

Supervised Learning - Regression (SLR)

$$y = mx + c \rightarrow c + mx$$

$$y = b_0 + b_1 x$$

$b_0 \rightarrow$ Constant (Intercept of best fit. line with y-axis)

ML - Machine identifying pattern without explicitly coding.

Supervised
Learning

Unsupervised
Learning

There is a Target variable
To be predicted

No target variable
To be predicted

Regression

Target variable is
(continuous) numeric

Target variable
is categorical

Understand if there
is a group in data

Classification

Clustering

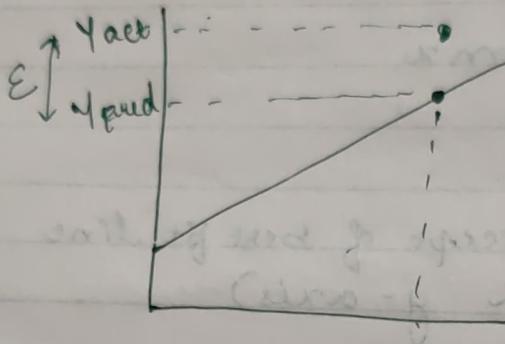
Target variable - dependent variable, - response variable
or y - all are same

Independent - Explanatory variables - x_i - all are
Predictor variable same

$$\hat{y} = b_0 + b_1 x_1$$

$\hat{y} \rightarrow y_{\text{pred}}$ → predicted $y \rightarrow y_{\text{pred}}$

y_{actual} b_0 & b_1 are parameters of Regress



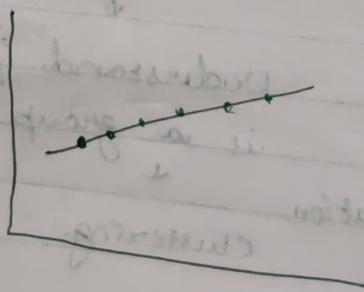
$$\text{Error: } Y_{\text{actual}} - \hat{Y}_{\text{(pred)}}$$

$$Y_{\text{actual}} = b_0 + b_1 x_1 + \varepsilon$$

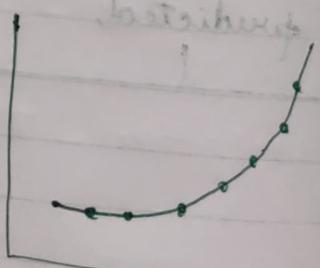
$$Y_{\text{actual}} = Y_{\text{predicted}} + \varepsilon$$

Types of Association

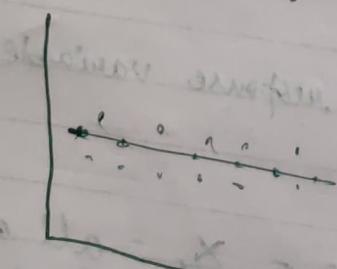
1) Linear



2) Non-Linear



3) No relationship



How to draw best fit line:-

1) Ordinary Least Squares Method.

Best fit line - line which gives us lowest sum of squares of Error.

$$\text{Min of } \sum \text{error}^2$$

→ Cost function or Loss function.

Solving partial derivatives

$$\frac{\partial \text{Cost func}}{\partial b_0} = 0 \quad \dots \quad (1)$$

$$\frac{\partial \text{Cost func}}{\partial b_1} = 0 \quad \dots \quad (2)$$

gives

$b_1 = \frac{\text{Cov}(x, y)}{\text{var}(x)}$
$b_0 = \bar{y} - b_1 \bar{x}$

Business interpretation of b_1 ,

For every unit increase in x , \hat{y} will increase by b_1 .

b_0 is not always interpreted. Depends on domain knowledge.

We use partial derivatives, as derivatives don't work on more than 2 unknowns!

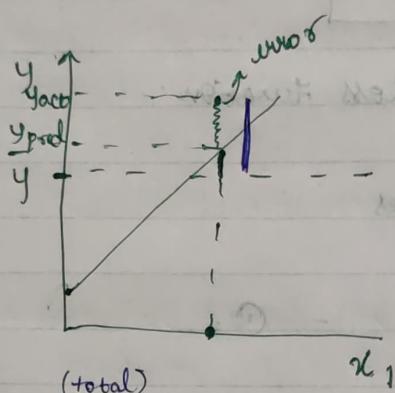
Performance Measure

How close the predicted value is compared with Actual.

- Higher the value of R^2 , better the prediction.

It's always b/w 0 & 1

- proportion of variation of y explained by the model is coefficient of Determination - R^2



$$SST = \sum (y_{act} - \bar{y})^2$$

R - Coeff of correlation

(regression)

$$SSR = \sum (y_{pred} - \bar{y})^2$$

(error)

$$SSE = \sum (y_{act} - y_{pred})^2$$

$$R^2 = \frac{SSR}{SST}$$

$$SST = SSR + SSE$$

$$SSR = SST - SSE$$

$$\frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

$$R^2 = 1 - \frac{SSE}{SST}$$

Errors is also known as Residual
 sometimes - Regression (green) is known as explained
 & variation explained by the model

(SSR)

(SSE)

DATE: _____
 PAGE: _____

** Problem with R^2 - R^2 goes up even if there is a useless column

In that case, we use adjusted R^2 .

$$\text{Adj } R^2 = 1 - \frac{(1 - R^2)(n-1)}{(n - k - 1)}$$

$n = \text{no of obs}$
 $k = \text{no of } X_s$
 (columns)

- Higher the better.

- If k goes up, Adj R^2 comes down

$$\text{Adj } R^2 = 1 - \frac{\text{SSE}/(n-k-1)}{\text{SST}/(n-1)}$$

is derived from

→ Used as Adj R^2 - when you are supposed to add/remove columns. (or decides if its useful or not).

new model
or

df of SSR = k - fitting of fit() - Trains the algorithm.

df of SST = $n - 1$

df of SSE = $n - k - 1$

- Reject zone for f-statistic is always on right side.
- Larger the f-statistics, we may end up in rejecting Null hypothesis.

$$f\text{-stat} = \frac{\text{SSR}/k}{\text{SSE}/(n-k-1)} \rightarrow \text{df of SSR}$$

To find p-value from f-stats
 scipy.stats → stats

stats.f. sf(f-stat, df of num=k, df of den=n-k)

- Testing if regression is valid/not - F test
- H₀: All coefficients = 0 → implies regression is not valid
- H_a: All coefficients $\neq 0$ → regression is valid
At least 1

To check if a column is useful/not :- T test.

- H₀: coefficient of that $X_8 = 0$ - implies that column is not significant.
- H_a: coefficient of that $X_8 \neq 0$ - column is significant

$$t_{\text{stat}} = \frac{\text{coeff}}{\text{std error}}$$

- look at pval for columns (not for const) and determine if H₀ is rejected/Accepted.
- if only 1 x, its simple linear Regression.
→ if more x, its multi linear regression

Modelname • resid → gives errors

Day 2

Assumptions of linear regression :-

- Tests before model building
- 1) Dependent variable must be numeric
- 2) Predictors must not show multicollinearity

Multicollinearity :- Correlation b/w X_8

A strong correlation b/w X_8 is not good for prediction

Multicollinearity detection.

- 1) Determinant of the correlation Matrix checks presence of multicollinearity
- 2) Condition Number
- 3) Let D be determinant of corr matrix

$D < 0$	high multicollinearity
$D \geq 1$	No "
- 4) $CN > 1000$ severe multicollinearity
 $100 \leq CN \leq 1000$ Moderate
 $CN < 100$ No "}
- 5) Correlation matrix which variables are involved
- 6) Variance Inflation Number (VIF) Factor in multicollinearity

VIF

Runs regression internally to check which variables have multicollinearity.

Target	Pred	R^2	$VIF = 1/(1-R^2)$
x_1	$x_2 x_3 x_4$		
x_2	$x_1 x_3 x_4$		
x_3	$x_1 x_2 x_4$		
x_4	$x_1 x_2 x_3$		

$$VIF = \frac{1}{1 - R^2}$$

If $VIF >= 5$, there is strong multicollinearity that affects the model

Auto correlation :-

Correlation with itself.

Is my record dependent on previous record

Durbin - Watson Test : (To check if its autocorrelated / not)

H_0 : The error terms are not autocorrelated.

H_a : The error terms are autocorrelated.

If test statistics = 2 - No autocorrelation
Practically exact & doesn't come. So take a range of 1.5 - 2.5

$0 < \text{statistics} < 2$ Positive autocorrelation

Statistics 0.0, 1.1, 2.2, 3.3

$2 < \text{statistics} < 4$ Negative autocorrelation

Homoscedasticity Assumption (errors should be homo)

Homoscedasticity

Same variance

of error term and outcome

Heteroscedasticity - Funnel type shape when plotted

We use Breusch Pagan Test

H_0 : The error terms are homoscedastic.

H_a : The error terms are heteroscedastic.

Tarque - Bera Test for Normality

H_0 : Data is normal

H_a : Data is not normal.

Based on Jarque-Bera statistic

Test performance is cpt.

DATE: / /
PAGE: _____

* Even if VIF $\gamma = 5$, do not drop categorical variables

Interaction Effect:-

→ particular variable's effect on combination of variables.

Eg:- Effect of experience alone on salary
" " Education "

Effect of both Experience & Education on Salary is even more high.

→ If intercept is needed, then only use sm. add-constant.

Else build the model without that.

$$\text{Mean abs percent error} = \frac{1}{n} \sum \left| \frac{\text{Act} - \text{Pred}}{\text{Act}} \right|$$

Overfitted:

• Significance difference b/w RMS of Test & Train implies overfitted. \rightarrow [Train is high Test is low]

Note:- Scaling does not affect/impact Regression

Feature Selection:

1) Forward selection

- Start with Null model and then go on adding columns one by one.

2) Backward selection

- Start with full model & then go on dropping columns

Stepwise Regression:

→ Combination of both fwd selection & backward step

Automated in SPSS

SFS - sequential feature selector (wrapper method)

RFE - Recursive feature eliminator

Validation

K fold cross validation

Total: 1000 rows

x_1, \dots, x_n	4 samples used for test
200	A was not in portfolio
200	B
200	C
200	D
200	E

Test function library with blind sets

Train: R^2

ABCD

ABC E

ABDE

ACDE

BCDE

(val is +ve)

Validate (test): R^2

E

D

C

B

A

Avg =

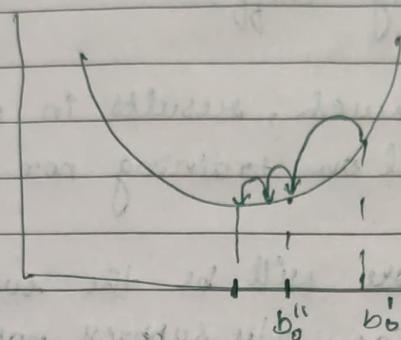
$y = \text{rate}$

variable - business

as per user has been used to predict -
as per user knows nothing

Leave one out cross Validation (LOOCV)

Gradient descent



Regularization:

(Eg. of teaching a child about dog (breeds)).
we know that,

LR :- Minimize $\sum (\text{Y}_{\text{act}} - b_0 - b_1 x_1)^2$,
it results/possibility of overfitting.

1) So,

$$\text{Minimize } \sum (\text{Y}_{\text{act}} - b_0 - b_1 x_1)^2 + \lambda \sum_{i=1}^n \beta_i^2$$

In python, $\lambda = \alpha$, β = coefficients = b_0, b_1

Algorithm → Ridge Regression (L2 regularization)
↳ Used when there is non-linear relationship
between the X_i .

2) Lasso Regression (L1 regularization)

$$\text{Minimize } \sum (\text{Y}_{\text{act}} - b_0 - b_1 x_1)^2 + \lambda \sum_{i=1}^n |\beta_i|$$

$\beta \rightarrow \text{coeff} = b_0, b_1, b_2, \dots$

$\lambda \rightarrow$ Hyperparameter

↳ we decide this b/c it is the parameter we pass through Algorithm.

- As λ increases, penalty increases.
effect - lines with small coeff gets selected (suppressing coeff).
- Increasing λ too much, results in underfitting. Neither works well on training nor testing.
- Larger the α/λ , more will be the suppression in Ridge. (Ridge can only suppress not drop)
- Even if there is multicollinearity in columns, Ridge and Lasso can be applied.
- If there is insignificant column, Lasso drops it automatically.

Elastic Net regression

$$\text{minimize } \sum (y_{\text{act}} - b_0 - b_1 x_1 - \dots) + \lambda_{\text{ridge}} \sum b_i^2 +$$

theoretically \rightarrow

$$\text{Practically } \lambda_b = \lambda_{\text{Ridge}} + \lambda_{\text{Lasso}}$$

$$\text{H ratio} = \frac{\lambda_{\text{Lasso}}}{\lambda_{\text{Ridge}} + \lambda_{\text{Lasso}}}$$

An algorithm used inside sklearn's linear regression is Gradient Descent.

Grid Search (Used to find tune Hyperparameter)

- If there are large no. of columns and you feel many are useless & ~~inter~~, go for Lasso and also interrelations.

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)

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

or $\Delta \beta \propto \frac{1}{\Delta \text{Error}}$

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

η = learning Rate

$$E = \frac{1}{2} (y - \hat{y})^2$$

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

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

diff w.r.t β_0 $- (y - \hat{y})$

$$\Delta \beta_0 = \frac{\partial E}{\partial \beta_0} = -(y - \hat{y})$$

$$\Delta \hat{y}_i = \frac{\partial E}{\partial \beta_i}$$

$$= \frac{1}{2} \cdot 2 \left(y - (\beta_0 + \beta_1 x) \right)$$

diff w.r.t β_i

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

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 (values) 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

$$\hat{y}_{\text{prob}} = \frac{(x_0 + 1) - 1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

$$(\hat{p} - p) = \frac{2B - 9}{95}$$

Also called as Soft Limit function

DATE: / /
PAGE:

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

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

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

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

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

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

Concept of Odd's

Mutually Exclusive \rightarrow Union (Addition)

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

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 effect

Coraxin

$$P_{SE} = 0.2$$

Covishield

$$P_{SE} = 0.1$$

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

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

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

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

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

$$O_{CS}(P_{SE})$$

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

$$y_{prob} = \frac{e^{-z}}{1 + e^{-z}} \rightarrow \text{let this be } x.$$

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

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

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

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

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

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

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

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

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

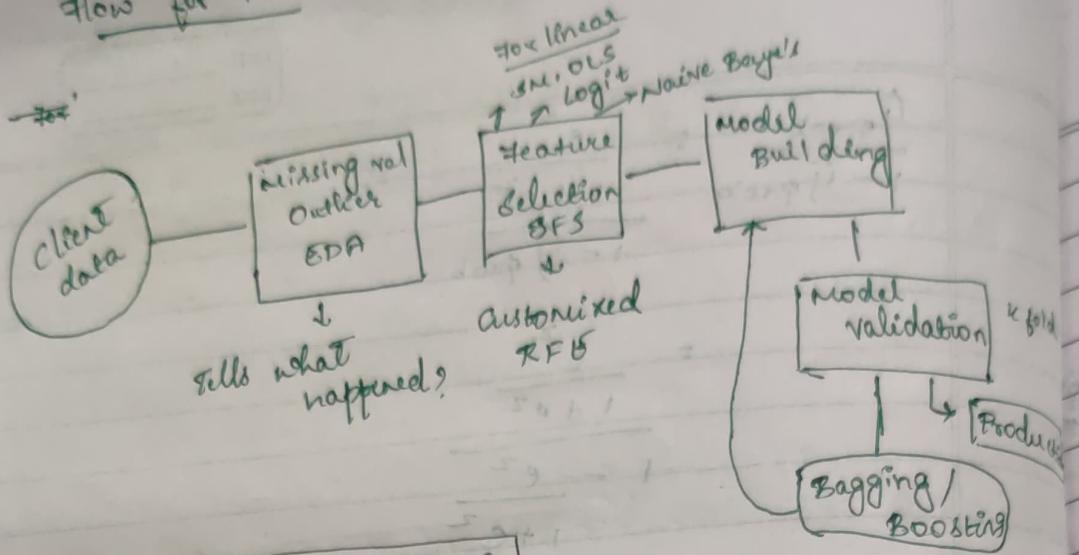
$$1 - \hat{y}_{\text{prob}} = (B_0 + Bx)$$

$e^{B_0 + Bx} \Rightarrow \text{odds}$

Raising coeff to $e^{(B_0 + Bx)}$ is called odds ratio.
 In logistic reg, we can't interpret slopes as yrs times x as we do in linear regression.

If odds > 1 , then odds are significant

Flow for ML



Bagging - \downarrow variance error

Boosting - \downarrow 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 RFE
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				

↓
Test

Shapes:

X_{Train}

4×3

X_{Test}

2×3

Y_{Train}

4×1

Y_{Test}

2×1

Scaling: Compulsory Models = KNN & Naive Baye's

Ex. add constant (λ)

Constant

1.00

1.00

1.00

This constant is the multiplication factor of β_0 (coeff of β_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 \ominus 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

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

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

$y_{\text{act}} \quad \hat{y}_{\text{prob}}$

0	0.2	$\hat{y}_{\text{prob}} = 1 - 0.2 = 0.8$
0	0.3	$\hat{y}_{\text{prob}} = 1 - 0.3 = 0.7$
1	0.8	$\hat{y}_{\text{prob}} = 1 - 0.8 = 0.2$
1	0.9	$\hat{y}_{\text{prob}} = 1 - 0.9 = 0.1$

What evaluation metric we choose in validation method matters.

Validation Method

Train/Test split - NO preferred as it is like (X, y) , testing model with only 1 test

✓ Cross-Val score

K-LOOCV - Not preferred as all rec are same

& only 1 rec is changed

(Test data is less & train data is repetitive)

for Linear Reg:-

```

kf = KFold(n_splits=5, shuffle=True, random_state=None)
rule = 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 T_{NVE} gives as better

But, in reality, lowest χ^2 is better
↳ Hence \neg is given. in scoring.

np. mean (np.abs(cause))

NP: std (or nse, df = 1)

for Logistic Reg:-

		Predicted		Yact	
		H	CVD	H	H
Yact	H	8	0	H	H
	(8)			H	H
CVD	H	2	0	H	H
	(2)			H	H
Predicted		CVD	H	CVD	H

N	True Negative specificity of model	False Positive - Type I Error α
P	False Negative Type 2 Error β	True Positive sensitivity of the model

★ Type I error is dangerous than Type II error

Confusion Matrix is performance metric similar to RMSE in Linear Reg.

Confusion Matrix - for all classification

methods to set threshold:

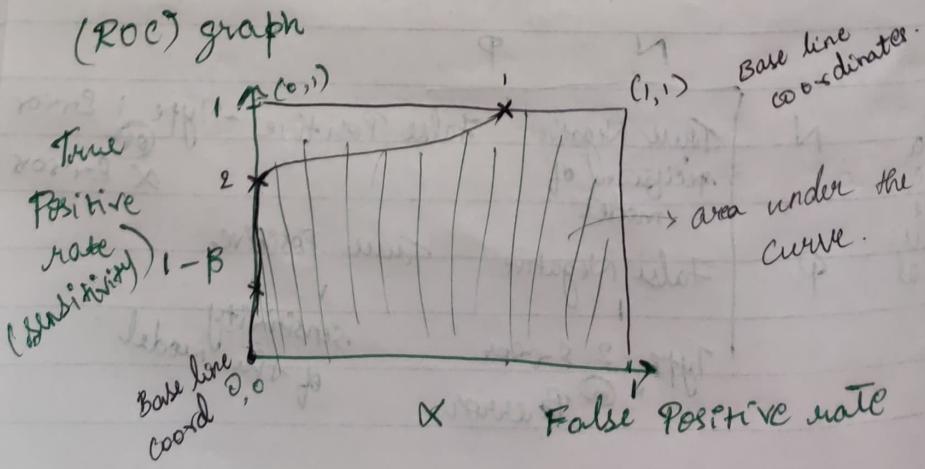
\hat{y}_{act}	\hat{y}_{prob}	\hat{y}_{class}	\hat{y}_{class}	\hat{y}_{class}
0 - H	0.35	1	0	0
H	0.2	0	0	0
H	0.4	(0.1, 0.9)	0	0
- CVD	0.68	(0.3, 0.7)	0	0
CVD	0.75	1	1	0
CVD	0.85	1	1	0

H		CVD		H		CVD		H		CVD	
H	2	1	2	H	3	0	H	3	0	H	3
3				3			3			3	
CVD	0	1	3	CVD	0	1	CVD	0	1	CVD	0
3				3			3			3	
overall Accuracy	$\frac{4}{6}$		$\frac{5}{6}$		$\frac{4}{6}$		$\frac{5}{6}$		$\frac{4}{6}$		

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 B confusion Matrix on ROC graph.

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

2nd $\alpha = 0\%$

sensitivity = $2/3 = 66.6\%$

3rd $\alpha = 0\%$

sensitivity = ~~2/3~~ / 3 = 0.33

Sklearn

model.predict(X)
 ↳ o/p class label

model.predict_proba(X)
 ↳ o/p prob of score

sm. logit

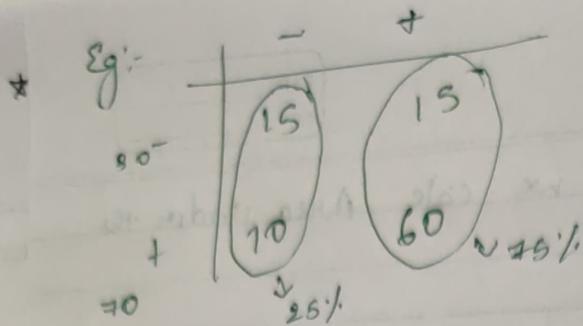
model.predict()

o/p prob score
 ↳ get class label ~~manually~~

Kappa score or Reliability score.

↳ Overall accuracy adjusted to Baseline Reference.

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



$$DA = 0.75$$

Expected by chance

$$75\% \text{ of } 70 + 25\% \text{ of } 30$$

$$52.5\% + 7.5\% \Rightarrow 60\%$$

$$\text{Kappa score} = \frac{0.75 - 0.6}{0.15} = 37.5\%$$

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

		G	S	Recall	Precision
		5	0	100% / 5/5	5/6 83%
		5	4	80% / 4/5	4/4 100%
5	G				
5	S				

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

(19) Which threshold makes TPR the highest & FPR the lowest we choose that.
Or Youden's Index

(20) Sum of α error & β error should be min
 ↗ Another way to choose threshold.
 ↗ Choose threshold which gives least total cost

(3) Simply choose Mean.

$$\min_{\alpha} (w_1 * FN + w_2 * FP)$$

3 ways ↑ to find threshold