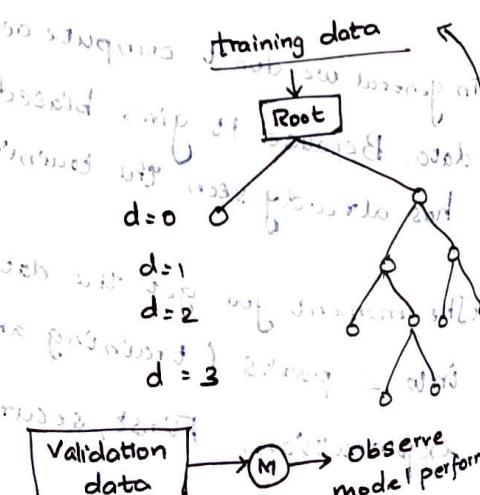
Step 1:

Wrote down data (sample)

→ Split into training and test

training validation test

keep it in complete darkness



- Step 2: Data exploration (on training data only).
- Step 3: Data cleaning (on training data only).
- Only after splitting data, we go with data exploration and then data cleaning.
- When sample size is 100, we first split as 80-20, $\frac{80}{100} = \frac{4}{5}$ test data, $\frac{20}{100} = \frac{1}{5}$ training data. When class distribution is uniform, we make use of SRS.
- For regression, we can use SRS, since there are no classes.

3. CROSS-VALIDATION:

(A) k-fold cross validation

k is an integer.

Let's say $k=5$, then

$$S = \boxed{S_1 \mid S_2 \mid S_3 \mid S_4 \mid S_5}$$

5 random and equal splits.

Iteration 1:

hold = S_1 (as validation data)

$$\text{training } (T) = S_2 \cup S_3 \cup S_4 \cup S_5$$

We train our model on T and validate the model on hold.

ALGORITHM:

1. Let S be the training dataset.

2. Randomly partition S into k equal parts.

S_1, S_2, \dots, S_k , such that $\bigcup_{i=1}^k S_i = S$ & $S_i \cap S_j = \emptyset$ if $i \neq j$.

3. for i from 1 to k do:

4. $h_i = S_i$ (holdout sample)

5. $t_i = \bigcup_{j=1, j \neq i}^k S_j$ (training sample)

6. m_i = model fitted on t_i

7. e_i = validation score of m_i on h_i

8. end for

9. c.v score = aggregated score of e_i (average).

Every time you take data randomly from population and use the data to fit into the regression (for each data, we have different parameters for the regression line (because the data changes) for every model we get different values of b_0 and b_1 , and \rightarrow similarly we can also say for classification decision tree. Each data results in different decision

trees. We compute the error of each models and compute the

CV score.

→ Lesser the value of size of training data, higher the value of k .

→ If training data = 20 and set $k = 20$.

I'll split such that during every iteration,

validation data = 1 observation and training data = 19 observations.

which means at a whole my model will be validated on each of the 20 observation (h_i where $i = 1 \text{ to } 20$)

→ If we have 4 models M_1 : logistic, M_2 : linear, M_3 : decision tree, M_4 : kNN, we'll perform cross validation on each of these models and determine the CV score of each model.

→ Suppose M_2 has lesser CV score, we choose M_2 over other models.

(B) LOOCV (Leave one out cross validation)

1. Let S be the training data of n examples.

2. Randomly partition S into n equal parts. O_1, O_2, \dots, O_n

such that $n(O_i) = 1, \forall i$, $\cup O_i = S$, $O_i \neq O_j \forall i \neq j$.
No duplicate observations

3. for $i \in [1 \text{ to } n]$ do :

4: $h_i = O_i$ (holdout data)

5. $t_i = \bigcup_{\substack{j=1 \\ j \neq i}}^n O_j$ (training data)

6. $m_i = \text{model fitted on } t_i$ {inside for loop}

7. $e_i = \text{validation score of } m_i \text{ on } h_i$

8. end for

9. C.V. score = aggregated score of e_i

Supervised Machine Learning:

→ we initiate the learning of the algorithm with labelled data

→ But the learning algorithms can be different.

You have a kid and show him square & rectangle. You say that what's length, breadth, if all length are same it's a square. However when lengths > breadth, then it's a rectangle.

→ We show the object to the kid and also we show the learning algorithm.

→ We are supervising the learning process with label.

→ We are supervising the learning process with label.

→ The child has been exposed to the several shapes along with supervision done using labels and learning algorithm.

→ Discriminate the shape of object.

For every value we're bring into the model, we supervise the learning using the algorithm.

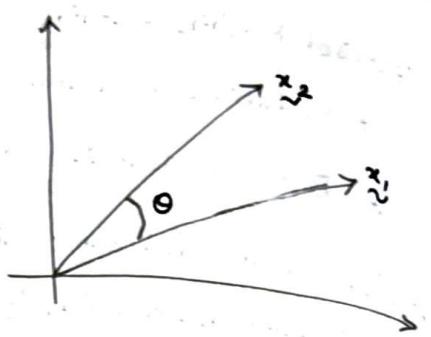
→ In KNN model Let's say we need to find the \hat{y} -value of y for a given vector $(100, 20, 10)$. In KNN we find the 5 (k) nearest neighbours by calculating the euclidean distance between x (vector) and all other datapoints represented as vectors and find the 5 nearest neighbours. There's no y variable involved while determining the nearest neighbours. KNN is controversial to be considered as supervised because we don't use the y variable to determine the neighbours.

→ We're not definitive what groups to created, there exists some groups and determine

→ Distance - measure of dissimilarity. If distance between 2 datapoints is larger, then they are more dissimilar.

→ Cosine - measure of similarity

Cosine between 2 vectors is
the measure of correlation between
them.



$$\cos \theta = \frac{\mathbf{x}_1^T \mathbf{x}_2}{\|\mathbf{x}_1\| \|\mathbf{x}_2\|}$$

$$= \frac{\sum_{i=1}^n x_{1i} \sum_{i=1}^n x_{2i}}{\sqrt{\sum_{i=1}^n x_{1i}^2} \sqrt{\sum_{i=1}^n x_{2i}^2}}$$

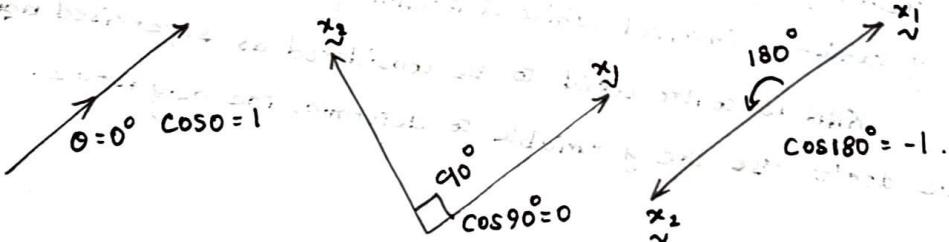
$$\cos \theta = \frac{\sum_{i=1}^n x_{1i} x_{2i}}{\sqrt{\sum_{i=1}^n x_{1i}^2} \sqrt{\sum_{i=1}^n x_{2i}^2}} \quad \text{--- (1)}$$

$$\sqrt{\sum_{i=1}^n x_{1i}^2} \sqrt{\sum_{i=1}^n x_{2i}^2}$$

$$\text{Cor}(\mathbf{x}_1, \mathbf{x}_2) = \frac{\sum (x_{1i} - \bar{x}_1)(x_{2i} - \bar{x}_2)}{\sqrt{\sum (x_{1i} - \bar{x}_1)^2} \sqrt{\sum (x_{2i} - \bar{x}_2)^2}} \rightarrow *$$

When $\theta = 0^\circ$, then $*$ becomes equal to (1).

Special case of correlation coefficient is cosine similarity



Cosine similarity \rightarrow -1 to +1. \rightarrow most similar.

\rightarrow Netflix recommendation systems
cosine similarity: most dissimilar

Customer A and B have the affinity

to buying products 1 and 2, hence
they are more similar. Group customers
based on similarities

Euclidean distance:

Customer A is light buyer and customer B is a heavy buyer.
Customer A buys much less than customer B. Distance between
them is high (\rightarrow) more dissimilar.

