

knowledge check 1.1! - how good

1) For each of Parts (a) through (d), indicate whether we would generally expect the performance of a flexible statistical learning method to be better or worse than an inflexible method. Justify your answer.

a) The sample size is extremely large, and the number of predictors is small.

Ans 1:-

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

Learning predictors 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 predictor and response is highly non-linear.

Ans:-

The performance of a flexible statistical learning method is better because of normal distribution.

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 problem, and indicate whether we are most interested in inference or prediction. Finally provide an analysis.

a) We collect a set of data for the top 500 firms in the US. For each firm we record Profit, number of Employees, industry and the CEO Salary. We are interested in understanding which factor affect CEO Salary.

Ans:-

* "Regression problem:-"

Because of this

Scenario, it's a continuous variable.

* Inference

* $n = 500, P = 4$

b) We are considering launching a new product and we want to know whether it will be a success or a failure. We collect data on 20 similar products that were previously

collect data on

launched - for each product that were previously launched. For whether it was success or failure, price charged for the product, marketing budget, competition, price and ten other variable.

Ans! -

* Classification Problem! -

Because of this scenario is a categorical Variable.

* Prediction

* $n=20, P=13$

(c) We are interested in predicting the y. change in the USD/Euro exchange rate in relation to the weekly changes in the world stock market, hence we collect weekly data for all of 2012. For each week we record the y. change in the USD/Euro, the y. change

in the US market, the % change in the British market, and the % change in the German market.

Ans:

* Regression Problem

* Prediction

* $n = 52, p = 4$

3) a) Describe three real-life application in which regression might be useful.

Describe the response, as well as the prediction. Is the goal of each application inference or prediction? Explain your answer.

Ans:

* Predicting the marks of a student

based on the number of hours he/she put into the preparation.

* Predicting and inferring whether

reports

are true or false.

Movie rating before directing the

at different regions.

b) Describe three real-life applications in which classification might be useful, describe the reason as well as the prediction. Is the goal of each application inference or prediction?

Ans:-

* Classify if it is spam or not (Email)

* The goal is to identify images of single digit 0-9 correctly (Handwritten digit recognition)

* Identify disease or tissue typed based on the gene expression levels.

Knowledge Check 1.2:-

1) We have revisited the bias - variance decomposition

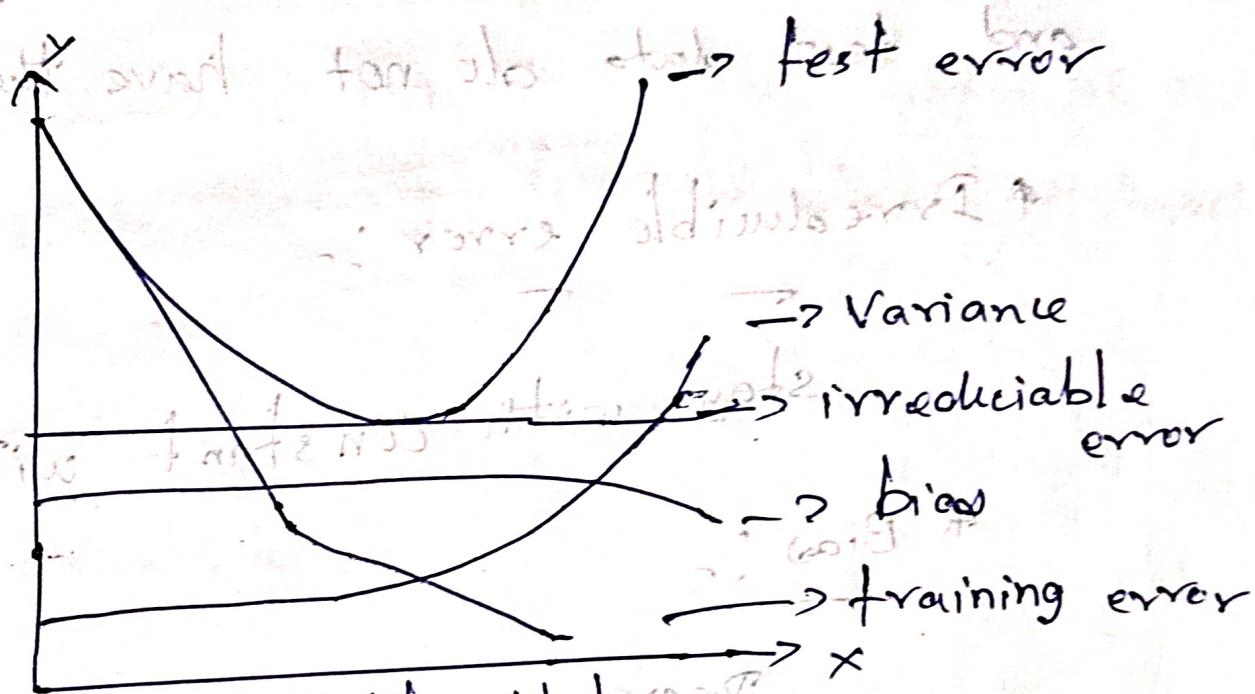
a) Provide a sketch of typical bias, variance training error, test error, and bayes error curves, on a single plot, as we go from flexible to statistical learning

method towards more flexible approach.

The x-axis should represent the amount of flexibility in the method.

and the y -axis should represent the values for each. There should be five make sure to label each one.

Ans :-



b) Explain why each of the five curves has the shape displayed in part(a):-

Ans:-

* Training Error! - (80% or 90%)

Decrease with flexibility - possible to better follow the data with more

* Test Error :- (20% or 30%)

Decreases and increases with flexibility, error increases because model is following of data in training set and test data do not have the same noise.

* Irreducible error :-

Irreducible error

Stays with constant with the method

* Bias :-

Decreases with flexibility

because more likely to appropriately fit the data.

* Variance :-

Increases with flexibility

more wobbly, follows the data more.

2) What are the advantages and disadvantages of a very flexible approach for regression or classification under what circumstances might a more flexible approach be preferred over a less flexible approach? When might a less flexible approach be preferred?

Ans :-

Advantage :-

- * Non-Linear data
- * Less bias
- * Variable Infractious

Disadvantage :-

- * Lack in interpretability
- * high in Variance.

Why take a more or less flexible option?

Chose more flexible when lots of data and many different

groups and chose less flexible when few data points for both regression and classification approaches.

3) Describe the differences between a parametric and a non-parametric statistical learning approach. What are the advantages 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 methods do not assume anything about the function when trying to estimate the fit of the data.

~~Advantages~~: associated with algorithms

~~Easy to fit with large multivariate data sets~~

* Linear

* easy to interpret

* easy to define

Disadvantages:

* Non-Linear

* No easy for interpret and defit.

- 4) The table below provides a training data set containing six observations, three predictors and one qualitative response variable.

	x_1	x_2	x_3	
1	0	3	0	Red
2	2	0	1	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.

$x_1 = x_2 = x_3 = 0$ using k-nearest neighbors

a) compute the Euclidean distance between each observation and the test point,

$$(x_1, x_2, x_3) \in \{0, 1, 0\}$$

Ans!.

$$x_0 = 0, 0, 0$$

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

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

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

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

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

$$6 \Rightarrow \sqrt{(-1-0)^2 + (0-0)^2 + (1-0)^2} = \sqrt{2} = 1.41$$

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

b) What is our prediction with $k=1$? why?

Ans!.

Green

c) What is our prediction with $k=3$? why?

Ans!.

d) If the Bayes decision boundary in this problem is highly non-Linear, then would we expect the best value for k to be large or small? why?

Ans:-

Small Value