

BDA

Chp -1

* Data Mining :- It's about extract the Data and discover some patterns & Model.

12 There are different Types of Models:-

- ① Statistical Modeling
- ② Machine learning
- ③ computational Approaches to Modeling
- ④ summarization
- ⑤ Feature extraction

1) Statistical Modeling :-

Statistical Modeling, that is an Underlying distribution from which the visible Data is drawn.

2) Machine learning :-

Machine learning practitioners use the data as a training set to train an algorithm of one of many types used by machine learning such as BayesNet support vector machines.

3) Computational Approaches to Modelling :-

To solve complex problem we use computational Model. There are many different Approach to modeling Data

4] Summarization:-

One of the most interesting forms of summarization is the Page Rank Idea, which made Google successful and which we will know. In the form of web mining the entire complex structure of the Web is summarized by a single number for each page.

5] Feature extraction:-

The typical feature based model looks for the most extreme examples of phenomenon and represent the data.

Some of the important kinds of feature extraction from large scale data that is:

- ① frequent Items
- ② similar items

* Statistical limits on Data Mining :-

① Total Information Awareness :-

- It was conjectured that the information needed to predict and foil the attack was available in data, but that there was then no way to examine the data and detect suspicious events.
- This response was a program called TIAA.

② Bonferroni's Principle :-

- A theorem of statistics known as the Bonferroni correction gives a statistically sound way to avoid most of these bugs.
- Bonferroni's principle help us avoid treating random occurrences as if they were real.

Example : There are one billion people who might be evil-doers.

③ Everyone goes to a hotel on day in 100.

④ we shall examine hotel records for 100 days.

* Hash Functions:-

- A hash function h takes a hash key value as an argument & produces a bucket number as a result.

- The bucket number is an integer.

- Hash key can be of any type.

- If hash key drawn randomly then h will send equal numbers of hash keys each of the buckets.

→ Uniform Distribution.

$$h(x) = e$$

- Σ of value & divide by Numbers of Bucket.

* Indexes:-

- It is an a data structure that makes it efficient objects given the value of one or more element of objects.

- It store the records. It allow to find record quickly through Index number.

* Secondary Storage :-

- It is Important, when dealing with large scale Data, that we have a good understanding of the difference in time taken to perform Computation.
- Disks are Organized into blocks.

January

JAN											
M	T	W	T	F	S	S	M	T	W	T	F
1	2	3	4	5	6	7	8	9	10	11	12
13	14	15	16	17	18	19	20	21	22	23	24
25	26	27	28	29	30	31					

Thursday

Wk 12

Chp-2

mmds.org

1 to 32 slide (Revise of BDA)

* Big Data Workflows :-

Understanding workflows & effect of big data is following :-

- ① Identify the big data sources you need to use.
- ② Map the big data types to your workflow data types.
- ③ Ensure that you have the processing speed & storage access to support your workflow.
- ④ Select the data store.

* SPARK :- It is a multi-language for executing data engineering, data science and machine learning on single node machine or clusters.

- It provides high level APIs in Java, Scala, Python & R and an optimized engine that supports general execution.

* RDD :- The main abstraction Spark provides is a resilient distributed dataset (RDD). which is collection of element partitioned across the nodes of the cluster that can be operated on in parallel.

- RDD automatically recover from node failures.

There are 2 ways to create RDD:-

① parallelizing an existing collection in your driver program.

② referencing a dataset in an external storage system.

January 20

13

Saturday

MS 352 WK 02

JAN											
M	T	W	T	F	S	S	M	T	W	T	F
1	2	3	4	5	6	7	8	9	10	11	12
13	14	15	16	17	18	19	20	21	22	23	24
25	26	27	28	29	30	31					

* RDD Operation in pySpark :-

The RDD supports two types of operations:-

① Transformations:-

- Transformations are the process which are used to create a new RDD.
- Few of transformations are given below:-

- ① map
- ② flatmap
- ③ filter
- ④ distinct
- ⑤ reduceByKey
- ⑥ mapPartitions
- ⑦ sortBy

14 Sunday



② Actions:-

- Actions are the processes which are applied on an RDD.

Monday 15

- Focus Actions and Following:

collect

collect AsNoP

reduce

countByKey

countByValue

take

first

* Narrow Operations in RDDs

① count()

reduce(f)

collect(f)

map

filter(f)

foreach(f)

join

cc

January 2018

JAN

M	T	W	T	F	S	S	M	T	W	T	F	S
1	2	3	4	5	6	7	8	9	10	11	12	13
14	15	16	17	18	19	20	21	22	23	24	25	26
27	28	29	30	31								

Tuesday

* Create RDD Using Spark Context parallelize()

Code:-

Create RDD from parallelize

data = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]

rdd = Spark.SparkContext.parallelize(data)

* Create RDD using SparkContext.textFile()

Create RDD from external Data

rdd2 = Spark.SparkContext.textFile("/Path/textfile.txt")