**UE18CS322 Lecture Notes : Topic – Introduction to Big Data**

### Question Bank for Lecture
- o T1 – Self assessment linked to LO 1.1 - #1, 4, 5
- o What are common pitfalls in analysis of data? Give examples
- o How is error estimation of model done in Big Data problems? How is this different from the traditional approach?
- o In the Peter Norvig video he defines the following function to compute the best segmentation for a sentence?


**UE18CS322 Lecture Notes : Topic – Introduction to Big Data 2**


### Question Bank for Intro2
- o T1: Self assessment for LO1.3 – 1, 2
- o T2: Multiple choice questions – 1-7, Short answer : 1-7
- o What are the 4Vs characteristics of Big Data?
- o What is the difference between volume and velocity of data?
- o What differences in data generation have resulted in higher volume?
- o Give some example of issues of veracity in Big Data?
- o What modules does a typical Big Data architecture have?
- o How Big Data and Hadoop are related to each other?
- o What are the steps one must adopt to deploy a big data solution?


**UE18CS322 Lecture Notes : Topic – Hadoop Distributed File System**

### Question Bank

- o T2 : MCQs 1-10.  Short answers 1-10
- o What is Hadoop HDFS cluster and in how many different modes can you configure a  Hadoop HDFS cluster?
- o What is the usage of HDFS for the read operation?
- o Described the use of Namenode? Why is it the most important process in HDFS
- o The namenode contains the metadata that maps the user file block to the datanodes that contain a copy of the block. It is the most important service as in its absence, the cluster will fail to function. As the namenode is running on the master node, any failure to the master node, renders the entire cluster unusable.
- o Can namenode and datanode be on commodity hardware?
- o What are the differences between Hadoop and RDBMS?
- o Explain the core components of Hadoop.
- o What is a block in HDFS and what is its default size in Hadoop 1 and Hadoop 2? Can we change the block size?
- o What is Commodity Hardware?
- o Why is HDFS only suitable for large data sets and not the correct tool to use for many small files?