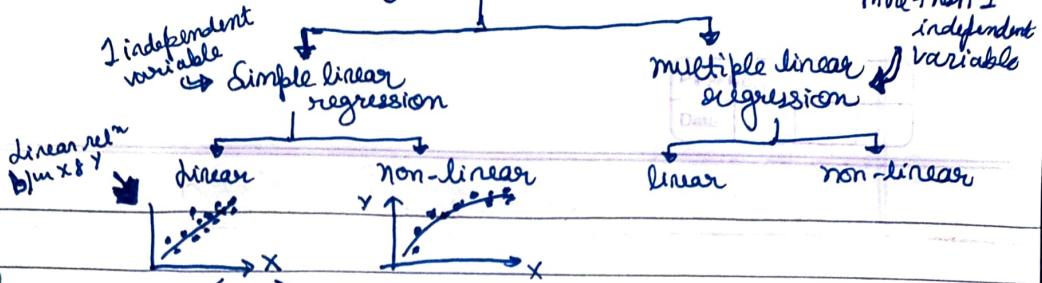


Regression Models



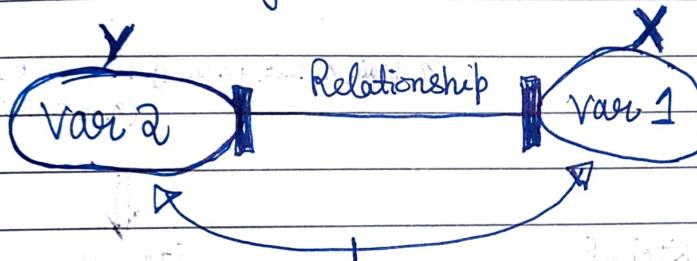
Unit - 2

Simple Linear Regression (SLR)

→ It is a statistical technique to find existence of an association relationship b/w a → 'dependent variable' and an 'independent variable'.

→ Here, there is only one independent variable in the model.

* Regression can only tell



does not mean that value of independent variable (Var 2 or Y) depends on value of Independent variable (Var 1 or X).

- Defⁿ ⇒ It is a statistical model in which there is only one independent variable and the ^{funcⁿ} reln b/w the dependent variable and regression coefficient is linear.

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

here,

Y_i = dependent variable in sample / response variable / outcome variable

X_i = independent variable in sample / predictor variable / input or explanatory variable

β_0 and β_1 = Regression parameter or
Regression coefficient

ϵ_i = random error (or residuals)

β_0 = intercept of regression line
(Y-intercept)

β_1 = slope of the regression line

Simple Linear Regression Model Building

Step 1

Collect / Extract

data

(on the dependent and independent variable. Time & resource consuming & expensive even with organisation having well-designed resource planning system (ERP))

Eg. adult = [.....]
Sale = [.....]
OR
df = pd.read_excel('---')
df.head()

Step 2

Pre-process Data

Ensure Quality of data

- Check for issues such as reliability, completeness, usefulness, accuracy etc.
- All the steps required to give clean data

Step 3

Dividing data into training and Validation datasets

→ Data divided into two sets - training data & Test (Validation data)

→ % of Training data is usually 70%, rest is testing data.

→ Training data used for developing model
Testing data to check & select the model

df_train, df_test, y_train, y_test = sk.train_test_split(df, y, test_size=0.2, random_state=42)

Step-4

Perform descriptive analysis of the data like boxplot, scatter plot to identify any outliers

→ to reveal col^n b/w the two variables
→ useful to determine functional reln b/w dependent and independent variable

Eg. df.boxplot(['column name'])
df.hist([])(df[''])
df[['column name']].median(),
mean
(R)
mean_value = np.mean(adult)
plt.boxplot(adult)

Step-5

Build the Model

→ Model is build using training data.

→ OLS (ordinary least square) method

used to estimate parameters in linear regression model.

→ The goal of OLS is to find the best fit line through a set of data points by minimizing the sum of squares of vertical distances (residual) b/w the observed and value predicted by the model.

[$\text{Residual}^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$ where, \hat{y}_i are predicted values.]

→ Using this we find values of P_0 and P_1 , and once these regression parameters are set \Rightarrow from a given x predicted value of y will be $\hat{y} = P_0 + P_1 \times x + E$. I.e. our model is made.

Eg. model = LinearRegression()
model.fit(x, y)
model.intercept_
model.coef_[0]

Step-6 → Perform model Diagnostic

- Important to run diagnostic test to see if the model runs correctly
- Important to ensure model runs consistently and somewhat accurately on the testing data as it did on training data.

Eg. y_predicted = model.predict(x)

Step - 7 : Validate the model and Measure model accuracy / Plot the model

⇒ Ensure model runs accurately on the testing data and perform various test to see if the models fit perfectly.

$$R^2 \text{-square} = \text{score}(y, y_{\text{predicted}})$$

$$\text{MSE} = \text{mean-squared-error}(y, y_{\text{predicted}})$$

Plot the model :-

plt.scatter(x, y, color='green')

plt.plot(x, y_predicted, color='red')

plt.grid

plt.show



Step - 8 : Decide on model deployment

⇒ Final step involves figuring out how to make strategic decisions like adjusting budgets, setting prices etc. Based on the model output. It basically help in decision making.

Regression Eqⁿ and its Analysis

$$Y_i = B_0 + B_1 X_i + E_i$$

y-intercept

y-slope

$$\text{So, if } :- \hat{Y}_i = 61655.3653 + 3076.1774 X_i \quad (\text{MBA stud vs % of marks in salary from class 10th})$$

Interpretation :-

If a student have

0% in class 10th still

his/her salary = ₹ 61,655 approx

It means for every 1% increase in class 10th score the salary increases by ₹ 3076.1774 ...

($X = 0$, then, $y = ?$ is
y-intercept)

* Estimation of Parameters (using OLS)

→ A crucial step in regression model building is estimation of regression parameters.

→ Given an dependent variable (y_i) and the corresponding independent variable values (X_i) and each subject to random error (ϵ_i), we have to find the best eqⁿ to represent the relationship b/w these two variables.

● Assumption of OLS-based Linear Regression model

① Linearity :- Regression model must be linear in its parameters meaning, β_0 and β_1 must be linear. [Variables X_i etc. can or cannot be linear]

② Variation in independent variable (X_i) :- There should be sufficient variation in X_i values then only model can be correct Eg. $Y = \text{consumption level}$ and $X = \text{income level}$. And the X is income taken of every person is $₹10,000 \rightarrow$ meaning no variation \Rightarrow so, model will not able to predict a \hat{Y} , based on a new X .

③ The explanatory variable X_i is assumed to be non-stochastic (ie. X is deterministic). [\Rightarrow Fixed value of independent variable]

⇒ This means, X is controlled and not subjected to random variation or influences \Rightarrow which will help predicting Y more straightforwardly.

Eg. $X = \text{Study hours}$? If, we take into account the variation of student feeling like $Y = \text{marks in exam}$ sometimes they wanna study so before 1 day they study 18 hrs but before they did not due to their mood \Rightarrow so it makes determining Y difficult. So we make ' X ' more controlled and consistent so its easy to determine Y .

④ Residuals (ϵ_i) follows a normal distribution and have zero mean.

→ This means, if we plot all residuals from our model, they would form the classic bell shaped curve of normal distribution.

→ Mean = 0, means on average the errors do not systematically overestimate or underestimate the dependent variable (y).

* Now, Based on above 4 assumption :- In OLS, our objective is to find optimal value of β_0 and β_1 such that we minimize the sum of Squared Errors (SSE) i.e.

$$SSE = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 X_i)^2$$

on partial deriving it w.r.t. β_0 and β_1 , we get :-

$$\beta_0 = Y_i - \beta_1 X_i, \quad \beta_1 = \frac{\text{cov}(X, Y)}{\text{var}(X)} \quad \text{or} \quad \beta_1 = \frac{\sigma_Y}{\sigma_X}$$

~~Jmp~~ calculate β_1 (σ_X, σ_Y are SD).
and put value here.

⑤ The variance (meaning how spread out the residuals are) of the residuals $\Rightarrow \text{Var}(\epsilon_i | X_i)$ is constant for all values of X_i .

→ When the variance of residual is constant for different values of X_i its $\text{Rfa} \Rightarrow$ Homoscedasticity.

→ But, when the variance of residual is non-constant its $\text{Rfa} \Rightarrow$ Heteroscedasticity

(we want Homoscedasticity ✓ and we don't want heteroscedasticity ✗)

Eg. Whether its a cold day (less ice cream sold) or a hot day (more ice cream sold), the amount to which our prediction is wrong should have some spread. If the variance changed with temp \Rightarrow then Heteroscedasticity which will violate our OLS assumption.

[Here, $X_i = \text{temp}$, $Y_i = \text{ice cream sold}$]

⑥ In case of time series data, residuals are uncorrelated that is $\text{cov}(\varepsilon_i, \varepsilon_j) = 0$ for all $i \neq j$.

→ In time series data \Rightarrow data is collected over time \Rightarrow residuals are uncorrelated means \Rightarrow residual/error at one point of time (e.g. error in predicting electricity consumption on Monday) does not tell anything about the residual (error in ... on Tuesday). They are independent of each other.

• Validation of Simple Linear Regression Model •

→ It's important to validate the regression model to ensure its validity and goodness of fit before it can be used for practical applications.

→ The following measures are used :- (+ Their interpretation on calculating them)

* Coeff. of Determination (R^2) [R-square]

→ Objective of regression is to explain the variation in Y (or predict) based on knowledge of X .

→ $R^2 \rightarrow$ measures the %age of variation in Y explained by the model ($\beta_0 + \beta_1 X_i$)

Simple Linear Regression model

↓
Broken into -

$$Y_{\text{(variation in } Y)} = \beta_0 + \beta_1 X_i \quad \text{(variation in } Y \text{ explained by the model)} + \varepsilon_i \quad \text{(variation in } Y \text{ not explained by the model)}$$

Explained variation Unexplained Variation

so, Total variation = Variation in \hat{Y} + Variation in Y
 in Y explained by model not explained by
 model.

$$\downarrow \quad Y_i - \bar{Y} = \hat{Y}_i - \bar{Y} + Y_i - \hat{Y}_i$$

Applying summation to this squares:-

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

$\underbrace{\sum_{i=1}^n (Y_i - \bar{Y})^2}_{SST}$
 (Sum of square of
 total variation)

$\underbrace{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}_{SSR}$
 (Sum of square of
 variation explained
 by regression model)

$\underbrace{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}_{SSE}$
 (Sum of square of
 variation)

(Sum of square of
 error or
 unexplained
 variation)

Coeff. of Determination = $R^2 = \frac{SSR}{SST}$ = Explained variation / Total variation

$$OR \quad 1 - \frac{SSE}{SST}$$

(since, $SSR = SST - SSE$)

Interpretation

\downarrow
 ① R^2 lies b/w 0 and 1

② Higher the value of R^2 (meaning closer to 1) the
 better fit the model.

③ R^2 less than 0.5 (or 50%) \Rightarrow not good fit
 (Independent var)

Eg: $R^2 = 0.1563 \Rightarrow$ means, grade 10 marks explain 15.63% of variation
 in the starting salary of MBA students

($X = \text{grade 10 marks}$, $Y = \text{salary}$) \downarrow (dependent var)

Ideally
 (not good fit) \rightarrow it should have
 been able to explain 100% variability.

Analysis of (ANOVA) (F-static)

some constant variance

* Significance Test (p-value)

→ In SLR we are trying to find a relⁿ b/w X and Y.

⇒ p-value is used to determine whether X has statistically significant relationship with Y (dependent var). (independent var)

How it works in SLR :-

- ① H_0 (null hypothesis) :- In SLR, null hypothesis states that there is "no" relⁿ b/w X and Y. ⇒ means, $\beta_1 = 0$ (β_1 = coeff. of X).
- ② H_A (alternative hypothesis) ⇒ There is relⁿ b/w two $\Rightarrow \beta_1 \neq 0$.
- ③ Significance level (α) ⇒ we choose a significance level ($\alpha = 0.05$ mostly) before running regression analysis.
- ④ p-value → Fit SLR model into data, we get a p-value from β_1 (coeff. of X).
- ⑤ Conclusion →

If, p-value $< \alpha \rightarrow$ Reject H_0 and hence there is relⁿ b/w two.
Accept H_A

If, p-value $> \alpha \rightarrow$ Accept H_0 , there is no relⁿ b/w two.
Reject H_A

Eg. data set of House prices (X) [trying to predict] based on size of house (X). ↓

⇒ perform SLR and find p-value of the ~~coeff~~ β_1 (coeff. of X), suppose it comes to be $= 0.03$.

⇒ our significance level $= 0.05$

Since, p-value $< 0.05 (\alpha) \rightarrow$ we reject H_0 and say there is relⁿ b/w house price and size of house.

p-value for either test \rightarrow t-test } same \rightarrow In simple linear regression.

so in other words,

p-value for t-test or f-test $< 0.05 \Rightarrow$ means ~~not~~ relⁿ b/w X and Y and model is statistically significant.

★ Interpretation

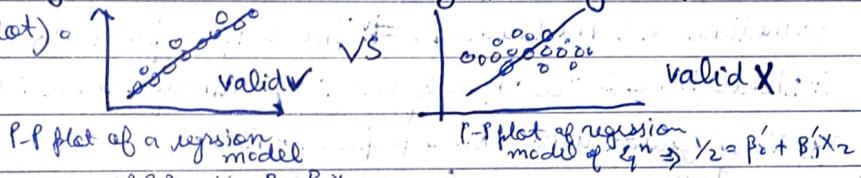
If, p-value less than 0.05 (α = level of significance) we reject H_0 meaning there is relationship b/w X (independent var) and Y (dependent var).

★ Residual Analysis

It is important to check if the assumptions of the SLR is being satisfied or not. It consists of checking the following four points:-

① The residual $(Y_i - \hat{Y}_i)$ is normally distributed or not.

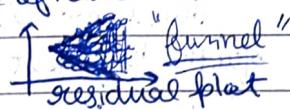
Easiest way to determine this is by plotting P-P plot (Probability-Probability plot).



dots close to diagonal \rightarrow so residual follow normal distribution. And model is valid.

② Test of Homoscedasticity

If there is heteroscedasticity (meaning variance of residual non-constant) there is funnel shape in residual plot.

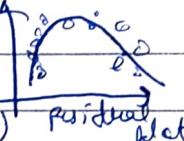


③ Testing functional form of Regression model.

Residual plot shouldn't have any pattern (a parabola).

if it is there then \rightarrow model is not valid.

(incorrect functional form used)



④ Outlier Analysis

Outlier are those values whose value show a large deviation from mean value. Presence of outlier can have influence on regression coeff \rightarrow changing the best fit line.



The following distance measures are used to identify influential observation for outliers:-

① Z-score \Rightarrow It is the distance of observation from its mean

value.

$$Z = \frac{Y_i - \bar{Y}}{\sigma_Y} \quad @ \sigma_Y \text{ and } \bar{Y} \text{ are SD and mean of dependent variable.}$$

* Interpretation

\Rightarrow Any observation with Z-score greater than 3 will be classified as outliers/influential obs. whose presence can alter regression coefficients.

② Mahalanobis distance \Rightarrow Its value greater than Chi-square critical value for an obs. mean that that obs is outlier.

③ Cook's distance :- Cook's distance value more than 1 indicates highly influential obs./outlier.

Extra point in R^2

that the model is

\hookrightarrow Not all times, high R^2 value mean better fit. It can sometimes be misleading. \Rightarrow This is K/a spurious regression :- meaning it shows there is relⁿ b/w X and Y even though there is n't.

Eg. No. of facebook users | No. of people left Finland.

2004

1 lakh

3.2 lakh

2005

2 lakh

2.9 lakh

etc.

\Rightarrow On calculating $R^2 = 0.8$ it means they both are related / model is correct. Our model eqⁿ is $\hat{Y} = 1.996 + 2.1X$

The is no relation \times
Bt the two, but R^2 is

\leftarrow means for every 1 person leaving finland 2.1 use

high \Rightarrow hence misleading so, R^2 is not always a correct measure to tell better fit.

Multiple Linear Regression

→ It is statistical technique used to find existence of relationship b/w a dependent variable (y) and several independent variables.

→ Functional form of MLR is given by:-

$$Y_i = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \beta_5 X_5 + \dots + \beta_n X_n + \epsilon$$

here, $\beta_1, \beta_2, \beta_3, \beta_4, \beta_5$ etc. are 'partial regression coeff'.

β_0 = y-intercept or constant term

ϵ = error term

Assumptions of MLR

→ same as of that of SLR actually!!

MLR eqⁿ and its Analysis :-

→ consider, two independent variables: promotion and rating point
dependent variable: advt. revenue

→ consider, dependent variable \Rightarrow House price
independent variable \Rightarrow size and Age of house

lets say,

$$Y = 41.238 + 5931.85 \times \text{Size} + -1000 \times \text{Age of house}$$

↓

① for every one unit ↑ in size \Rightarrow price of house
 \rightarrow say ≈ 5931.85 (keeping age const.)

② for every one unit ↑ in age of house \Rightarrow price of house
 \rightarrow say ≈ -1000 (keeping size const.).

→ In MLR, it may not be possible to control a variable (i.e. keeping it constant) in many situations.

Those which represent characteristic like person gender, marital status. It can take numerical value. E.g. 1 for male, 2 for female etc.

Working with categorical variables

E.g. $y = \text{Salary}$, $x_1 = \text{Highschool}$, $x_2 = \text{Bachelor}$, $x_3 = \text{Master}$, $x_4 = \text{Phd}$.

We take x_1 as reference and calculate others in reference to it.

Multi-collinearity and VIF

→ It refers to presence of high correlation b/w two or more independent variables in a multiple linear regression model.

→ It can cause issue with the estimation of coeff. and effect the model.

→ It refers to ~~high~~ presence of perfect linear relationship b/w few or all independent variables in a regression model.

Eg. There are $n \rightarrow$ independent variable then

$$\alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n = 0 \quad (\text{where } \alpha_i \text{ are constant})$$

↳ Perfect relⁿ with each other.

→ One of the Assumptions of CLRM (classical linear regression model) is that there should be no multicollinearity among the independent variable.

→ As the no. of independent variable $\geq 2 \Rightarrow$ problem of multicollinearity can also \geq .

→ Therefore we need to detect such multicollinearity in data and one such way is by calculating VIF.

Variance Inflation Factor

→ It is used to measure identifying existence of multi-collinearity.

For Eg. consider a regression model,

$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots$ and their R^2 value be R_{12}^2 then,

$$\text{VIF value} = \frac{1}{1 - R_{12}^2}$$

Formula

\Rightarrow VIF greater than 5 requires investigation.

Interpretation *

$VIF = 1 \Rightarrow$ no correlation b/w independent variables

$1 < VIF < 5 \Rightarrow$ moderate multicollinearity

$VIF > 5 \Rightarrow$ High degree multicollinearity, problematic

\rightarrow To fix high multicollinearity we can do :-

(1) Increase Sample Size.

(2) Drop non-essential ^{independent} variables.

Outlier Analysis \Rightarrow Discussed in SLR.

* Autocorrelation

\hookrightarrow also k/a serial correlation. It happens when there is a relⁿ b/w current obs. and one or more past obs. It occurs when the residual of a regression model exhibit a pattern.

Eg. Measuring temp. every day \Rightarrow If since last 2-3 days temp is high \Rightarrow then this increases chance that temp. tomorrow will also be high. In other words, autocorrelation occurs when past values in a data series influence future values.

\rightarrow Durbin-Watson test is a statistical test used to detect the presence of autocorrelation. It is based on estimated residuals.

* Interpretation : Its Range 0 to 4

(1) A value of 2.0 \Rightarrow no autocorrelation

(2) Value less than 2.0 \Rightarrow +ve autocorrelation (High temp yesterday \rightarrow High temp tomorrow)

(3) Value greater than 2.0 \Rightarrow -ve \dots (rare, Eg. value = 3.5 means sales this quarter can go down if last quarter were up. Sale High (Past) \rightarrow Low now.)

Validation of MLR model

To validate a MLR model we can follow the given measures and perform below test to check it.

① coeff. of multiple determination (R-square) and adjusted R-square

It explain how well

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

(R^2 = portion of variance in y explained by the model)

Your model explain variation in the dependent variable.

→ we use adjusted R^2 more bcz in MLR, normal R^2 can be manipulated as it can be increased artificially by adding more variables, but adjusted R^2 takes care of that aspect.

⇒ Adjusted R^2 considered more reliable than R^2 .

★ Interpretation

② $R^2 \rightarrow R^2$ range 0 to 1. $R^2 = 0 \rightarrow$ means model doesn't explain any variability.

→ more its near 1, the better fit the model with data is.

e.g. $R^2 = 70\%$ or 0.70 means 70% of variation in dependent variable (y) is explained by the independent variable.

③ Adjusted R-square → ④ Also, b/w 0 to 1

⑤ less than or equal to R-square

⑥ More near to 1 the better fit the model is.

⑦ T-test, F-test

$$\begin{aligned} H_0 &= \text{null hypothesis} \\ H_A &= \text{alt hypothesis} \end{aligned}$$

calculate p-value → If, $p\text{ value} < 0.05 \Rightarrow$ H₀ reject, rel "b/w H₀ accept" $\times y \checkmark$

$p\text{ value} > 0.05 \Rightarrow$ opposite of above.

- ③ Check for multicollinearity, if present take appropriate steps.
- ④ And, check for autocorrelation (in case of timeseries data).
- ⑤ Residual Analysis (same as that of SLR \Rightarrow check if Homoscedasticity is followed or not).

Unit-3

* Logistic and Multinomial Regression

* Logistic Regression

(Y)

→ It is used when the dependent variable has two possible outcomes i.e. binary

→ Consider it answering a Yes or No question based on some info.

Eg. A doctor says a person has a disease (Yes or No) based on their symptoms (like cough and cold). ^{simplifies}

→ Logistic Regression will help in calculating probability the patient has disease based on system. It outputs value 0 and 1, which can be interpreted as probabilities.

* Multinomial Regression

→ It is extension of Logistic Regression, but it is used when dependent variable has more than ~~two~~ two categories.

Eg: In above case based on symptoms (like cold, cough, fever) a newly appointed doctor said that the patient \Rightarrow Yes (Has this normal fever), No (he doesn't have ^{it's simply due to weather}), Maybe.

→ So, your dependent variable is no longer binary and has more categories.

Eg: Based on customer age, region, income level \Rightarrow The ice cream company has three flavours (= Y (dependent variable)) : Choclate, Vanilla, Strawberry.

↓

Instead of binary outcome, the outcome/dependent variable has multiple categories.

★ Logistic Fn x (Sigmoid Fnx)

→ The logistic regression model (or binary regression model) is given by this sigmoid fnx :-

$$\frac{f(x)}{1+f(x)}$$

$$P(Y=1) = \frac{e^z}{1+e^z}$$

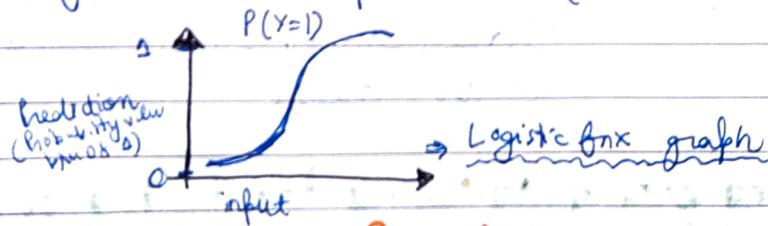
where,

$P(Y=1)$ - probability of obs labelled as 1

$$z = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots$$

↳ X_1, X_2, X_3 are independent var.

⇒ The logistic fnx is a S-shaped curve (as it is a probability fnx).



Estimation of Parameters in Logistic Regression

→ one of the major assumption of MLR was, residuals follow normal distribution, but here, in logistic regression residuals don't follow normal distribution, so it is difficult to follow a normal distribution → Thus we cannot use OLS to calculate parameters fine. They are calculated using maximum likelihood estimator (MLE).

using which β_0 & β_1 are calculated.

Interpretation of Parameters of Logistic Reg.

→ An easy interpretation of the regression coeff, β_i , is that:-

- ① +ve coeff ($\beta_i = +ve$) → means probability of ~~happier~~ event occurring ($P(Y=1)$) also rises
- ② -ve coeff ($\beta_i = -ve$) → means probability of event occurring also falls.

Eg. Predict student pass ($Y=1$) or fail ($Y=0$) based on no. of hours they study (X)

→ If, $\beta_1 = 0.8 \Rightarrow$ means, each additional hour of him/her studying will rise the probability of them passing.

$\beta_1 = -0.5 \Rightarrow$ means, each additional hour of them studying will fall the probability of them failing. (counter-intuitive consider $X = \text{stress}$).

Estimating Probability using Logistic Regression

→ It involves fitting data into the model, and then using this model to predict the probability of an event occurring.

Example :- Estimating Probability of getting a Loan Approved.

↳ Consider you are working for a bank and want to predict whether a person gets loan approved or not X based on their age & income.

Step-1 :- Collect Data

3IM

① Gather past loan applicants data including their age, income & whether they got their loan approved or not.

Step-2 :- Fit into Regression model.

- ④ Use collected data to fit a regression model. The regression model would have an eqn:-

$$\log\left(\frac{P(Y=1)}{1-P(Y=1)}\right) = \beta_0 + \beta_1 \times \text{Income} + \beta_2 \times \text{CreditScore} \quad \text{--- (1)}$$

⇒ here, $P(Y=1)$:- probability of loan approval, β_0 is intercept, β_1 and β_2 are coeff of income and credit score.

Step-3 :- Use model for prediction

⇒ once model is fitted, we can predict whether applicant loan will be approved or not.

Eg. If age = 25 and income = ≥ 50,000 ^{* new data}

we put into the above eqn (1) to get log odds.

then convert that log odd to get the probability using logistic fnx:-

$$P(Y=1) = \frac{e^{\text{log-odds}}}{1 + e^{\text{log-odds}}}$$

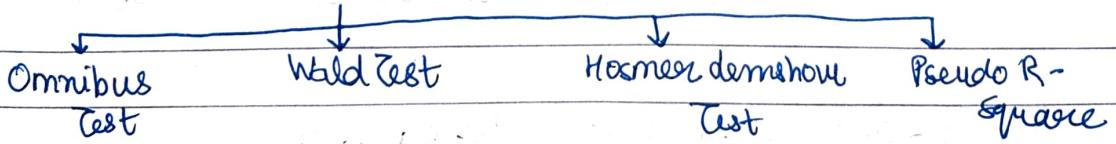
Step-4 :- Interpret the results

→ The result will be b/w 0 to 1 Eg. If its comes 0.7 or 70%. Meaning there is 70% chance that his/her loan will be approved.

Logistic Model Diagnostic

↳ We perform diagnostic test before a binary / logistic regression model is accepted for deployment.

4 test :-



Omnibus :-

- ① To check overall significance of a logistic regression.
- ② Checks whether models with predictors (if X is an independent var) is significantly better at explaining the variance than a model with no predictors (null model).

Formula \Rightarrow omnibus = $\frac{\text{deviance reduced} - \text{deviance full}}{\text{deviance reduced}}$ \geq we obtain these value by SPSS (Statistical soft.)

Interpretation :-

\Rightarrow After performing this test \rightarrow then doing hypothesis test \rightarrow If, we get $p\text{-value} < 0.05 \Rightarrow$ means, our predictors (X_1, X_2 etc.) were able to explain variance and have a significant impact on Y .

Eg. Model :- predict admission ✓ or X based on GPA and test-score. This test will show whether they both together significantly predict admission chances than a model that predicts admission without these X_1 & X_2 or predictors.

Wald Test :-

- ① Check significance of individual coeff in the model (by checking significance of B_1, B_2 etc.)

Formula \Rightarrow Wald test = $W = \frac{(\text{estimated coeff})^2}{(\text{standard error of coeff})^2}$

Interpretation :-

\hookrightarrow If the t -value calculated after Wald test < 0.05 then, the Wald test statistic for a coeff is a significant contributor to model.

\hookrightarrow $p\text{-value for } B_1 \text{ of } X_1 = 0.205 > 0.05 \Rightarrow$ this means X_1 may not have a significant impact on Y .

Hosmer Lemeshow Test

- ① It is a goodness-of-fit test for logistic regression model.
- ② It evaluate the models performance in predicting outcome based on observed vs predicted probabilities.

Interpretation :-

\hookrightarrow If $p\text{-value high} \Rightarrow$ model a good fit

$p\text{-value low} \Rightarrow$ mode not a good fit

Pseudo R-square :-

- ① Similar to what R^2 measure in LSR, but not directly comparable to it since R^2 is 0-1.
- ② Pseudo R-square [since model fit is less here in logistic]

③ It provides indication of how well predictor variable (X) explain the variance in dependent variable (y) in logistic regression.

Interpretation

There are 3 types of pseudo R-square

- McFadden
- Cox and Snell
- Nagelkerke \Rightarrow Range 0 to 1 closer to 1 the better

↳ Generally,

Higher values indicate a better fit.

Eg. Cox and Snell R-square = 0.5961 \Rightarrow 59% \Rightarrow meaning that around 59% of variance in the dependent/outcome variable (X) is explained by the model which is moderate here.

Model Performance \rightarrow The following things help in determining the model's performance.

* Classification Table (or confusion matrix)



Helps to determine the performance of a model.

		Actual (Reality matrix)	
		+ve	-ve
Predicted	+ve	TP <small>(True Positive)</small>	FP <small>(Type I error)</small>
	-ve	FN <small>(Type II error)</small>	TN

(True = correctly)
False = incorrectly

① Sensitivity / True Positive Rate / Recall

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

$$\text{Sensitivity} = \frac{\text{Predict bhi kya raha hai}}{\text{actual mai +ve}} \\ = \frac{TP}{TP + FN}$$

② Specificity / True Negative Rate

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

$$\text{Specificity} = \frac{\text{Actual -ve ka predict gya hai}}{\text{actual mai -ve hai}} \\ = \frac{TN}{FP + TN}$$

③ Accuracy :- Accurately diagnosed to the total no. of patients

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

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

→ we predicted correctly

④ Precision :- correctly predicted

$$\text{Precision} = \frac{TP}{TP + FP}$$

→ actual mai bhi
→ predicted +ve

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

⑤ F-score (or F1 score) :- It is the harmonic mean of Recall and Precision,

$$F\text{-score} = \frac{1}{\frac{1}{\text{Recall}} + \frac{1}{\text{Precision}}}$$

$$\frac{1}{\text{Recall}} + \frac{1}{\text{Precision}}$$

$$F\text{-score} = \frac{2 \times \text{Precision} \times \text{Recall}}{(\text{Precision} + \text{Recall})}$$

		observed		
		Correct	Incorrect	
Prediction	Correct	17	3	FP
	Incorrect	0	4	TN

① Precision = $\frac{17}{20}$, ② Recall = $\frac{17}{17+0} = 1$, ③ F score = $2 \times \frac{17 \times 1}{20}$

④ Specificity = $\frac{4}{4+3} = \frac{4}{7}$

Gini Coeff / Gini Index

- It has a value b/w 0 and 1.
- performance of the model improves with ↑ in gini coeff. (more near to 1 ⇒ better the model)
- Gini coeff can be calculated from AUC of ROC.

$$\text{Gini coefficient} = 2 \times \text{AUC} - 1$$

ROC (Receiver operating characteristic)

→ The performance of logistic regression model can be assessed graphically by:-

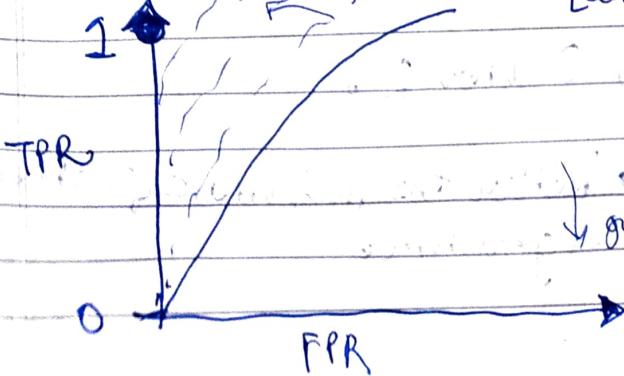
ROC → It is the graph b/w TPR and FPR

graph more near to 1
better the model

(True positive rate) (False positive rate)

TP	FP
FN	TN

[Sensitivity] [1 - False positive rate (i.e. specificity)]



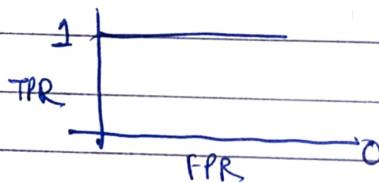
graph more near to FPR, worst the model.

→ Basically we plot TPR and FPR at various threshold points then \Rightarrow ROC curve is formed.

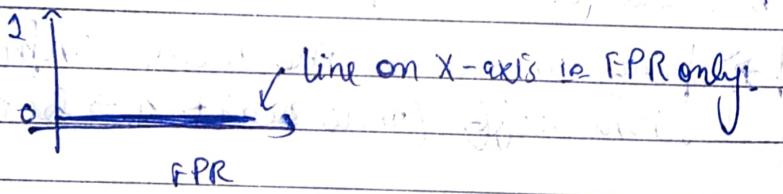
$$TPR = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{FP + TN}$$

② With a $TPR = 1$ and $FPR = 0$, a perfect classifier so, ROC curve would look like:-



③ With $TPR = 0$ and $FPR = 1$, worst model, ROC curve look like:-



★ Area Under Curve (AUC)

→ When making a ROC \Rightarrow AUC is used to access its effectiveness.

→ Calculating AUC means calculating the area under ROC curve
→ no formula \Rightarrow done using software.

\Rightarrow AUC tells (more the AUC) \Rightarrow tells that model is able to ~~explain all the~~ give better performance.

\Rightarrow AUC values lie b/w 0 and 1.

↓
More the AUC \Rightarrow better the performance of the algorithm.