

Machine Learning

Simple Linear Regression

Supervised ML $\begin{cases} \text{Regression} \\ \text{Classification} \end{cases}$

O/p \rightarrow Continuous

O/p \rightarrow Binary, multiclass categories

Dataset

Weight	Height
74	170
80	180

Train

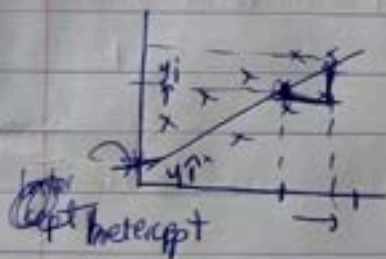
New Weight \rightarrow

model

New height
actual

\rightarrow Predicted

Error $\rightarrow \{y_i - \hat{y}_i\}$



$$\hat{y} = mx + c$$

$$h_0(x) = \theta_0 + \theta_1 x$$

θ_0 = Intercept

θ_1 = Slope (or) Coefficient

Cost function

$$J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_0(x))^2 \Rightarrow \text{Mean squared error}$$

y_i = actual value

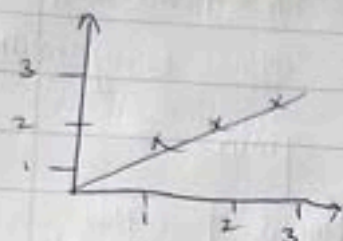
$h_0(x)$ = predictor value

n = no. of datapoint

(To get best fit line)

Final aim \rightarrow Minimize Cost function $J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_0(x))^2$

Optimization: {minimize the cost function}



$$h_0(x) = \theta_0 + \theta_1 x$$

$$\theta_0 = 0$$

$$h_0(x) = \theta_1 x$$

$$\text{Let } \theta_1 = 1$$

$$x=1 \quad h_0(x) = 0 + 1(1)$$

$$h_0(x) = 1$$

$$x=2 \quad h_0(x) = 0 + 1(2)$$

$$h_0(x) = 2$$

$$x \quad y \quad h_0(x)$$

$$1 \quad 1 \quad 1$$

$$2 \quad 2 \quad 2$$

$$3 \quad 3 \quad 3$$

$$\text{Let } \theta_1 = 0.5$$

$$h_0(x) = 0 + 1(0.5)$$

$$= 0.5$$

$$h_0(x) = 0 + 0.5(2)$$

$$= 1$$

$$= 0 + 0.5(3)$$

$$= 1.5$$

Cost function

$$J(\theta_1) = \frac{1}{n} \sum_{i=1}^n \overbrace{(y_i - h_0(x_i))^2}^{\text{error}}$$

$$n=3$$

$$= \frac{1}{3} [(1-1)^2 + (2-2)^2 + (3-3)^2]$$

$$J(\theta_1) = 0$$

$$\theta_1 = 0.5$$

$$= \frac{1}{3} [(1-0.5)^2 + (2-1)^2 + (3-1.5)^2]$$

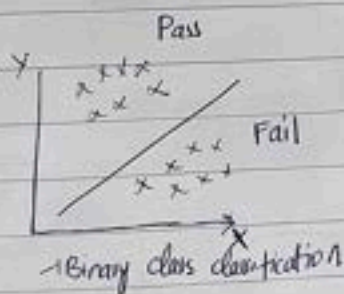
$$= \frac{1}{3} [0.25 + 1 + 2.25]$$

Logistic Regression

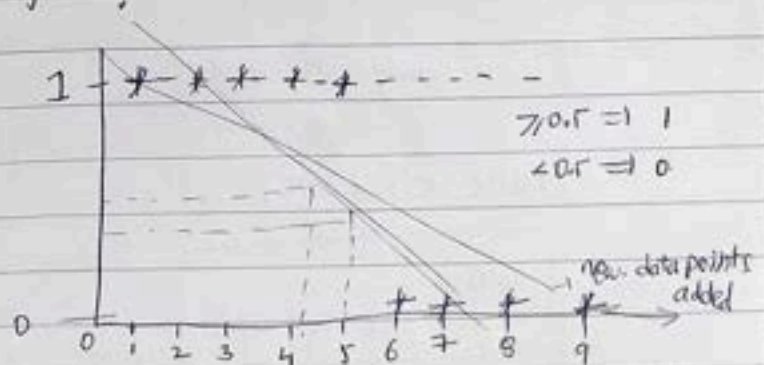
To solve classification

- Binary classification \rightarrow O/P \rightarrow 2 categories
- Multiclass classification \rightarrow O/P \rightarrow > 2 categories

Dataset	Independent feature	O/P (or) dependent feature
No. of play hours	Pass/Fail (y)	
9	0 Fail	<div> <div>new data</div> <div>Model</div> <div>P/Fail</div> </div>
8	0 Fail	
7	0 Fail	
6	0 Fail	
5	1 Pass	Accuracy
4	1 Pass	



Can we solve this classification problem using Regression?



if new data points added (or) having outliers
our predict are going wrong.

Why we cannot use linear Regression for classification.

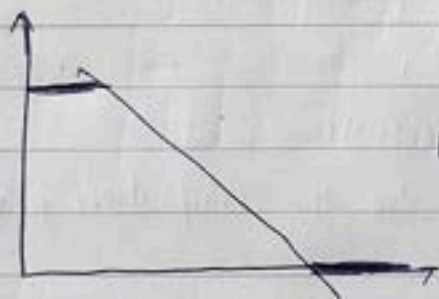
① Best fit line changes because of outliers \rightarrow prediction goes wrong

② The outcomes come > 1 and < 0

\Downarrow
Logistic Regression

\downarrow
0 to 1 \rightarrow Squashing Technique

How Logistic Regression solves classification problem



$h_{\theta}(x) = \theta_0 + \theta_1 x_1$ \rightarrow Best-fit line
Prediction

\Downarrow
Sigmoid Activation function

\downarrow
0 to 1

$$\sigma = \frac{1}{1 + e^{-z}} \Rightarrow 0 \text{ to } 1$$

Squashing the best-fit line

$$h_{\theta}(x) = \sigma(\theta_0 + \theta_1 x_1)$$

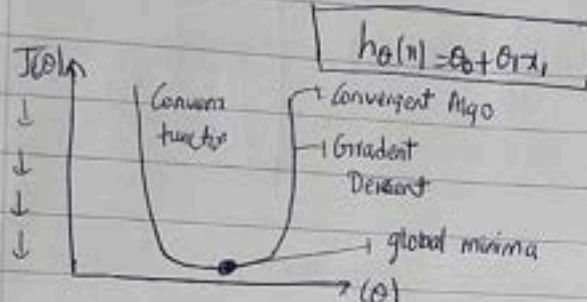
$$h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1)}}$$

Cost function

$\hat{y}_i \rightarrow$ predicted

Linear Regression cost function

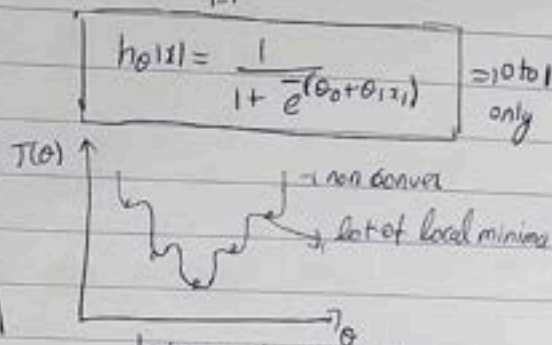
$$J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2$$



$$h_{\theta}(x) = \theta_0 + \theta_1 x_1$$

Logistic Regression cost function

$$J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2$$



we won't use this one of cost function

$$h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1)}} \Rightarrow 0 \text{ to } 1 \text{ only}$$

Log loss function

$h_{\theta}(x)_i =$ predicted value $y_i =$ actual data

$$J(\theta_0, \theta_1) = -y_i \log(h_{\theta}(x)_i) - (1 - y_i) \log(1 - h_{\theta}(x)_i)$$

Cost function

If $y_i = 1$

$$J(\theta_0, \theta_1) = \begin{cases} -\log(h_{\theta}(x)_i) & \text{if } y=1 \\ -\log(1 - h_{\theta}(x)_i) & \text{if } y=0 \end{cases}$$

Final Aim Minimize cost function $J(\theta_0, \theta_1)$ by changing θ_0 & θ_1

Convergence Algorithm

$$\begin{array}{l} \text{Repeat until convergence} \\ \{ \\ \theta_j : \theta_j - \alpha \frac{\partial}{\partial \theta} J(\theta_0, \theta_1) \\ \} \end{array}$$

Logistic Regression with Regularization Parameter

Cost function

$$J(\theta_0, \theta_1) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

$$h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}}$$

$$J(\theta_0, \theta_1) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y=1 \\ -\log(1-h_{\theta}(x)) & \text{if } y=0 \end{cases}$$

$$J(\theta_0, \theta_1) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x)) + L_2 \text{ Regularization}$$

Ridge
Reduce Overfitting

$$J(\theta_0, \theta_1) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x)) + L_1 \text{ Regularization}$$

L_1 feature selection
- lasso

$$J(\theta_0, \theta_1) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x)) + L_2 \text{ reg} + L_1 \text{ reg} \rightarrow \text{elastic}$$

L₂ Regularization → Reduce overfitting

λ = hyperparameter

$$T(\theta_0, \theta_1) = -y \log(h\theta(x)) - (1-y) \log(1-h\theta(x)) + \lambda \sum_{i=1}^n (\text{slope})^2$$

L₁ Regularization

$$T(\theta_0, \theta_1) = -y \log(h\theta(x)) - (1-y) \log(1-h\theta(x)) + \lambda \sum_{i=1}^n |\text{slope}|$$

Elastic Net

$$T(\theta_0, \theta_1) = -y \log(h\theta(x)) - (1-y) \log(1-h\theta(x)) + \lambda \sum_{i=1}^n (\text{slope})^2 + \lambda \sum_{i=1}^n |\text{slope}|$$

$$C \propto \frac{1}{\lambda}$$

$$C = 2.0 \quad \lambda = \frac{1}{2}$$

from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=2)

from sklearn.preprocessing import StandardScaler → z-score mean divide

scaler = StandardScaler()

X_scaled = scaler.fit_transform(X_train)

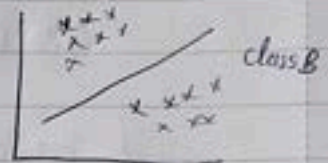
from sklearn.linear_model import LogisticRegression

log_reg = LogisticRegression(random_state=0).fit(X_train_scaled, y_train)

Performance Metrics : Accuracy, Precision, Recall & F-Beta

Topics to covered

- ① Confusion Matrix
- ② Accuracy
- ③ Precision
- ④ Recall
- ⑤ F-Beta Score



DATASET		Actual value		predict
x_1	x_2	y	\hat{y}	
—	—	0	1	→ Wrong Prediction
—	—	1	1	→ correct
—	—	0	0	→ correct
—	—	1	1	→ correct
—	—	1	1	→ correct
—	—	0	1	→ Wrong
—	—	1	0	

① Confusion matrix

		1	0	→ actual values (y)
1		2+1	1+1=2	
0		1	1	
	predicted			

accurate data

		1	0	Actual
1		TP	FP	
0		FN	TN	
	Predicted			

$$\text{Model Acc} = \frac{TP + TN}{TP + FP + FN + TN}$$

$$= \frac{3+1}{3+2+1+1}$$

$$= \frac{4}{7} = 57.1\%$$

True positive

False positive

False negative

True Negative

③ Precision

DATASET \rightarrow Imbalance dataset

1000 datapoints $\left\{ \begin{array}{l} 900 \rightarrow 1 \\ 100 \rightarrow 0 \end{array} \right\}$ Imbalanced dataset

B. Dumb model $\rightarrow 1$

Acc = 90%

getting wrong

	1	0 actual
1 Predict	900	100
0	0	0

So are another 3.

③ Precision = $\frac{TP}{TP+FP}$ } Out of all the actual value how many are correctly predicted

	1	0 Actual
1 Predict	TP	FP
0	FN	TN

FP is important \rightarrow To reduce FP $\downarrow\downarrow$

④ Recall = $\frac{TP}{TP+FN}$ } Out of all the predicted value how many are predicted with actual values

	1	0
1	TP	FP
0	FN	TN

FN is Important \rightarrow To reduce FN $\downarrow\downarrow$

Use case 1 \rightarrow Spam Classification

	1	0
1	TP	FP
0	FN	TN

Text = Model \Rightarrow Spam / Not Spam

TP \leftarrow Mail \rightarrow Spam
Model \rightarrow Spam

TN \leftarrow Mail \rightarrow not a spam
Model \rightarrow Not a spam

Important \leftarrow FP \leftarrow mail \rightarrow Not a spam⁰
 Model \leftarrow spam¹ } Wrong predictions
 { Blunders }

FN \rightarrow mail \rightarrow spam¹
 Model \leftarrow Not a spam⁰ } Wrong predictions

\rightarrow we use precision

\rightarrow Disease \rightarrow FN is Important

To predict whether a person has diabetes (or) not.

	Diabetic ¹	No. Diabetic ⁰
Diabetic ¹	TP	FP
No. Diabetic ⁰	FN	TN

Recall

$$\frac{TP}{TP + FN}$$

TP \leftarrow Actual \rightarrow Diabetes
 Model \rightarrow Diabetes } Correct

TN \leftarrow Actual \rightarrow no Diabetes
 Model \rightarrow no Diabetes } Correct

FP \leftarrow Actual \rightarrow no Diabetes⁰
 Model \rightarrow Diabetes¹ } Wrong Predictions
 but

FN \rightarrow Actual \rightarrow Diabetes¹
 Model \rightarrow No Diabetes⁰ } Blunders

Assignment - Tomorrow the stock market will crash or not

1	0
11	FP 11

Actual + it will crash.
Predict it will not crash.
is big problem 0.

$$FP \downarrow \rightarrow Precision = \frac{TP}{TP+FP}$$

⑤ F-beta score = $(1 + \beta^2) \frac{Precision + Recall}{Precision + Recall}$

① If FP and FN are both are important

$$\beta = 1$$

$$F_1 \text{ score} = (1+1) \frac{P+R}{P+R} \quad \left. \vphantom{\frac{P+R}{P+R}} \right\} \text{ Harmonic mean}$$

② If FP is more important than FN.

$$\beta = 0.5$$

$$F_{0.5} \text{ score} = (1+0.25) \cdot \frac{P+R}{P+R}$$

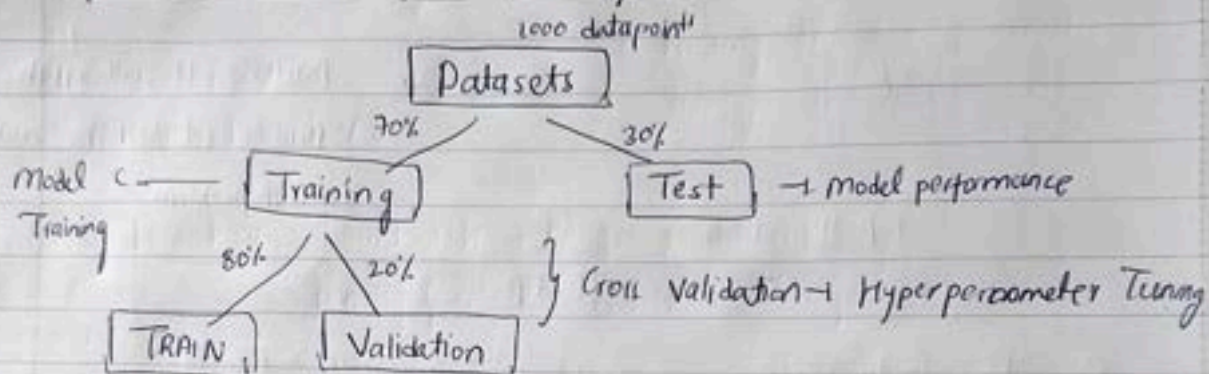
③ If FN is more important than FP

$$\beta = 2$$

$$F_2 \text{ score} = (1+4) \frac{P+R}{P+R}$$

Cross Validation and Its Types

this is used further in hyperparameter tuning



Random state

85%, 75%, 70%, 60% } to find Average accuracy we do cross validation

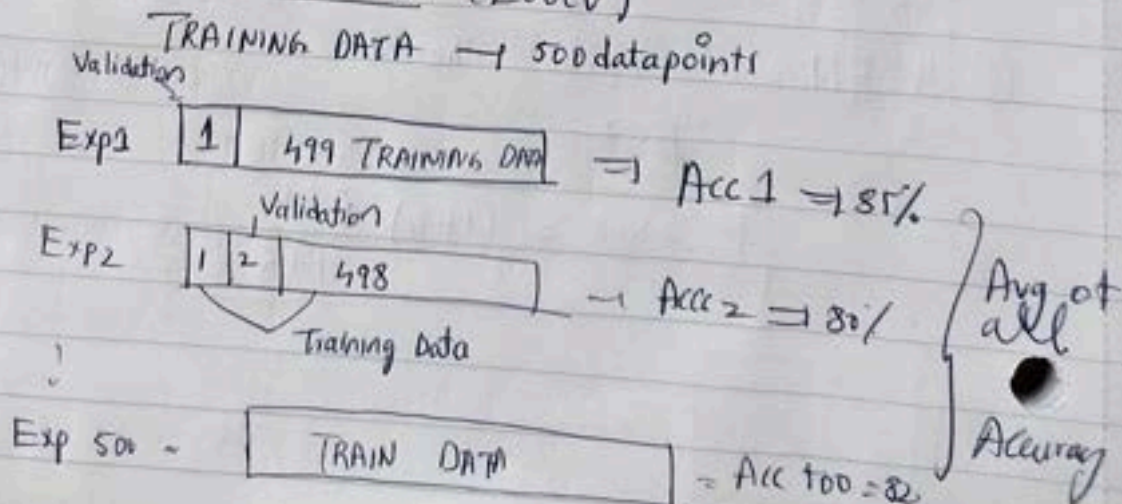
TRAINING \Rightarrow CV = 5

Exp1	TRAIN	VALIDATION	Acc
Exp2	TRAIN	VALIDATION	Acc
Exp3	TRAIN	VALIDATION	Acc
Exp4	TRAIN	VALIDATION	Acc
Exp5	TRAIN	VALIDATION	Acc

Mean

Types of cross validation

① Leave One Out Cross Validation (LOOCV)

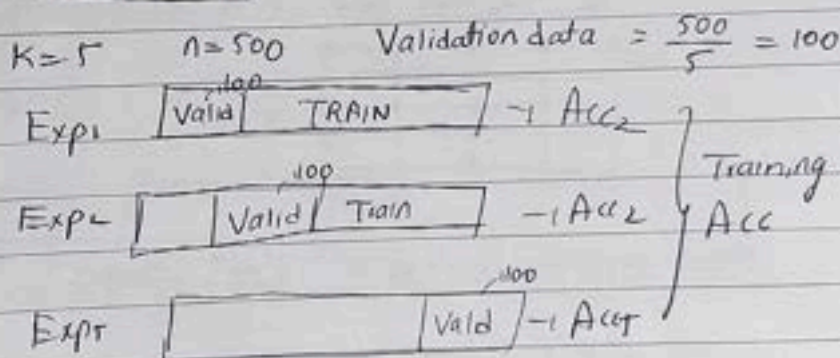


Disadvantage

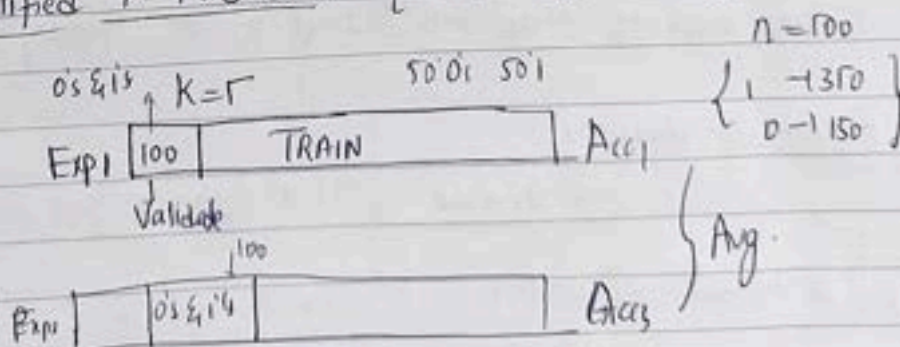
- ① Time complexity is huge for Training dataset
- ② Model overfit \Rightarrow Training Acc $\uparrow\uparrow$
New data \rightarrow Acc $\downarrow\downarrow$

② 1 fold CV $p=10, 20, 100$
 $p = \text{hyperparameter}$

③ K fold cross validation



④ Stratified K Fold C.V { Imbalanced Dataset }



Hyper parameter Tuning with Cross Validation

↳ Finding the best parameter while training the model.

① Grid search cv

② Randomized search cv

① Grid search cv [Grid search + Cross validation] [2-combination]

↓
This says which parameter (combination) is useful.

penalty { 'L1', 'L2', 'elasticnet', 'None' }

Solver { 'lbfgs', 'liblinear', 'newton-cg', 'newton-cholestry', 'sag', 'saga' }

	L1	80% L2	80% EL	None
lbfgs	•	•	•	•
liblinear	•	•	•	•
	•	•	•	•
	•	•	•	•
	•	•	•	•

Logistic Regression (L, Solver)

↓
Cross validation [k fold cv]

[Val | Training] Acc1 } Avg Acc

Disadvantage

i) Time complexity increase with ^{huge} dataset is for Training the model

② Randomized search cv

[n_iter = 100] + [CV = 5]

10 different combination + CV = 5 ⇒ 50

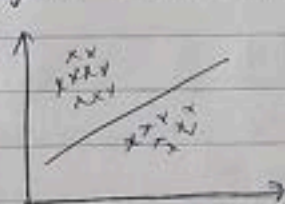
↓

Select the best parameter

Advantage: The time complexity is decreased

Now to remove outlier $Z_{22} = 1.5 + 0.5$

Logistic Regression For Multiclass Classification



each 2 combine each 2

- ① OVR \rightarrow One versus Rest
- ② Multinomial

One versus Rest

t_1	t_2	t_3	O/p
-	-	-	1
-	-	-	0
-	-	-	2
-	-	-	2
-	-	-	0
-	-	-	1

one hot encoding

m_1	m_2	m_3
O_1	O_2	O_3
0	1	0
1	0	0
0	0	1
0	0	1
0	0	0
0	1	0

t_1	t_2	t_3	O/p
-	-	-	1
-	-	-	0
-	-	-	2
-	-	-	2
-	-	-	0
-	-	-	1

m_1 model \rightarrow $I/p = t_1, t_2, t_3$ $O/p = 0_1$
 \parallel
 Binary classification

m_2 model \rightarrow $I/p = t_1, t_2, t_3$ $O/p = 0_2$

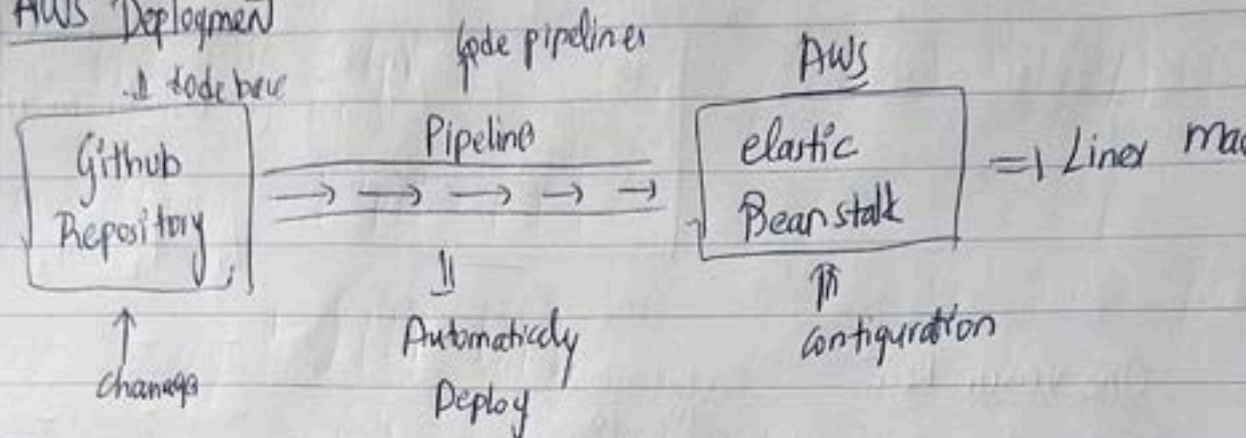
m_3 model \rightarrow $I/p = t_1, t_2, t_3$ $O/p = 0_3$

New test data $\rightarrow [m_1, m_2, m_3]$

$[0.25, 0.25, 0.5]$

New DATA $\rightarrow [0_2]$

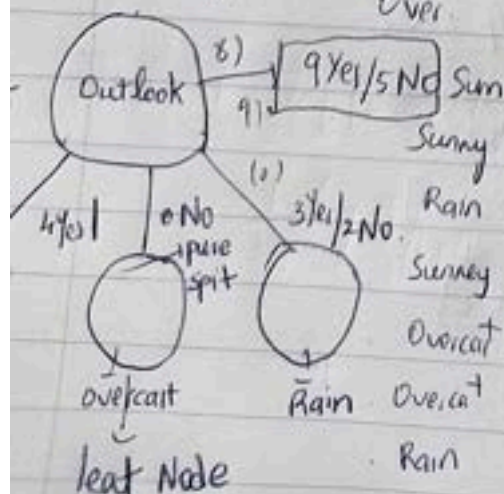
AWS Deployment



Decision Tree

Dataset

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
1)	Sunny	Hot	High	Weak	No
2)	Sunny	Hot	High	Strong	No
3)	Overcast	hot	High	Weak	Yes
4)	Rain	mild	High	Weak	Yes
5)	Rain				Yes
6)	Rain				No
7)	Over				Yes
8)					No



① Purity split check = Pure split (or) Impure split

$\left. \begin{array}{l} \text{Entropy} \\ \text{Gini Impurity} \end{array} \right\} \text{measure of purity}$

② What feature you need to select to start the split -> Information gain

① Purity Check

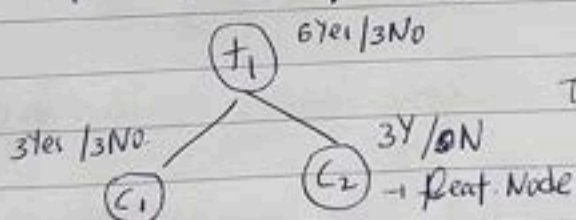
Binary Classification

① Entropy

$$H(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

p_+ = probability of positive category

p_- = Probability of negative category



$$H(C_1) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

$$= -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$

$$\boxed{H(C_1) = 1} \quad \text{This is impure}$$

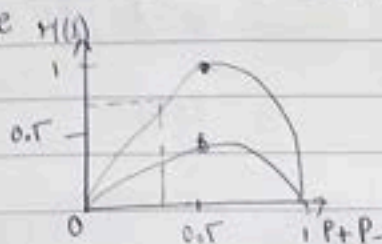
② Gini Impurity

$$G.I = 1 - \sum_{i=1}^h (p_i)^2$$

$$\begin{aligned}
 p_+ &= \frac{\text{Total positive count}}{\text{Total count}} = \frac{3}{6} = \frac{1}{2} \\
 p_- &= \frac{\text{Total negative count}}{\text{Total count}} = \frac{3}{6} = \frac{1}{2}
 \end{aligned}$$

$$H(C_2) = -\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3}$$

$$\boxed{H(C_2) = 0} \quad \text{pure split}$$



② Gini Impurity - max = 0.5

3/4 1/4

$$G.I = 1 - \sum_{i=1}^n (p_i)^2$$

$$G.I(C_2) = 1 - \left[\left(\frac{3}{3} \right)^2 + \left(\frac{0}{3} \right)^2 \right]$$

$$Gini(C_1) = 1 - \left[(p_+)^2 + (p_-)^2 \right]$$

$$= 1 - \left[\left(\frac{1}{2} \right)^2 + \left(\frac{1}{2} \right)^2 \right]$$

$$= 1 - \frac{1}{2} = 0.5 \Rightarrow \text{Impure split}$$

$$= 1 - 1$$

$C_2 = 0$ - Purity split

Multiclass classification Problem : 3 categories in o/p

$$H(S) = -p_{C_1} \log_2 p_{C_1} - p_{C_2} \log_2 p_{C_2} - p_{C_3} \log_2 p_{C_3}$$

$$G.I = 1 - \left[(p_{C_1})^2 + (p_{C_2})^2 + (p_{C_3})^2 \right]$$

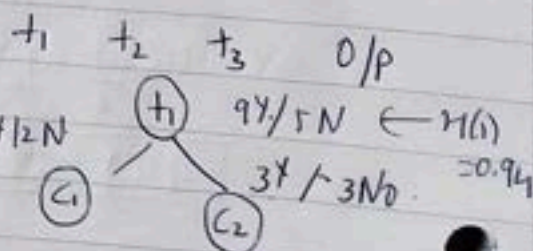
② Information Gain : -> Which feature to select to start the split?

$$\text{Gain}(S, t_1) = H(S) - \sum_{v \in \text{val}} \frac{|S_v|}{|S|} H(S_v) \rightarrow \text{Entropy of child node}$$

$$H(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

$$= -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \left(\frac{5}{14} \right)$$

$$\approx 0.94$$



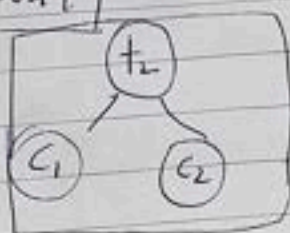
$$H(C_1) = -\frac{6}{8} \log_2 \left(\frac{6}{8}\right) - \frac{2}{8} \log_2 \left(\frac{2}{8}\right) \approx 0.81$$

$$H(C_2) = -\frac{3}{6} \log_2 \left(\frac{3}{6}\right) - \frac{3}{6} \log_2 \left(\frac{3}{6}\right) = 1$$

$$\text{Gain}(S, f_1) = H(S) - \sum_{\text{value}} \frac{|S_v|}{|S|} H(S_v) \quad \text{Entropy of categories}$$

$$= 0.94 - \left[\underbrace{\frac{8}{14}}_{\text{Total goat}} \times \underbrace{0.81}_{H(C_1)} + \underbrace{\frac{6+1}{14}}_{\text{Total cow}} \right]$$

$$\boxed{\text{Gain}(S, f_1) = 0.049}$$



$\Rightarrow \text{Information loss} = 0.051$

$$\boxed{\text{Gain}(S, t_2) = 0.051} \Rightarrow \boxed{\text{Gain}(S, t_1) = 0.049}$$

We need to start splitting using f_2 feature

Entropy

Entropy vs Gini Impurity

When dataset is small \rightarrow Entropy [log formula]

When dataset is huge \rightarrow Gini Impurity [simple math]

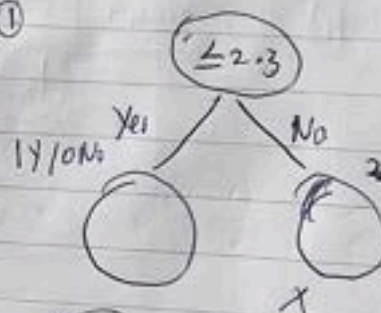
What if my feature is continuous

t_1	O/p
2.3	Yes
3.6	Yes
4	No
5.2	No
6.7	Yes
7.8	No

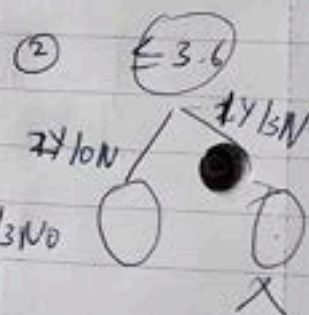
① Sort the feature t_1

① Threshold = 2.3

①



②



③



Time complexity is

huge when dataset

is huge from sklearn.tree import DecisionTreeClassifier.

classifier = DecisionTreeClassifier()

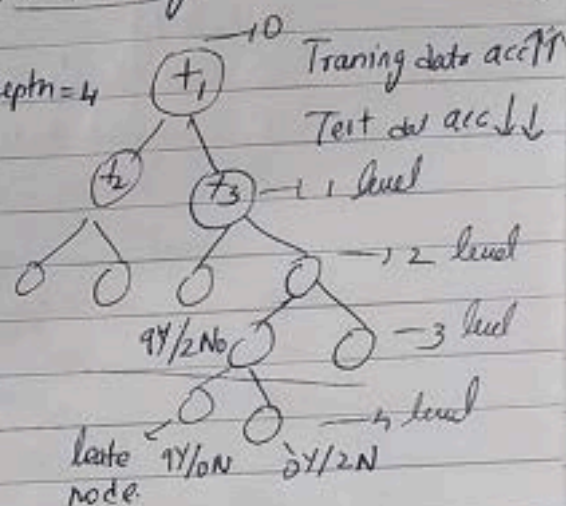
Decision Tree post pruning and pre pruning [Reduce overfitting]

Training data

$t_1 \ t_2 \ t_3 \ O/P$

Generalized model
Reduce overfitting

Max-depth=4



① Post Pruning

- ① Construct the entire decision Tree to complete leaf node
- ② Pruning the decision
- ③ ~~for~~ suitable for smaller Dataset

② Pre pruning

- ① No Hyperparameter Tuning to select Best parameters
(grid cv, Random CV)

Decision Tree Classifier

- ① Entropy
- ② Gini (Gini)
- ③ Information Gain

Decision Tree Regressor

- ① Variance Reduction
- ② Variance

Dataset

exp	career gap	Salary
2	Yes	40k
2.5	Yes	42k
3	No	52k
4	No	60k
4.5	Yes	56k

O/p → continuous.

Salary

40k

42k

52k

60k

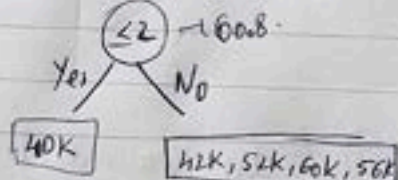
56k

$\bar{y} = 50k$

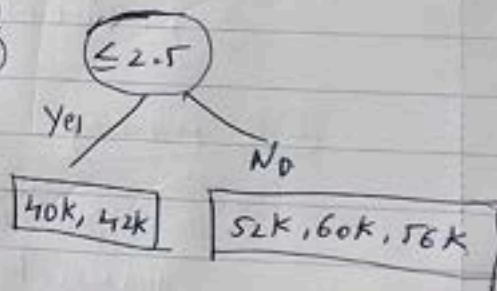
← avg. salary (mean)

Variance Reduction

① [40k, 42k, 52k, 60k, 56k]



②



Final aim → Variance reduction

$$\text{Variance error} = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \quad [\text{mean squared error}]$$

$$\begin{aligned} \text{Variance of Root} &= \frac{1}{5} [(40-50)^2 + (42-50)^2 + (52-50)^2 + (60-50)^2 + (56-50)^2] \\ &= \frac{1}{5} [100 + 64 + 4 + 100 + 36] \\ &= 60.8 \end{aligned}$$

$$\text{Variance of left child} = \frac{1}{1} [(40-10)^2]$$

$$= 100 //$$

$$\text{Variance of right} = \frac{1}{4} [(42-50)^2 + (52-50)^2 + (60-50)^2 + (50-50)^2]$$

$$= 51 //$$

$$\text{Variance of Reduction} = \text{Var}(\text{Root}) - \sum w_i \text{Var}(\text{child})$$

$$= 60.8 - \left[\frac{1}{5} * 100 + \frac{4}{5} * 51 \right]$$

Let side having only one element by total elements

$$= 0$$

$$\text{Variance Reduction} = 0$$

Select the greatest Variance Reduction