

## Variance and Bias:

\* The error induced because of estimating the population is known as bias.

$d \rightarrow$  estimate

$\theta \rightarrow$  parameter

$(d - \theta)$  is error

\*  $d$  is a random variable so expected value of  $d$  is  $E(d)$ .

$E[d] - \theta$  is defined as bias

The difference between  $E[d] - \theta$  is error

\* The expected value of statistics is equal to population parameters then this statistic is equal to unbiased estimate population parameter

$\bar{x}$  is unbiased estimate of  $\mu$

$$\mu = \bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

## Assignment:

1, Prove that  $\bar{x}$  is unbiased estimate of  $\mu$

$$\mu = \bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

2. Prove that  $\bar{\pi}$  is not unbiased estimate of  $\sigma^2$

$$\sum_{i=1}^n (\pi_i - \bar{\pi})^2$$

3. In the process of proving 2<sup>nd</sup> proof. Prove that

$$\sigma^2 = \frac{\sum_{i=1}^n (\pi_i - \bar{\pi})^2}{n-1}$$

Evaluating an Estimator: Bias and Variance

Let  $X$  be a sample from a population.

$$X = \{x^t, y^t\}$$

$\downarrow$  Input       $\downarrow$  output

$g$  is an estimate of  $v$

$g \sim v$ .  $v$  is a population parameter

$E[(x-g)^2]$  is Mean square error (MSE)

$$MSE = E[(x - E(x)) + (E(x) - g)]^2$$

$$E[(x - E(x))^2 + (E(x) - g)^2 + 2(x - E(x))(E(x) - g)]$$

$$= \underbrace{E[(x - E(x))^2]}_{\text{variance}} + \underbrace{E[(E(x) - g)^2]}_{\text{Bias}} + 0$$

ref notes

$$\underbrace{2[E(x) - g][E(x) - E(x)]}_{\text{Constant}}$$

$$\Rightarrow 2(E(x) + \lambda)(E(x) - E(x))$$

$$= 0$$

Notes:

- a, Variance is coming from data, It is scattered around mean.
- b, the one which we can control is Bias
- c, the one which we cannot control is variance
- d, Statistic is estimation of the population parameter

\* Variance: Variance measures how much, on average,  $x_i$  vary around the expected value (going from one dataset to another),

\* Bias: Bias measures how much, the expected value varies from the correct value  $g$

\* Error: error is the sum of the variance and the square of the bias.

Decision tree:

A decision tree is a hierarchical model for supervised learning whereby the local region is identified in a sequence of recursive splits in a smaller number of steps. A decision tree composed of internal decision nodes and terminal leaves

It is an efficient nonparametric method, which can be used for classification and regression.

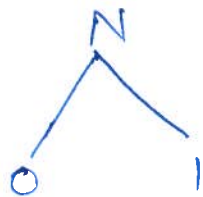


It is a binary tree.

Classification problem:

- \* Complex decision based on several parameters
- \* Basic decision stage

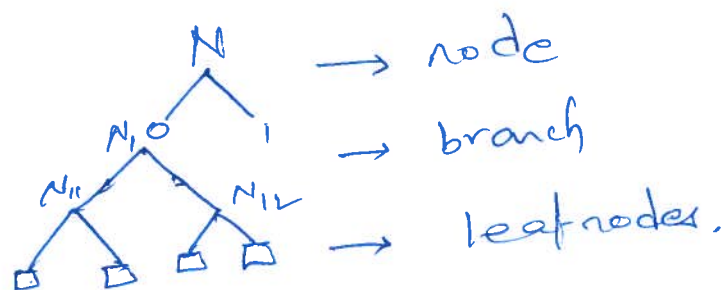
\* Let's consider  $N$  as property called price, used to decide whether a car is family car or not for a number of samples



If price  $< 500000$  No represents 0 leaf in tree.

If price  $> 10,00000$  Yes represents 1 leaf in tree.

The above problem can be extended to multiple parameters example price & engine power.



Each decision node  $N$  implements a test function  $f_N(x)$  with discrete outcomes labelling the branches. Given an input, at each node, a test is applied and one

5

of the branches is taken depending on the outcomes.

This process starts at the root and is repeated recursively until leaf node is hit, at which point the value written in the leaf constitutes the output.

- \* Decision tree can manage discrete data very well.
- \* Continuous variable cannot be handled by decision tree, continuous variable has to be discretize.

### Drawbacks of Decision tree:

- It gives best result when number of parameters are limited.
- It gives better result for very big amount of data which takes care of all possible values atleast once.
- Theoretically each node should have atleast 30 points.
- Decision tree is specific to sample.

### Advantages:

- No need to eliminate outliers
- error induced because of outliers is negligible

