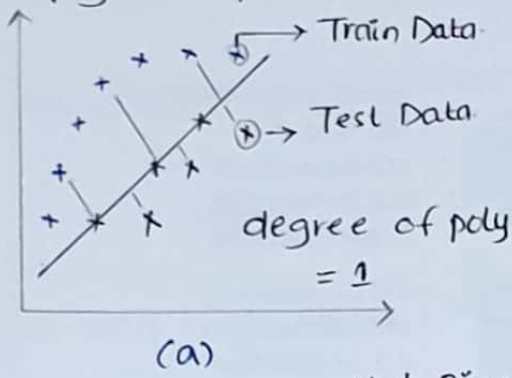


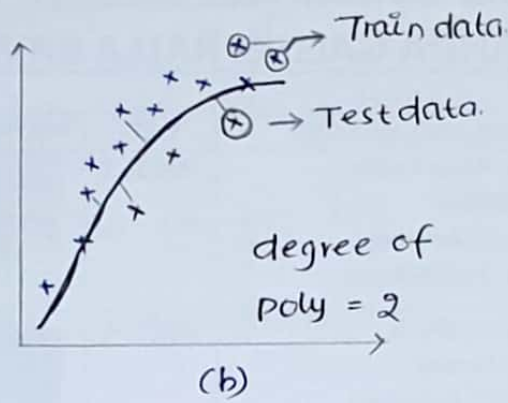
# Overfitting & Underfitting :-

①

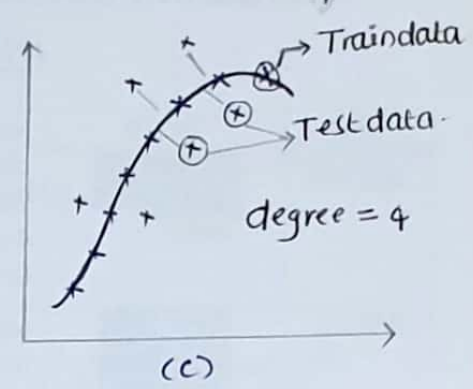
Regression :-



Under-fitting → High Bias  
→ High Variance



Best Model → low B.  
→ low V.



Overfitting → low Bias  
→ High Variance

Case - I :- from the above fig we draw a best fit line or Curve Using linear regression where we use polynomial degree, in fig (a) we apply degree = 1 its act like a Simple linear regr. and draw a line on training data points when we use Scatter plot graph, after it we calculate  $R^2$ -Error, if we observe the error will be high because the distance between the actual value and predicted value is more i.e. error rate will be more.

Note :- In case of training data itself the model gives us high error. i.e. "high bias" is called as "Underfitting" issue.

Case - II :- if we observe graph (c), where we applying degree = 4 its form a Curve which is best-fitted to all the training datapoints. So from this we can say On training data there is no error because Curve fitted good., if we have a new test data points come in the graph. then that makes some distance from Curve i.e. errors formed in model. here Our aim is to fit the Curve best, from that recorded low error from this we can say Accuracy of train data is more, but for test data.

is low because error-form in test data is called "Overfitting" ②

In case - i the Accuracy of training data is low (huge errors) as well the Accuracy of a test data also low (errors).

Here our aim is to be get high Accuracy in case of train data as well test data (ie prevent the errors)

if we observe the figure (b) i think its a best-fitted curve at Polynomial degree = 2 and where has no errors ie minute errors in train data as well test data. So i prefer this model, because this gives high Accuracy in train data, test data and error rate should be low as compare to other two models where overfitting & underfitting problem effected to the model.

if we observe model 2 that gives us low bias and low Variance. if in Underfitting case that gives high bias and high Variance ie that the error rate will be high in train data called as high bias. Similarly, if error rate will high in test data called as high Variance.

In case of Overfitting that gives low bias and high Variance ie the error rate in train data is low but error rate in test data is high.

Variance :- if your model changes a lot with a small change in

the training data this means the model has "High Variance"

the Issue with high Variance is that the model becomes very complex.

to avoid Complexity we have to choose Simple model.

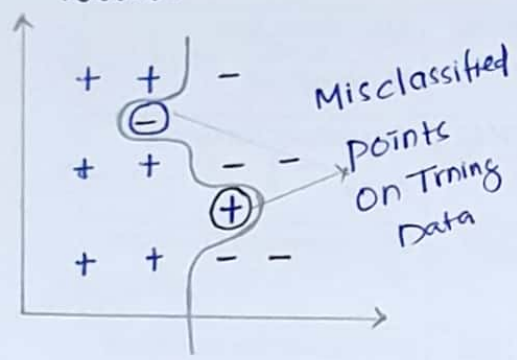


Occam's Razor theorem Says that try to learn Simple models but not too Simillar models.

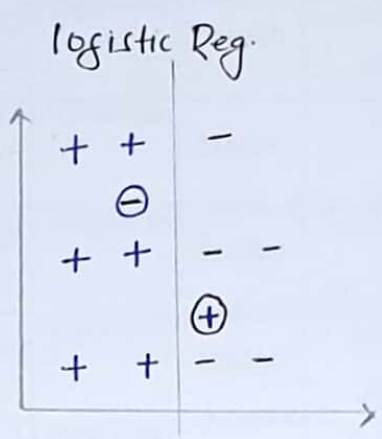
Reason behind choose Simple model is that going to be more "generalize" On test data (ie generalizable).

Note :- In Regression we classify the model as overfit or Under-fit based on "Errors" in training data as well in test data, Similarly In classification we classify based On misclassified data points in training data as well test data.

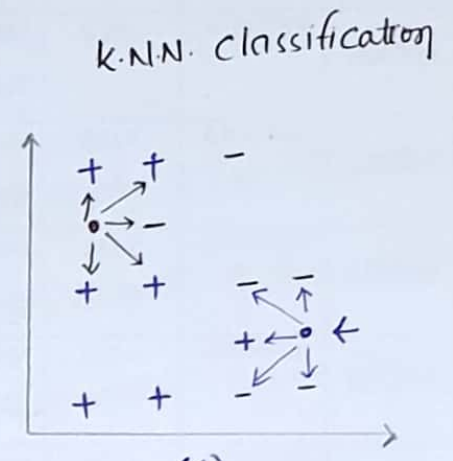
Classification :-



Overfitting { low bias  
                  high variance



Best model { low bias  
                  low var.



Underfitting { high bias  
                  high var.

from the above figure (a) if we plot a line which is best classify or Separate a data points is Consider as classification task if we plot a line across this data. we get this two misclassified points on graph. for Suppose if we have a new data point (or) test data Comes to In this this will get more change in data set ie when train data. we have less classifier in correct point, when ever we consider test data points the change in classifier is more then we can say that the model is "Overfitting"

(4)  
This overfit model containing low bias and high variance.

Case - II :- if we observe a graph (b), that line separates the data points even in that we get some misclassified points where we applying logistic Regression classifier. if we get low misclassified points on separate data in training data as well when we apply test data to this model then also if we get some minute misclassified points in test data. that are changeable. then there is no more impact of test data. (ie not classified more points) then we can say that is the "Best model", this model containing low bias and low variance.

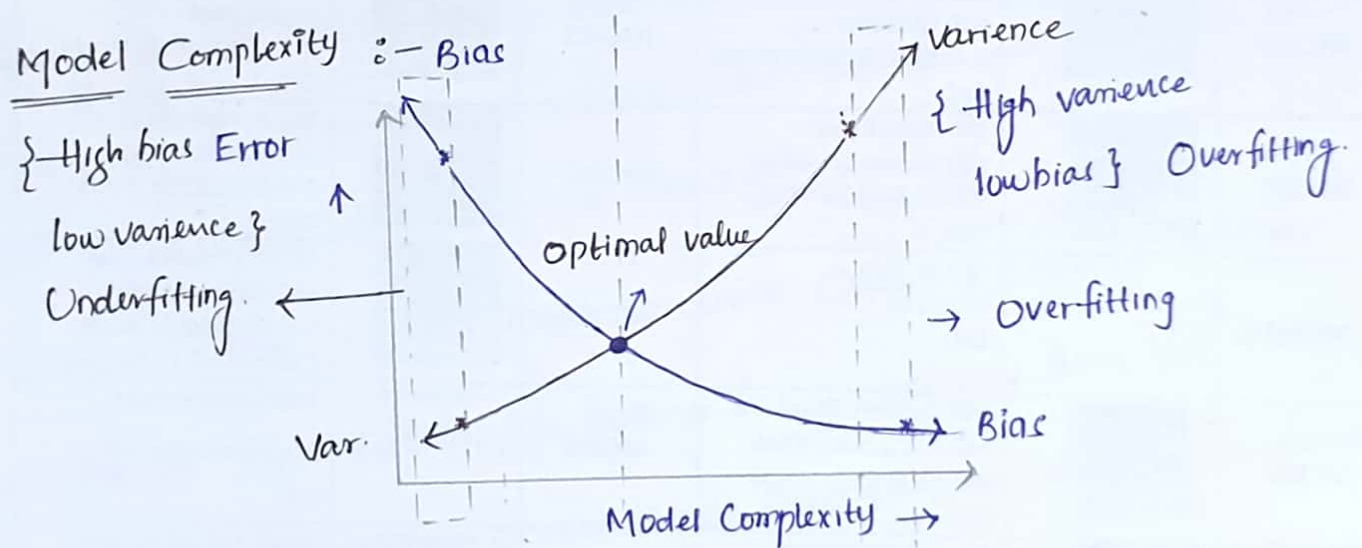
Case - III :- if we observe figure (c), there we apply KNN classification algorithm to classify a data point ie belonging to which category of that point. In KNN the term say that K-nearest points observation where the nearest points are which group or belonging to which category are obtained more that point should be that category ie the test point have also similar properties. here depends on 'K' value the category of point is classified.

if we observe the graph (c) if we take any new test data point on graph from the value if we consider "K = 4" then compare all data points for suppose we get (3+ve) (1-ve) point then we can say the new point is "tve" point but if we know that the point is -ve (ie the test point) but in K nearest points all are positive then we can say that the point is misclassified. here 'K' value is more important.



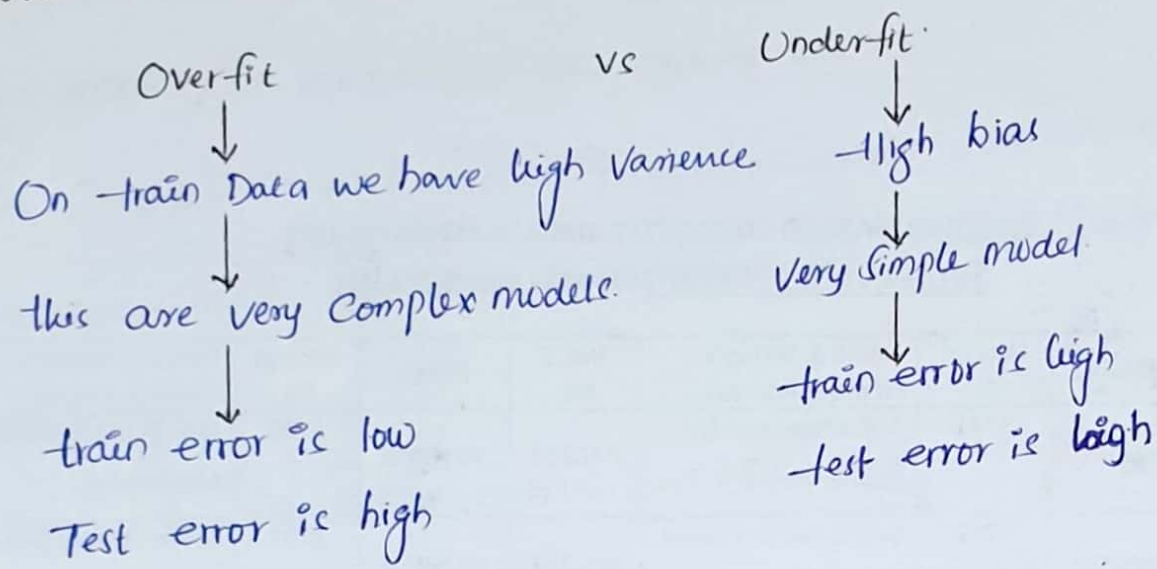
Let Consider we have total training points are '13' and we mentioned ' $k=13$ ' that means all points comes under Observation. then we find the nearest points distance of all 13 data points. after that we sort ~~the~~ all data points. but in this we consider  $k$  value as total training data points. So no need of Sorting if we classify the all data points here +ve's are 7 and -ve's are 6. So the new data point should be consider as a +ve point this is big problem because the model is dumb model (ie because of highest value as "+ve" class, So every time whenever new data point comes we classify that as +ve point.) where high misclassified data points comes in training data as well test data. So then we can easily say that this model should be overfitting model.

Overfit model consist of ~~Low~~ <sup>high</sup> variance as well high bias.



So from this graph we can say that On training data In Overfitting contain high variance, Similarly Underfitting contain ~~high variance~~ low variance. For test data In Overfitting contain low bias, Similarly On Underfitting contain high bias as shown graph.

## Short Notes On Overfit vs Underfit :-



Bias Variance Tradeoff :- it is also called as "Optimal value point" as shown in figure "model Complexity", actually what this point do is when ever this problems come Underfitting and Overfitting. we unnecessarily have high variance and low bias in Complex models, where we have high bias and low variance in Simple models. this both high variance and high bias going to kill model (ie large error rate) and it never be generalizable on test data. In order to reduce (ie maintain low variance, low bias) we use this Optimal value point. (or) trade off.

Regularization :- In order to treat Overfit & Underfit problems we ~~use~~ do Regularization, is nothing but a technique used for tuning the function by adding on additional Penalty term in the error function. In Other words, this is a form of regression, that regularize or shrinks the coefficient estimate towards zero and this discourages learning a more Complex or Simple model, So to avoid the risk of Overfitting.



Note :- All Machine Learning Algorithms will have some regularization steps inbuilt. (ie we have to take care of this terms)

Let Consider linear Regression & logistic regression here we regularize some terms in this algorithm that prevent the issue of Overfitting and Underfitting.

Consider the linear reg eqn -  $m^*, c^* = \arg_{m, c} \text{Min} \sum (y_i - (mx + c))^2$

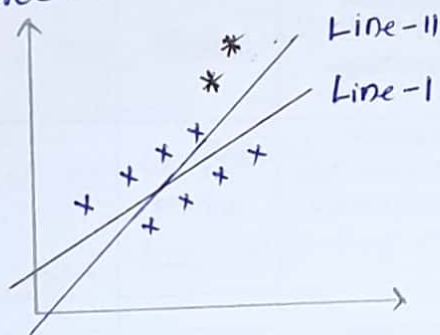
from the eqn we replace  $m$  with  $w$ , - we get eqn -  $w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_0 \end{bmatrix}$

$$w^* = \arg_w \text{Min} \sum (y_i - (w^T x_i + w_0))^2$$

Cost function

Here we minimize the Error from the eqn

So in regression process we fit a line that best fit the data points.



Overfitting Model

from the above graph assume initially there is no lines, we plot data points on graph and making a line-I that best fit the all data points. for Suppose if we consider two more points (test points) which are away from all the points (ie Outliers). In this situation line-I is not a best fit line so we have to again plot a line that best fit the data points. let consider line-II that best fit all the data points. (ie minimize the error) from this line making model will be change if we observe this

We can say for a small change in train data, the entire model will be changed (ie high variance) effected. Simply we can say our linear regression model has overfitting problem.

In order to prevent this problem we use Regularization, this mainly effect on errors. from the linear eqn we have to Cost-function. So our aim is to minimize Cost-function. So simply we can represent -

$$w^* = \arg_w \text{Min} \{ \text{Cost function} \}$$

In order to reduce cost-function in linear regression we have to use 'Gradient Descent' parameter.

So now we have to represent regularized term in order to min error

$$w^* = \arg_w \text{Min} \{ \text{Cost function} + \text{Reg.}(w) \}$$

where  $w \rightarrow$  coefficient (m,c) in linear Regression.

So now here we used two diff regularized terms for linear Regression

they are 1) L1 Regularized also called as lasso Regularizer

2) L2 Regulariser also know as Ridge Regularizer.

these both terms comes from  $(L_1, L_2)$  Manhattan, Euclidean distance

In L1 Regularizer we do -  $\sum |w_i|$  absolute.

In L2 Regularizer we do -  $\sum (w_i)^2$  because M.S.E

when ever we use these regularizers we have to multiply this term with lambda denoted as  $(\lambda)$ , it is a Inbuilt function in

Regularizers. like  $(L_1 \& L_2)$



Now finally Eq's will be -

linear reg eqn with  $L_1$  Reglizer -  $w^* = \arg_w \text{Min} \left\{ \sum (y_i - (w^T x_i))^2 + \lambda \sum_{i=1}^n |w_i| \right\}$

linear reg with  $L_2$  -  $w^* = \arg_w \text{Min} \left\{ \sum (y_i - (w^T x_i))^2 + \lambda \sum_{i=1}^n (w_i)^2 \right\}$  Lasso Reglizer

Logistic Regression :-

Ridge Regularizer

Now we consider logistic eqn -  $m^*, c^* = \arg_{m,c} \text{Max} \left\{ \sum (y_i * (mx_i + c)) \right\}$

In this eqn we have a problem with outlier, because of this we get incorrect result; to avoid this we use Sigmoid function, when apply this function the eqn will be -  $m^*, c^* = \arg_{m,c} \text{Max} \left\{ \sum \frac{1}{1+e^{-x}} \right\}$

Now we consider the  $(y_i * \hat{y})$  as a function - 'f' then -

$$m^*, c^* = \arg_{m,c} \text{Max} \left\{ \sum \frac{1}{1+e^{-f}} \right\} \Rightarrow \text{Conbe} \Rightarrow \exp \{-f\}$$

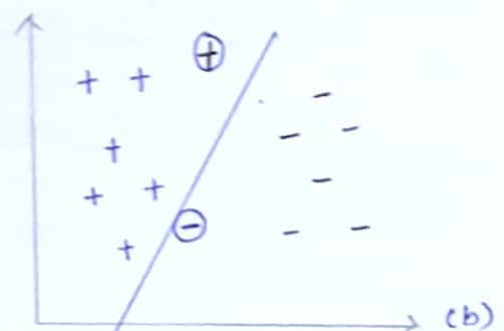
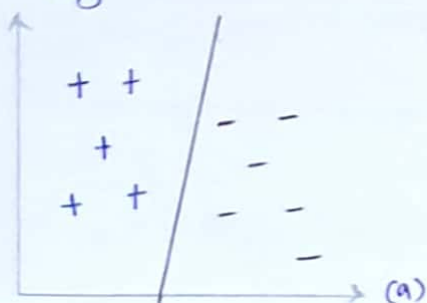
Now from eqn put 'f' value  $\Rightarrow m^*, c^* = \arg_{m,c} \text{Max} \left\{ \sum \frac{1}{1+\exp \{-y_i * (mx_i + c)\}} \right\}$

So when we minimize the value of 'f' then only we get Max result

In Order to do it we do inverse of it -  $\max \frac{1}{f} = \min f$

finally we get -  $m^*, c^* = \arg_{m,c} \text{Min} \sum_{i=1}^n \left\{ \frac{1}{1+\exp \{-y_i * (mx_i + c)\}} \right\}$  cost func.

This eqn represent logistic Regression after treatment of Outliers with use of Sigmoid function, which tries to classify the positives to negatives



from figure (a) we know that the line which is best

Separate the data from positives to negatives. but what happens when new data point (test data) comes, then we have to drastically line formation changes according to new data points if we observe fig (b) the line angle will be change as compare with fig (a). and in fig (c) we correctly classify the data points initially which are misclassified data points.

"Our aim is to reduce the misclassified points in logistic Regression"  
So that the best separated line fig (c) reduce the misclassified data points from this separation line making model will be changed. if we observe this we can say for a small change in train data the entire model will be changed (ie high variance) effected so we can say this model has overfitting issue. to solve it we use regularizer methods.

In logistic Regression we have a Cost-function, so our aim is to correctly classify the data points, the eqn will be represented as -

$$w^* = \arg_w \text{Min} \{ \text{Cost-function} \}$$

So now we represent regularized ~~term~~ term in order to current classify.

$$w^* = \arg_w \text{Min} \{ \text{Cost-function} + \text{Reg}(w) \}$$

In logistic also we used  $L_1$  &  $L_2$  Regularizers. -

Logistic reg. with  $L_1$  Regulr.  $w^* = \arg_w \text{Min} \{ \sum (1 + \exp\{-y_i * (mx_i + c)\}) + \lambda \sum_{i=1}^n |w_i| \}$   $\rightarrow$  lasso Reglrr

$L_2$  Reglrr -  $w^* = \arg_w \text{Min} \{ \sum (1 + \exp\{-y_i * (mx_i + c)\}) + \lambda \sum_{i=1}^n (w_i)^2 \}$  Ridge