

Exam Roll - CSE 218056

Class Roll - 0017AOS01029

BCSE-IV 2nd Semester.

Subject :- Big Data Analytics.

1.

i) The characteristics of big data are as follows:-

1. Volume :- The amount of data handled by big data systems are very huge. Such large amounts cannot be handled by conventional systems.

2. Velocity :- Since the data involved in big data is very huge, conventional systems are not able to provide results in reasonable time.

Thus providing results on demand at a faster pace is a major aspect of Big Data.

3. Variety :- Big data is generated in multiple varieties. Compared to traditional data like texts, phonenumbers, addresses, the latest trend in data is also in the form of images, video, audio etc. A high volume of the data is unstructured.

4. Veracity :- Veracity refers to the degree of reliability the data has to offer. Big Data algorithms are required to filter out huge amounts of data and translate them so that they are relevant and reliable to the particular requirement.

Ques

1.

ii) The different types of analytics are :-

1. Descriptive :- It involves evaluating new descriptions of entities and their properties.
It basically is the summarization of past data in an easily readable form.
2. Diagnostic :- It involves determining whether something has happened based on past data.
3. Predictive :- It involves forecasting events that have not happened yet.
4. ~~Cognitive~~ Prescriptive :- It involves finding out the best course of action in a given situation.
5. Cognitive :- It involves combining analytics with AI and Machine Learning.

1.

iii) Data ingestion is a process by which data is moved from one or more sources to a destination where it can be stored and further analysed.

Data is collected from different sources in different formats. These non-homogeneous data needs to ~~be~~ be transformed in a way so that ~~is~~ it is possible to analyse it together.

Since Big Data involved data in very high volume, with high variety and ~~and~~ velocity, ~~as~~ conventional systems fail. Thus Big Data ~~ingest~~ ingestion ~~therefore~~ tools like Apache Kafka, Wavefront, Amazon Kinesis etc. are required.

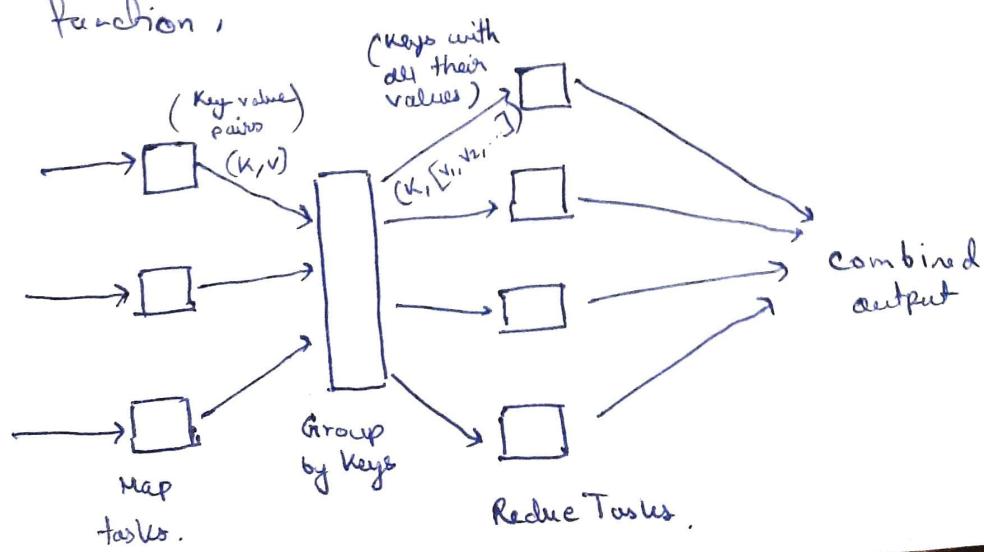
v). Map Reduce is a programming processing technique that can handle large-scale parallel computation in a way that is tolerant to hardware failures.

The user needs to write two functions Map & Reduce. The execution occurs as follows:-

a) Some number of Map tasks are given one or more chunks from a distributed file system. The Map function converts the chunks in a ~~series~~ of Key-Value pairs.

b) The 'key-value' pairs from each Map task are collected by a master-controller and sorted by key. They keys are sorted divided among all the Reduce tasks, so all key-value pairs with the same key wind up at the same Reduce task.

c) The Reduce tasks work on one key at a time, and combine all the values associated with that key in ~~a~~ a way specified by the Reduce function,



Q.

iv) The broad features of Hadoop distributed file system are:

- Can manage very large distributed file system.
It can scale upto thousands of nodes and handle ~~million~~ large amounts of data.
- It assumes commodity hardware. Even though a large number of nodes are handled, the nodes themselves are not expected to be very powerful. Also nodes are expected to fail. Files are replicated to handle hardware failure. The architecture can detect failures and recover them.
- It is optimized for batch processing. The data locations are exposed so that computations can move to where the data resides. This provides very high aggregate bandwidth.
- A single namespace is used for the entire cluster.
- Hadoop has a very robust ecosystem that is ~~well~~ well suited to meet the analytical needs of developers and small to large scale organizations.

7.

i) Attribute Value Frequency (AVF) Algorithm is a simple & fast algorithm to detect outliers in categorical data.

It is based on the intuition that outliers are those points which are infrequent in the dataset, and that the ideal outlier point in a categorical dataset is one whose each & every value is extremely irregular or infrequent.

Given a dataset where each point x has m attributes

$x_i = [x_{i1}, x_{i2}, \dots, x_{im}]$ where x_{ij} is the j th attribute of x_i .

$$\therefore \text{AVFScore}(x_i) = \frac{1}{m} \sum_{j=0}^m f(x_{ij})$$

Outliers will have infrequent ~~categories~~ attributes, thus will have a lower AVF Score.

Once the AVF Score of all points are calculated we can designate the points with very low scores as outliers.

Using MapReduce architecture we can parallelize AVF and thus can operate on large datasets.

The MR-AVF is comprised of 2-stages.

Stage 1:-

Map:- The first map function associates attribute's each distinct value to the map's output key.

Reduce:- The reduce function computes the attribute freq. count of each attribute value.

~~Sohor~~

Stage 2 :-

Map:- AVF score of each data point is calculated.

Reduce:- sorting operation the calculate AVF scores.

PSUDO CODE :-

Input: Database D (n points x m attributes)

Output: K-detected outliers.

HashTable H;

```

map  $k_1 = i, v_1 = D_i = x_i, i = 1 \dots n$  begin
  | for each  $l \in x_i; l = 1 \dots m$  do
    |   collect  $(x_{il}, 1)$ ;
  | end
end.
  
```

Stage 1

```

reduce  $k_2 = x_{il}, v_2$  begin
  |  $H(x_{il}) += v_2$ ;
end
  
```

```

map  $k_1 = i, v_1 = D_i = x_i$  begin
  |  $AVF = \sum_{l=1}^m H(x_{il})$ ;
  | collect  $(k_1, AVF)$ ;
end
  
```

Stage 2.

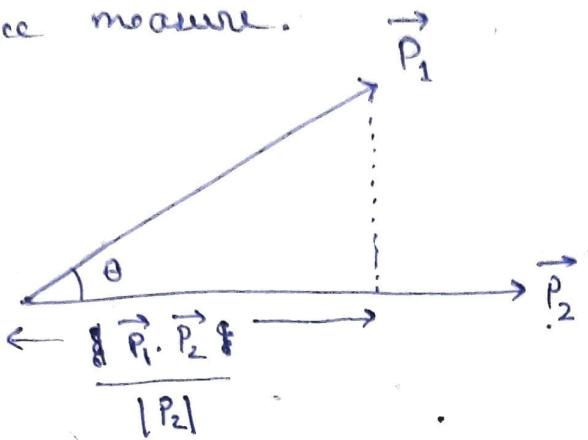
reduce $k_2 = AVF_1, v_1 = i$;

7

ii) Cosine distance :-

Assume ~~open~~ points as vectors from the origin to its location.

Cosine distance simply refers to the angle between such two vectors. It is a non-Euclidean distance measure.



$$d(P_1, P_2) = \theta = \arccos(\vec{P}_1 \cdot \vec{P}_2 / |\vec{P}_1| \cdot |\vec{P}_2|)$$

We can prove that ~~as~~ here θ is a distance measure:-

$$\rightarrow d(x, x) = 0 \text{ because } \arccos(1) = 0$$

$\rightarrow d(x, y) \geq 0$ because any two intersecting vectors make an angle in the range 0° to 180° .

$$\rightarrow d(x, y) = d(y, x) \text{ by symmetry.}$$

\rightarrow Triangle inequality: If we rotate an angle x to z and then from z to y , we can't ~~not~~ rotate less than from x to y .

$$\therefore d(x, z) + d(z, y) \geq d(x, y)$$

\therefore Cosine distance is a distance measure.

2. The main motive behind search autocomplete is to reduce typing time of the users, thus providing relevant results faster.

In most search engines like Google & Yahoo, search engine autocomplete is based on past searches.

The search engine sorts through trillions of past search queries to find those matching phrases to provide suggestions.

One way of implementing this auto complete feature is through fuzzy search.

Fuzzy search is implemented by breaking the search query into ~~into~~ K-shingles.

After Shingling the search queries are converted into bit vectors and then using MinHashing their hashes are stored.

When a user searches a new query, its minHash is compared with other pre present searches in the database.

An efficient method of comparing the ~~two~~ hashes is using Locality Sensitive Hashing. ⑩

The MinHash signatures of the search ~~and~~ queries are broken down into a number of bands. These bands are then again hashed into ~~multiple buckets~~, to a hash table with k buckets. Queries that hash to the same bucket are likely to be similar.

Once we have the list of all similar ~~past~~ queries another algorithm is used to ~~filter~~ filter those set in ~~order~~ order to show ~~the~~ only the most relevant queries to the user.

The queries may be sorted by relevance and the most relevant queries are shown to the user.

In some search engines like Google, if the user is logged in, the engine first few suggestions may be the user's own ~~past~~ past search queries.

Queries are also sorted by relevance of a query. Some may be affected by various ~~re~~ factors. Some of these factors include trending searches, location, language used to search etc.

In the screenshots we can see that the search suggestions of Yahoo & Google are quite different, although there are ~~are~~ a few common suggestion in case of the "ce".

This difference may arise due to ~~due~~ difference in filtering algorithms of ~~the~~ both the search engines.

Also due to difference in user base, the search query history database of both the search engines ~~are~~ is expected to be different.

If the user is logged in on the search engine (has an account made on the platform) user data like past searches & other preferences can also play a role in these suggestions.

3. Advertisers bid for query strings on search engines.

The search engines shows adds on the search results page that maximizes their ~~profits~~ revenue.

The expected revenue by displaying an add is

calculated as :- $\text{Bid} * \text{CTR}$

where Bid is the amount that the advertiser bid on the search query, and CTR is Click Through Rate i.e. the probability that once an ad is shown a user will click on the ~~other~~ ad.

Based on this $\text{Bid} * \text{CTR}$ value a search engine can choose which ~~old~~ ads to show. Various algorithms ~~are~~ like Greedy approach or Balance algorithm can be used in order to select the ads.

~~In the first screenshot~~

If we look into the screenshots of ~~of~~ of the first query string "t-shirt for men":-

→ Google displays adds from "The House" & "Bewakoof.com".

→ Yahoo displays adds from "Nykaa" & "Shoppers Stop". Along with that Yahoo also displays websites on the right.

In the second query "pulse oximeter":-

- Google displays various ads for Oximeter from various companies like "Groqit technologies", "Dr. Trust", "Flipkart" etc.
- In case of yahoo all the ads on the top are from "Healthmag.com". And there are other related websites that are shown on the right.

The difference in ads ~~may be due~~ in both search engines may be due to different bids from the companies. Like in the 2nd query, there is only one company in the Yahoo page results. The difference may also arise from different CTR on the different search engines.

The search engines may also use different strategies for optimizing their revenue, thus the order of ads appearing can also be different.

4. Shingles used by Google :-

- Central government
- discussions
- Pfizer
- American Pharmaceutical Company
- provide its COVID-19 vaccine possibly
- July, informed Dr. VK Patel, chair
- National Expert Group
- vaccine Administration (NEG-VAC)

Shingles used by Yahoo :-

- Pfizer
- pharmaceutical
- its COVID-19 vaccine.
- vaccine
- Nation Expert Group Vaccine Administration (NEG-VAC)

In Google, most of the articles on the first sequence of websites ~~almost~~ ~~about~~ have ~~comple~~ almost match with the paragraph. As we gradually go down, the ~~match~~ matches become ~~scarse~~ ~~scarse~~ more and more not less often..

On the other side in Yahoo, even the first website doesn't have a complete match with the paragraph. (This may be because the ~~new~~ ~~web~~ web pages containing the news article are not indexed yet by Yahoo).

But here also, as we ~~do~~ go down, relevance of the sites decreases.

The webpages ~~are~~ on the ~~internet~~ internet are indexed by the search engines. In order to match context measure relevance of content the content of the websites are converted to K-shingles. These shingles are then ~~are~~ hashed and ~~so~~, the hashes are stored.

When a search query is made the search string is broken into shingles ~~is~~ and converted into a bit array. This bit array is hashed & compared with ~~in~~ the hashes of indexed webpages. ~~is~~

B) Different search engines use different comparison algorithms. One such algorithm can be LSH (Locality Sensitive Hashing).

A) ~~Various~~ Various search engines use various algorithms in order to rank the web pages that have ~~the~~ ~~are~~ content similar to the ~~search~~ search query. Some of the factors that contribute to the ranking may include, ~~dates~~ region, degree of match with the search string, data, location of the search, user preferences etc.

5.

(A) → Authority (H) → Hub

a) "Kalyani University" in Google

1. www.kly.univ.ac.in → ~~A~~(A)
2. collegeuniv.com → (H)
3. en.wikipedia.com → (H)
4. West-Bengal.Result91.com → (H)
5. shiksha.com → (H)
6. educationuniv.com → (H)
7. exametc.com → (H)
8. career360.com → (H)
9. CollegeAdmissions.in → (H)
10. dl.acm.org → (H)
11. ~~india~~ www.indiatimes.com → (A)
12. timesofindia.indiatimes.com → (A)
13. researchgate.net → (H)
14. pincode.net.in → (H)
15. iitkalyani.ac.in → (A)

"Kalyani University" in Yahoo.

1. klyuniv.ac.in → (A)
2. collegeuniv.com → (H)
3. career360.com → ~~(A)~~(H)
4. ~~klyuniv.ac.in/conf~~
5. www.bruniv.ac.in → (A)
6. in.seekweb.com → (H)
7. in.zapmeta.search.icon → (H)
8. www.americancollegeSpain.com → (H)
9. www.agoda.com → (H)
10. www.icbse.com → (H)
11. educationiconnect.com → (H)
12. www.Sarvgyan.com → (H)
13. educationuniv.com → (H)
14. www.fresherslive.com → (H)
15. instituto.sdaaglasem.com → (H)

6). "Gurusaday Museum" in Google

1. en.wikipedia.org → (H)
2. www.gurusadaymuseum.org → (A)
3. scroll.in → (H)
4. www.tripadvisor.in → (H)
5. facebook.com → (H)
6. www.telegraphindia.com → (A)
7. www.museums-of-india.com → (H)
8. www.getkeral.com → (H)
9. www.ixigo.com → (H)
10. www.sahapedia.org → (H)
11. www.inspirack.com → (H)
12. www.anirudh.in → (A)
13. www.Kolkataonwheels.com → (H)
14. lbb.in → (H)
15. indiantribalheritage.org → (H)

"Gurusaday Museum" in Yahoo

1. ~~www.~~ museums-of-india.org → (H)
2. www.gurusadaymuseum.org → (A)
3. indianetzone.com → (H)
4. lbb.in → (H)
5. www.tripadvisor.in → (H)
6. www.~~republic~~republicworld.com → (H)
7. www.linehistoryIndia.com → (H)
8. transindiatravels.com → (H)
9. www.tripadvisor.com → (H)
10. www.ketto.org → (H)
11. ~~www.~~ museums-of-india.org → (H)
12. yahoo.com → (H)
13. ~~Ketto~~ www.Kolkatatonwheels.org → (H)
14. en.wikipedia.org → (H)
15. www.indianetzone.com → (H)

Good quality hubs have higher hub scores & good quality of content.

Hubs are pages that link to authorities.

Good hubs are those that have appropriate content and links to authorities. Ex - wikipedia.

Sites that have content but very few links to other authorities are called bad hubs.

Sites that have links to authorities but improper or unverified content are bad hubs, ex - social media sites.

The hubness 'h' and authoritativeness 'a' of a web page can be computed using HITS algorithm.

HITS algorithm

Each page i has 2 scores

authority score : a_i

Hub score : h_i

HITS algorithm :-

→ Initialize : $a_i^{(0)} = 1/\sqrt{N}$, $h_j^{(0)} = 1/\sqrt{N}$ for $j = 1 \dots N$;

→ Keep iterating till convergence?

$$\forall i : \text{Authority} : a_i^{(t+1)} = \sum_{j \rightarrow i} h_j^{(t)}$$

$$\forall i : \text{Hub} : h_i^{(t+1)} = \sum_{i \rightarrow j} a_j^{(t)}$$

$\forall i$: Normalize :

$$\sum_i (a_i^{(t+1)})^2 = 1, \sum_j (h_j^{(t+1)})^2 = 1$$

Here $h_i \rightarrow$ denotes the hubness of the web page

6. Both ~~on~~ Amazon and flipkart list a lot of books.

6. In Amazon's search results almost all the ~~re~~ books in the first page has the ~~words~~ phrase "They all fall down" in the book title.

But in flipkart all the books in the first page has the ~~names~~ ~~feet~~ phrase "fall down" ~~in~~

⁶ ~~All~~ "all fall down" in the books.

Books titled "They all fall down" appear only on the 6th row of the search ~~re~~ results.

Like most search engines, Amazon & Flipkart also perform searches based on key words.

Based on the above results one may ~~refer that~~ conclude that the difference in ~~output~~ search results of the two may be because flipkart is considering the word "they" as unimportant while searching for ~~at~~ keywords.

Both ~~on~~ Amazon & flipkart after filtering out books with the key words given in the search query, sort the books based on relevance (by default) relevance can ~~be~~ consist of various factors, like degree of match with the search ~~query~~ query, ratings, number of people who bought it, availability on a the location etc.

When we select a book on ~~@ Amazon~~

There were mostly 3 types of recommendation :-

1. "Frequently bought together" (Amazon) :-

These recommendations contain items that are very closely related to each other. Like other books in the series. These recommendation can be generated ~~as~~ using frequent item set mining algorithms. To find a frequent item set, we compare baskets which had selected books and ~~as~~ the frequent items in those baskets can be used as a recommendation. This can be done by ~~algo~~ algorithms like ~~Apro~~ A-priori Algorithm or PCY algorithm.

2. "Customers who bought this item also bought" (Amazon) :-

These recommendations are ~~as~~ generated by the help of user-user collaborative filtering. Here we filter the user who also ~~has~~ bought the selected books and based on their purchase history, display recommendations ~~the~~ products that are frequently bought.

3. "Products related to the item" (Amazon)

or

"Similar products" (Flipkart) :-

These recommendations are generated based on key words loosely related to the ~~current~~ book/product.

It can be done using searching based on key words or even content based ~~label~~ to ~~catalog~~ labeling.

7.

iii) Future of big data :-

- a) Business Business strategies :- take with the exponential increase in ~~existing~~ population & ~~now~~. The customers for businesses are also increasing. Thus is increasing data for businesses to handle. More and more businesses will understand the importance of big data analytics in providing a more tuned and streamlined products & services to the customers.
- Big data tools will become more & more common & easily ~~are~~ accessible even to ~~is~~ new small startups.

b) Releasess

Calamity Avoidance :- ~~The~~ The COVID-19 pandemic is not the first pandemic that the world has faced and ~~now~~ nor will it be the ~~is~~ last. Big data ~~can play~~ as ~~an~~ ~~over~~ crisis and other related distributed application can play a very crucial role in ~~avoiding~~ preventing and recovering from these natural calamities.

We may see more and more predictive analysis taking place and if there are more ~~as~~ pandemics in the future ~~now~~ prescriptive analysis can help.

Q) User Privacy & Security:- With more & more (22)
~~Big data~~ ~~Big data~~ applications
collecting user ~~per~~ data & ~~prefers~~
~~per~~ preferences using big data analytics
• data breaches are becoming
more & more common. ~~The~~ ~~eg~~ Various
organizations ~~like~~ like the EU & ~~society goes~~
~~govern~~ country governments & ~~with~~
are stepping in and making ~~bust~~ laws
on how user data is being
handled.