

— X — X —

Simple Linear Regression

(1) Regression

→ return to a less developed state
i.e. coming down to a single
variable

→ technique to investigate relⁿ
b/w dependant & independent
variables

Eg: Historical data

	Area	Dist	Bedrooms	Price
H ₁				
H ₂				
⋮				
H ₃₀₀				

↓
ID variable

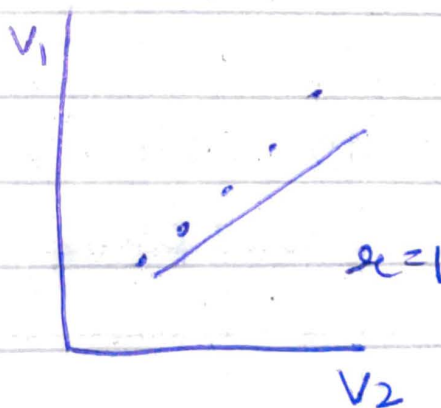
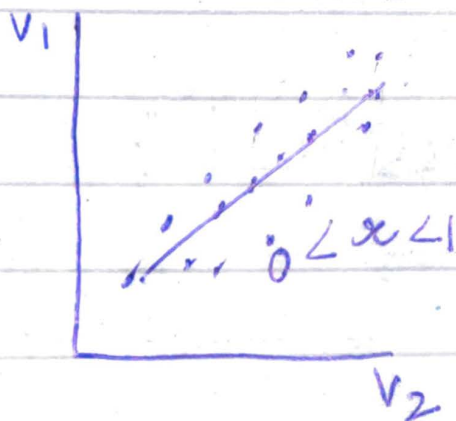
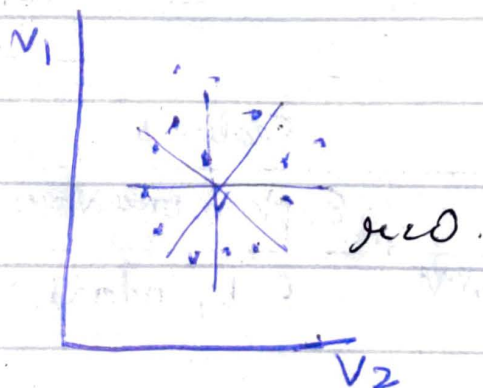
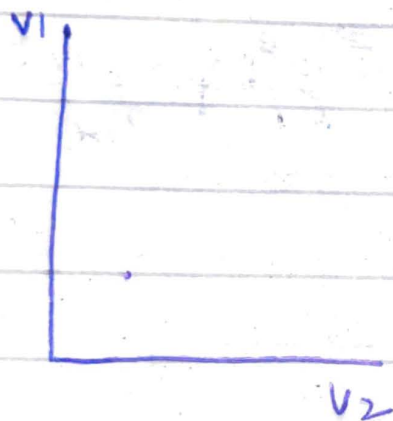
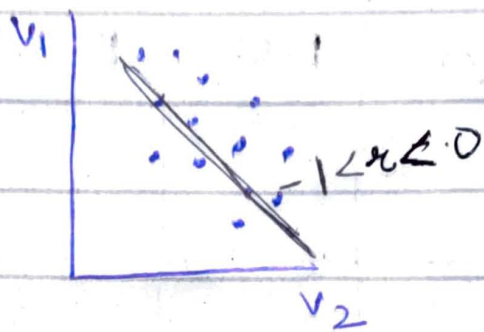
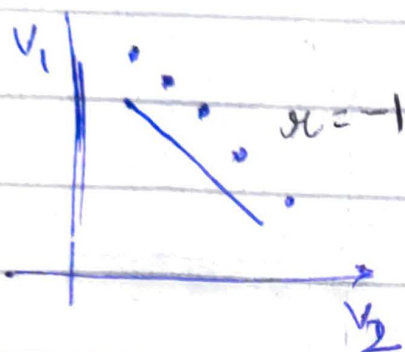
Features/
Explanatory Variables/
Independent Variables

Target/Dependent vari.

RULES:

→ ID variables are not used for predictions

Correlation

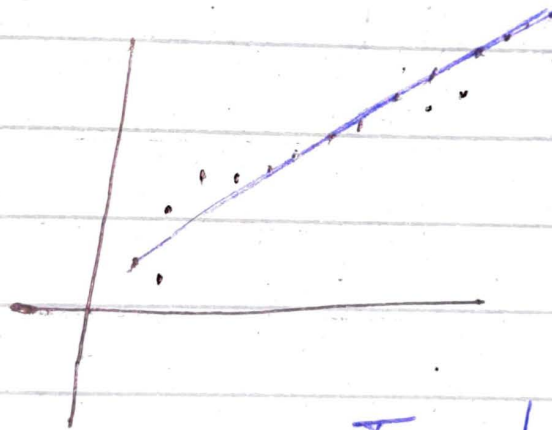


Difference b/w

Corr = 2 variables ;

reg - Variables, hypothesised,
- advance understanding of system

Ice Cream:



Temp	IC	IC = 2T	Error
given	given	predi.	Given - Pred
			<div style="border: 1px solid black; width: 50px; height: 20px; margin: 0 auto;"></div>
			minimise

General Hypothesis

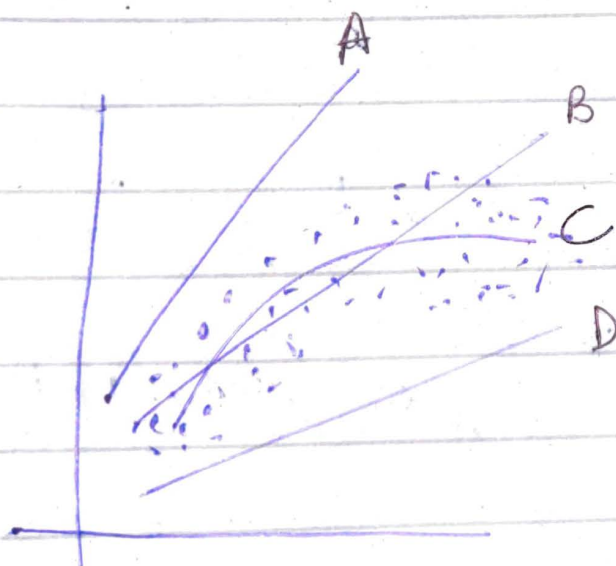
$$I.C = a + b \text{Temp} \quad a, b = \text{unknown}$$

(Ice cream)

Find a & b for which error is min.

a, b are called parameters

\therefore we need to find parameters to minimize the error.



out of ABD
B is best

out of ABCD
C is best

out of B & C,
C is best as
it minimizes the error.

The best line also has errors. If we

= unknown

non

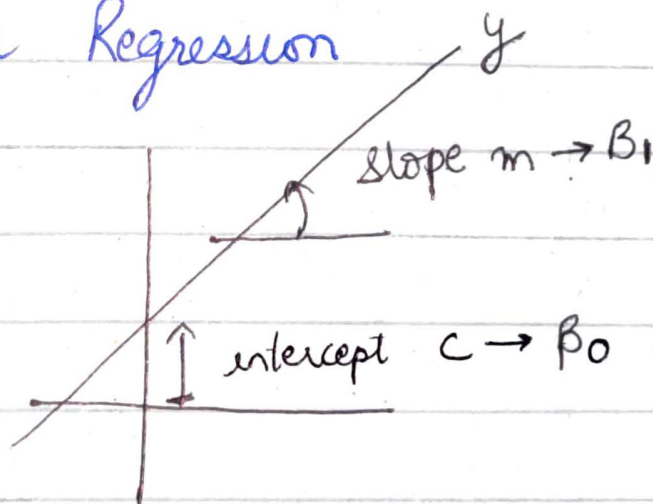
reach an error that is good enough
we reach line B. Further if
the error is still significant we
move further to line C.

→ defendable / fixed

Causation is not Regression
Regression is not correlation

BD

Linear Regression Line



$$y = mx + c$$

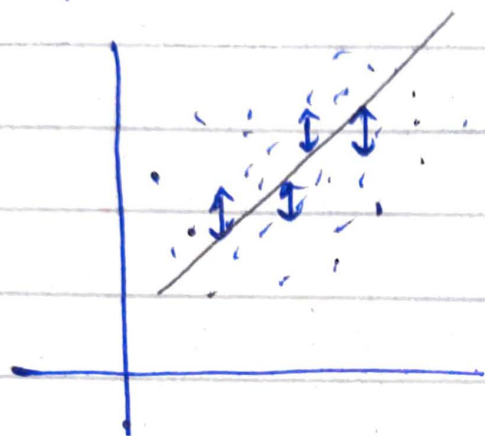
$$y = B_0 + B_1x$$

n dim \rightarrow $(n-1)$ dimensional quantity

2d \rightarrow line (1D)

3d \rightarrow plane (2d)

Linear regn. is a linear approach to find the relation b/w a dependant & independent variables



B_0 & B_1 such that ERROR is min.

$$\text{ERROR} = Y_{\text{actual}} - Y_{\text{predicted}}$$

$$\text{I} \quad \frac{1}{n} \sum_{i=1}^n (Y_a - Y_p)$$

$$\text{II} \quad \frac{1}{n} \sum_{i=1}^n |Y_a - Y_p|$$

$$\text{RMSE III} \quad \sqrt{\frac{1}{n^2} \sum_{i=1}^n (Y_a - Y_p)^2}$$

I is not used since positive & -ve errors cancel to give a false 0 error.

II \rightarrow not used as mod is used;
the mod value makes the error
non-differentiable; hence inc.
difficulty in computⁿ.

III \rightarrow used since all values are
considered & hence computation also
possible.

The sq. magnifies the error but
the sq root balances.

\rightarrow To make III, further simple

the SSR is used. $SSR = \text{Sum Square Residual.}$

$$SSR = \frac{1}{n} \sum_{i=1}^n (y_a - y_p)^2$$

Hence B_0 & B_1 are such that

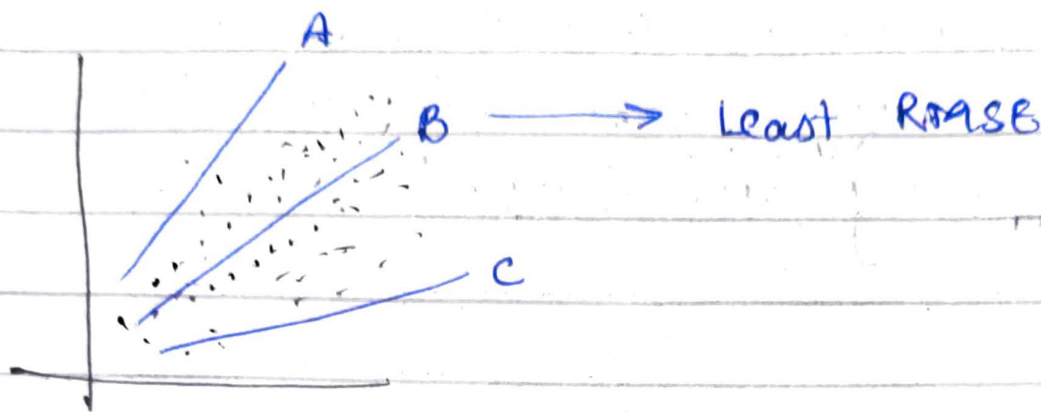
SSR is min.

$$RMSE = \sqrt{SSR}$$

Model Evaluation Metrics

① RMSE

② R-Square

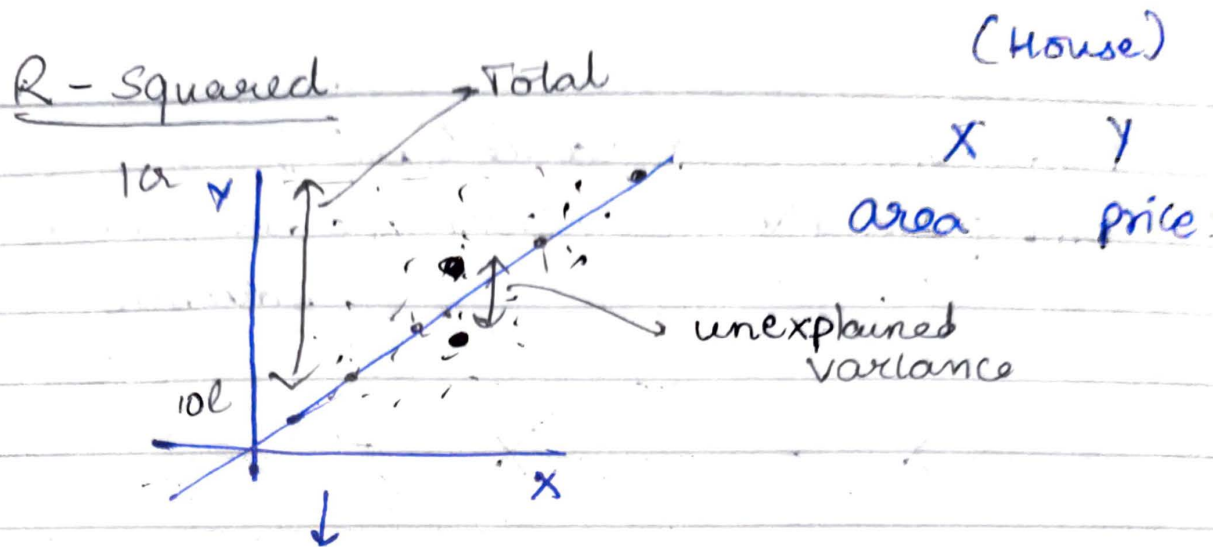


$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_p - y_a)^2}$$

$$\text{SSR} = \frac{1}{n} \sum_{i=1}^n (y_p - y_a)^2$$

A best fit line is not a perfect line

RMSE & R-sq are universal for all continuous techniques.



In the absence of an explanatory variable ; the spread of y is different

\therefore variance of y is explained by x

$$R^2 = \frac{\text{Explained Variance}}{\text{Total Variance}}$$

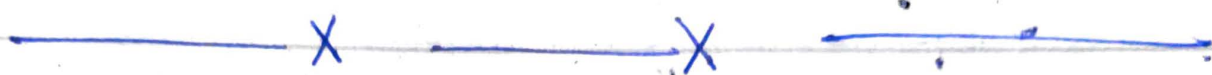
↳

$$R^2 = \frac{\text{Total Variance} - \text{Unexp. Variance}}{\text{Total Variance}}$$

$\therefore R^2$ is the % of variance explained

Total Var = Variance (in Y)

Unexplained Var = Variance in error
column



Multiple Dimensions

p explanatory variables

$X_1 \quad X_2 \quad Y \Rightarrow p$ explanatory
values

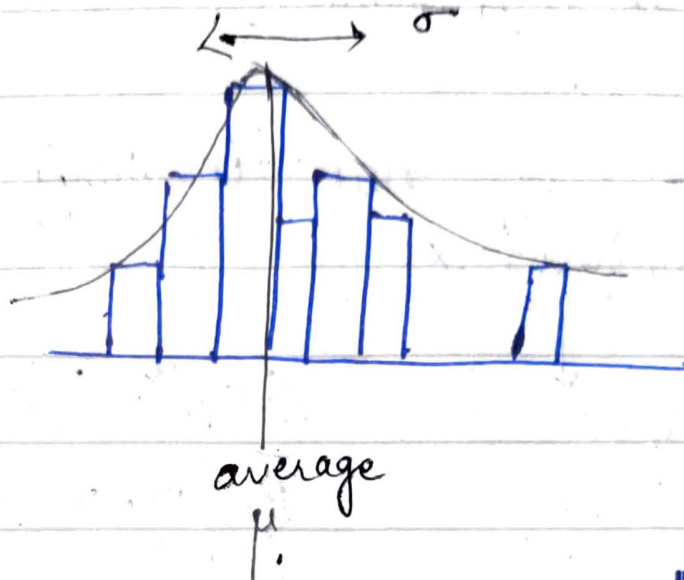
$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

↓
"Hyper plane"

'Error' = || to y axis & not
|| to the plane.

ASSUMPTIONS

① Normal Distribution



$$\mu \pm \sigma = 66\%$$

$$\mu \pm 3\sigma = 99\%$$

① The target Variable should be normal distributed

no skewness.

if not normal, transform y

$\left. \begin{array}{l} \sqrt{y} \\ \log(y) \end{array} \right\}$

$$x \sim \log y$$

$$y = e^x$$

② Linear relationship

→ TV & Independent Variable should be linearly related

- In case of non-linearity, we can transform x

① x^2

② \sqrt{x}

③ $\log x$

④ e^x

⑤ $\sin x$

④

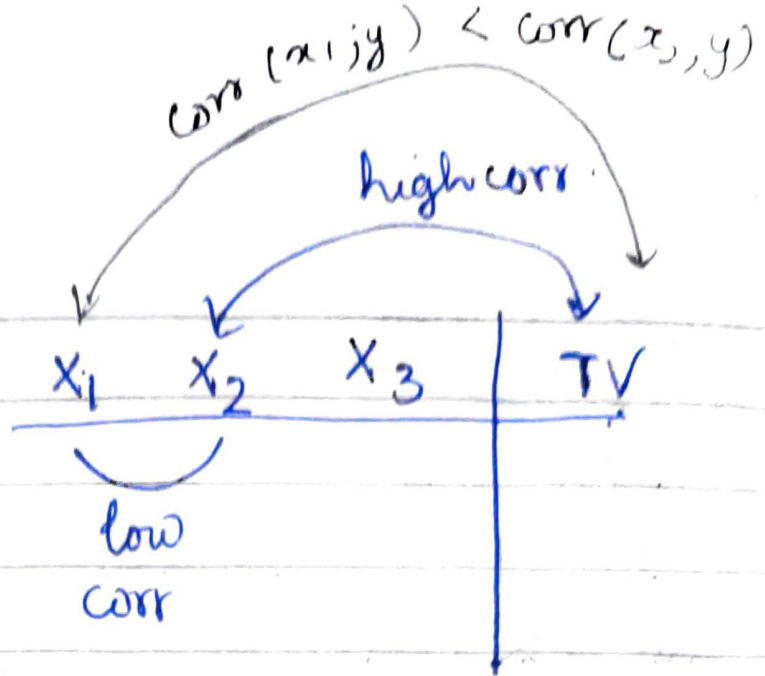
③ Independent Variables should not be correlated

no corr" & truly independent

x_1 & $x_2 \rightarrow$ low correlation

$$|r| < 0.75$$

if 2 columns are highly correlated hence 1 is dropped.



X_1 is dropped,
even if

④ Variance in error should be
Constant \rightarrow called as homoscedasticity

① done post analysis ; as evalⁿ
is done on error (after predⁿ)

② If variance / ^{in error} is not constant,
the model is defective.
The variance ^{in error} is not
constant only if the previous
assumptions are violated ~~or~~ or
transform of variables is not
correct.