

## Sb + Day 1 Sum up

what is machine learning?

→ is about How machine learning Certain things.

→ Three phase (two phases)

i.e. patterns

(e.g. ↑ no. of  
cus  $\rightarrow$  ↑ sales)

1 Training phase : learns patterns

2 Testing phase : replicate the learned patterns (i.e. predict)

||

Then that machine (that learned some pattern)  
is called model.

Model

↳ what? learns patterns  
↳ how? learn patterns [train]  
↳ purpose? replicate patterns [testing]  
(Optional: e.g. VSL)  
↳ How? (learns patterns)

↳ check the performance  
(i.e. % of patterns it  
replicate successfully).

large amount of data

|| need to

learn patterns.

Data

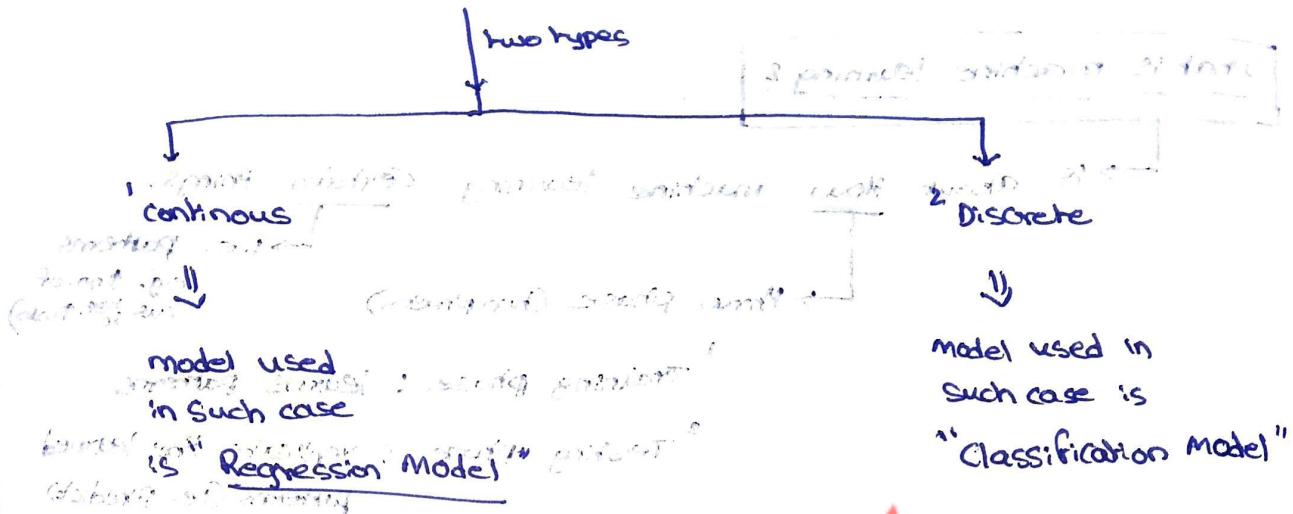
$x, y$

Independent variable  
(feature / cause.)

Target variable / dependent variable  
outcome / effect.

↳ y: is pain area of the client, problem of client.  
[is our objective, get from directly from client]  
↳ i.e. requirement.

↳ y / target variable / pain area



Model and formal model methods for what?

→ Model

$\rightarrow$  **Regression model**

→ classification model

```

graph LR
    A["based on training data"] --> B["(Given/not)"]
    B --> C["Supervised Learning Model  
X, y given."]
    B --> D["Unsupervised Learning Model  
X given; y not given."]

```

The diagram shows a flowchart starting with "based on training data". An arrow points to a box labeled "(Given/not)". From this box, two arrows branch out: one to the left leading to "Supervised Learning Model X, y given.", and one to the right leading to "Unsupervised Learning Model X given; y not given."

## Estimation vs prediction

## estimation

$$y = mx + c$$

$$y = b_0 + b_1 x$$

(See *Index*.) *Also see*



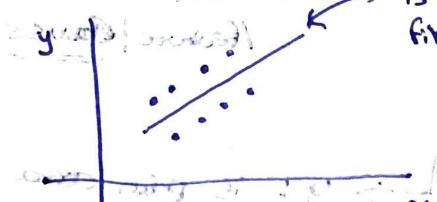
~~Hornbeam~~ p. 16-26

## Prediction

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

$$\hat{y} = b_0 + b_1 x + e$$

## Widowed women, etc.



is the best fit line.

3

linear regression  
line /  
prediction  
line.

at It is not prediction.

\* It is not estimation.

## Error

$\rightarrow \text{error} = \text{Actual} - \text{predicted}$

$$\rightarrow \text{So, } \text{Error} = \frac{1}{\text{no. of observations}} \sum_{i=1}^n (\text{Actual} - \text{predicted}) = \sum_{i=1}^n (y - \hat{y})$$

$$\text{from above all in eqn } y = b_0 + b_1 x + e.$$

but in reality, it is  
Sum of Squared Error (SSE)

$$\text{i.e. } \sum_{i=1}^n (y - \hat{y})^2$$

Error/SSE can be negative, 0, positive.

## How to find the best fit line?

(Build the Model)

using: minimum SSE (min. Sum of Squared Error)

finds where the line with min. SSE.

$\hat{y}_i = b_0 + b_1 x_i$   
method used (Build the Model)

① Statistical way      ② Machine learning way



use OLS  
(Ordinary least squared)

Method used

use Gradient  
descent.

$$d = 0$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

$$\hat{y}_i = \bar{y}$$

$$b_1 = \frac{\text{cov}(x, y)}{\text{cov}(x)} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

purpose/when to use:

- To test if there is a linear relationship between variables.
- To find significant variable.
- To check if regression assumptions are met.
- To know performance of the model.
- (Most widely used)

purpose/when to use:

Test performance  $\Rightarrow$  Math use by  $R^2$ .

→ baseline / benchmark: Average model.

$\hookrightarrow$  The model with no independent variables

$\hookrightarrow$  so it is the worst model you can have ever.

White noise or random error

Using mnemonic:

$\hookrightarrow$  3 variations

1 Sum of squared errors (SSE)

2 Sum of Squared Regression (SSR)

3 Sum of Squared Total (SST)

More benefit from 3rd, add information to model

$\hookrightarrow$  Now, we use Stats way, formulaic already derived:

$$b_0 = \bar{y} - b_1 \bar{x}$$

$\bar{x}, \bar{y} \rightarrow$  are means

$$b_1 = \frac{\text{Cov}(x, y)}{\text{Var}(y)}$$

From formula: Average Model  $\rightarrow b_1 = 0 \Rightarrow$

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

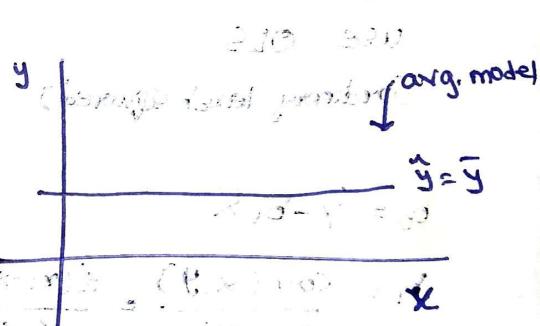
$$\hat{y} = b_0 + (0)x$$

$$\hat{y} = b_0$$

$$\hat{y} = \bar{y} - b_1 \bar{x}$$

$$\hat{y} = \bar{y} - (0)\bar{x}$$

$$\boxed{\hat{y} = \bar{y}}$$



Mean est method (squares)

no. of observations

$$SSE = \sum_{i=1}^{n_{\text{obs}}} (\text{Actual} - \text{predicted})^2 \quad [\text{model fails to explain}]$$

due to randomness, noise, influence beyond model scope

based on  $(y_i, \hat{y}_i)$  no. of observations

$$SSR = \sum_{i=1}^{n_{\text{obs}}} (\text{predicted} - \text{Avg. Model})^2 \quad [\text{how the independent variables contribute to predicted dependent variable}]$$

$$SST = \sum_{i=1}^{n_{\text{obs}}} (\text{Actual} - \text{Avg. Model})^2 \quad [\text{variation of } y]$$

$\hookrightarrow$   $SST = SSR + SSE$   
 total variance  $\rightarrow$  explained by known factor  
 unexplained variance  $\rightarrow$  unknown factor

$$R^2 = \frac{SSR}{SST} \quad \text{and} \quad R^2 = 1 - \frac{SSE}{SST}$$

$R^2$  of total variation  
 explainable by  
 Explainable Variation  
 Caused by regression.

$R^2 \rightarrow$  Square of Pearson Coeff. ( $r$ )  
 i.e.  $y_{\text{pred}}$  and  $y_{\text{actual}}$

$\hookrightarrow R^2 \in [0, 1]$  as  $-1 \leq r \leq 1$

①  $R^2 = 0$  when  $SSR = 0 \rightarrow SST = SSE$ .

↓ worst model [also note:  $\hat{y} = \bar{y}$  (Avg. model)] as predicted = Actual for  $SST = SSE$

②  $R^2 = 1$  when  $SSR = SST \rightarrow SSE = 0$

↓ best model

is that so? i.e.  $SSE = 0$

↳ impossible in reality as error  $\neq 0$ .

[comes concept of overfitting (good fit, underfit)]

One of assumption of linear

regression  $\rightarrow$  strong correlation b/w  $X$  &  $Y$

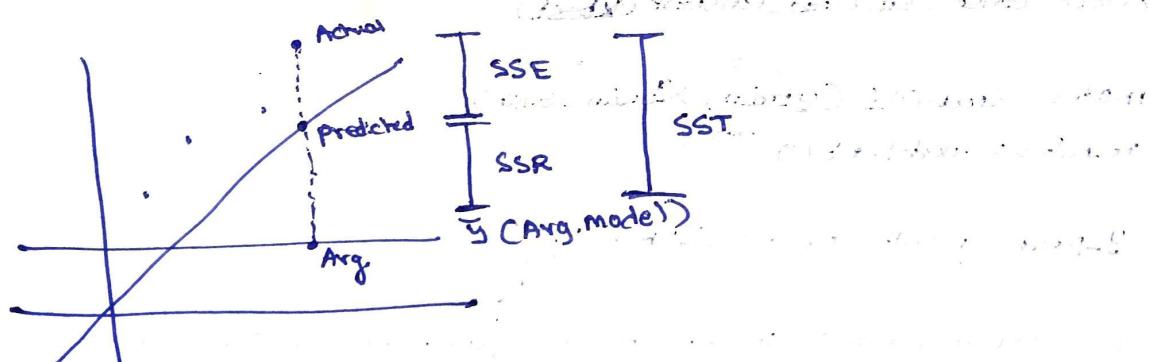
(any factor) (strong) (else bring down the model performance)

inline with  $R^2=0$  as worst model

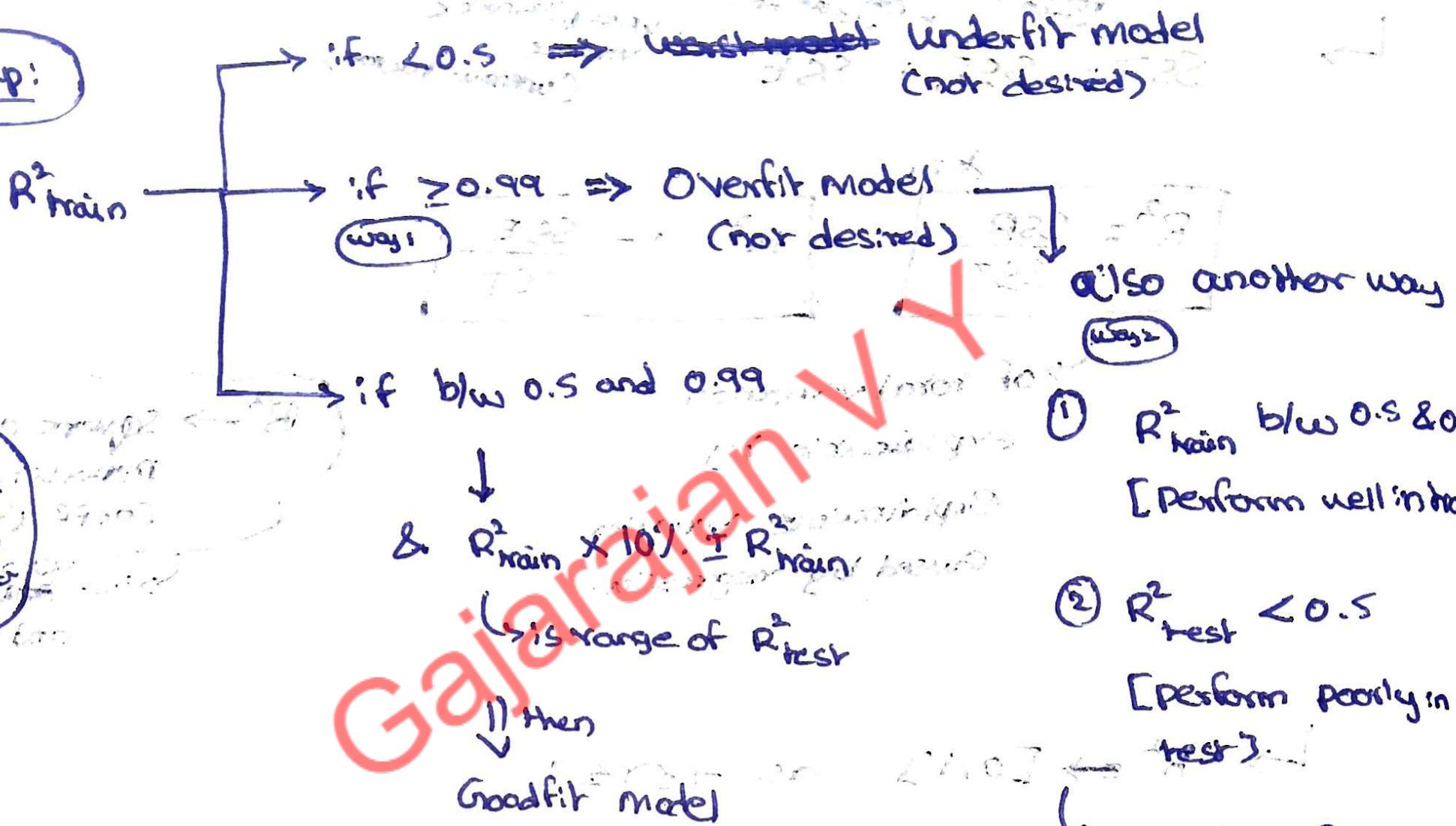
and  $R^2=1$  as 'best' model (actually overfit)

↑ use  $R^2_{\text{train}}, R^2_{\text{test}}$  & comparison b/w them.

↑ other way also exists using RMSE:



recap:



$$R^2_{\text{train}} = 0.5$$

$$R^2_{\text{test}} = 0.49$$

## SL-1 Day 2 Sum up

### Model building - Six Step Process

1. Split  $X$  (independent variable) and  $y$  (dependent variable)
  - ↓ then
2. Split each  $X$  and  $y$  as train and test
  - ↓ then
3. Initialize Model / Algorithm (e.g. Lr, Randomforest, etc)
  - ↓ then
4. Train Model / Learn patterns
  - ↓ then
5. Test Model / predict (replicate patterns)
6. Check Performance [use RMSE and/or  $R^2$  CAPE]
  - ↓ then

### performance checking

1. Error =  $(\text{Actual} - \text{predicted})$

2. Squared error =  $(\text{Actual} - \text{predicted})^2$

3. Mean Squared Error =  $\frac{1}{n} \left( \sum_{i=1}^n (\text{Actual} - \text{predicted})^2 \right)$  [ $n \rightarrow \text{no. of observations}$ ]

4. Root mean Squared Error =  $\sqrt{\frac{1}{n} \left( \sum_{i=1}^n (\text{Actual} - \text{predicted})^2 \right)}$   
[RMSE]

Note:

$R^2 \rightarrow \text{range } [0,1] \Rightarrow \text{so gives strength useful in saying good fit / overfit / underfit.}$

$\text{RMSE} \rightarrow \text{range } [0, \infty) \Rightarrow \text{so not gives strength so, hardly used in saying good fit / overfit / underfit.}$

$\text{SMAPE} \rightarrow \text{any arbitrary value}$  used in saying good fit / overfit / underfit.

$\text{MAPE} \rightarrow \text{any arbitrary value}$  used in saying good fit / overfit / underfit.

$R^2$  & RMSE  $\Rightarrow$

$RMSE \xrightarrow{\text{tend to } 0} 0 \Rightarrow \text{error} \xrightarrow{\text{tend to } 0} 0 \Rightarrow R^2 \xrightarrow{\text{tend to } 1} 1$

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

Good fit vs Overfit vs Underfit

$R^2_{\text{train}} \& R^2_{\text{test}}$

$R^2$ value	fit type of model	additional criteria / guideline (not rule)
$R^2_{\text{train}} = 1$	Overfit	$R^2_{\text{train}} - R^2_{\text{test}} \leq 10\%$
$R^2_{\text{train}} = 0$	Underfit	$R^2_{\text{train}} - R^2_{\text{test}} \geq 10\%$
$R^2_{\text{train}} \in [0, 0.5]$	Underfit	$R^2_{\text{train}} - R^2_{\text{test}} \geq 10\%$
$R^2_{\text{train}} \in [0.5, 0.99]$	Good fit	$R^2_{\text{train}} - R^2_{\text{test}} \leq 10\%$

Difference b/w  $R^2$  value of train prediction and test prediction is ~~is~~ not more than 10%.

Model fit type	Definition
Overfit	<ul style="list-style-type: none"> <li>* Model perform very well in training, but very poor in testing.</li> <li>* 2 conditions           <ul style="list-style-type: none"> <li><math>R^2_{\text{train}}, R^2_{\text{test}} \approx 1</math></li> <li>model perform very well in training data but very poor in testing data.</li> </ul> </li> </ul>
Good fit	<ul style="list-style-type: none"> <li>* Model perform very well in training data and also perform very well in testing data.</li> </ul>
Underfit	<ul style="list-style-type: none"> <li>* Model perform very poor in training data and also perform very poor in testing data.</li> </ul>

$R^2_{\text{train}}$  Condition check

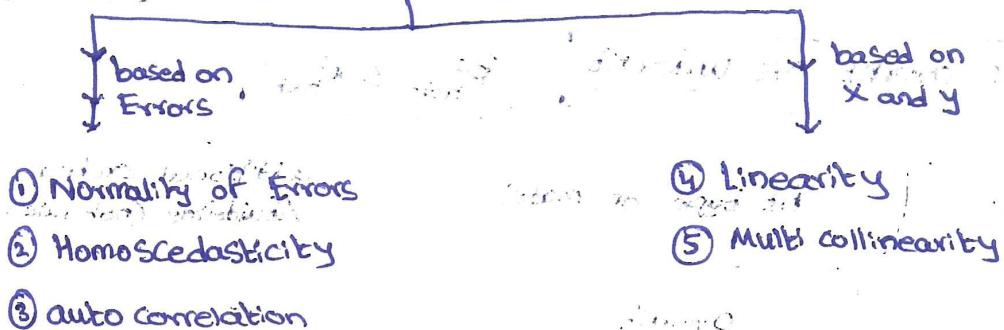
if only for  
 $R^2_{\text{train}} \in [0, 0.99]$

$R^2_{\text{test}}$  Condition of 10% check

Good enough  
Good enough

## Assumptions of linear Regression

### Types of Assumptions



NHA  
LM

	$x_1$	$x_2$	...	$y_{actual}$ (A)	$y_{predicted}$ (P)
R <sub>1</sub>	-	-	-	-	-
R <sub>2</sub>	-	-	-	-	-
R <sub>3</sub>	-	-	-	-	-

1) Normality of errors  
 2) Homoscedasticity  
 3) auto Correlation  
 4) Linearity  
 5) Multi collinearity

Normality

Homoscedasticity

auto Correlation

Linearity

Multi collinearity

error (A-P)

e<sub>1</sub>

e<sub>2</sub>

e<sub>3</sub>

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

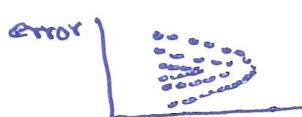
;

;

;

## \* Heteroscedasticity

↳ i.e. unequal variance for errors. [i.e. pattern forming]



fitted/predicted value

fig. unequal variance  
of convergent pattern



fitted/predicted value

fig. unequal variance  
of divergent pattern

## ③ Auto Correlation

↳ i.e. dependency/correlation b/w errors.

[i.e. follows certain patterns]

$e_1 = 100$   
 $e_2 = 50$   
 $e_3 = 0$   
 $e_4 = -50$

$\rightarrow$  So, pattern  $e_i = e_{i-1} - 50$



Example: If we follow same data set, predicted value will be different

↳ Test:

Rainbow test, Durbin-Watson

Key idea for assumptions based on error

How Model Learns to Predict? (using i.e. generalized linear model)

Let's say (y)  $\rightarrow$  based on pattern (in) data, output (y-hat) is sent to function (f), which will be used (and passed) to

so, then if pattern in error, it learns it

↳ So, pattern  
in error(s) is  
not preferred  
or must be  
avoided

$\downarrow$  lead to

reduce errors (i.e. most of time Actual = predicted)

$\downarrow$  case of

Overfit [not learn pattern, just by hands]

## ④ Linearity

↳ i.e. High relationship/correlation / linearly related

blue X, and Y.  
(Independent Variable)  
(Dependent Variable)

↳ Test: Rainbow test

Other test: ADF, LM, Breusch-Pagan test, etc.

## ⑤ Multicollinearity

↳ i.e. High correlated / dependency / related b/w

Independent variable(s).

↳ e.g.

$$X_3 = X_1 + X_2$$

$$\text{or} \quad X_3 = X_4 - X_5$$

$$X_3 = X_6/X_7$$

↳ Test:

COND NO. (Condition number)

VIF test

COND NO.	
>1000	High Multi Collinearity
blw 1000 and 100	Moderate Multi collinearity
<100	Low multi. collinearity

for this only using  $(N-1)$  dummies encoding

\* If High Multi collinearity (or in general collinearity), then,  $\beta$  value of the dependent variables (in form of 'independent variable') [e.g.  $X_3$ ] is very high than of the actual independent variable [e.g.  $X_1$ ].



↳ Consider if means  $X_3$  is given more importance (as contributor to predict) but it is not the actual case ( $X_1$  &  $X_2$  need to be given more importance instead as least importance by the Model built with Multicollinearity).

extra points

\*  $R^2_{\text{test}} \geq R^2_{\text{train}}$  (at least, possible cases by chance)

↳ when the best data pattern is the most learned by the Model in training phase from train data.

\* 1<sup>st</sup> check Good fit / not then only 2<sup>nd</sup> check model performance

\* OLS given linear regression with min. SSE can be worst (as it's just min SSE but not may good model), so compare with Avg. model using  $R^2$ .

- \* Null hypothesis (usual thing / status quo) is mutually exclusive to Alternative hypothesis.

Gajarajan ✓ ✓

## SL-1 Day 3 Sum UP

- 1 Durbin Watson
- 2 VIF test,
- 3 Q-Q plot
- 4 Error
- 5 Bias error
- 6 Variance error
- 7 error = 0 ?
- 8 BE & VE
- 9 BE, NE, overfit, underfit.
- 10 Regularization
- 11 2 types of Reg.

### Durbin Watson

- test for autocorrelation.
- Range [0, 4]
- If in bw 1.5 and 2.5  $\Rightarrow$  Autocorrelation exists.

### VIF test

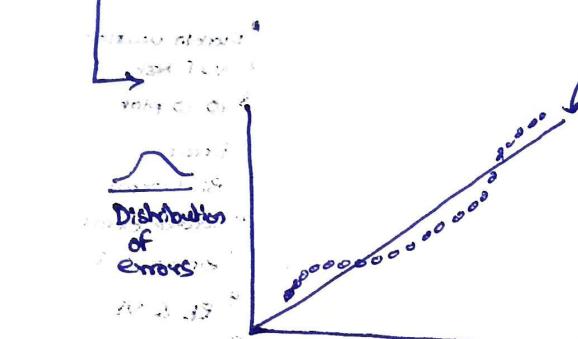
#### 1. Multi Collinearity

- If exist in dataset  $\Rightarrow$  check with help of COND. NO.
- If exists in which variables / features  $\Rightarrow$  check with help of VIF test.
- | VIF    | Status                     |
|--------|----------------------------|
| $< 5$  | low multicollinearity      |
| $5-10$ | Moderate Multicollinearity |
| $> 10$ | Severe Multicollinearity   |

Next step  
Keep & good to go  
remove that variable  
feature & again  
Check VIF test.
- $$VIF = \frac{1}{1-R^2}$$
- High multicollinearity  $\Rightarrow R^2 \rightarrow 0$   $\Rightarrow VIF \rightarrow \infty$ .
- VIF is "Variation Inflation Factor"
- only for numeric columns.

## Q-Q plot

→ Tests the normality of errors/residuals.



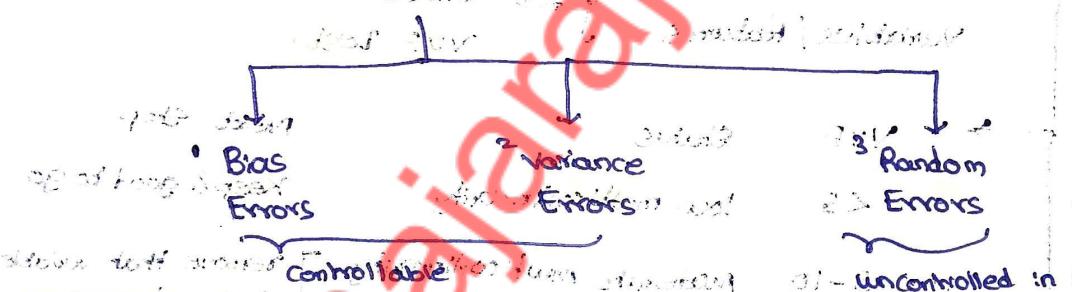
This line represents both are same.  
And any deviation from it indicates the dist. of errors is not normally dist.

↳ Quantile - Quantile plot.

## Errors

→ is (Actual - predicted)

→ 3 types (Mainly)



## Bias Errors

↳ is the inability of the model to predict correctly right data

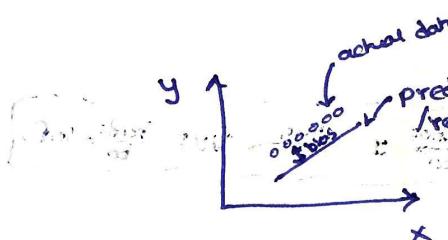


fig: low bias

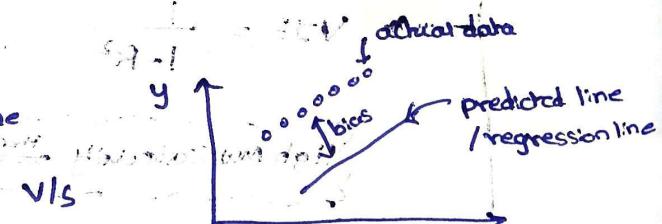


fig: High bias

↳ Sources of bias:

1 Collection of data (e.g. sample, may miss values)

2 Feature Selection (e.g. remove w/o domain knowledge)

3 Feature engineering (e.g. intro. noise may be)

eg. missing values, outliers

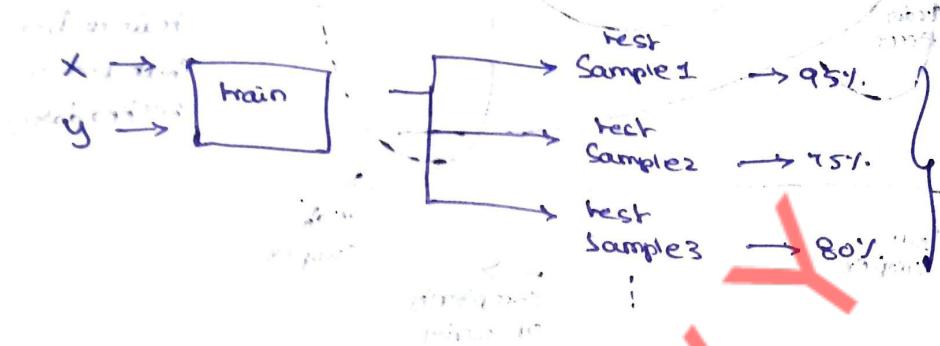
4 Introduced by model  
(if more data is on)

one category → learn pattern well  
only to that category

5 Introduced by Data科學 (EDA)

## Variance Error

- is the difference in the performance of the model for different samples.
- So inconsistency in performance.
- model with low variance is preferred. \* (bias comes in pic.)



Variance error  
(as diff.  
exists in  
performance  
across diff.  
samples)

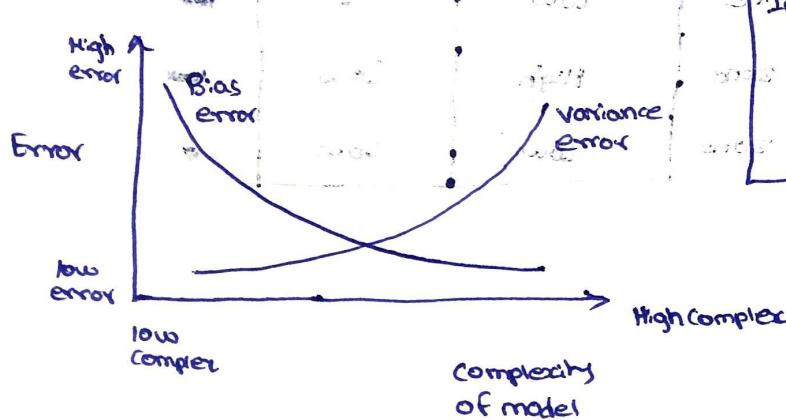
## Random Error

- is of random in nature. [So, we are not known about it]

## Can Errors = 0?

- Errors = Bias Error + Variance Error + Random Error
  - ↳ can be reduce, but not 0.
  - ↳ can be reduced, but not 0.
  - ↳ uncontrolled.
- So Errors ≠ 0

## Bias Error and Variance Error



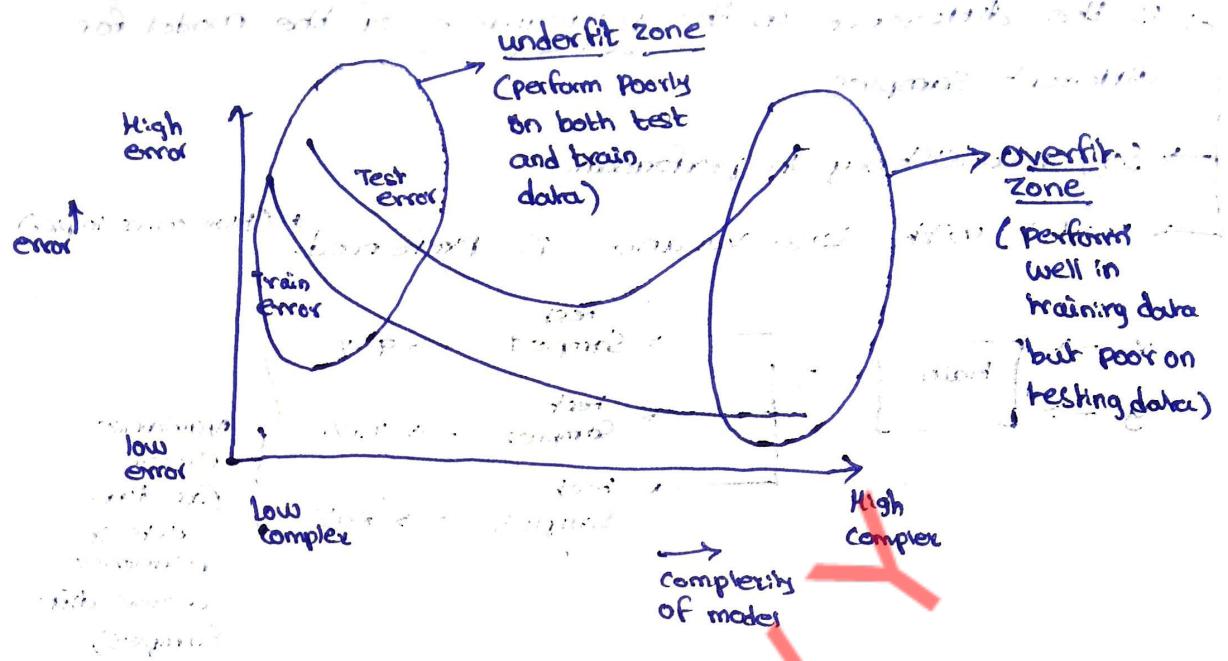
### Inference 1:

- ↑ complexity  $\Rightarrow$  ↑ Variance error
- ↑ complexity  $\Rightarrow$  ↓ Bias error

- Inference 2: seasons Y/n.
- ↑ Bias error  $\Rightarrow$  ↓ Variance error
- ↓ Bias error  $\Rightarrow$  ↑ Variance error

## Overfit and underfit

(based on 'Train error & Test error')



## Bias Error, Variance Error, Overfit, underfit

is base / foundation for 'Regularization'.

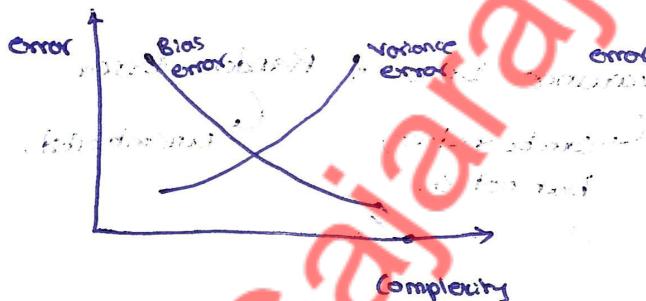


Fig 1

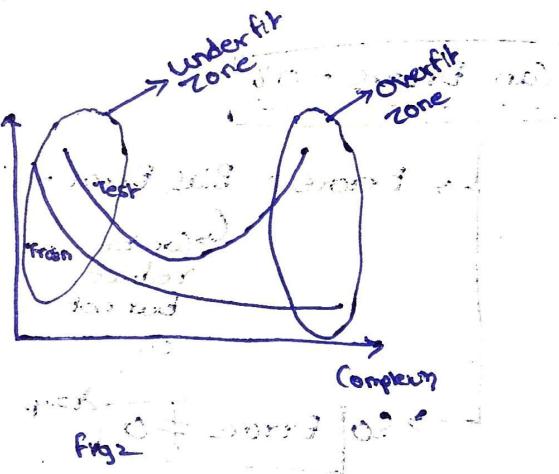


Fig 2

From fig 1 & fig 2,

fit zone	Bias	Variance
Overfit zone	Low	High
underfit zone	High	Low
Cross fit zone	Low	Low

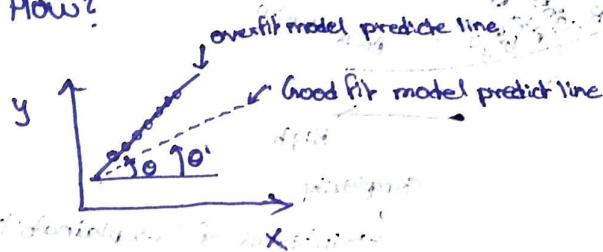
## Regularization

→ is overfitted Model (Low Bias, High Variance)

↓  
to

Good fit model

→ How?



Key thing of  
Reg. is

↑ bias & ↓ variance  
error  
by ↓  $\beta$

\* So, it is done by introducing some bias (i.e.

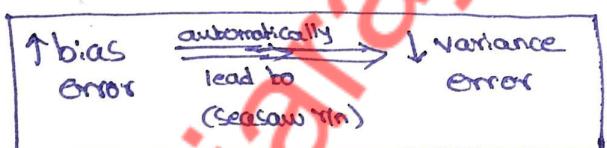
The inability of model to predict correctly (right data)

↓ How to introduce bias, then?

↓ By shrinking/reducing the coefficient of the independent variables

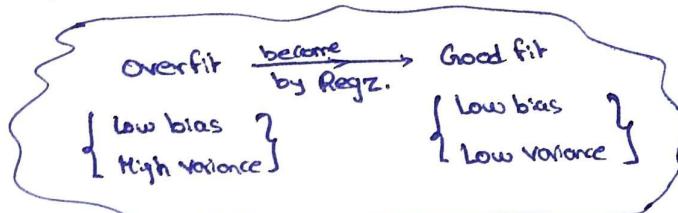
By Shrinking / reducing the coefficient of the independent variables

Flow:



By Shrinking / reducing Coefficients

Regularization → is all about Shrinking coefficients.



## Types of Regularization

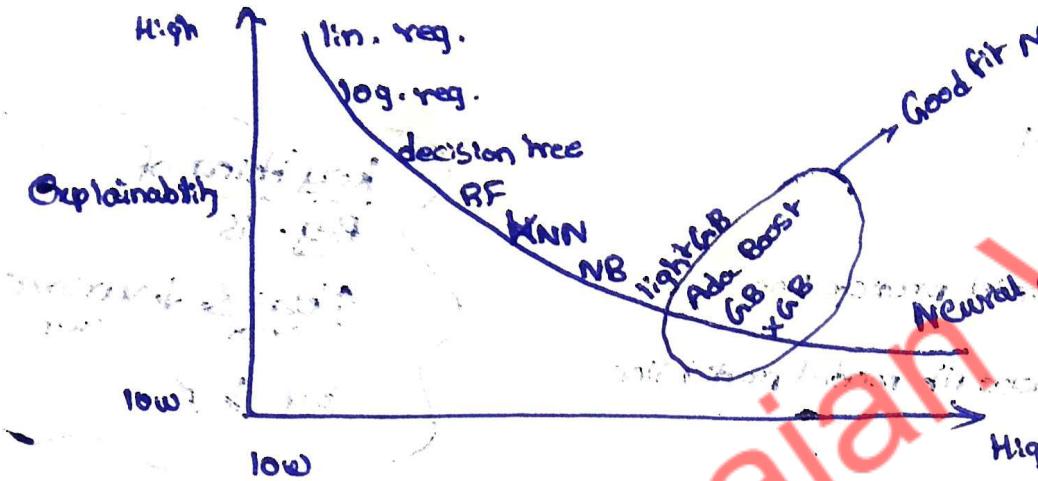
3 Main type

Ridge (Some  $\beta$  becomes near to 0)

Lasso ( $\frac{\text{some}}{\beta}$  becomes 0)

Elastic Net (Some  $\beta$  becomes 0 and some  $\beta$  becomes near to 0)

## Model and complexity and explainability



Good fit model {low bias & low variance}

(Simple models)

Complex models

Complexity

↳ in terms of 'explainability'

$\rightarrow^2$  No. of parameters

Extra

① Scaling only for model based on 'weights' &  $\sigma^2$  distance

② Scaling and Transformation can be done of target variable.

③ DON'T SCALE for target variable

DON'T Impute missing value for target variable

## SL-1 Day 4 Sum UP

### Regularization

→ By shrinking  $\beta$  using penalty factor.  
also trade off b/w SSE &  $\beta$ s balanced\*.

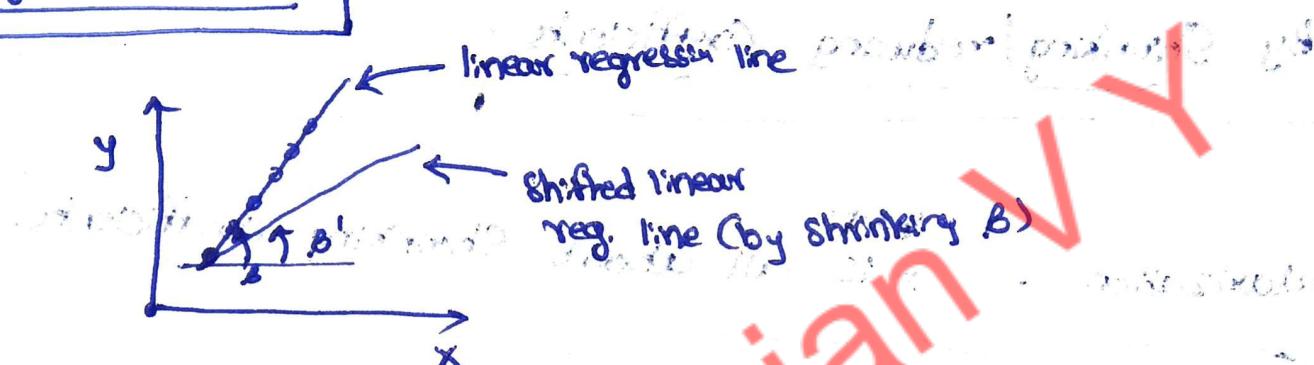


Fig:

\* Regularization: [aim: overfit model to Good fit (if done excessively then it may become underfit)]

Shrink/reduce  $\beta$

↓ so

T Bias (indirectly ↓ Variance)

↓ so

↓ Standardized  $\beta$  of each  $x$

(Goodfit line and  $\beta$  don't depend on  $x$ )

what?

$x^{up}$

Prerequisite:

Scaling of independent variables.

Cost function	Type of regularization	
$\sum (y_i - \hat{y}_i)^2 + \lambda (b_0^2 + b_1^2 + \dots)$	Ridge regularization	near zero <sup>*</sup> and never zero
$\sum (y_i - \hat{y}_i)^2 + \lambda  b_0  +  b_1  + \dots$	Lasso regularization	zero/not zero
$\sum (y_i - \hat{y}_i)^2 + \lambda_p (b_0^2 + b_1^2 + \dots) + \lambda_L (b_0 + b_1 + \dots)$	Elastic Net regularization	$L_1 \text{ ratio} = \frac{\text{llasso}}{\text{llasso} + \text{ridge}}$
		some zero some near to zero

\* So, any regularization,

Trade off b/w ① SSE (sum of square error) [need to be min.]  
 ②  $\beta/b$  (coeff) [need to be reduced]

↓ so

min(SSE) + min(penalty) → so, min. cost function.  
 [aim of Regularization]

$(y - \hat{y})^2 + \lambda (\text{betas})$  → penalty factor / regularization factor / alpha

↓ the first part: sum of square errors

why trade off exists?

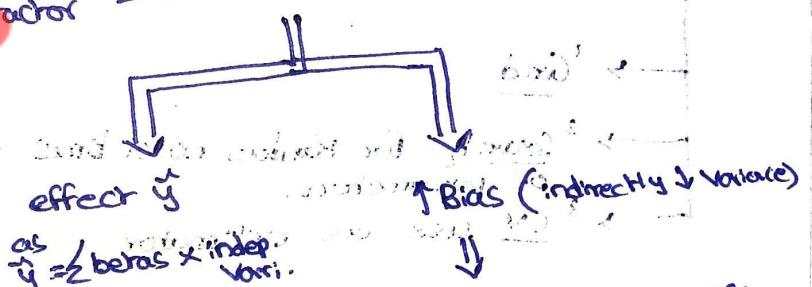
not having enough info ↑ alpha/penalty factor

↓ Betas

### Tradeoff in SLR

① Bias & Variance

② SSE &  $\beta^T \beta$

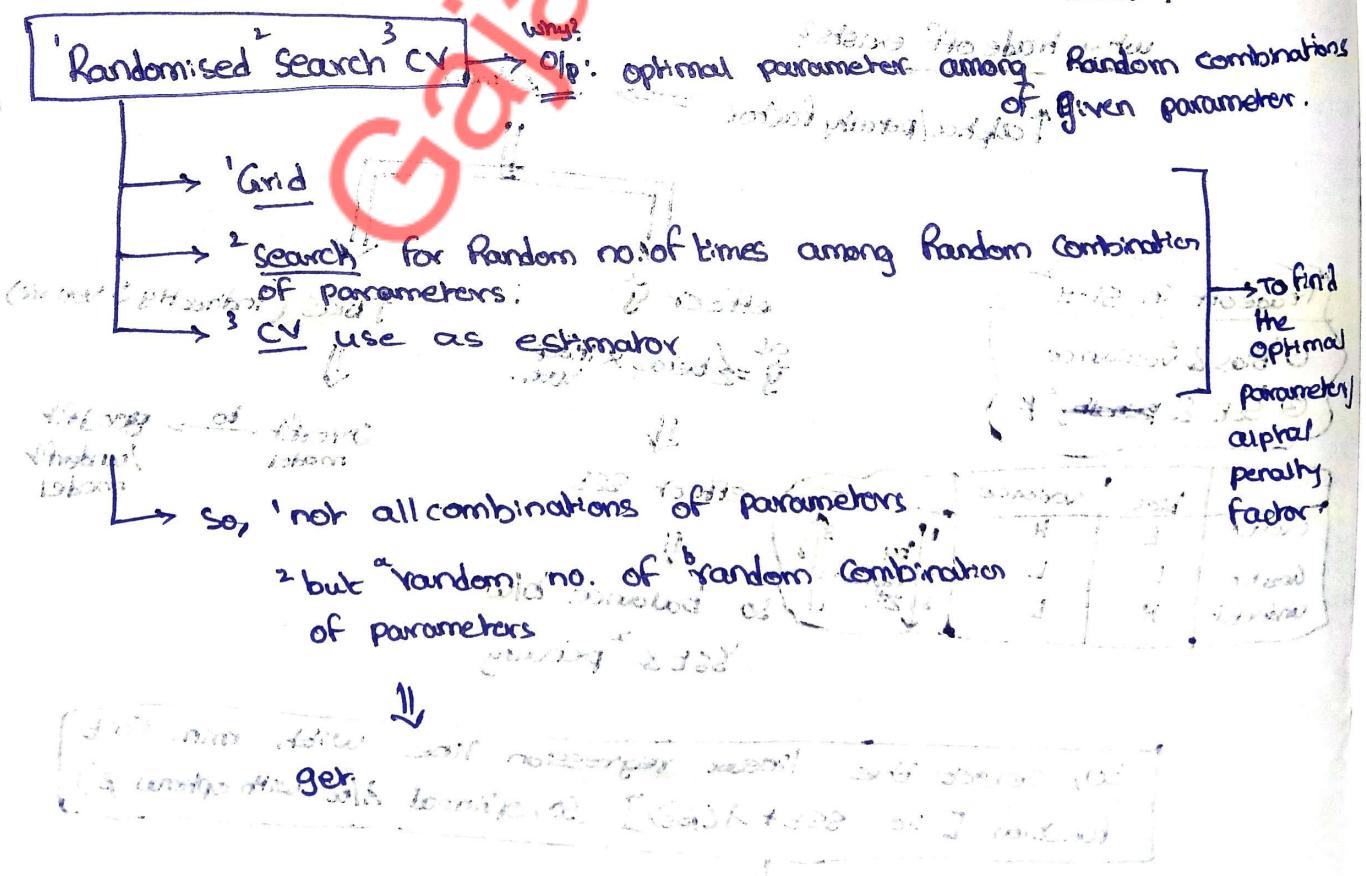
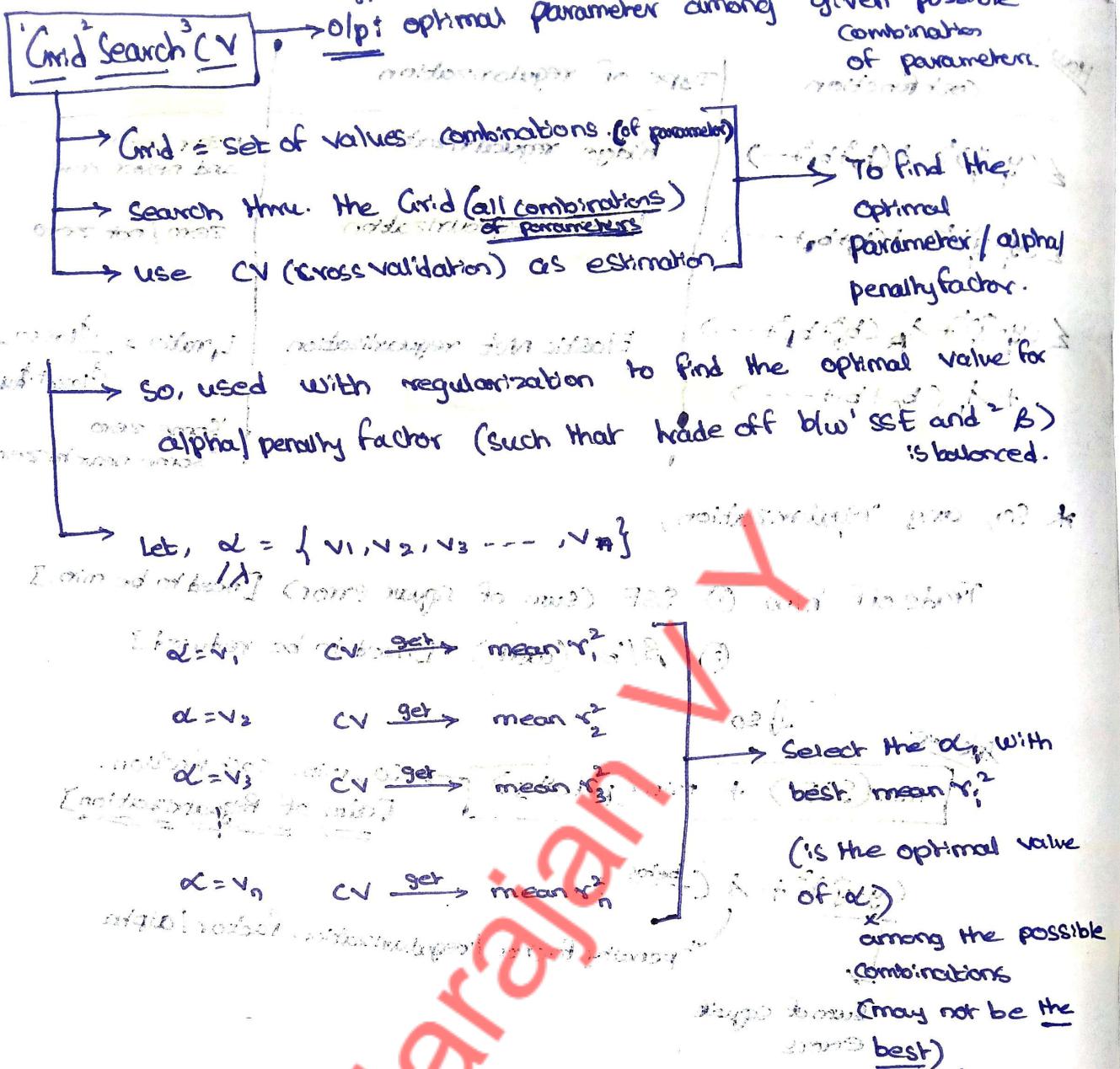


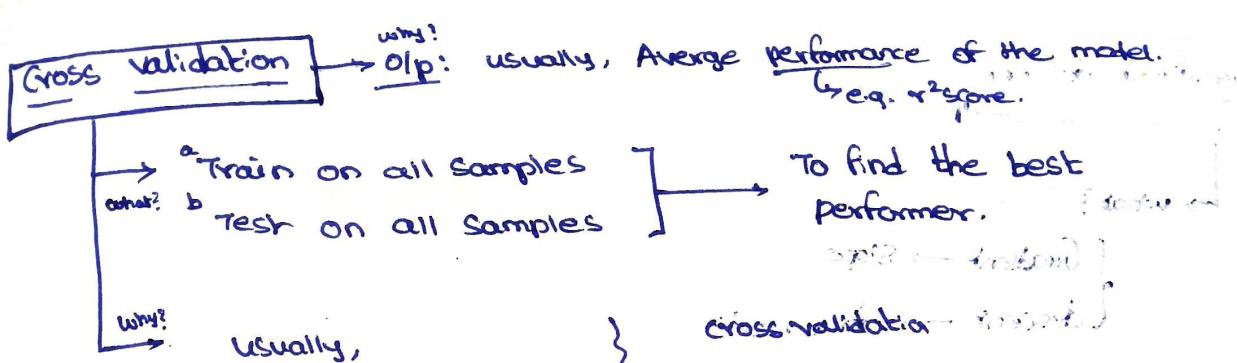
	bias	Variance	
Overfit	L	H	↑ by regression
Good fit	L	L	2.1
Underfit	H	L	3.2

so balance b/w them and SSE &  $\beta^T \beta$

So, select the linear regression line with min. cost function [ $= SSE + \lambda(\beta^T \beta)$ ] (so, optional d/d with optimal  $\beta$ )

Then How?





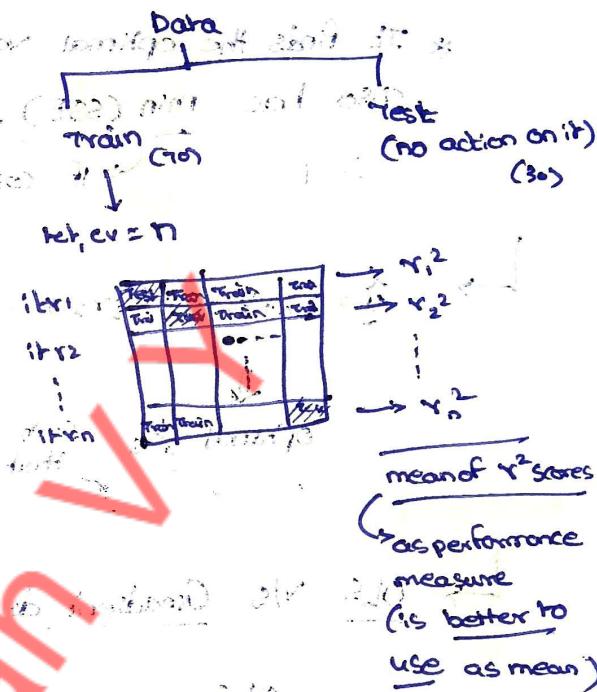
Data

Train (70)

Test (30)

model built on  
it may miss  
to learn some  
patterns (due to  
not much data on  
it)

instead in  
 performance  
 measure.



Similarly, Cross-val-score, kfold cross validation.

→ i.e.  $CV = K$  (300)  $\times 100$

→ i.e. split data into  $K$  folds.

For each iteration, if Shuffle = True,

→ then Shuffle

→ then Split

→ then

## Regularization use cases

- 1 Overfit model  $\xrightarrow{\text{to}}$  Good fit / underfit model.
  - 2 reduce multicollinearity as reducing  $\beta_{(S)}$  values.
  - 3 use Lasso regularization for feature Selection. (as  $\beta_{(S)}$  can be zero)

## Gradient descent

↓ what?  
 Gradient → Slope  
 descent → downwards

\* It finds the optimal value for  $\beta(s)$  / coeff. Such that it also has  $\min(\text{SSE})$ .

→ if cost func. =  $\min(\text{SSE})$  [w/o regularization]

↓ if no reg.

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2 + \dots$$

Optimal Bias

such that

$\min(\text{SSE}) \Rightarrow$  Best fit line

↓ OLS vs Gradient descent

### OLS

1. find optimal  $\beta(s)$  such that  $\min(\text{SSE})$

↓ slow & less accurate

why  
GD over OLS?

1. works best only for small dataset.
2. slower in processing.

3. Ancient / Old Statistical way. (before GD)

↓  $b_0 =$

$$b_0 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

↳ ↑ cost<sub>x</sub> for larger dataset  
 with more no. of rows & columns.

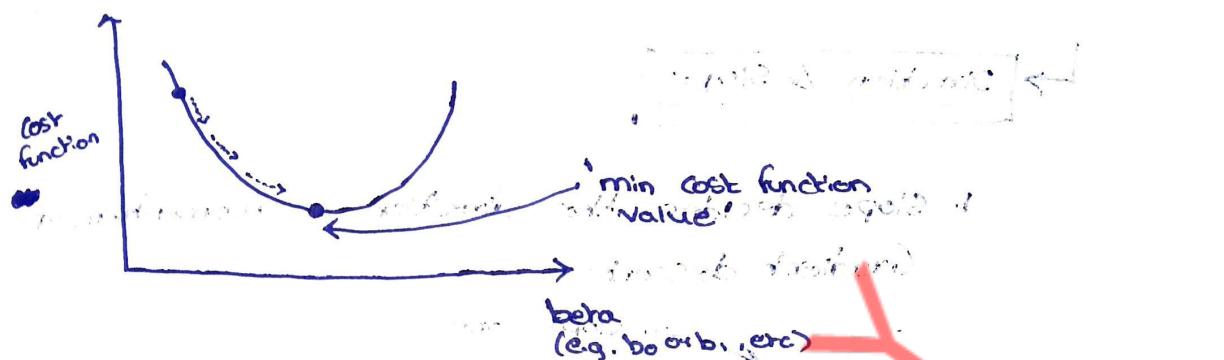
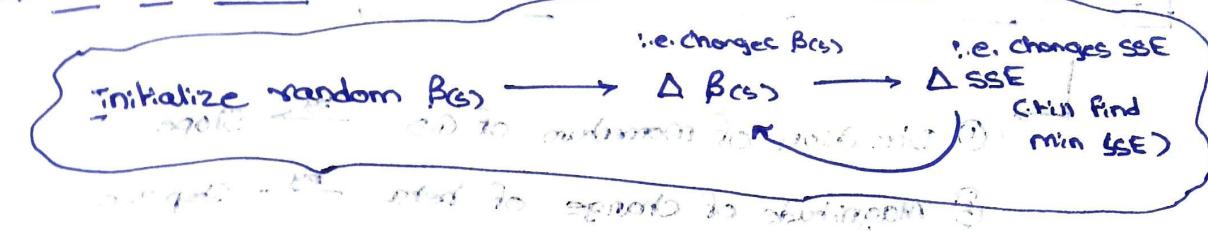
### GD

1. find optimal  $\beta(s)$  such that  $\min(\text{SSE})$  [Both OLS, GD same function]

1. works best even for large dataset. (as converges fast)

2. faster in processing.

↳ How it works?



Note:

1. cost function can be  $\sum (y - \hat{y})^2$  [i.e. SSE] or  $\sum (y - \hat{y})^2 + \lambda (\text{beta})$  [i.e. SSE + penalty for beta]

Penalty  $\rightarrow$  linear regression with regularization (Lasso/Ridge/Elasticnet)

2. cost function never zero.

as Error = Bias error + variance error + Random error

Error  $\neq 0$ .

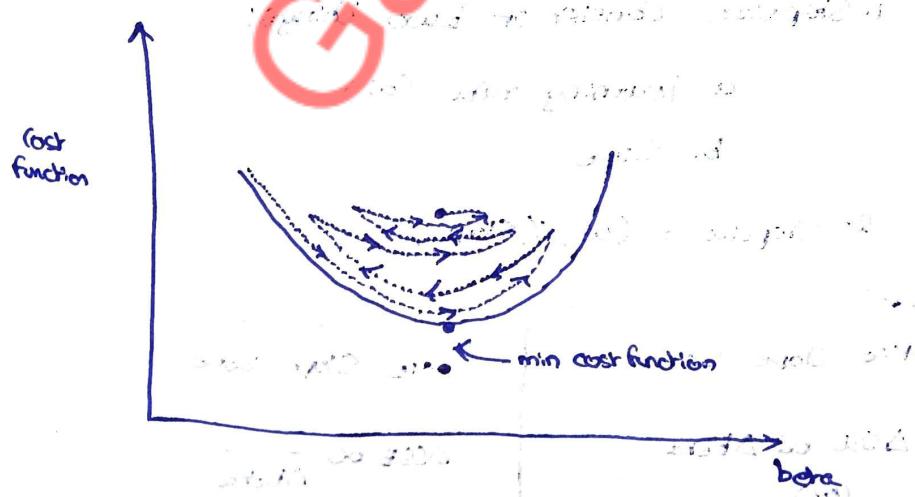


fig: process of working of GD (in 2 variables)

start with some  $\beta$ s

is similar to the ball oscillate in a 'V' shape back & front and finally settles.

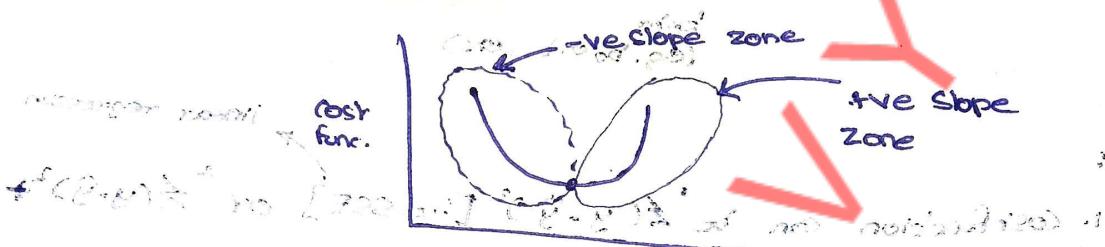
## Factors responsible for working of GD

Cost function, learning rate, step size

- ① Direction of momentum of GD  $\rightarrow$  Slope
- ② Magnitude of change of beta  $\rightarrow$  Step size

### Direction & Slope

1. Slope decides the direction of momentum of Gradient descent.



2. How?

-ve slope  $\xrightarrow{\text{indicates}}$  direction of reducing error/cost func.

tive slope  $\xrightarrow{\text{indicates}}$  direction of increasing error/cost func.

step size of momentum is constant & moving left = moving up & moving right = moving down

### Stepsize

1. Step size consists of two things:

- a. learning rate ( $\alpha$ )
- b. Slope

$$2. \text{Step size} = (\alpha) \times (\text{Slope})$$

+ve slope zone

-ve slope zone

$$\Delta SSE \propto \Delta \theta_0 \\ (\text{or})$$

$$\Delta SSE \propto \frac{1}{\Delta \theta_0}$$

$\Delta(\text{cost func.}) \propto \Delta \theta_0$  to prevent to jumping into local minima

$$\Delta SSE \propto -\Delta \theta_0$$

(or)

$$\Delta(\text{cost func.}) \propto -\Delta \theta_0$$

↳ So, to find optimal beta value

$$\text{new beta} = \text{old beta} - \alpha \times \text{Slope}$$

e.g. new  $b_0 = \text{old } b_0 - \alpha \times \text{Slope}$   
 new  $b_1 = \text{old } b_1 - \alpha \times \text{Slope}$

logic behind it working:

-ve slope zone  $\Rightarrow$

Slope is -ve  $\Rightarrow$

$T_{\text{newbeta}}$

④ Ave Slope zone  $\Rightarrow$

Slope is ave  $\Rightarrow$

$\downarrow \text{newbeta}$

Adip vrtsp  
Formula:

$$b_0 = -(y - \hat{y})$$

$$b_1 = -(x)(y - \hat{y})$$

↳ learning rate ( $\alpha$ )

Early 2014

After 2014

① fixed  $\alpha$

① dynamic  $\alpha$  based

② Need manual intervention for precision.

on min. cost function / min. SSE

if

e.g. Adams algorithm

↳ what does 'back and front glide' indicates?

↳ precision is the answer.

cost func.

cost

## Types of Gradient descent (Mainly)

### ① Batch Gradient descent

→ 'Compute Gradient' and 'update Parameters'  
(slope) (beta, etc)

global  
min.  
only one  
correct  
problem

based on entire dataset



useful in small dataset & convex problem.

### ② Stochastic Gradient descent

→ 'Compute Gradient' and 'update Parameters'  
(slope) (beta, etc)

escape  
local  
min.  
in non-  
convex

based on each datapoint in the dataset.

useful in large dataset, & non-convex problem.

### ③ Mini-Batch Gradient descent

→ 'Compute Gradient' and 'update Parameters'  
(slope) (beta, etc)

based on subset of datapoints in the dataset

useful in stability & speed balancing

Everage parallel Computing

### When to Stop?

based on any one of two criteria achieving

1. max. no. of iterations  
(usually, 100)  
2. tolerance level.  
(i.e. precision based)

Where GD is used

Linear regression → cost func. is BASE

Regularization (Lasso, Ridge, elastic net) → cost func. is BASE + L1 or L2 term

Ensembling techn. (xGB, etc)

Neural network ( $\beta$ 's  $\rightarrow$  weights)

**B or P-value**

**B or Explanatory power**

→ it says for ~~per~~ 1 unit of change in  $x_i$  (feature) brings  $\beta$  units of change in  $y$  (dependent variable)

**P-value**

→ it tells about whether variable is significant or not.

$H_0$ : insignificant variable

$H_1$ : Significant variable

Significance is the contribution by that independent variable to target variable.

**Interaction effect**

→ is independent variable / features ( $x_1, x_2$ ) interact with each other (i.e. multiplication operator) to bring change in target variable ( $y$ ).

$H_0$ : No interaction effect

$H_1$ : interaction effect

→ i.e. based on p-value of the new created interaction variable  $[x_1 * x_2]$

→  $\neq$  Multicollinearity (i.e. High Correlation b/w independent variables)

**Measures**

**MAPE**

→ Mean Absolute percentage error (MAPE)

$$\text{Error} = A - P$$

$$\text{percentage Error (PE)} = \left( \frac{A - P}{A} \right) \times 100$$

$$\text{Absolute percentage Error} = \left| \frac{A - P}{A} \right| \times 100$$

$$\text{Mean absolute percentage error} = \frac{1}{n} \sum \left| \frac{A - P}{A} \right| \times 100$$

↳ RMSE is preferred over MAPE (as MAPE gets inflated with min. A and max. P)

### MAE

→ MAE or mean absolute error is sum of absolute errors.

→ Mean absolute error (MAE) =  $\frac{1}{n} \sum |A_i - P_i|$

→ less preferred to use.

### Adjusted R<sup>2</sup>

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

$n \rightarrow$  no. of observations / rows

$k \rightarrow$  no. of features

$$\text{Adj. } R^2 \leq R^2$$

↳ why Adj. R<sup>2</sup>?

↳ because of drawback of R<sup>2</sup>:

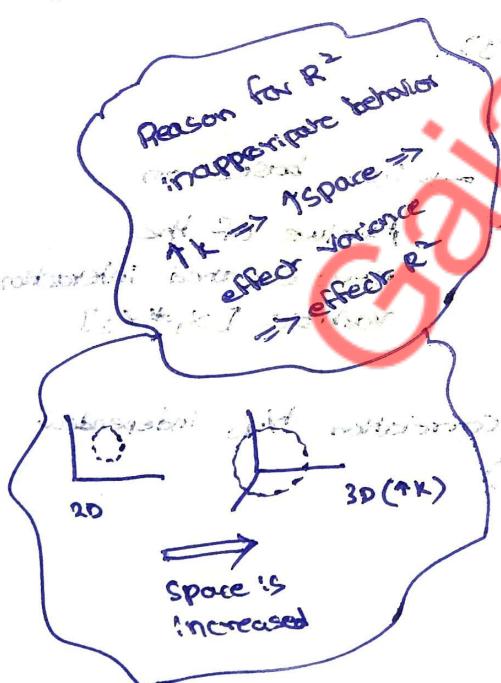
forward and backward selection

As  $n \gg k$  (no. of features), irrespective of which variable is important R<sup>2</sup> value increases telling that variable is important (but in reality may not be the case).

Reason for R<sup>2</sup> inappropriate behavior:

- 1D  $\Rightarrow$  1 space  $\Rightarrow$  effect variance
- 2D  $\Rightarrow$  2 effect variance
- 3D ( $\gg k$ )  $\Rightarrow$  effect variance

So, Adj. R<sup>2</sup> penalizes the value by  $(n-k-1)$  for increasing  $k$ .



(Forward) and (backward) selection

Y-axis standard

$\text{OOS}(\frac{S_{\text{test}}}{S_{\text{train}}})$  test and training

$\text{OOS}(\frac{S_{\text{test}}}{S_{\text{train}}})$  - model selection and validation