USN ☐☐☐☐☐☐☐☐☐☐



### Internal Assessment Test 3 – May 2019

| Sub: | Big Data Analytics | | | | | Sub Code: | 15CS82 | Branch: | CSE | | |
|------|------|------|------|------|------|------|------|------|------|------|------|
| Date: | 14/05/2019 | Duration: | 90 min's | Max Marks: | 50 | Sem / Sec: | | 8th A,B,C | | OBE | |

| | Answer any FIVE FULL Questions | MARKS | CO | RBT |
|------|------|------|------|------|
| 1 | What is Naïve Bayes technique? Explain the model with a simple classification example. | [3+7] | CO5 | L3 |
| 2 (a) | With a neat diagram explain the text mining process. How it is different from Data Mining. | [4+2] | CO5 | L2 |
| (b) | What is term document matrix? Explain with example. | [4] | CO5 | L2 |
| 3 | What is SVM? Explain the SVM model. Write advantages and disadvantages of SVMs. | [3+3+4] | CO5 | L2 |
| 4 | Describe the Influence Flow Model to compute the importance of a node and compute the rank values for the nodes of the following network.  | [4+2+4] | CO5 | L3 |
| 5 | Describe the three different types of web mining. | [10] | CO5 | L2 |
| 6 | Explain the process of Data import and export in Sqoop with neat diagrams. Compare the Version 1 and Version 2 of Sqoop | [7+3] | CO2 | L2 |
| 7. | What is Apache Flume? Explain the components of a Flume agent with neat diagrams. | [10] | CO2 | L2 |
| 8 | With a neat diagram explain the Hadoop version 2 ecosystem. | [10] | CO2 | L2 |

| Course Outcomes | | Modules covered | PO1 | PO2 | PO3 | PO4 | PO5 | PO6 | PO7 | PO8 | PO9 | PO10 | PO11 | PO12 | PSO1 | PSO2 | PSO3 | PSO4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CO 1 | Master the concepts of HDFS and MapReduce framework | 1 | 1 | 3 | 2 | 3 | 3 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 2 | 3 |
| CO 2 | Investigate Hadoop related tools for Big Data Analytics and perform basic Hadoop Administration | 1, 2 | 0 | 3 | 2 | 2 | 3 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 2 | 3 |
| CO 3 | Recognize the role of Business Intelligence, Data warehousing and Visualization in decision making | 3 | 1 | 2 | 1 | 2 | 1 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 2 |
| CO 4 | Infer the importance of core data mining techniques for data analytics | 3, 4 | 2 | 3 | 3 | 3 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 2 |
| CO 5 | Analyze Data Mining Techniques | 3, 5 | 2 | 2 | 3 | 3 | 1 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 2 | 3 |

**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING**

Date: 28th July 2018

**CO's to PO's & PSO's mapping**

| | | | | |
|---|---|---|---|---|
| Name of the course | : | **Big data analytics** | Sub Code : | 15CS82 |
| Name of the Faculty/s | : | Mrs Poonam | Sem & Sec : | 8th A,B,C |

SCHEME

| Question # | Description | Marks Distribution | | Max Marks |
|---|---|---|---|---|
| 1 | Naïve Bayes technique<br>Explanation of the model<br>Classification example | 3M<br>5M<br>2M | 10M | 10M |
| 2 a | The text mining process<br>Difference from Data Mining | 3M<br>2M | 5M | 5M |
| 2 b | Term document matrix<br>Example | 3M<br>1M | 4M | 4M |
| 3 | SVM<br>Explanation of the SVM model.<br>advantages and disadvantages of SVMs. | 3M<br>3M<br>4M | 10M | 10M |
| 4 | Describe the Influence Flow Model<br>compute the rank values for the nodes | 4M<br>2M<br>4M | 10M | 10M |
| 5 | the three different types of web mining | 10M | 10M | 10M |
| 6 | Process of Data import and export in Sqoop with neat diagrams. | 7M<br>3M | 10M | 10M |

| | | | | |
|---|---|---|---|---|
| | Comparison of Sqoop the Version 1 and Version 2 | | | |
| 7 | Apache Flume<br>Components of a Flume agent with neat diagrams. | 2M<br>8M | 10M | 10M |
| 8 | Hadoop version 2 ecosystem.<br>neat diagram | 8M<br>2M | 10M | 10M |

**Solution**

**Q1> What is Naïve Bayes technique?**

It is a classification technique based on Bayes' Theorem with an assumption of independence among predictors. In simple terms, a Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature. For example, a fruit may be considered to be an apple if it is red, round, and about 3 inches in diameter. Even if these features depend on each other or upon the existence of the other features, all of these properties independently contribute to the probability that this fruit is an apple and that is why it is known as 'Naive'.

Naive Bayes model is easy to build and particularly useful for very large data sets. Along with simplicity, Naive Bayes is known to outperform even highly sophisticated classification methods.

Bayes theorem provides a way of calculating posterior probability $P(c|x)$ from $P(c)$, $P(x)$ and $P(x|c)$. Look at the equation below:

Bayes_rule-300x172Above,

$P(c|x)$ is the posterior probability of class (c, target) given predictor (x, attributes).

$P(c)$ is the prior probability of class.

$P(x|c)$ is the likelihood which is the probability of predictor given class.

$P(x)$ is the prior probability of predictor.

How Naive Bayes algorithm works?

Let's understand it using an example. Below I have a training data set of weather and corresponding target variable 'Play' (suggesting possibilities of playing). Now, we need to classify whether players will play or not based on weather condition. Let's follow the below steps to perform it.

Step 1: Convert the data set into a frequency table

Step 2: Create Likelihood table by finding the probabilities like Overcast probability = 0.29 and probability of playing is 0.64.

Bayes_4

Step 3: Now, use Naive Bayesian equation to calculate the posterior probability for each class. The class with the highest posterior probability is the outcome of prediction.

Problem: Players will play if weather is sunny. Is this statement is correct?

We can solve it using above discussed method of posterior probability.

$P(Yes | Sunny) = P( Sunny | Yes) * P(Yes) / P (Sunny)$

Here we have P (Sunny |Yes) = 3/9 = 0.33, P(Sunny) = 5/14 = 0.36, P( Yes)= 9/14 = 0.64

Now, P (Yes | Sunny) = 0.33 * 0.64 / 0.36 = 0.60, which has higher probability.

Naive Bayes uses a similar method to predict the probability of different class based on various attributes. This algorithm is mostly used in text classification and with problems having multiple classes.

**Q2> With a neat diagram explain the text mining process. How it is different from Data Mining.**

Text Mining is a rapidly evolving area of research. As the amount of social media and other text data grows, there is need for efficient abstraction and categorization of meaningful information from the text.

The first level of analysis is identifying frequent words. This creates a bag of important words. Texts – documents or smaller messages – can then be ranked on how they match to a particular bag-of-words. However, there are challenges with this approach. For example, the words may be spelled a little differently. Or there may be different words with similar meanings. The next level is at the level of identifying meaningful phrases from words. Thus _ice' and _cream' will be two different key words that often come together. However, there is a more meaningful phrase by combining the two words into _ice cream'. There might be similarly meaningful phrases like _Apple Pie'.

The next higher level is that of Topics. Multiple phrases could be combined into Topic area. Thus the two phrases above could be put into a common basket, and this bucket could be called _Desserts'.

Text mining is a semi-automated process. Text data needs to be gathered, structured, and then mined, in a 3-step process.
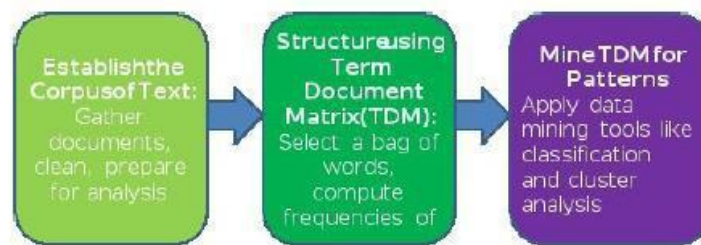


Fig : Text Mining Architecture

**Q2> What is term document matrix? Explain with example.**

A document-term matrix or term-document matrix is a mathematical matrix that describes the frequency of terms that occur in a collection of documents. In a document-term matrix, rows correspond to documents in the collection and columns correspond to terms. There are various schemes for determining the value that each entry in the matrix should take. One such scheme is tf-idf. They are useful in the field of natural language processing.

When creating a database of terms that appear in a set of documents the document-term matrix contains rows corresponding to the documents and columns corresponding to the terms. For instance if one has the following two (short) documents:

- D1 = "I like databases"

- D2 = "I hate databases",
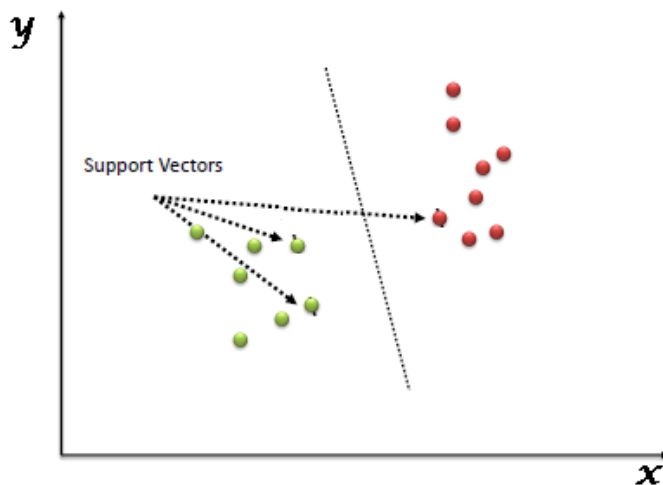
then the document-term matrix would be:

|     | I | like | hate | databases |
| --- | --- | --- | --- | --- |
| **D1** | 1 | 1 | 0 | 1 |
| **D2** | 1 | 0 | 1 | 1 |

which shows which documents contain which terms and how many times they appear.

Note that more sophisticated weights can be used; one typical example, among others, would be tf-idf.

**Q3> What is SVM? Explain the SVM model. Write advantages and disadvantages of SVMs.**

"Support Vector Machine" (SVM) is a supervised machine learning algorithm which can be used for both classification or regression challenges. However, it is mostly used in classification problems. In this algorithm, we plot each data item as a point in n-dimensional space (where n is number of features you have) with the value of each feature being the value of a particular coordinate. Then, we perform classification by finding the hyper-plane that differentiate the two classes very well (look at the below snapshot).



Support Vectors are simply the co-ordinates of individual observation. Support Vector Machine is a frontier which best segregates the two classes (hyper-plane/ line).
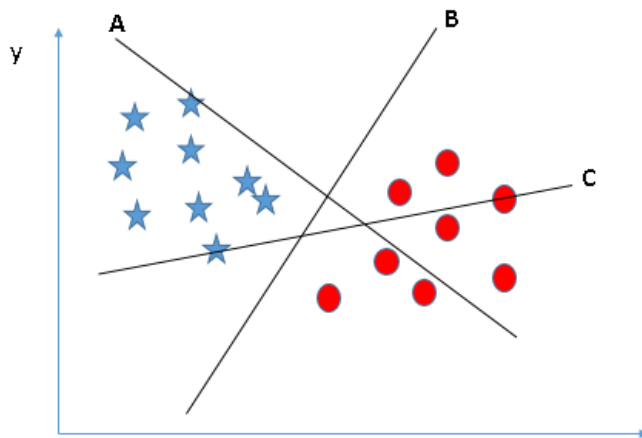
You can look at support vector machines and a few examples of its working here.

How does it work?

Above, we got accustomed to the process of segregating the two classes with a hyper-plane. Now the burning question is "How can we identify the right hyper-plane?". Don't worry, it's not as hard as you think!
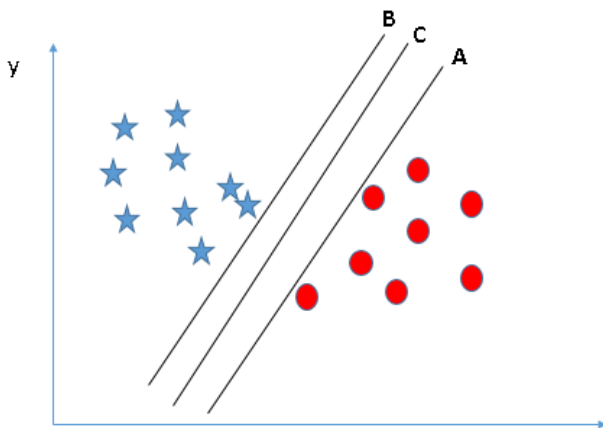
Let's understand:

- Identify the right hyper-plane (Scenario-1): Here, we have three hyper-planes (A, B and C). Now, identify the right hyper-plane to classify star and circle.
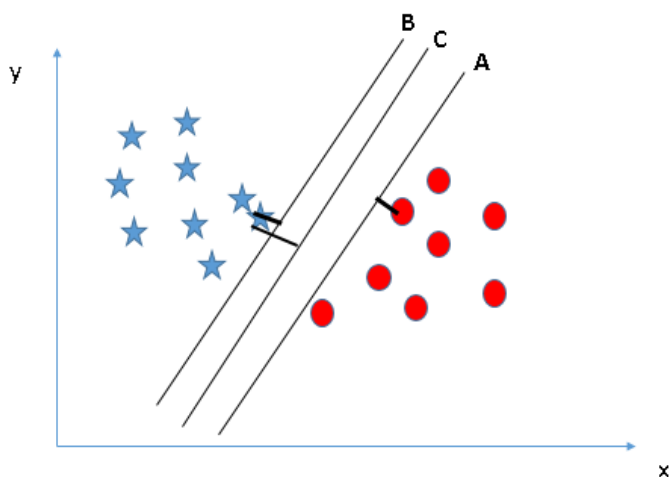


×  You need to remember a thumb rule to identify the right hyper-plane: "Select the hyper-plane which segregates the two classes better". In this scenario, hyper-plane "B" has excellently performed this job.

- Identify the right hyper-plane (Scenario-2): Here, we have three hyper-planes (A, B and C) and all are segregating the classes well. Now, How can we identify the right hyper-plane?



×  Here, maximizing the distances between nearest data point (either class) and hyper-plane will help us to decide the right hyper-plane. This distance is called as Margin. Let's look at the below snapshot:
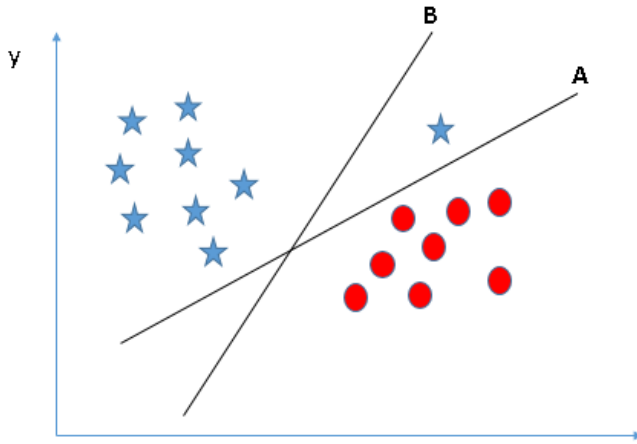


Above, you can see that the margin for hyper-plane C is high as compared to both A and B. Hence,
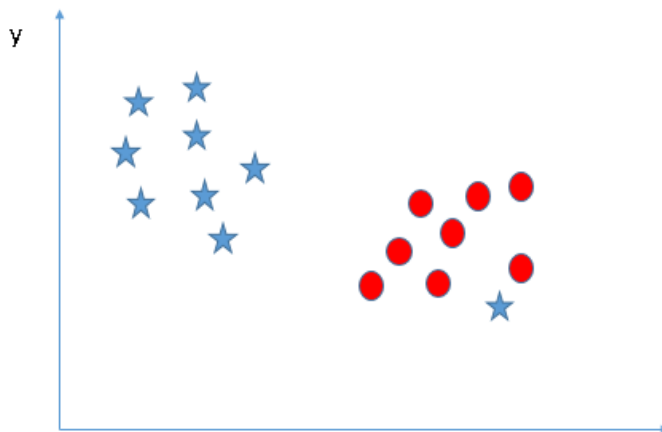
we name the right hyper-plane as C. Another lightning reason for selecting the hyper-plane with higher margin is robustness. If we select a hyper-plane having low margin then there is high chance of miss-classification.

- Identify the right hyper-plane (Scenario-3):Hint: Use the rules as discussed in previous section to identify the right hyper-plane
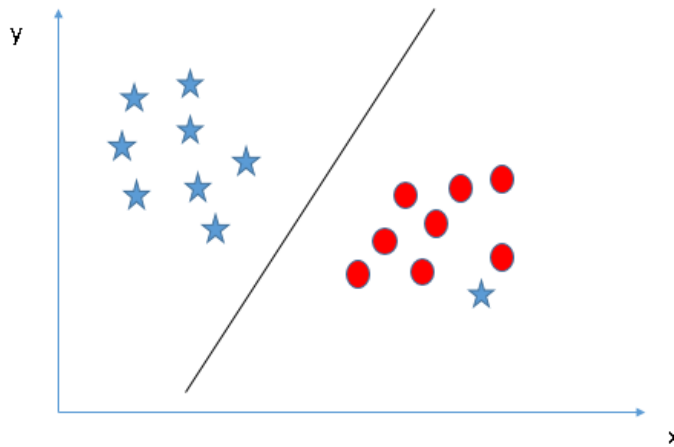


Some of you may have selected the hyper-plane B as it has higher margin compared to A. But, here is the catch, SVM selects the hyper-plane which classifies the classes accurately prior to maximizing margin. Here, hyper-plane B has a classification error and A has classified all correctly. Therefore, the right hyper-plane is A.

- Can we classify two classes (Scenario-4)?: Below, I am unable to segregate the two classes using a straight line, as one of star lies in the territory of other(circle) class as an outlier.
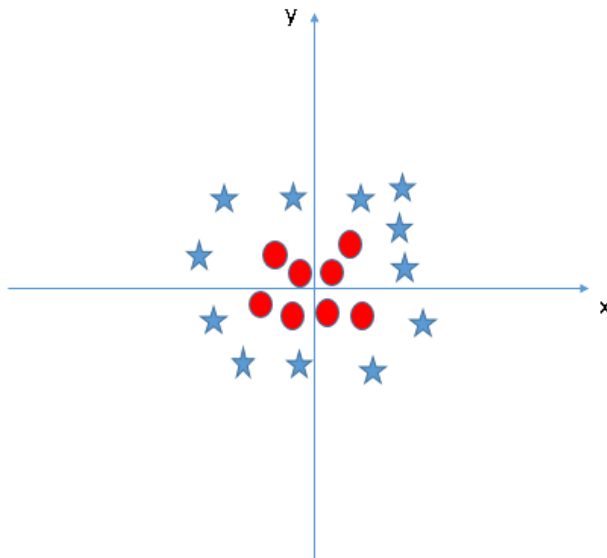


As I have already mentioned, one star at other end is like an outlier for star class. SVM has a feature to ignore outliers and find the hyper-plane that has maximum margin. Hence, we can say, SVM is robust to
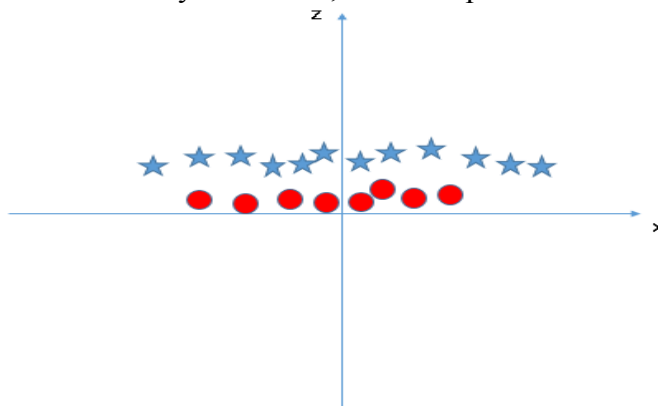
outliers.



- Find the hyper-plane to segregate to classes (Scenario-5): In the scenario below, we can't have linear hyper-plane between the two classes, so how does SVM classify these two classes? Till now, we have only looked at the linear hyper-plane.



SVM can solve this problem. Easily! It solves this problem by introducing additional feature. Here, we will add a new feature $z = x^2 + y^2$. Now, let's plot the data points on axis x and z:
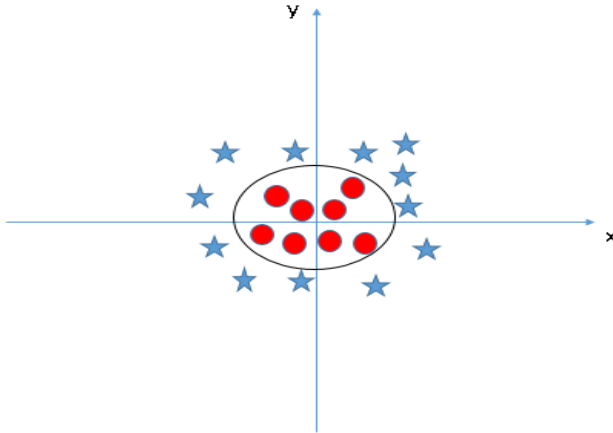


In above plot, points to consider are:

- All values for z would be positive always because z is the squared sum of both x and y

o   In the original plot, red circles appear close to the origin of x and y axes, leading to lower value of z and star relatively away from the origin result to higher value of z.

In SVM, it is easy to have a linear hyper-plane between these two classes. But, another burning question which arises is, should we need to add this feature manually to have a hyper-plane. No, SVM has a technique called the kernel trick. These are functions which takes low dimensional input space and transform it to a higher dimensional space i.e. it converts not separable problem to separable problem, these functions are called kernels. It is mostly useful in non-linear separation problem. Simply put, it does some extremely complex data transformations, then find out the process to separate the data based on the labels or outputs you've defined.

When we look at the hyper-plane in original input space it looks like a circle:



Now, let's look at the methods to apply SVM algorithm in a data science challenge.

SVM Advantages

- SVM's are very good when we have no idea on the data.

- Works well with even unstructured and semi structured data like text, Images and trees.

- The kernel trick is real strength of SVM. With an appropriate kernel function, we can solve any complex problem.

- Unlike in neural networks, SVM is not solved for local optima.

- It scales relatively well to high dimensional data.

- SVM models have generalization in practice, the risk of over-fitting is less in SVM.

- SVM is always compared with ANN. When compared to ANN models, SVMs give better results.

SVM Disadvantages

- Choosing a "good" kernel function is not easy.

- Long training time for large datasets.

- Difficult to understand and interpret the final model, variable weights and individual impact.

- Since the final model is not so easy to see, we can not do small calibrations to the model hence its tough to incorporate our business logic.

- The SVM hyper parameters are Cost -C and gamma. It is not that easy to fine-tune these hyper-parameters. It is hard to visualize their impact

Q4>

**Q5>**

### Web Mining

Web mining is the art and science of discovering patterns and insights from the World-wide web so as to improve it. The world-wide web is at the heart of the digital revolution. More data is posted on the web every day than was there on the whole web just 20 years ago. Billions of users are using it every day for a variety of purposes. The web is used for electronic commerce, business communication, and many other applications. Web mining analyzes data from the web and helps find insights that could optimize the web content and improve the user experience. Data for web mining is collected via Web crawlers, web logs, and other means. Here are some characteristics of optimized websites:

1. *Appearance*: Aesthetic design. Well-formatted content, easy to scan and navigate. Good color contrasts.
2. *Content*: Well planned information architecture with useful content. Fresh content. Search-engine optimized. Links to other good sites.
3. *Functionality*: Accessible to all authorized users. Fast loading times. Usable forms. Mobile enabled.

This type of content and its structure is of interest to ensure the web is easy to use. The analysis of web usage provides feedback on the web content, and also the consumer's browsing habits. This data can be of immense use for commercial advertising, and even for social engineering. The web could be analyzed for its structure as well as content. The usage pattern of web pages could also be analyzed. Depending upon objectives, web mining can be divided into three different types: Web usage mining, Web content mining and Web structure mining .


Fig: 5.2 Web Mining structure

### Web content mining

A website is designed in the form of pages with a distinct URL (universal resource locator). A large website may contain thousands of pages. These pages and their content is managed using specialized software systems called Content Management Systems. Every page can have text, graphics, audio, video, forms, applications, and more kinds of content including user generated content.

The websites keep a record of all requests received for its page/URLs, including the requester information using _cookies'. The log of these requests could be analyzed to gauge the popularity of those pages among different segments of the population. The text and application content on the pages could be analyzed for its usage by visit counts. The pages on a website themselves could be analyzed for quality of content that attracts most users.

Thus the unwanted or unpopular pages could be weeded out, or they can be transformed with different content and style. Similarly, more resources could be assigned to keep the more popular pages more fresh and inviting.

## Web structure mining

The Web works through a system of hyperlinks using the hypertext protocol (http). Any page can create a hyperlink to any other page, it can be linked to by another page. The intertwined or self-referral nature of web lends itself to some unique network analytical algorithms. The structure of Web pages could also be analyzed to examine the pattern of hyperlinks among pages. There are two basic strategic models for successful websites: Hubs and Authorities.

1. *Hubs*: These are pages with a large number of interesting links. They serve as a hub, or a gathering point, where people visit to access a variety of information. Media sites like Yahoo.com, or government sites would serve that purpose. More focused sites like Traveladvisor.com and yelp.com could aspire to becoming hubs for new emerging areas.

2. *Authorities*: Ultimately, people would gravitate towards pages that provide the most complete and authoritative information on a particular subject. This could be factual information, news, advice, user reviews etc. These websites would have the most number of inbound links from other websites. Thus Mayoclinic.com would serve as an authoritative page for expert medical opinion. NYtimes.com would serve as an authoritative page for daily news.

## Web usage mining

As a user clicks anywhere on a webpage or application, the action is recorded by many entities in many locations. The browser at the client machine will record the click, and the web server providing the content would also make a record of the pages served and the user activity on those pages. The entities between the client and the server, such as the router, proxy server, or ad server, too would record that click.

The goal of web usage mining is to extract useful information and patterns from data generated through Web page visits and transactions. The activity data comes from data stored in server access logs, referrer logs, agent logs, and client-side cookies. The user characteristics and usage profiles are also gathered directly, or indirectly, through syndicated data. Further, metadata, such as page attributes, content attributes, and usage data are also gathered.

The web content could be analyzed at multiple levels

1. The *server side analysis* would show the relative popularity of the web pages accessed. Those websites could be hubs and authorities.

2. The *client side analysis* could focus on the usage pattern or the actual content consumed and created by users.

    1. Usage pattern could be analyzed using ‗clickstream' analysis, i.e.analyzing web activity for patterns of sequence of clicks, and the location and duration of visits on websites. Clickstream analysis can be useful for web activity analysis, software testing, market research, and analyzing employee productivity.

    3. Textual information accessed on the pages retrieved by users could be analyzed using text mining techniques. The text would be gathered and structured using the bag-of-words technique to build a Term-document matrix. This matrix could then be mined using

cluster analysis and association rules for patterns such as popular topics, user segmentation, and sentiment analysis.



Fig: 5.3 Web Usage Mining architecture

Web usage mining has many business applications. It can help predict user behavior based on previously learned rules and users' profiles, and can help determine lifetime value of clients. It can also help design cross-marketing strategies across products, by observing association rules among the pages on the website. Web usage can help evaluate promotional campaigns and see if the users were attracted to the website and used the pages relevant to the campaign. Web usage mining could be used to present dynamic information to users based on their interests and profiles. This includes targeted online ads and coupons at user groups based on user access patterns.

Q6> Explain the process of Data import and export in Sqoop with neat diagrams. Compare the Version 1 and Version 2 of Sqoop

| Feature | Sqoop 1 | Sqoop 2 |
|---|---|---|
| Connectors for all major RDBMS | Supported. | Not supported.<br>**Workaround**: Use the generic JDBC Connector which has been tested on the following databases: Microsoft SQL Server, PostgreSQL, MySQL and Oracle.<br>This connector should work on any other JDBC compliant database. However, performance might not be comparable to that of specialized connectors in Sqoop. |
| Kerberos Security | Supported. | Supported. |

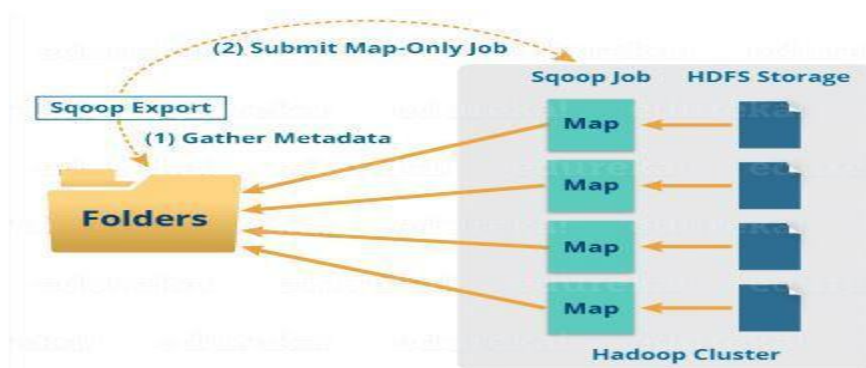| Feature | Sqoop 1 | Sqoop 2 |
|---|---|---|
| Integration | | |
| Data transfer from RDBMS to Hive or HBase | Supported. | Not supported.<br><br>**Workaround:** Follow this two-step approach.<br><br>1. Import data from RDBMS into HDFS<br>2. Load data into Hive or HBase manually using appropriate tools and commands such as the LOAD DATA statement in Hive |
| Data transfer from Hive or HBase to RDBMS | Not supported.<br><br>**Workaround:** Follow this two-step approach.<br><br>1. Extract data from Hive or HBase into HDFS (either as a text or Avro file)<br>2. Use Sqoop to export output of previous step to RDBMS | Not supported.<br><br>Follow the same workaround as for Sqoop 1. |

*Fig : Sqoop WorkFlow*

Sqoop is a tool designed to transfer data between Hadoop and relational databases. You can use Sqoop to import data from a relational database management system (RDBMS) into the Hadoop Distributed File System (HDFS), transform the data in Hadoop, and then export the data back into an RDBMS.

Sqoop can be used with any Java Database Connectivity (JDBC)–compliant database and has been tested on Microsoft SQL Server, PostgresSQL, MySQL, and Oracle

When we submit Sqoop command, our main task gets divided into sub tasks which is handled by individual Map Task internally. Map Task is the sub task, which imports part of data to the Hadoop Ecosystem. Collectively, all Map tasks imports the whole data

Export also works in a similar manner.

When we submit our Job, it is mapped into Map Tasks which brings the chunk of data from HDFS. These chunks are exported to a structured data destination. Combining all these exported chunks of data, we receive the whole data at the destination, which in most of the cases is an RDBMS (MYSQL/Oracle/SQL Server).

**Q7. What is Apache Flume? Explain the components of a Flume agent with neat diagrams.**

Apache Flume is an independent agent designed to collect, transport, and store data into HDFS. Often data transport involves a number of Flume agents that may traverse a series of machines and locations. Flume is often used for log files, social media-generated data, email messages, and just about any continuous data source.

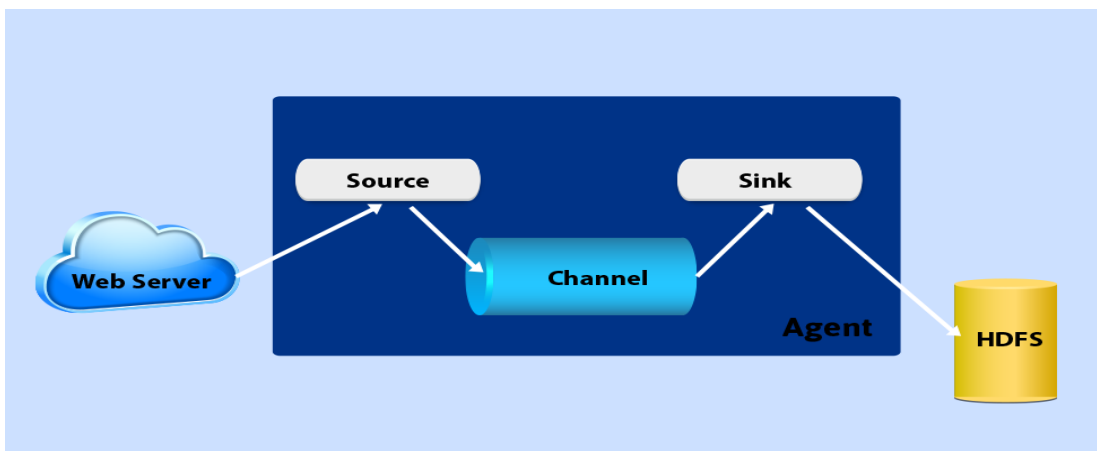As shown in Figure 2.3, a Flume agent is composed of three components.



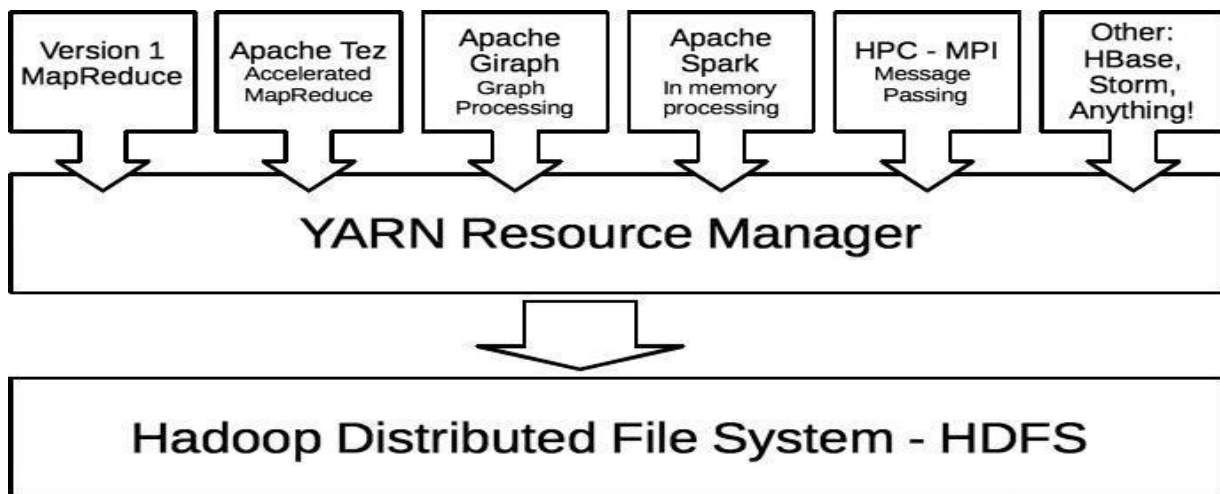Fig 2.3. Flume agent with source, channel, and sink

• Source. The source component receives data and sends it to a channel. It can send the data to more than one channel. The input data can be from a real-time source (e.g., weblog) or another Flume agent.

• Channel. A channel is a data queue that forwards the source data to the sink destination. It can be thought of as a buffer that manages input (source) and output (sink) flow rates.

• Sink. The sink delivers data to destination such as HDFS, a local file, or another Flume agent.

A Flume agent must have all three of these components defined. A Flume agent can have several sources, channels, and sinks. Sources can write to multiple channels, but a sink can take data from only a single channel. Data written to a channel remain in the channel until a sink removes the data. By default, the data in a channel are kept in memory but may be optionally stored on disk to prevent data loss in the event of a network failure.

**Q8> With a neat diagram explain the Hadoop version 2 ecosystem.**

YARN presents a resource management platform, which provides services such as scheduling, fault monitoring, data locality, and more to MapReduce and other frameworks. Figure 7 illustrates some of the various frameworks that will run under YARN. Note that the Hadoop version 1 applications (e.g., Pig and Hive) run under the MapReduce framework.



*Fig 2.7 Example of the Hadoop version 2 ecosystem.*

**Distributed-Shell**

As described earlier in this chapter, Distributed-Shell is an example application included with the Hadoop core components that demonstrates how to write applications on top of YARN. It provides a simple method for running shell commands and scripts in containers in parallel on a Hadoop YARN cluster.

**Hadoop MapReduce**

MapReduce was the first YARN framework and drove many of YARN's requirements. It is integrated tightly with the rest of the Hadoop ecosystem projects, such as Apache Pig, Apache Hive, and Apache Oozie.

**Apache Tez**

One great example of a new YARN framework is Apache Tez. Many Hadoop jobs involve the execution of a complex directed acyclic graph (DAG) of tasks using separate MapReduce stages. Apache Tez generalizes this process and enables these tasks to be spread across stages so that they can be run as a single, all-encompassing job. Tez can be used as a MapReduce replacement for

projects such as Apache Hive and Apache Pig. No changes are needed to the Hive or Pig applications.

**Apache Giraph**

Apache Giraph is an iterative graph processing system built for high scalability. Facebook, Twitter, and LinkedIn use it to create social graphs of users. Giraph was originally written to run on standard Hadoop V1 using the MapReduce framework, but that approach proved inefficient and totally unnatural for various reasons.. In addition, using the flexibility of YARN, the Giraph developers plan on implementing their own web interface to monitor job progress.

**Hoya: HBase on YARN**

The Hoya project creates dynamic and elastic Apache HBase clusters on top of YARN. A client application creates the persistent configuration files, sets up the HBase cluster XML files, and then asks YARN to create an ApplicationMaster. YARN copies all files listed in the client's application-launch request from HDFS into the local file system of the chosen server, and then executes the command to start the Hoya ApplicationMaster. Hoya also asks YARN for the number of containers matching the number of HBase region servers it needs.

**Apache Spark**

Spark was initially developed for applications in which keeping data in memory improves performance, such as iterative algorithms, which are common in machine learning, and interactive data mining. Spark differs from classic MapReduce in two important ways. First, Spark holds intermediate results in memory, rather than writing them to disk. Second, Spark supports more than just MapReduce functions; that is, it greatly expands the set of possible analyses that can be executed over HDFS data stores. It also provides APIs in Scala, Java, and Python.

Since 2013, Spark has been running on production YARN clusters at Yahoo!. The advantage of porting and running Spark on top of YARN is the common resource management and a single underlying file system

**Apache Storm**

Traditional MapReduce jobs are expected to eventually finish, but Apache Storm continuously processes messages until it is stopped. This framework is designed to process unbounded streams of data in real time. It can be used in any programming language. The basic Storm use-cases include real-time analytics, online machine learning, continuous computation, distributed RPC (remote procedure calls), ETL (extract, transform, and load), and more. Storm provides fast performance, is scalable, is fault tolerant, and provides processing guarantees. It works directly under YARN and takes advantage of the common data and resource management substrate.