# Machine learning.

knowledge check 1:1 :-

1. For each of parts (a) through (d), indicate whether wether we would generally expect the performance of a flexible statistical learning method. Justify your answer.

a) The sample size n is extremely large, and the number of prrediction p is small.

Ans:-

In this case the better performance is flexible statistical learning comparing inflexible statistical learning method.

b) The number of preditors p is extremely large, and the number of observations n is small.

Ans:-

The performance of a flexible. statistical learning method is

worst because of overfitting is very high.

c) The relationship between the predictors and response is highly non-linear.

Ans:-

The performance of a flexible statistical learning method is better because of normal distribu-

d) The variance of the error terms i.e. $\sigma^2 = \text{var}(\epsilon)$ is extremely high.

Ans:-

The performance of a flexible statistical learning is worse when the variance term is very high.

2) Explain whether each scenario is a classification or Regression problems, and indicate whether use are most intrested in Intorence or prediction. Finally provide $n$ and $P$.

a) we collect a set of data on the top 500 trims In the us for each Ferms we Record, profits, number of Employes, and the CEO salary. we are intrested in understanding which factor attect CEO salary.

Ans:-

Regression problems. Because of this scenario is a continious variable.

* Interence

* $n = 500$, $P = 4$

b) we are considaring launching a new product and wish to know whether it will be success of failure. we collect data on 20 similar products

The were previash launched for each product we have beorded whether it was a success or failure, price charged for the product, marketting budget conperio price and the ten other variabels

Ans:-

* clasitication problem, because of this scenario is a categorical variable.

* Prediction.

* $n = 20$, $P = 13$.

c) we are intrested in the predicting %. chang in the USD/Euro exchange rate in relation to the weekly changes in the world stock markets. Hence we collectly weekly data for all of 2012 For each week we record the % in the USD/Euro

chang in the use Market, and
the % change in the German
market.

Ans :-

&ast; Regression problem.

&ast; prediction

&ast; $n = 52$. $D = 4$

3) a) Discribe 3 real - life Application
in which Regression might be
usetul. Describe The Response as
well as the predictors is the
goal ot each application
intorence or predication Explain
your answer.

Ans:

&ast; predicting the marks of a
student Based on the number
ot hourse he/she put into
preperation.

* predicting intererence the intereree
Reports.

* Movie Rating betore ask, Directi
The Ditterence region.

b) Discribe the real-lite applicatic
In which classitication might
be usetul. Discribe the respons
as well as predictors. Is th
goal ot each application
interence or prediction? Explai
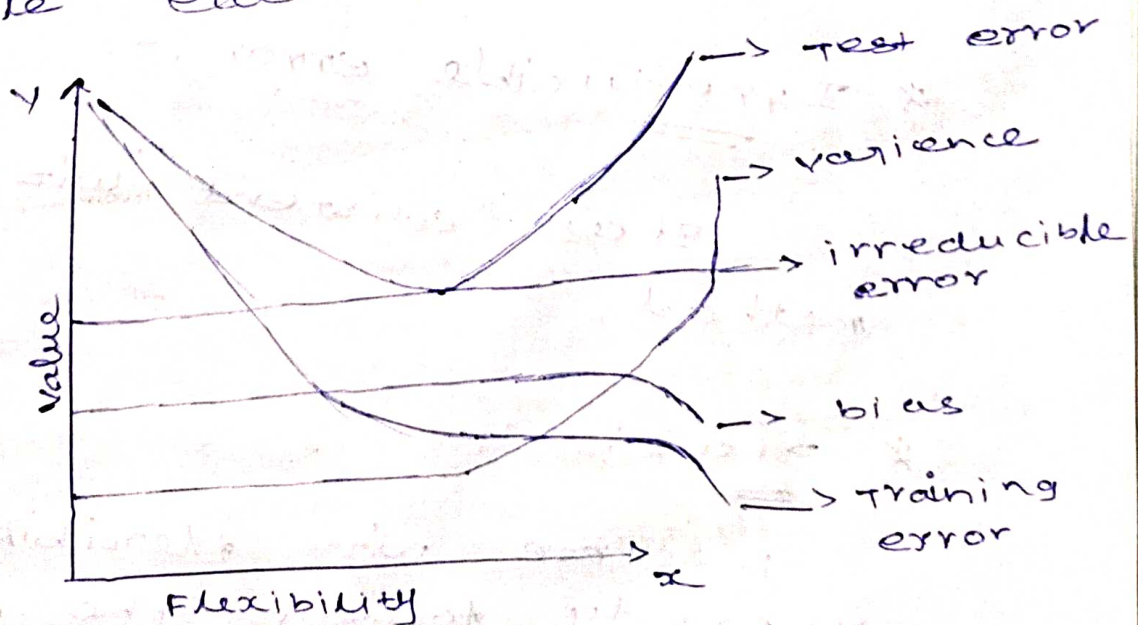your answer.

Ans :-

* classity it is spam or not
Email.

* The goal is Identity image
ot single digit 0.9 correctly
( hand written digit Recogonitia

* Identity disease other tissue
type based on the goan
expression Levels.

knowledge check 1.2

1. We now revisit the bias-variance decomposition.

a) provide a sketch of typical (squariel) bias, variance, training error, test error, and payers errror curves, on a single plot, as we go from less flexible statistical learining method towards more flexible approches the x_ axis should represent the amound of flexible in the method, and the y-axis should resprent the values for each curves. there should be five curves, makes sure to lable each on.

b) Explain why each of the five corves has the shape displayed in part (a).

Ans:-

* Trcuing Error :- (80% or 70%)

Decreass with flexibity - possible to be better follow the data with more flexible more.

* Test error : (20% or 30%)

Decreases end then increas with flexibility error increase because model is following nois of data in traning set and test data po not have the same noise.

* Irreducible error :-

stay constent with the method.

* Bias :-

Decreee with flexibity mor the data

more appropriately fit the data.

* Increase with flexibility more unsteaddy. follower the data more.

2. what are the advantages and disadvantages of a very flexible approch for regression or classification under what cricumstancess hight a more flexible approch be preferred to a less flexible approch? when might a less flexible approch be preferred?

Ans:-

Advantage:-

* Non Linear data

* Less bias

* variable intractions.

Disadvantage:-

* Lack in interpretability.

* high in varience.

Why take more or less flexible option?

more flexible when bits of data and many different groups and chose less flexible when few data Doing for Both Regression and classification approaches

3) Describe the Difference between a parametric and a non-parametric statistical learning approach what are the advantage of a parametric approach to regression or classification? what are its disadvantages?

Ans :-

Parametric method :-

This method make an assumption about the function of the model and that it is linear.

# Non - parametric method :-

This method do not assume anything about to funection when trying to estimate the fit of the data.

## Advantages :-

* Linear
* easy to intrepret
* easy to dotit

## Disadvantage :-

* non - Linear
* no easy to intrepret and ditit.

4) The table below provides a traning data set continning six observation, three predictors and one Qualitative response variable.

| | $x_1$ | $x_2$ | $x_3$ | Y | |
|---|---|---|---|---|---|
| 1 | 0 | 3 | 0 | Red | |
| 2 | 2 | 0 | 0 | Red | |
| 3 | 0 | 1 | 3 | Red | |
| 4 | 0 | 1 | 2 | Green | |
| 5 | -1 | 0 | 1 | Green | |
| 6 | 1 | 1 | 1 | Red | |

suppose we wish to use this data set to make a prediction for y when $x_1 = x_2 = x_3 = 0$ using k - nearest neighbors

a) compute the Evelidean disadva between each observation and the test point, $(x_1, x_2, x_3) = (0,0,0$

·Ans :-

$$1 = \sqrt{(0-0)^2 + (8-0)^2 + (0-0)^2} = \sqrt{9} = 3$$

$$2 = \sqrt{(0-0)^2 + (0-0) + (0-0)^2} = \sqrt{11} = 2$$

$$3 = \sqrt{(0-0)^2 + (1-0)^2 + (2-0)^2} = \sqrt{10} = 3.16$$

$$4 = \sqrt{(0-0)^2 + (1-0)^2 + (2-0)^2} = \sqrt{5} = 2.23$$

$$5 = \sqrt{(-1-0)^2 + (0-0)^2 + (1-0)^2} = \sqrt{2} = 1.41$$

$$6 = \sqrt{(1-0)^2 + (1-0)^2 + (1-0)^2} = \sqrt{3} = 1.73$$

b) what is out prediction with

k = 1 ? why ?

Ans'.-

  Green

c) what is out prediction with

k = 3 ? why ?

Ans :-

  Red .

d) It the Buyers decision boundary in this problem is highly non linear, then would we expect the best value for k to be Large on small? why ?

Ans :

  small value.