

## Regression Notes Strong / perfect Regression line.

↓  
its explanatory variable.  
مستغنى يكونوا مرتبطتين ببعض  
Independent Correlated  
=  
Multi Correniality.

Auto Correlation  $\Rightarrow$  Relationship Between variable & itself.

Muli Correlation  $\Rightarrow$  Relationship Between explanatory variable of Regression line.

Correlation  $\Rightarrow$  Relationship Between two independent variables.

\* the Regression line "perfect"  $\Rightarrow$  Should build on Independent Variables.

\* لو فيه 5 او 6 وفيه علاقة بينهم لازم نسيبها الى عامل مشترك قبل ما  
ابدأ اعمل Regression Model

\* عشان افسر اناشوف العلاقة بين ال Variables اعمل Scatterplot وتكون الى  
ليحدلي "المفروض يطلع" Randomized "دكة معناه ان مفروض علاقة

\* Scatterplot  $\Rightarrow$  use to find Relationship Between variables.  
 $\Rightarrow$  use to Detect the outliers.

\* اناشوف لو انا فلتا مؤثرين من ليقنع اعمل ال Regression Model وها موجودين.  
\* يعني ال Model الاول مرة بينهم و مرة من غيرهم و اناشوف لو فيه  
اختلاف في ال output يبقى كسبراهم لو كوكو يبقى كسبراهم  
وها كدة من مؤثرين.

Nonlinear. صفت على

## Correlations Linear Relationships:

لدى قوة الطلاقة يرتفع

Correlation

NS

Covariance.

Strength of Linear Relationships Bet. 2 variables.

Measure of association. Between two variables. (numerator.)

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{S_x S_y (n-1)}$$

لما ابي افان بين اكر متباين  
استخدم ال Correlation حسن

\* التنبؤ رى يعنى بيقدر ال Strength  
بسن مقرر شى اعتمد على ال Covariance  
وانا بفان بين pairs of variables

\* Simple Linear Regression: this line quantifies Relationships.

Between two variables  $x \rightarrow y$   
explanatory  $\leftarrow$  dependent

Scatterplot. عن طريق ال

يحدد نقطة من النقطة  
line

\* تحاول ادور على line بكل Min ال Deviations بقا ال Datapoints.  
(x) مش ان تحاول اجيب line يقع عليه كل او مظم ال Datapoints

Fitted value Magnitude value of point if this point on the straight line.  
( $\hat{y}$ ) بحسب هوى لو كانت على الخط كانت شيق فيمتا كام.

( $\hat{y} - y$ )  $\rightarrow$  Residual. (الباق) if (-ve)  $\therefore$  point under the line.  
if (+ve)  $\therefore$  point over the line.

Least Squares Estimation:  $\sum (-ve)$  وال  $(+ve)$  Errors ال Square بعمل  
 $\sum (-)^2$

\* the Best Fitting line  $\rightarrow$  Line with the Smallest Sum of Squared Residuals

Imp

$$\sum e_i^2$$



output (dependent)  $\leftarrow$

Regression eq.  $y = a + b x$   $\rightarrow$  Feature (explanatory).

Intercept  $\downarrow$  Slope  $\downarrow$

Intercept  $\rightarrow$  (+ve), (-ve)

$b = r_{xy} \frac{s_y}{s_x}$

$a = \bar{y} - b\bar{x}$

Imp for نظری

observed value = Fitted value + Residual.

Residual  $\Rightarrow$  Difference Between Actual & Fitted value of dependent variable.  $(y - \hat{y})$

Standard error of estimate  $\Rightarrow$  Magnitude of Residuals. i.e. Standard deviation of Residuals.

$R^2 \Rightarrow$  Imp to Measure goodness of fit between (0 - 1).

$$S_e = \sqrt{\frac{\sum e_i^2}{n-2}}$$

no Relationship  $\leftarrow$  perfect.

$$R^2 = 1 - \frac{\sum e_i^2}{\sum (y_i - \bar{y})^2}$$

note: In Simple Linear Regression  $R^2$  is Square of Correlation Between Dependent var (y) and explanatory var (x).

Residual = (Actual - predicted)

Standardized Residual =  $\frac{\text{residual}}{\text{Standard deviation}}$

$\sim$  if  $> (\pm 2 \text{ or } \pm 3)$   
 $\therefore$  it is outliers.

# Checking Assumptions

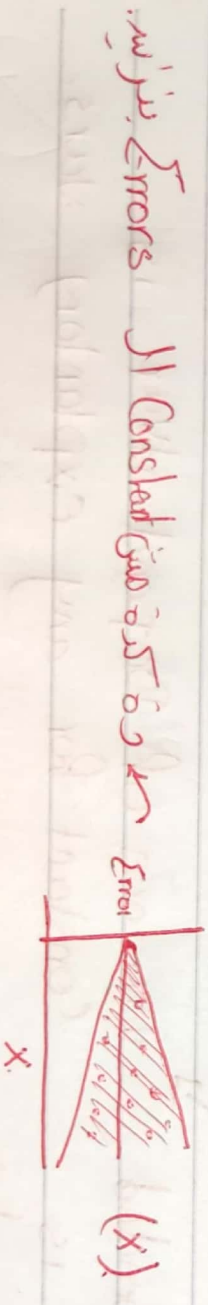
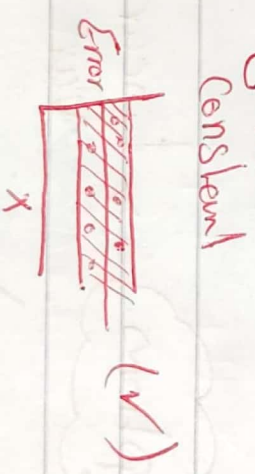
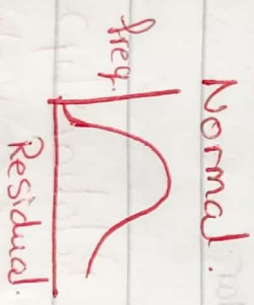
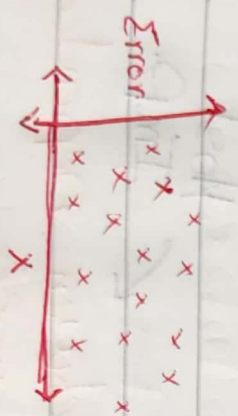
لترى نيل Check على كذا حاجة في الـ Regression Model. بيد انك الـ

Linearity Normality Homoscedasticity Independence of Errors. Successive observations should be not Related.

Linear Errors. Variation about Regression line Imp when Independent variable is time.

1) Scatterplot. Histogram. is constant

2) Residual plot. Residuals. plot. Constant





**Multiple Regression**  $\Rightarrow$  Fitting a plane to Data. in 3D. space.  
 $\Rightarrow$  Slope term for each explanatory variable in equation but interpretation of these terms are different.  
 $\Rightarrow R^2$  Same as Simple Regression.

$$y = a + b_1 x_1 + b_2 x_2 + \dots + b_k x_k$$

$$S_e = \sqrt{\frac{\sum e_i^2}{n - k - 1}}$$

$\hookrightarrow$  no. of explanatory var.

Most Imp.

**Check Assumptions**  $\rightarrow$  Linear Relationship Between Dependent var & explanatory vars.

Errors prob. Independent  
 Dependent var is Normally Distributed.  
 variance of Dependent var is constant for any explanatory vars.

**Note** More Common problem is. Multicollinearity, where explanatory vars. are highly correlated.

VIF greater than 10.  
 $\therefore$  have Correlations

**(from table)** VIF  $\Rightarrow$  (1) No correlation  
 $\Rightarrow$  (1-5) Moderate correlation.

Cross Sectional Data

Time Series Data

usually taken for granted.

Assumption often violated.

$\rightarrow$  Bec. auto correlation.

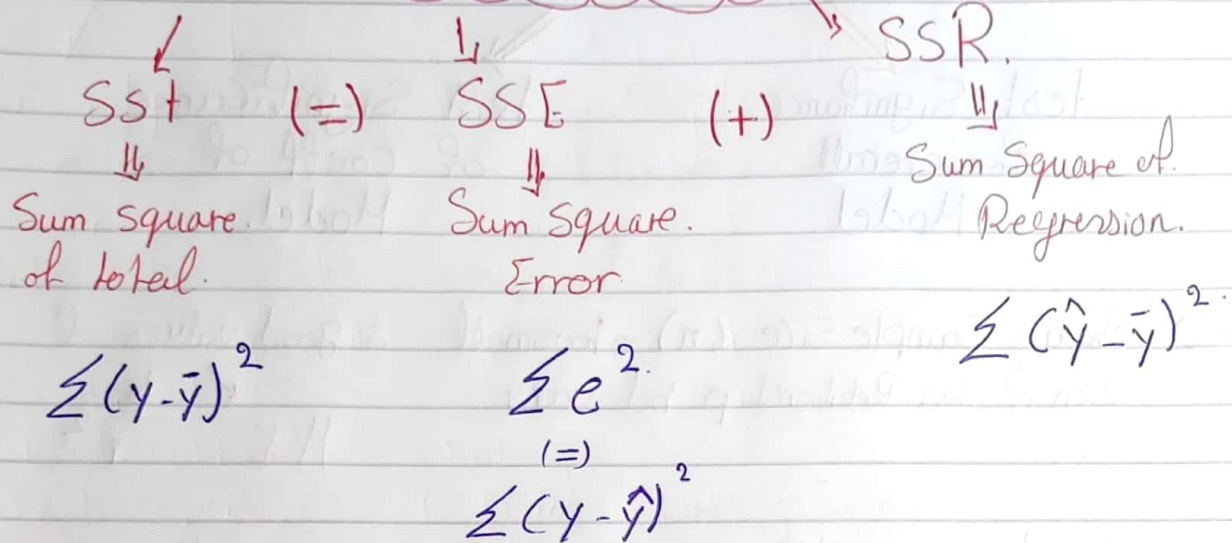
**Note** Durbin-Watson (DW) statistic. (Measure auto correlation). (always bet (0 - 4))

2  $\Rightarrow$  No auto correlation.

0 - 2  $\Rightarrow$  Positive correlation.

2 - 4  $\Rightarrow$  negative correlation

## Measure the Fit of Regression Model



**Note**  $\hat{y} \Rightarrow$  value which I get. (from Reg equation).

$\bar{y} \Rightarrow$  Average of all values

$y \Rightarrow$  Real value I have.

$$r^2 = \frac{SSR}{SST} \quad \text{or} \quad 1 - \frac{SSE}{SST}$$

## Anova table

	df	SS	MS	F	Significance.
Regression.	K	SSR	$\frac{SSR}{K}$	$\frac{MSR}{MSE}$	$p(F) \frac{MSR}{MSE}$
Residual	$n-K-1$	SSE	$\frac{SSE}{n-K-1}$		
total.	$n-1$	SST			

$\downarrow$  p value for calculated F value from table.



we do hypothesis test  
F-test from Anova table.  
testing Model for Significance.  
help Determine if the values are Meaningful.

test Significance for overall Model.

test Significance of coeff of Model.

Notes when Sample Size (n) is too small  $\therefore$  good values of even if No Relationship bet vars  
MSE  $\swarrow \searrow$   $r^2$

F-test  $\rightarrow H_0: B_1 = 0 \therefore$  No Relationship Bet x, y.  
 $\rightarrow H_1: B_1 \neq 0 \therefore$  Here is linear Relationship.

if very little Error  $\therefore$  MSE (small) & F (large)

if F (large)  $\therefore$  p-value (low)

$\rightarrow$  we can Reject Null hypothesis & Accept linear Relationship Bet x-y & MSE &  $r^2$  are Meaningful.

$\therefore$  Steps of Hypothesis test  $\therefore$

1.  $H_0: B_1 = 0$   
 $H_1: B_1 \neq 0$

2. Select Level of Significance ( $\alpha$ ) (given).

3. Calculate F-value =  $\frac{MSR}{MSE}$

4. Reject if F calculated  $> F_{\alpha, df_1, df_2} \rightarrow n-1$   
or Reject if p value  $< \alpha$ .

Choose any Method

**Multi Collinearity** \* when there is Strong Linear Relationships.  
Among set of explanatory variables.

\* في الحالة هنا ال  $(B_i)$  الى لوال Stop متو الوحيه الى يدعرب  
لدى تعلق  $(x)$  بال  $(y)$  لذن فيه  $(x)$  تانيه ظهرت بتأثر فال  $(y)$

\* there are various degrees of Multi Collinearity, but in each of them there is Linear.

**Include / Exclude Decisions**  $\Rightarrow$  we use t-value. & p-value.

||> Finding Best  $x_i$  to include them  
in Regression equation.

**Notes** \* t value & p value  $\Rightarrow$  if p-value Above  $(\alpha)$ .  
 $\therefore$  Candidate for exculsion.

if t value less than or ~~greater than~~  $(\alpha)$   
 $\therefore$  excluded from eq.

**Stepwise Regression**  $\rightarrow$  Stepwise.

Forward.

||

begins with  
No explanatory  
var then Add.  
one By one.  
until No Remain.  
variables found.

Backward.

||

begins with all.  
Variables. &  
Delete them  
one By one.  
until further  
Delete no harm good.

||>  
Same the Forward.  
except also considers  
possible deletions.  
along the way.