

Grid Search (Used to find tune hyperparameters)

- 1) If there are large no. of columns and you feel many are useless & ~~interact~~, go for Lasso and also interactions.

2) Dropouts like to follow the rule of thumb 10% - 20%.

Supervised Learning:

Classification

- Dependent / Target is categorical variable.
- If means are 1.96 SD away from each other, then it is considered to be significant.
- For Non linear model, there are no assumptions.

Gradient Descent: (Iterative Approach)

$$\Delta \beta \propto -\Delta \text{Error}$$

$$\textcircled{2} \quad \Delta \beta \propto \frac{1}{\Delta \text{Error}}$$

$$\Delta \beta = -\eta \Delta \text{Error}$$

η - learning Rate

$$\textcircled{3} \quad \sigma^2 = \frac{1}{2} (y - \hat{y})^2$$

$$\frac{1}{2} (y - (\beta_0 + \beta_1 x))^2$$

$$\text{differentiate : } \frac{1}{2} \cdot 2 (y - (\beta_0 + \beta_1 x))$$

$$\text{diff w.r.t } \beta_0 \quad -(y - \hat{y})$$

$$\Delta \beta_0 = \frac{\partial \sigma^2}{\partial \beta_0} = -(y - \hat{y})$$

$$\Delta \hat{P}_i = \frac{\partial E}{\partial P_i}$$

$$(\text{Error})^2 = \frac{1}{2} \cdot 2 \left(y - (P_0 + P_1 x) \right)^2$$

dry diff w.r.t. P_i will be equal to zero if

$$\Delta P_i = -(y - \hat{y}) x$$

As error gradient is smaller, it will approach the best fit line

- Once the SSE reaches saturation - 2 consecutive iteration's SSE difference is insignificant (stopping criteria)
- Once the SSE starts increasing, we get to know that, we are going away from best fit line (value there we stop)

Drawback

- Tends to underfit - full batch gradient - high bias error
- Tends to overfit - SGD - high variance error

Mini batch learning \rightarrow Non linear model.

Logistic Regression: (Sigmoid)

\hat{y} of classification should give probability score

$$\hat{y}_{\text{prob}} = \frac{(x_0 + 1) - 1}{1 + e^{-(P_0 + P_1 x)}}$$

$$(P - \hat{P}) = 1 - \frac{e^{-P}}{1 + e^{-P}}$$

Also called as Soft Limit Function

DATE: / /
PAGE: _____

$$\text{Error Gradient} = \frac{1}{2} \left(y - \frac{1}{1 + e^{-(B_0 + B_1 x_1)}} \right)$$

$$\frac{\partial E}{\partial B_0} = \frac{1}{2} \cdot 2 \left(y - \frac{1}{1 + e^{-(B_0 + B_1 x_1)}} \right)$$

$$= \left(0 - \left(\frac{1 + e^{-(B_0 + B_1 x_1)}}{1 + e^{-(B_0 + B_1 x_1)}} \right) \right) - (1) \left(0 + e^{-(B_0 + B_1 x_1)} \right)$$

$$= -e^{-(B_0 + B_1 x_1)} \frac{(1 + e^{-(B_0 + B_1 x_1)})^2}{(1 + e^{-(B_0 + B_1 x_1)})^2}$$

$$\boxed{\frac{\partial E}{\partial B_0} = -\hat{y} (1 - \hat{y}) (y - \hat{y})}$$

$$\frac{\partial E}{\partial B_1} = -(y - \hat{y}) \hat{y} (1 - \hat{y}) x$$

If probability score $> 0.5 \Rightarrow \text{class} = 1$
 $< 0.5 \Rightarrow \text{class} = 0$

Concept of Odd's

Mutually Exclusive \rightarrow Union (Addition)

$$\boxed{\text{odd's} = \frac{P(\text{occurring})}{P(\text{non occurring})}}$$

odd's ratio is Relative Risk.

biased coin $P(H) = 0.7$

$$O(H) = \frac{0.7}{0.3} = 2.33$$

Fair coin $P(H) = 0.5$

$$O(H) = \frac{0.5}{0.5} = 1$$

odd's of head in biased coin is 2.33 times the odd's of fair coin.

SE - side effects

Coraxin

$$P_{SE} = 0.2$$

Covishield

$$P_{SE} = 0.1$$

$$O(P_{SE}) = \frac{0.2}{0.8}$$

$$O(P_{SE}) = \frac{0.1}{0.9}$$

$$O(P_{SE}) = 0.25$$

$$O(P_{SE}) = 0.11$$

$$\frac{O_{CE}(P_{SE})}{O_{CS}(P_{SE})} = 2.3$$

Odds of side effects of Coraxin is 2.3 times the odds of side effects of Covishield.

$$\hat{y}_{Prob} = \frac{e^{-z}}{1 + e^{-z}} \rightarrow \text{let this be } x$$

$$= \frac{1}{1 + e^{-z}} \times \frac{e^z}{e^z} \quad \begin{matrix} (x \\ e^{-z}) \text{ as } e^z \end{matrix}$$

$$= \frac{e^z}{1 + e^z}$$

$$\frac{\hat{y}_{\text{prob}}}{1 - \hat{y}_{\text{prob}}} = \frac{e^z}{1 + e^z}$$

$$= \frac{e^z}{1 + e^z} \cdot \frac{1 - e^z}{1 - e^z}$$

$$= \frac{e^z}{1 + e^z} \cdot \frac{1 + e^z - e^z}{1 + e^z}$$

$$\hat{y}_{\text{prob}} = e^z$$

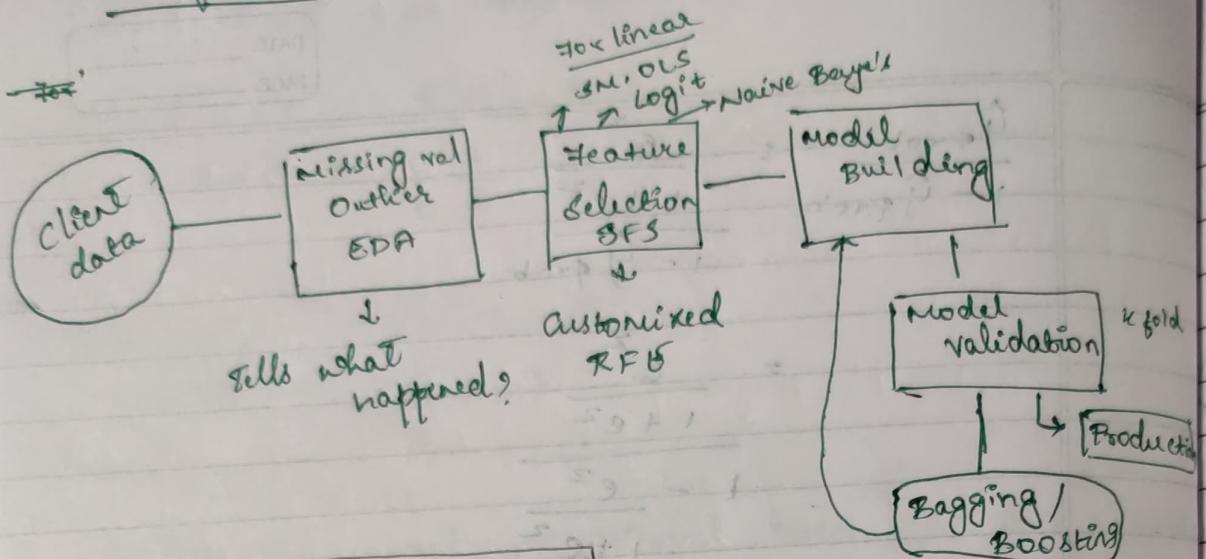
$$1 - \hat{y}_{\text{prob}} = e^{(B_0 + B_1 x)}$$

$e^{(B_0 + B_1 x)} \Rightarrow \text{odds}$

Raising coeff to e^z is called odds ratio.
In logistic reg, we can't interpret slopes as plus times x as we do in linear regression.

If odds > 1 , then odds are significant.

Flow for ML



Bagging - \downarrow des variance error.

Boosting - \downarrow des bias error.

→ Can be applied
to any model

drawback of SFS

- Doesn't see value of columns but looks into only R^2
- Drawback of RFB
User has to specify no of features

Linear regression - strictly used as Regressor.

Logistic regression - strictly used as classifier

Naive Bayes - strictly used as classifier

Linear
Models

K-NN - can be used as Regressor or as classifiers - It is Non linear Model

Decision Tree - Regressor & classifier

Random Forest - Regressor & classifier

Ensemble Models - Non Linear Bagging & Boosting.

1 sample mean test / 1 sample proportion Test is used to validate the client data.

	x_1	x_2	x_3	y
1				
2				
3				
4				
5				
6				

Shapes:

x_{Train}

4×3

x_{Test}

2×3

y_{Train}

4×1

y_{Test}

2×1

Scaling Compulsory Models = KNN & Naive Bayes

Ex. add constant (x)

constant

1.00

1.00

1.00

this constant is the multiplication factor of B_0 (coeff of B_0)

In sklearn, model.fit(X, Y)

But in statsmodels, logit(Y, X).fit()

$$\text{Pseudo } R^2 = 1 - \frac{\text{log likelihood}}{\text{LL Null}}$$

↳ practical range of $R^2 = 0.2 \text{ to } 0.6$

loss func of classifiers = log loss \circledast entropy

$$\text{log loss} = -(y_{\text{act}} \cdot \log_2 \hat{y}_{\text{prob}} + (1 - y_{\text{act}}) \cdot \log_2(1 - \hat{y}_{\text{prob}}))$$

① for log loss to be minimal,

\hat{y}_{prob} for $y_{\text{act}} = 0$ should near 0 &

\hat{y}_{prob} for $y_{\text{act}} = 1$ should near 1

② Maximum likelihood \rightarrow log loss will be minimal

$$\hat{y}_{\text{prob}} = \max$$

$$1 - \hat{y}_{\text{prob}} = \max$$

$y_{\text{act}} \quad y_{\text{prob}}$

$$0.2 \quad 1 - \hat{y}_{\text{prob}} = \hat{y}_{\text{prob}}$$

$$0 \quad 0.3 \quad 1 - \hat{y}_{\text{prob}} = \hat{y}$$

$$1 \quad 0.8 \quad \hat{y}_{\text{prob}} = \hat{y}$$

$$1 \quad 0.9 \quad \hat{y}_{\text{prob}} = \hat{y}$$

What evaluation metric we choose in validation method matters.

Validation Method

Train/Test split

- NO preferred as it is like assessing model with only 1 test

✓ Cross val score

KOOCV

- Not preferred as all exec are same & only 1 exec is changed.
(Test data is less & train data is repetitive)

for Linear Reg:-

$kf = \text{kFold}(\text{n_split} = 5, \text{shuffle} = \text{True}, \text{random} = \text{True})$
 score = cross_val_score(LR, X, Y, cv=KF)
 scoring = 'neg_root_mean_squared_error')

↳

by default, $R^2 \rightarrow$ higher the R^2 better.
 higher R^2 gives as better.

But, in reality, lowest R^2 is better
 ↳ Hence 'neg' is given. in scoring.

np.mean(np.abs(score))

np.std(score, df=1)

for Logistic Reg:-

Yact \hat{Y}_{pred}

Yact \hat{Y}_{pred}

Confusion Matrix: (Performance Metric)
 PREDICTED.

Actual

PREDICTED

H CVD

H H

H H

H H

CVD

2 0

CVD

H

H

Predicted

CVD H

N P

		N	P	
		True Negative specificity of model	False Positive	-Type I Error or Errors
Actual	P	False Negative	True Positive	Sensitivity of the model
		Type II Error (B error)	Type I Error (A error)	

** Type I error is dangerous than Type II.

Confusion Matrix is performance metrics

similar to RMSE in Linear Reg.

Confusion Matrix - for all classification model

Methods to set threshold:

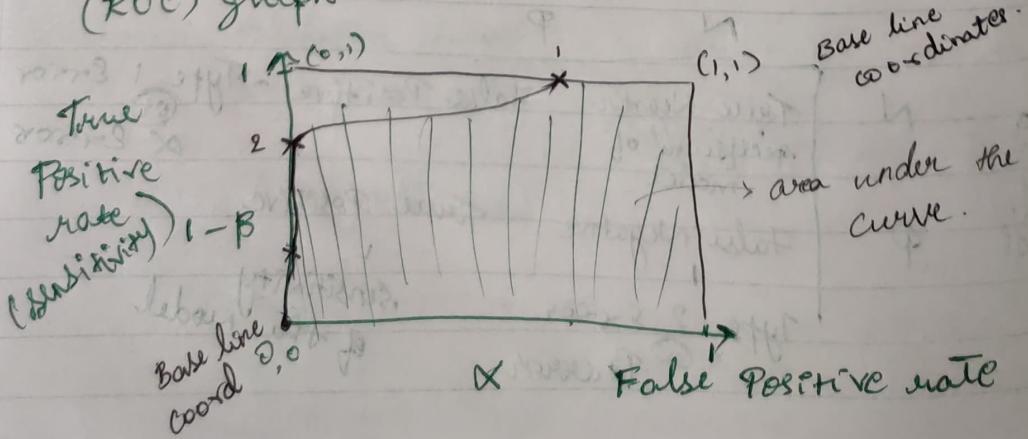
	\hat{y}_{prob}	\hat{y}_{class} $T_1 = 0.25$	\hat{y}_{class} $T_2 = 0.4$	\hat{y}_{class} $T_3 = 0.8$
H	0.35	1	0	0
H	0.2	0	0	0
H	0.4	(comes 1.00. qn)	0	0
CVD	0.68	(1 = H, 0 = CVD)	0	0
CVD	0.75	1	1	0
CVD	0.85	1	1	1

	H	CVD		H	CVD		H	CVD
H	2	1		3	0		3	0
CVD	0	3		CVD	2		CVD	2
Overall Accuracy	$\frac{4}{6}$			acc = $\frac{5}{6}$			acc = $\frac{4}{6}$	

Accuracy is not sensitive enough to capture the performance - hence plot.

Receiver operating characteristics (ROC)

(ROC) graph



Using this graph, we calc Area under the curve (AUC)

map 3 confusion Matrix on ROC graph.

$$\text{For } 1^{\text{st}} \text{ one} = \alpha = 2/3 = 66.6\% \\ \text{sensitivity} = 3/3 = 100\%$$

$$2^{\text{nd}} \alpha = 0\% \\ \text{sensitivity} = 2/3 = 66.6\%$$

$$3^{\text{rd}} \alpha = 0\% \\ \text{sensitivity} = 1/3 = 0.33$$

Sklearn

\downarrow
model.predict(X)
 \hookrightarrow o/p class label

\downarrow
model.predict_proba(X)
 \hookrightarrow o/p prob score

lm. logit

\downarrow
model.predict()

\downarrow
o/p prob score

\downarrow
get class label ~~manually~~
manually

kappa score or Reliability score.

\hookrightarrow overall accuracy adjusted to Baseline Reference.

$$\text{kappa score} = \frac{\text{OA} - \text{Pchance}}{1 - \text{Pchance}}$$

Eg:-

		-	+
-	15	15	
	10	60	
+			25%
Total	25	75	

Exp predict by chance = $\frac{25}{100} + \frac{75}{100}$

$$= 50\% + 75\% = 125\%$$

$$50\% + 75\% = 125\%$$

$$\text{Kappa score} = \frac{0.75 - 0.6}{1 - 0.6} = 37.5\%$$

Multicollinearity

corr() if < 0.3 OK (no coll)

$0.3 \text{ to } 0.5$ moderate coll

> 0.5 High coll

Recall & precision:-

G G G G G S S S S S

G G G G G G S S S S S

		G	S	ways	Recall	Precision
G	S	5	0	1	$\leftarrow G$	$100\%, 5/5$
		1	4	5/6	83%	$5/4$ 100%

$$F1 \text{ score} = \frac{2(R \times P)}{R + P}$$

- (1) Which threshold makes TPR the highest & FPR the lowest we choose that.

Or Youden's Index

- (2) Sum of α error & β error should be min
 ↳ Another way to choose threshold.
 ↳ Choose threshold which gives least total cost
 $\min\left(\frac{w_1 * FN}{P} + \frac{w_2 * FP}{\bar{A}}\right)$
- (3) Simply choose Mean.

3 ways ↑ to find threshold

NAIVE BAYE'S

Mutually Exclusive Variables

$$P(H \text{ or } T) = P(H) + P(T)$$

$$= 0.5 + 0.5 = 1$$

$$P(K \text{ or } Q) = P(K) + P(Q)$$

$$= \frac{4}{52} + \frac{4}{52} = \frac{8}{52}$$

$$P(A \cap B) = 0$$

Mutually Non Exclusive Variables

$$P(A \cap B) \neq 0$$

$$P(K \text{ or } \Delta) = P(K) + P(\Delta) - P(K \cap \Delta)$$

$$= \frac{4}{52} + \frac{13}{52} - \frac{1}{52}$$

Mutually Independent Variables

$$P(A \text{ and } B) = P(A) \cdot P(B)$$

conditional Probability
Mutually Non Ind Var

$$P(A \text{ and } B) = P(A) \cdot P\left(\frac{B}{A}\right)$$

A is affecting B

Baye's theorem

or- Union
and- Intersection

$$P(B) P(A|B) = P(A) P(B|A)$$

$$P(A|B) = \frac{P(A) P(B|A)}{P(B)}$$

Naive Baye's classifier

color Model Train Booking Y/N

(P(Booking = 'Y') / Red, Asta, MT) =

$$\frac{P(\text{Red}/\text{yes}) * P(\text{Asta}/\text{yes}) * P(\text{MT}/\text{yes})}{P(\text{Red}) * P(\text{Asta}) * P(\text{MT})}$$

P(Booking = 'No') / Red, Asta, MT) =

$$P(\text{Red}/\text{No}) * P(\text{Asta}/\text{No}) * P(\text{MT}/\text{No})$$

$$P(\text{Red}) * P(\text{Asta}) * P(\text{MT})$$

classifier classifies based on probability which is higher (Denominator is same & doesn't matter)

for Normal distribution

$$\frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

if $x \in A$

$$-(\alpha)^2 + (\beta)^2 = (\alpha - \beta)^2$$

$$(\alpha - \beta)^2$$

Never scale y.

In classifiers, bias error is the balance or 100% - auc-roc / wt

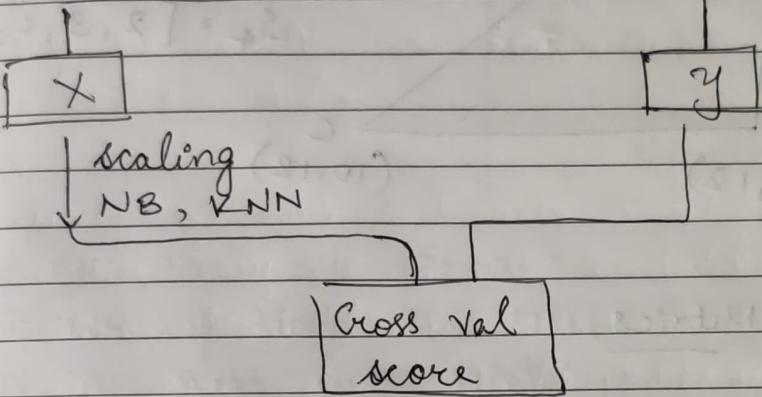
$Z = \frac{x - \bar{x}}{\sigma}$	1 = 2 × 3.2
DATE:	1 / 38 4
PAGE:	25 / 30 3 / 82
	29 / 12

Data frame



Dropna

1



cross_val_score(model, X_scaled, y, cv=kfold, score=

bi_class = 'roc-auc')

multi_class = 'f1-score'

log neg = 'neg-mean-square'

lin neg = 'r square'

No feature selection in cross val score happens

→ Bias Error ~~2.4~~ 4.190

Target price in \$

Var Error ~~2.0~~ 0.481

mean ~~3.772~~ (8)

Inference :-

$$0.481 \pm 4.190$$

$$4.190 \pm 0.481$$

× 1000

$$4190 \$ \pm 481 \$$$

Pred Avg car price is

Avg difference b/w the actual price & predicted price will lie b/w 3772 & 4608.

K - Nearest Neighbour (KNN)

(Classification Alg ⚡)

(5, 16)

A

(5, 12)

B

(10, 12)

C

similar

$$S_1 = [4, 7, 3, 2]$$

$$S_2 = [1, 2, 5, 8]$$

$$S_3 = [5, 6, 2, 3]$$

$$S_4 = [2, 3, 6, 9]$$

Distance Metrics

- 1) Euclidian Distance [Can travel diagonally]

$$\text{Distance b/w 2 points} = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

A & C

- 2) Manhattan Distance [Can't travel diagonally]
 Has to travel straight

$$|x_1 - x_2| + |y_1 - y_2|$$

- 3) Cosine Distance [Commonly used in Recommendation System].

First compute Cosine Similarity Score (CS)

$$\text{Cosine distance} = (1 - CS)$$

$$CS = \frac{\vec{A} \cdot \vec{C}}{\|\vec{A}\| \cdot \|\vec{C}\|}$$

$$CS = \frac{(5, 16) \cdot (10, 12)}{\|(5, 16)\| \cdot \|(10, 12)\|}$$

$$= \frac{50 + 192}{\sqrt{5^2 + 16^2} \cdot \sqrt{10^2 + 12^2}}$$

	ED	MD	CD (1-CS)
(Close) S ₁ -S ₃	2.0	4	CS = 0.97 CD = 0.03
(far) S ₁ -S ₂	8.6	16	CS = 0.57 CD = 0.428

For (S₁, S₃) & (S₁, S₂)

$$S_1, S_2 \approx \sqrt{(4-1)^2 + (7-2)^2 + (3-5)^2 + (2-8)^2}$$

$$= \sqrt{3^2 + 5^2 + 2^2 + 6^2} = 8.6.$$

$$S_1, S_3 \approx \sqrt{(4-5)^2 + (7-6)^2 + (3-2)^2 + (2-3)^2}$$

$$= 2.$$

$$S_1-S_3 = \frac{20 + 42 + 6 + 16}{\sqrt{4^2 + 7^2 + 3^2 + 2^2} \cdot \sqrt{5^2 + 6^2 + 2^2 + 3^2}} = 84$$

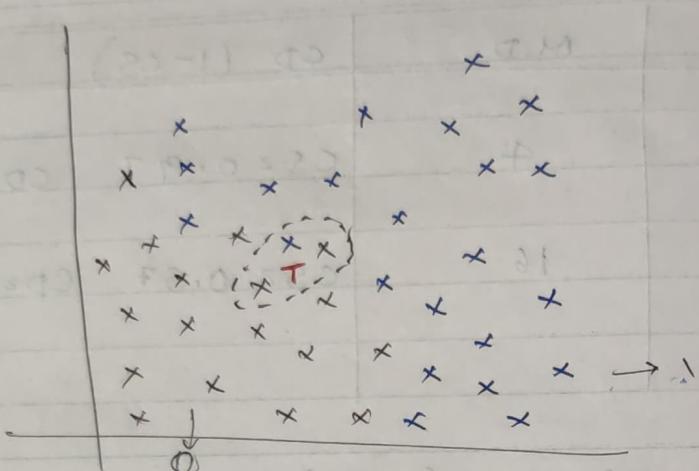
$$S_1, S_2 = \frac{4 + 14 + 15 + 16}{\sqrt{4^2 + 7^2 + 3^2 + 2^2} \cdot \sqrt{1^2 + 2^2 + 5^2 + 8^2}} = 49$$

Default mode of opⁿ | Unweighted or Uniform mode

Working of KNN

$n=3$.

optimal n -neighbour
2-5% of datarize



Blue - 1 KNN will wait for the Test records

Black - 0 Once they come in it will

T - Test rec place it on value and calculate

the distance from that test value to all the points and arranges it in ascending order & picks the first 3 and predicts it to be the max class of 3 (here it is black-0)

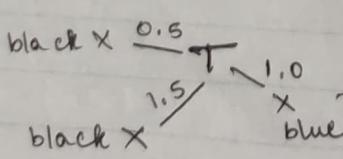
2nd Hyper parameter

K Nearest Neighbor (n -neighbor, weights = 'distance')

by default, weights = 'uniform'

~~Weighted Voting~~

$$NV = \frac{1}{0.5^2} + \frac{1}{1.5^2}$$



Weighted Voting

$$NV = \frac{1}{\text{dist}^2} = \frac{1}{1^2} = 1$$

$$NV = \frac{1}{1.5^2} + \frac{1}{0.5^2} = 4.44$$

KNN cannot be Boosted

DATE: / /
PAGE: _____

- model predicts the one which has wv more.
weighted mode - makes the model overfit
→ still we prefer this mode only, as we can bag it later. But never take risk of underfitting by doing unweighted / uniform
→ More refined results rather than going with the crowd.

can also be used as Regressor.

Both weighted & Unweighted can be done.
Avg of 3 as predicted.

180 mg
190 mg
200 mg

If unweighted avg glu lvl = 190 mg,
weighted avg will be < 190.

$$\frac{1}{0.5^2} * 180 + \frac{1}{1^2} * 190 + \frac{1}{1.5^2} * 200 = < 190$$

New mean will get attracted to the one with least distance.

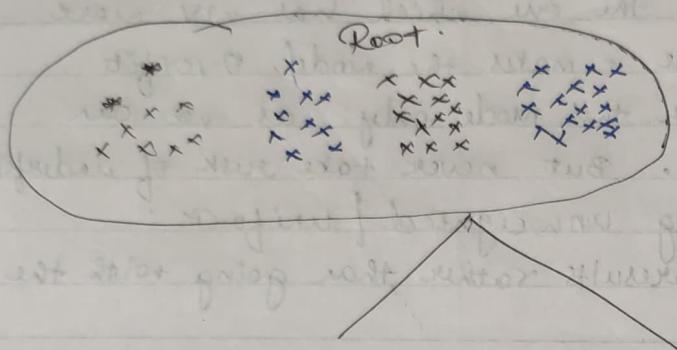
Creditab.

Last in Q16 Day 3	Model	Bias Error	Variance Error
N B	11.3		0.0115
KNN	6.4 $(100 - 93.6)$ $(100 - \text{KNN mean})$	0.0074	direct SD is taken

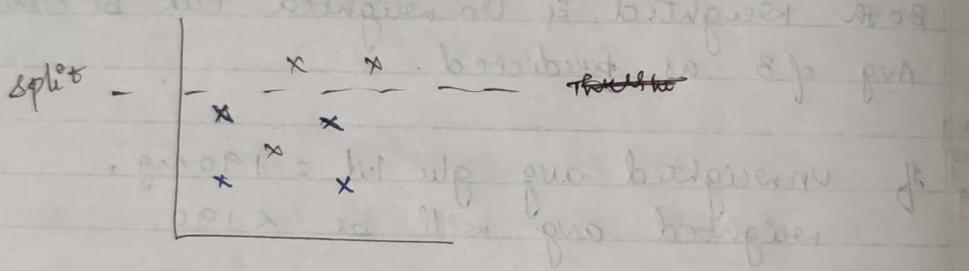
As BE & VE both are low in KNN,
KNN is a better model

Complexity

Decision Tree Working



Forms Axis line (horizontal / vertical) but not inclined



Which feature, what threshold can get best

Entropy or Gini Score

- Measure of Uncertainty

At root level, before splitting.

$$-(3/7) * \log_2(3/7) + (4/7) * \log_2(4/7) \text{ - entropy}$$

Entropy range - 0 to 1

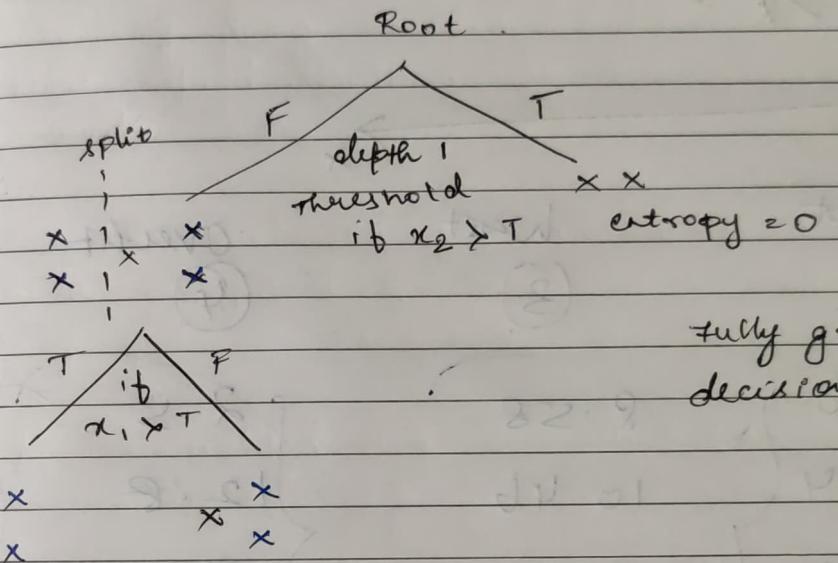
Gini score - 0 to 0.5

$$1 - ((3/7) * 2 + (4/7) * 2) \text{ - Gini score}$$

Aim of DT is to minimize uncertainty

Weighted entropy = entropy \times how many points / Total points

Difference b/w root Entropy and weighted entropy of depth 1 is known as information gain.



Entropy = 0

DT as Regressor

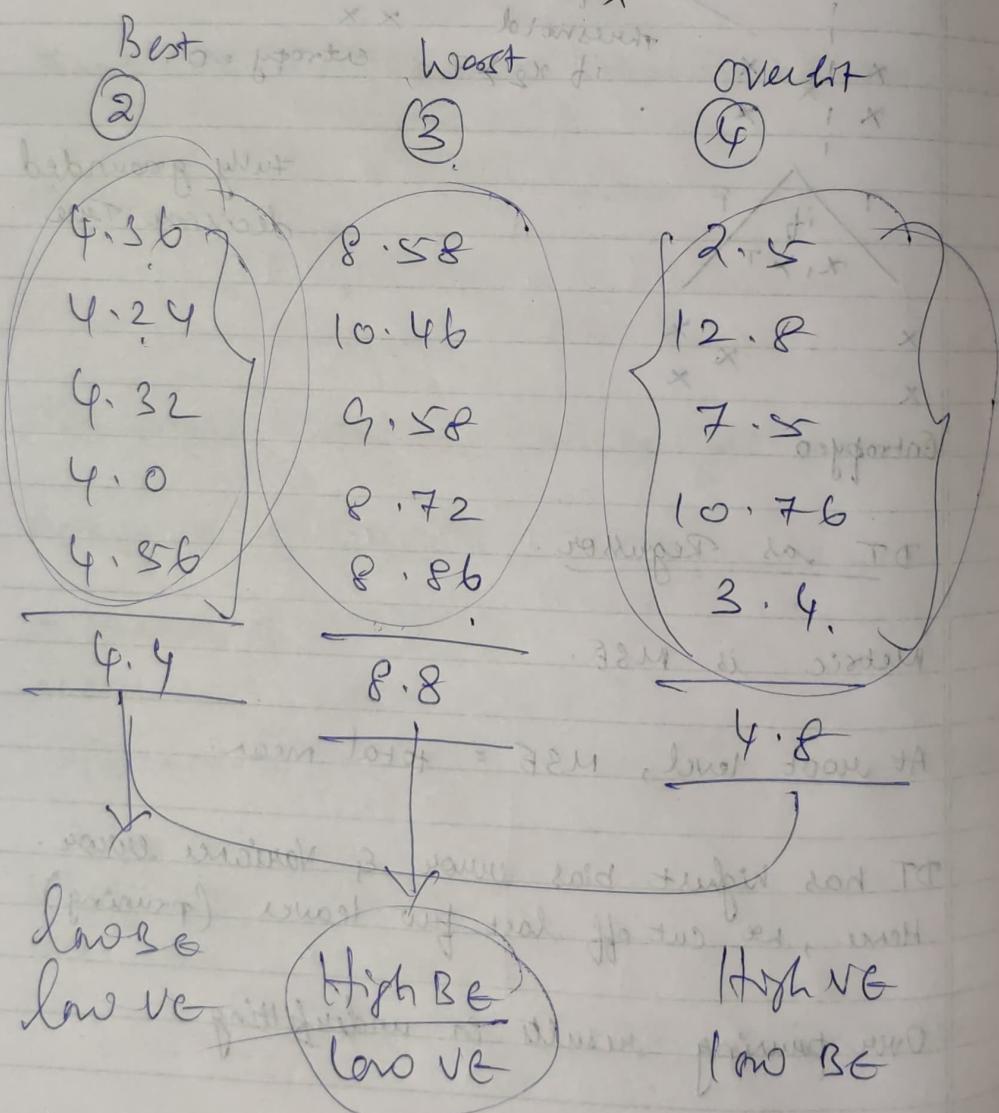
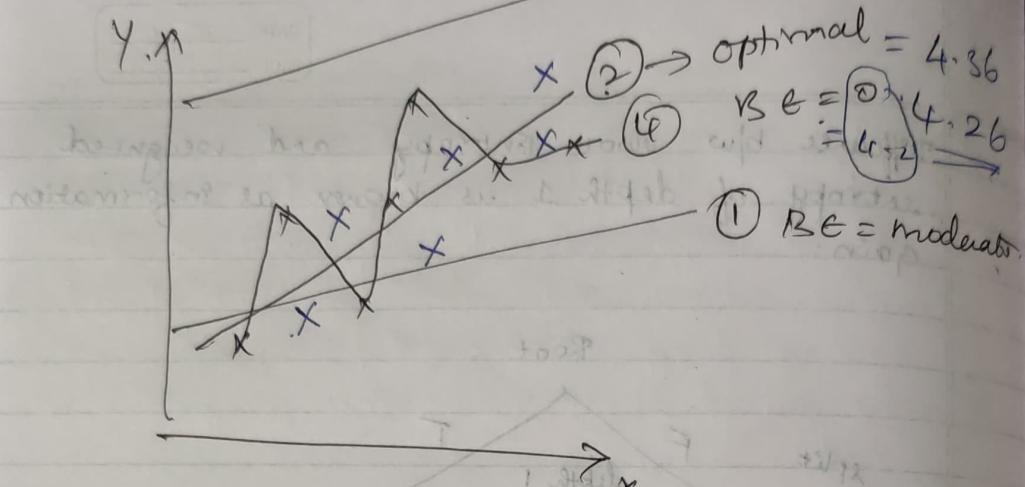
Metric is MSE.

At root level, MSE = total mean.

DT has highest bias error & Variance error.
Hence, we cut off last few leaves (pruning)

Over pruning results in underfitting

$$\textcircled{3} \quad BE = VH$$



90

800
=

71

72

73

74

75

51

90

92

86

85

92

→ ve ↓

52

95

70

58

92

60

→ ve

55

35

34

40

45

42

↓
40

ve ↓. Be A

Be 60%.

ve

\$7600

✓
13400