

Unit -1

Q. What is data and Characteristics of data.

Data is like the raw material of the digital world—bits and pieces of information that, when organized and analyzed, can reveal meaningful insights. It can be anything from numbers and text to images and videos. Now, when it comes to characteristics, here are a few key ones:

1. Volume: There's a lot of it! Big data is called "big" for a reason. We're talking about massive amounts of information.
2. Velocity: Data comes at us fast and furious. It's not just about the quantity, but also the speed at which it's generated and needs to be processed.
3. Variety: It's not all 1s and 0s. Data comes in different formats—structured (like databases), unstructured (social media posts), and semi-structured (XML files). It's a diverse bunch.
4. Veracity: Data quality matters. With so much information out there, ensuring accuracy and reliability is a challenge.
5. Value: The ultimate goal is to extract meaningful insights from data. It's not just about having a lot of information; it's about turning that information into valuable knowledge.

So, data is like a treasure trove, but you've got to navigate through the complexities to find the gold.

Q Types of Digital Data.

****Structured Data:****

This is the well-behaved, organized type. Think of it like a spreadsheet or a database. It's neatly arranged in tables with rows and columns, making it easy to search, sort, and analyze. Each piece of information has a designated place, and it follows a clear and predefined structure.

****Unstructured Data:****

Now, imagine the rebel of the data world. Unstructured data doesn't conform to a fixed data model. It's like a wild jungle of information—texts, images, videos, social media posts. It's messy, but it's also rich with insights. The challenge is taming this wild side to extract meaningful patterns.

****Semi-Structured Data:****

BDA

Semi-structured data is like the middle child. It has some structure, but it's not as rigid as structured data. It may have tags, markers, or keys that help organize the information, but it doesn't follow a strict schema. Examples include XML files or JSON documents. They provide a bit of order to the chaos while allowing for flexibility.

In the grand data scheme, each type plays a unique role. Structured data is the neat and tidy librarian, unstructured data is the creative artist, and semi-structured data is the hybrid, finding balance between order and chaos.

Q. Sources of Data

Data originates from a multitude of sources, creating a rich landscape of information. Here's a brief overview:

- **Databases:** Structured data is housed in databases, organized in tables for efficient storage and retrieval.
 - **Social Media:** Platforms like Facebook, Twitter, and Instagram generate vast amounts of unstructured data through posts, comments, and multimedia content.
 - **Websites:** User interactions, clicks, and searches on websites contribute valuable data.
 - **Sensors and IoT Devices:** Smart gadgets and sensors continuously collect and transmit data, providing insights into various environments.
 - **Surveys and Forms:** Gathering information through feedback forms, surveys, and questionnaires.
 - **Government Records:** Census data, public records, and official databases serve as comprehensive sources for analysis.
 - **Machine-generated Data:** Automated systems, logs, and machine-to-machine communication contribute to the data stream.
 - **Personal Devices:** Smartphones, tablets, and computers generate data, including location data, app usage, and more.
 - **Audio and Video Sources:** Multimedia platforms like podcasts, YouTube, and streaming services contribute audio and video data.
 - **Financial Transactions:** Credit card purchases, online transactions, and financial records offer insights into economic activities.
-

Q. Handling Unstructured Data

BDA

Handling unstructured data requires a strategic approach to extract meaningful insights. Here are some methods:

1. Text Mining/Natural Language Processing (NLP): Use NLP techniques to analyze and understand patterns in textual data. This includes sentiment analysis, keyword extraction, and language modeling.
 2. Machine Learning Algorithms: Train machine learning models to categorize and classify unstructured data. This could involve image recognition, speech-to-text, or even content categorization.
 3. Data Tagging and Annotation: Manually or automatically tag and annotate unstructured data to add structure. This helps in organizing and categorizing information for future analysis.
 4. Content Analysis: Break down unstructured data into components for deeper understanding. This could involve examining the elements within text, images, or videos to derive insights.
 5. Clustering and Segmentation: Group similar unstructured data together based on inherent similarities. This helps in identifying patterns and trends within the data.
 6. Data Lake or Data Hub: Create a centralized repository for storing unstructured data, allowing for easy access and analysis. This is especially useful when dealing with large volumes of diverse information.
 7. Graph Databases: Use graph databases to represent relationships within unstructured data. This is valuable for scenarios where understanding connections between different data points is crucial.
 8. Semantic Analysis: Understand the meaning and context behind unstructured data by employing semantic analysis techniques. This helps in deciphering the intended meaning of words or phrases.
-

Q. What is Big Data?

Big Data refers to the massive volume of structured and unstructured data that is generated at an unprecedented rate in our digital world. The term is not just about the size of the data but also encompasses the challenges and opportunities of processing, analyzing, and deriving meaningful insights from such vast and complex datasets.

The characteristics of Big Data are often described using the three Vs:

BDA

1. ****Volume:**** Big Data involves large amounts of data, often ranging from terabytes to petabytes and beyond. This sheer volume requires specialized processing capabilities.
 2. ****Velocity:**** Data is generated at high speed in real-time or near-real-time. This could include social media interactions, sensor data, and other sources that produce a constant flow of information.
 3. ****Variety:**** Data comes in various formats, including structured data (like databases), unstructured data (such as text and images), and semi-structured data (like XML files). Managing this diversity poses a unique challenge.
-

Evolution of Big Data:

If we see the last few decades, we can analyze that Big Data technology has gained so much growth. There are a lot of milestones in the evolution of Big Data which are described below:

1. **Data Warehousing:**

In the 1990s, data warehousing emerged as a solution to store and analyze large volumes of structured data.

2. **Hadoop:**

Hadoop was introduced in 2006 by Doug Cutting and Mike Cafarella. Distributed storage medium and large data processing are provided by Hadoop, and it is an open-source framework.

3. **NoSQL Databases:**

In 2009, NoSQL databases were introduced, which provide a flexible way to store and retrieve unstructured data.

4. **Cloud Computing:**

Cloud Computing technology helps companies to store their important data in data centers that are remote, and it saves their infrastructure cost and maintenance costs.

5. **Machine Learning:**

Machine Learning algorithms are those algorithms that work on large data, and analysis is done on a huge amount of data to get meaningful insights from it. This has led to the development of artificial intelligence (AI) applications.

6. **Data Streaming:**

Data Streaming technology has emerged as a solution to process large volumes of data in real time.

7. **Edge Computing:**

Edge Computing is a kind of distributed computing paradigm that allows data processing to be done at the edge or the corner of the network, closer to the source of the data.

Overall, big data technology has come a long way since the early days of data warehousing. The introduction of Hadoop, NoSQL databases, cloud computing, machine learning, data streaming, and edge computing has revolutionized how we store, process, and analyze large volumes of data. As technology evolves, we can expect Big Data to play a very important role in various industries.

Importance of Big data

Big Data importance doesn't revolve around the amount of data a company has. Its importance lies in the fact that how the company utilizes the gathered data.

Every company uses its collected data in its own way. More effectively the company uses its data, more rapidly it grows.

The companies in the present market need to collect it and analyze it because:

1. Cost Savings

Big Data tools like Apache Hadoop, Spark, etc. bring cost-saving benefits to businesses when they have to store large amounts of data. These tools help organizations in identifying more effective ways of doing business.

2. Time-Saving

Real-time in-memory analytics helps companies to collect data from various sources. Tools like Hadoop help them to analyze data immediately thus helping in making quick decisions based on the learnings.

3. Understand the market conditions

Big Data analysis helps businesses to get a better understanding of market situations.

For example, analysis of customer purchasing behavior helps companies to identify the products sold most and thus produces those products accordingly. This helps companies to get ahead of their competitors.

4. Social Media Listening

Companies can perform sentiment analysis using Big Data tools. These enable them to get feedback about their company, that is, who is saying what about the company.

Companies can use Big data tools to improve their online presence.

5. Boost Customer Acquisition and Retention

Customers are a vital asset on which any business depends on. No single business can achieve its success without building a robust customer base. But even with a solid customer base, the companies can't ignore the competition in the market.

If we don't know what our customers want then it will degrade companies' success. It will result in the loss of clientele which creates an adverse effect on business growth.

Big data analytics helps businesses to identify customer related trends and patterns. Customer behavior analysis leads to a profitable business.

6. Solve Advertisers Problem and Offer Marketing Insights

Big data analytics shapes all business operations. It enables companies to fulfill customer expectations. Big data analytics helps in changing the company's product line. It ensures powerful marketing campaigns.

Q. Challenges of Big Data

Big Data, while immensely valuable, comes with its fair share of challenges. Here are some key challenges associated with handling and analyzing large datasets:

BDA

1. ****Volume Overload:****

- Dealing with the sheer volume of data can be overwhelming. Storage, processing, and managing the vast amounts of information pose significant infrastructure challenges.

2. ****Velocity of Data:****

- The speed at which data is generated, processed, and analyzed in real-time can be challenging. Traditional systems may struggle to keep up with the constant flow of data.

3. ****Variety of Data:****

- Big Data encompasses various data types, including structured, unstructured, and semi-structured data. Integrating and analyzing this diverse range can be complex.

4. ****Veracity - Data Quality:****

- Ensuring the quality and reliability of data is a significant challenge. Inaccurate or incomplete data can lead to faulty insights and decisions.

5. ****Data Security and Privacy:****

- With the abundance of sensitive information, maintaining data security and ensuring privacy compliance become critical. Unauthorized access and data breaches are constant concerns.

6. ****Lack of Skilled Professionals:****

- There's a shortage of professionals with the necessary skills to manage and analyze Big Data. The rapidly evolving nature of technology adds to the challenge of keeping skills up-to-date.

7. ****Integration Across Platforms:****

- Integrating Big Data technologies with existing systems and platforms can be complex. Ensuring seamless communication and compatibility is a continual challenge.

8. ****Cost Management:****

- The infrastructure and tools required for Big Data can be expensive. Managing and optimizing costs while ensuring performance is a balancing act.

Q. Traditional Data vS Big Data

BDA

Traditional Data	Big Data
Traditional data is generated in enterprise level.	Big data is generated outside the enterprise level.
Its volume ranges from Gigabytes to Terabytes.	Its volume ranges from Petabytes to Zettabytes or Exabytes.
Traditional database system deals with structured data.	Big data system deals with structured, semi-structured, database, and unstructured data.
Traditional data is generated per hour or per day or more.	But big data is generated more frequently mainly per seconds.
Traditional data source is centralized and it is managed in centralized form.	Big data source is distributed and it is managed in distributed form.
Data integration is very easy.	Data integration is very difficult.
Normal system configuration is capable to process traditional data.	High system configuration is required to process big data.
The size of the data is very small.	The size is more than the traditional data size.
Traditional data base tools are required to perform any data base operation.	Special kind of data base tools are required to perform any databaseschema-based operation.

Q. Big Data Analysis Life Cycle

The Big Data Analytics Life cycle is divided into nine phases, named as :

1. Business Case/Problem Definition
2. Data Identification

BDA

3. Data Acquisition and filtration
4. Data Extraction
5. Data Munging(Validation and Cleaning)
6. Data Aggregation & Representation(Storage)
7. Exploratory Data Analysis
8. Data Visualization(Preparation for Modeling and Assessment)
9. Utilization of analysis results.

Let us discuss each phase :

- **Phase I Business Problem Definition –**

In this stage, the team learns about the business domain, which presents the motivation and goals for carrying out the analysis. In this stage, the problem is identified, and assumptions are made that how much potential gain a company will make after carrying out the analysis. Important activities in this step include framing the business problem as an analytics challenge that can be addressed in subsequent phases. It helps the decision-makers understand the business resources that will be required to be utilized thereby determining the underlying budget required to carry out the project.

Moreover, it can be determined, whether the problem identified, is a Big Data problem or not, based on the business requirements in the business case. To qualify as a big data problem, the business case should be directly related to one(or more) of the characteristics of volume, velocity, or variety.

- **Phase II Data Definition –**

Once the business case is identified, now it's time to find the appropriate datasets to work with. In this stage, analysis is done to see what other companies have done for a similar case.

Depending on the business case and the scope of analysis of the project being addressed, the sources of datasets can be either external or internal to the company. In the case of internal datasets, the datasets can include data collected from internal sources, such as feedback forms, from existing software, On the other hand, for external datasets, the list includes datasets from third-party providers.

BDA

- **Phase III Data Acquisition and filtration –**

Once the source of data is identified, now it is time to gather the data from such sources. This kind of data is mostly unstructured. Then it is subjected to filtration, such as removal of the corrupt data or irrelevant data, which is of no scope to the analysis objective. Here corrupt data means data that may have missing records, or the ones, which include incompatible data types.

After filtration, a copy of the filtered data is stored and compressed, as it can be of use in the future, for some other analysis.

- **Phase IV Data Extraction –**

Now the data is filtered, but there might be a possibility that some of the entries of the data might be incompatible, to rectify this issue, a separate phase is created, known as the data extraction phase. In this phase, the data, which don't match with the underlying scope of the analysis, are extracted and transformed in such a form.

- **Phase V Data Munging –**

As mentioned in phase III, the data is collected from various sources, which results in the data being unstructured. There might be a possibility, that the data might have constraints, that are unsuitable, which can lead to false results. Hence there is a need to clean and validate the data.

It includes removing any invalid data and establishing complex validation rules. There are many ways to validate and clean the data. For example, a dataset might contain few rows, with null entries. If a similar dataset is present, then those entries are copied from that dataset, else those rows are dropped.

- **Phase VI Data Aggregation & Representation –**

The data is cleansed and validates, against certain rules set by the enterprise. But the data might be spread across multiple datasets, and it is not advisable to work with multiple datasets. Hence, the datasets are joined together. For example: If there are two datasets, namely that of a Student Academic section and Student Personal Details section, then both can be joined together via

common fields, i.e. roll number.

This phase calls for intensive operation since the amount of data can be very large. Automation can be brought into consideration, so that these things are executed, without any human intervention.

- **Phase VII Exploratory Data Analysis –**

Here comes the actual step, the analysis task. Depending on the nature of the big data problem, analysis is carried out. Data analysis can be classified as Confirmatory analysis and Exploratory analysis. In confirmatory analysis, the cause of a phenomenon is analyzed before. The assumption is called the hypothesis. The data is analyzed to approve or disapprove the hypothesis.

This kind of analysis provides definitive answers to some specific questions and confirms whether an assumption was true or not. In an exploratory analysis, the data is explored to obtain information, why a phenomenon occurred. This type of analysis answers “why” a phenomenon occurred. This kind of analysis doesn’t provide definitive, meanwhile, it provides discovery of patterns.

- **Phase VIII Data Visualization –**

Now we have the answer to some questions, using the information from the data in the datasets. But these answers are still in a form that can’t be presented to business users. A sort of representation is required to obtain value or some conclusion from the analysis. Hence, various tools are used to visualize the data in graphic form, which can easily be interpreted by business users.

Visualization is said to influence the interpretation of the results.

Moreover, it allows the users to discover answers to questions that are yet to be formulated.

- **Phase IX Utilization of analysis results –**

The analysis is done, the results are visualized, now it’s time for the business users to make decisions to utilize the results. The results can be used for optimization, to refine the business process. It can also be used as an input for the systems to enhance performance.

Unit – 2

Q. Need of Big Data Analytics?

1. Risk Management

Use Case: Banco de Oro, a Phillippine banking company, uses Big Data analytics to identify fraudulent activities and discrepancies. The organization leverages it to narrow down a list of suspects or root causes of problems.

2. Product Development and Innovations

Use Case: Rolls-Royce, one of the largest manufacturers of jet engines for airlines and armed forces across the globe, uses Big Data analytics to analyze how efficient the engine designs are and if there is any need for improvements.

3. Quicker and Better Decision Making Within Organizations

Use Case: Starbucks uses Big Data analytics to make strategic decisions. For example, the company leverages it to decide if a particular location would be suitable for a new outlet or not. They will analyze several different factors, such as population, demographics, accessibility of the location, and more.

4. Improve Customer Experience

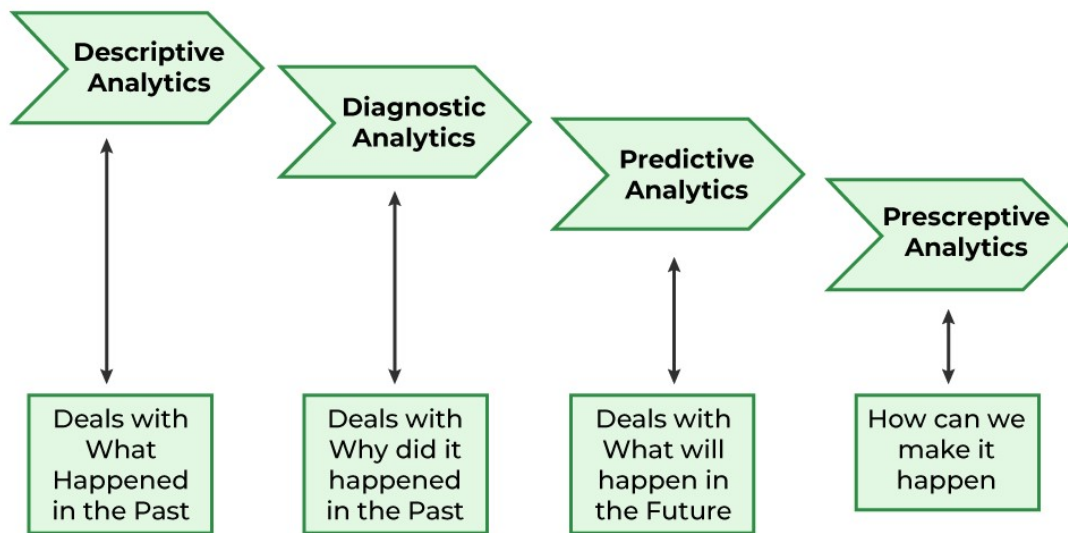
Use Case: Delta Air Lines uses Big Data analysis to improve customer experiences. They monitor tweets to find out their customers' experience regarding their journeys, delays, and so on. The airline identifies negative tweets and does what's necessary to remedy the situation. By publicly addressing these issues and offering solutions, it helps the airline build good customer relations.

Q. Types of Big Data Analytics

Types of Data Analytics

There are four major types of data analytics:

1. Predictive (forecasting)
2. Descriptive (business intelligence and data mining)
3. Prescriptive (optimization and simulation)
4. Diagnostic analytics



Data Analytics and its Types

Predictive Analytics

Predictive analytics turn the data into valuable, actionable information. predictive analytics uses data to determine the probable outcome of an event or a likelihood of a situation occurring. Predictive analytics holds a variety of statistical techniques from modeling, [machine learning](#), [data mining](#), and [game theory](#) that analyze current and historical facts to make predictions about a future event. Techniques that are used for predictive analytics are:

- [Linear Regression](#)
- [Time Series Analysis and Forecasting](#)
- [Data Mining](#)

Basic Corner Stones of Predictive Analytics

- Predictive modeling
- Decision Analysis and optimization
- Transaction profiling

Descriptive Analytics

Descriptive analytics looks at data and analyze past event for insight as to how to approach future events. It looks at past performance and understands the performance by mining historical data to understand the cause of success or failure in the past. Almost all management reporting such as sales, marketing, operations, and finance uses this type of analysis.

The descriptive model quantifies relationships in data in a way that is often used to classify customers or prospects into groups. Unlike a predictive

BDA

model that focuses on predicting the behavior of a single customer, Descriptive analytics identifies many different relationships between customer and product.

Common examples of Descriptive analytics are company reports that provide historic reviews like:

- Data Queries
- Reports
- Descriptive Statistics
- Data dashboard

Prescriptive Analytics

Prescriptive Analytics automatically synthesize big data, mathematical science, business rule, and machine learning to make a prediction and then suggests a decision option to take advantage of the prediction.

Prescriptive analytics goes beyond predicting future outcomes by also suggesting action benefits from the predictions and showing the decision maker the implication of each decision option. Prescriptive Analytics not only anticipates what will happen and when to happen but also why it will happen. Further, Prescriptive Analytics can suggest decision options on how to take advantage of a future opportunity or mitigate a future risk and illustrate the implication of each decision option.

For example, Prescriptive Analytics can benefit healthcare strategic planning by using analytics to leverage operational and usage data combined with data of external factors such as economic data, population demography, etc.

Diagnostic Analytics

In this analysis, we generally use historical data over other data to answer any question or for the solution of any problem. We try to find any dependency and pattern in the historical data of the particular problem.

For example, companies go for this analysis because it gives a great insight into a problem, and they also keep detailed information about their disposal otherwise data collection may turn out individual for every problem and it will be very time-consuming. Common techniques used for Diagnostic Analytics are:

- Data discovery
- Data mining
- Correlations

BDA

Q. Challenges to Big Data Analytics

Storage

With vast amounts of data generated daily, the greatest challenge is storage (especially when the data is in different formats) within legacy systems.

Unstructured data cannot be stored in traditional databases.

Processing

Processing big data refers to the reading, transforming, extraction, and formatting of useful information from raw information. The input and output of information in unified formats continue to present difficulties.

Security

Security is a big concern for organizations. Non-encrypted information is at risk of theft or damage by cyber-criminals. Therefore, data security professionals must balance access to data against maintaining strict security protocols.

Finding and Fixing Data Quality Issues

Many of you are probably dealing with challenges related to poor data quality, but solutions are available. The following are four approaches to fixing data problems:

- Correct information in the original database.
- Repairing the original data source is necessary to resolve any data inaccuracies.
- You must use highly accurate methods of determining who someone is.
- Scaling Big Data Systems

Database sharding, memory caching, moving to the cloud and separating read-only and write-active databases are all effective scaling methods. While each one of those approaches is fantastic on its own, combining them will lead you to the next level.

Evaluating and Selecting Big Data Technologies

BDA

Companies are spending millions on new big data technologies, and the market for such tools is expanding rapidly. In recent years, however, the IT industry has caught on to big data and analytics potential. The trending technologies include the following:

- Hadoop Ecosystem
- Apache Spark
- NoSQL Databases
- R Software
- Predictive Analytics
- Prescriptive Analytics
- Big Data Environments

In an extensive data set, data is constantly being ingested from various sources, making it more dynamic than a data warehouse. The people in charge of the big data environment will fast forget where and what each data collection came from.

Real-Time Insights

The term "real-time analytics" describes the practice of performing analyses on data as a system is collecting it. Decisions may be made more efficiently and with more accurate information thanks to real-time analytics tools, which use logic and mathematics to deliver insights on this data quickly.

Data Validation

Before using data in a business process, its integrity, accuracy, and structure must be validated. The output of a data validation procedure can be used for further analysis, BI, or even to train a machine learning model.

Q. Association rule

Association rule mining finds interesting associations and relationships among large sets of data items. This rule shows how frequently a itemset occurs in a transaction. A typical example is a Market Based Analysis.

BDA

Market Based Analysis is one of the key techniques used by large relations to show associations between items. It allows retailers to identify relationships between the items that people buy together frequently.

Given a set of transactions, we can find rules that will predict the occurrence of an item based on the occurrences of other items in the transaction.

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Before we start defining the rule, let us first see the basic definitions.

Support Count() – Frequency of occurrence of a itemset.

Here $(\{\text{Milk, Bread, Diaper}\})=2$

Frequent Itemset – An itemset whose support is greater than or equal to minsup threshold.

Association Rule – An implication expression of the form $X \rightarrow Y$, where X and Y are any 2 itemsets.

Example: $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$

Rule Evaluation Metrics –

- **Support(s)** –

The number of transactions that include items in the {X} and {Y} parts of the rule as a percentage of the total number of transaction. It is a measure of how frequently the collection of items occur together as a percentage of all transactions.

- **Support = $(X+Y)$ total** –

It is interpreted as fraction of transactions that contain both X and Y.

BDA

- **Confidence(c)** –
It is the ratio of the no of transactions that includes all items in {B} as well as the no of transactions that includes all items in {A} to the no of transactions that includes all items in {A}.
- **Conf(X=>Y) = Supp(X ∪ Y) / Supp(X)** –
It measures how often each item in Y appears in transactions that contains items in X also.
- **Lift(l)** –
The lift of the rule X=>Y is the confidence of the rule divided by the expected confidence, assuming that the itemsets X and Y are independent of each other. The expected confidence is the confidence divided by the frequency of {Y}.
- **Lift(X=>Y) = Conf(X=>Y) / Supp(Y)** –
Lift value near 1 indicates X and Y almost often appear together as expected, greater than 1 means they appear together more than expected and less than 1 means they appear less than expected. Greater lift values indicate stronger association.

Example – From the above table, {Milk, Diaper}=>{Beer}

$$s = \frac{(\{\text{Milk, Diaper, Beer}\})}{|T|}$$

$$= 2/5$$

$$= 0.4$$

$$c = \frac{(\text{Milk, Diaper, Beer})}{(\text{Milk, Diaper})}$$

$$= 2/3$$

$$= 0.67$$

$$l = \frac{\text{Supp}(\{\text{Milk, Diaper, Beer}\})}{\text{Supp}(\{\text{Milk, Diaper}\}) * \text{Supp}(\{\text{Beer}\})}$$

$$= 0.4 / (0.6 * 0.6)$$

$$= 1.11$$

The Association rule is very useful in analyzing datasets. The data is collected using bar-code scanners in supermarkets. Such databases consists of a large number of transaction records which list all items bought by a customer on a single purchase. So the manager could know if certain groups of items are consistently purchased together and use this data for adjusting store layouts, cross-selling, promotions based on statistics.

BDA

Q. Clustering Vs Classification

Classification	Clustering
Classification is a supervised learning approach where a specific label is provided to the machine to classify new observations. Here the machine needs proper testing and training for the label verification.	Clustering is an unsupervised learning approach where grouping is done on similarities basis.
Supervised learning approach.	Unsupervised learning approach.
It uses a training dataset.	It does not use a training dataset.
It uses algorithms to categorize the new data as per the observations of the training set.	It uses statistical concepts in which the data set is divided into subsets with the same features.
In classification, there are labels for training data.	In clustering, there are no labels for training data.
Its objective is to find which class a new object belongs to form the set of predefined classes.	Its objective is to group a set of objects to find whether there is any relationship between them.
It is more complex as compared to clustering.	It is less complex as compared to clustering.

Q. Mahout

We are living in a day and age where information is available in abundance. The information overload has scaled to such heights that sometimes it becomes difficult to manage our little mailboxes! Imagine the volume of data and records some of the popular websites (the likes of Facebook, Twitter, and Youtube) have to collect and manage on a daily basis. It is not uncommon even for lesser known websites to receive huge amounts of information in bulk.

Normally we fall back on data mining algorithms to analyze bulk data to identify trends and draw conclusions. However, no data mining algorithm can be efficient enough to process very large datasets and provide outcomes in quick time, unless the

BDA

computational tasks are run on multiple machines distributed over the cloud.

We now have new frameworks that allow us to break down a computation task into multiple segments and run each segment on a different machine. **Mahout** is such a data mining framework that normally runs coupled with the Hadoop infrastructure at its background to manage huge volumes of data.

What is Apache Mahout?

A *mahout* is one who drives an elephant as its master. The name comes from its close association with Apache Hadoop which uses an elephant as its logo.

Hadoop is an open-source framework from Apache that allows to store and process big data in a distributed environment across clusters of computers using simple programming models.

Apache **Mahout** is an open source project that is primarily used for creating scalable machine learning algorithms. It implements popular machine learning techniques such as:

- Recommendation
- Classification
- Clustering

Apache Mahout started as a sub-project of Apache's Lucene in 2008. In 2010, Mahout became a top level project of Apache.

Features of Mahout

The primitive features of Apache Mahout are listed below.

- The algorithms of Mahout are written on top of Hadoop, so it works well in distributed environment. Mahout uses the Apache Hadoop library to scale effectively in the cloud.
- Mahout offers the coder a ready-to-use framework for doing data mining tasks on large volumes of data.
- Mahout lets applications to analyze large sets of data effectively and in quick time.

BDA

- Includes several MapReduce enabled clustering implementations such as k-means, fuzzy k-means, Canopy, Dirichlet, and Mean-Shift.
- Supports Distributed Naive Bayes and Complementary Naive Bayes classification implementations.
- Comes with distributed fitness function capabilities for evolutionary programming.
- Includes matrix and vector libraries.

Applications of Mahout

- Companies such as Adobe, Facebook, LinkedIn, Foursquare, Twitter, and Yahoo use Mahout internally.
 - Foursquare helps you in finding out places, food, and entertainment available in a particular area. It uses the recommender engine of Mahout.
 - Twitter uses Mahout for user interest modelling.
 - Yahoo! uses Mahout for pattern mining.
-

Unit-3

Q. Hadoop

INTRODUCTION:

Hadoop is an open-source software framework that is used for storing and processing large amounts of data in a distributed computing environment. It is designed to handle big data and is based on the MapReduce programming model, which allows for the parallel processing of large datasets.

What is Hadoop?

Hadoop is an open source software programming framework for storing a large amount of data and performing the computation. Its framework is based on Java programming with some native code in C and shell scripts.

Hadoop is an open-source software framework that is used for storing and processing large amounts of data in a distributed computing environment. It is designed to handle big data and is based on the MapReduce programming model, which allows for the parallel processing of large datasets.

Hadoop has two main components:

- HDFS (Hadoop Distributed File System): This is the storage component of Hadoop, which allows for the storage of large amounts of data across multiple machines. It is designed to work with commodity hardware, which makes it cost-effective.
- YARN (Yet Another Resource Negotiator): This is the resource management component of Hadoop, which manages the allocation of resources (such as CPU and memory) for processing the data stored in HDFS.
- Hadoop also includes several additional modules that provide additional functionality, such as Hive (a SQL-like query language), Pig (a high-level platform for creating MapReduce programs), and HBase (a non-relational, distributed database).
- Hadoop is commonly used in big data scenarios such as data warehousing, business intelligence, and machine learning. It's also used for data processing, data analysis, and data mining. It enables the distributed processing of large data sets across clusters of computers using a simple programming model.

History of Hadoop

Apache Software Foundation is the developers of Hadoop, and it's co-founders are **Doug Cutting** and **Mike Cafarella**. It's co-founder Doug Cutting named it on his son's toy elephant. In October 2003 the first paper release was Google File System. In January 2006, MapReduce development started on the Apache Nutch which consisted of around 6000 lines coding for it and around 5000 lines coding for HDFS. In April 2006 Hadoop 0.1.0 was released.

Hadoop is an open-source software framework for storing and processing big data. It was created by Apache Software Foundation in 2006, based on a white paper written by Google in 2003 that described the Google File System (GFS) and the MapReduce programming model. The Hadoop framework allows for the distributed processing of large data sets across clusters of computers using simple programming models. It is designed to scale up from single servers to thousands of machines, each offering local computation and storage. It is used by many organizations, including Yahoo, Facebook, and IBM, for a variety of purposes such as data warehousing, log processing, and research. Hadoop has been widely adopted in the industry and has become a key technology for big data processing.

Features of hadoop:

1. it is fault tolerance.
2. it is highly available.
3. it's programming is easy.
4. it have huge flexible storage.
5. it is low cost.

Hadoop has several key features that make it well-suited for big data processing:

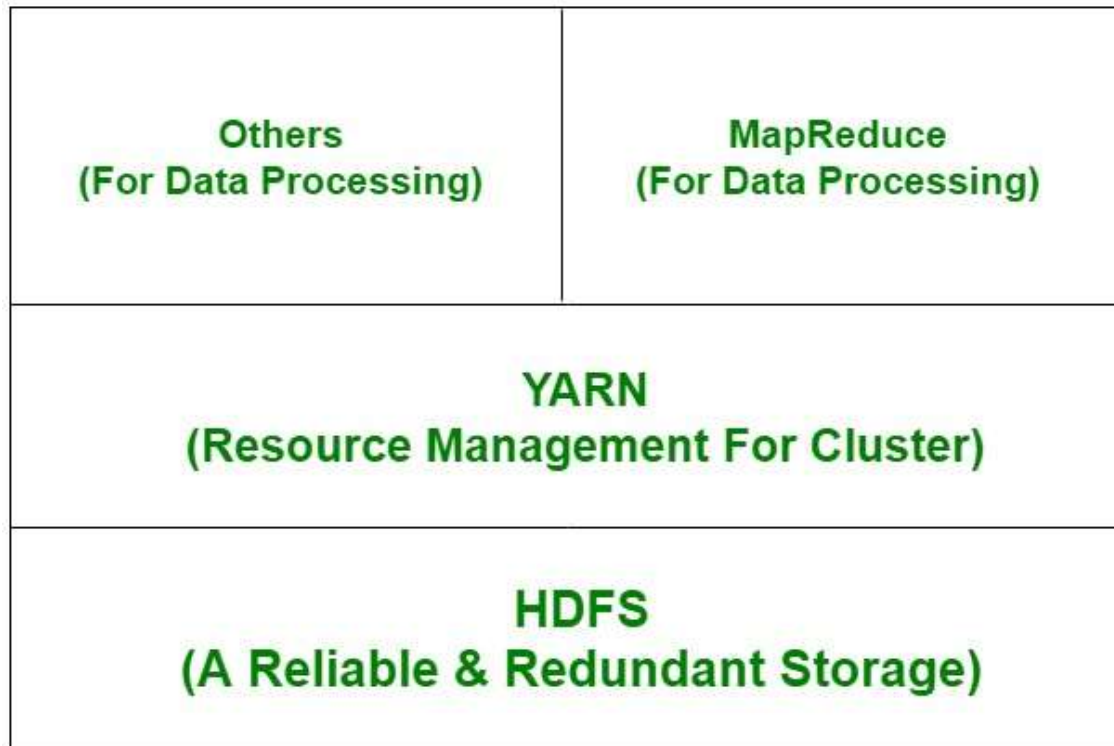
- **Distributed Storage:** Hadoop stores large data sets across multiple machines, allowing for the storage and processing of extremely large amounts of data.
- **Scalability:** Hadoop can scale from a single server to thousands of machines, making it easy to add more capacity as needed.
- **Fault-Tolerance:** Hadoop is designed to be highly fault-tolerant, meaning it can continue to operate even in the presence of hardware failures.

BDA

- Data locality: Hadoop provides data locality feature, where the data is stored on the same node where it will be processed, this feature helps to reduce the network traffic and improve the performance
- High Availability: Hadoop provides High Availability feature, which helps to make sure that the data is always available and is not lost.
- Flexible Data Processing: Hadoop's MapReduce programming model allows for the processing of data in a distributed fashion, making it easy to implement a wide variety of data processing tasks.
- Data Integrity: Hadoop provides built-in checksum feature, which helps to ensure that the data stored is consistent and correct.
- Data Replication: Hadoop provides data replication feature, which helps to replicate the data across the cluster for fault tolerance.
- Data Compression: Hadoop provides built-in data compression feature, which helps to reduce the storage space and improve the performance.
- YARN: A resource management platform that allows multiple data processing engines like real-time streaming, batch processing, and interactive SQL, to run and process data stored in HDFS.

Hadoop Distributed File System

It has distributed file system known as HDFS and this HDFS splits files into blocks and sends them across various nodes in form of large clusters. Also in case of a node failure, the system operates and data transfer takes place between the nodes which are facilitated by HDFS.



HDFS

Advantages of HDFS: It is inexpensive, immutable in nature, stores data reliably, ability to tolerate faults, scalable, block structured, can process a large amount of data simultaneously and many more. **Disadvantages of HDFS:** It's the biggest disadvantage is that it is not fit for small quantities of data. Also, it has issues related to potential stability, restrictive and rough in nature. Hadoop also supports a wide range of software packages such as Apache Flumes, Apache Oozie, Apache HBase, Apache Sqoop, Apache Spark, Apache Storm, Apache Pig, Apache Hive, Apache Phoenix, Cloudera Impala.

Some common frameworks of Hadoop

1. Hive- It uses HiveQL for data structuring and for writing complicated MapReduce in HDFS.
2. Drill- It consists of user-defined functions and is used for data exploration.
3. Storm- It allows real-time processing and streaming of data.
4. Spark- It contains a Machine Learning Library(MLlib) for providing enhanced machine learning and is widely used for data processing. It also supports Java, Python, and Scala.
5. Pig- It has Pig Latin, a SQL-Like language and performs data transformation of unstructured data.

BDA

6. Tez- It reduces the complexities of Hive and Pig and helps in the running of their codes faster.

Hadoop framework is made up of the following modules:

1. Hadoop MapReduce- a MapReduce programming model for handling and processing large data.
2. Hadoop Distributed File System- distributed files in clusters among nodes.
3. Hadoop YARN- a platform which manages computing resources.
4. Hadoop Common- it contains packages and libraries which are used for other modules.

Advantages and Disadvantages of Hadoop

Advantages:

- Ability to store a large amount of data.
- High flexibility.
- Cost effective.
- High computational power.
- Tasks are independent.
- Linear scaling.

Hadoop has several advantages that make it a popular choice for big data processing:

- Scalability: Hadoop can easily scale to handle large amounts of data by adding more nodes to the cluster.
- Cost-effective: Hadoop is designed to work with commodity hardware, which makes it a cost-effective option for storing and processing large amounts of data.
- Fault-tolerance: Hadoop's distributed architecture provides built-in fault-tolerance, which means that if one node in the cluster goes down, the data can still be processed by the other nodes.
- Flexibility: Hadoop can process structured, semi-structured, and unstructured data, which makes it a versatile option for a wide range of big data scenarios.
- Open-source: Hadoop is open-source software, which means that it is free to use and modify. This also allows developers to access the source code and make improvements or add new features.

BDA

- Large community: Hadoop has a large and active community of developers and users who contribute to the development of the software, provide support, and share best practices.
- Integration: Hadoop is designed to work with other big data technologies such as Spark, Storm, and Flink, which allows for integration with a wide range of data processing and analysis tools.

Disadvantages:

- Not very effective for small data.
- Hard cluster management.
- Has stability issues.
- Security concerns.
- Complexity: Hadoop can be complex to set up and maintain, especially for organizations without a dedicated team of experts.
- Latency: Hadoop is not well-suited for low-latency workloads and may not be the best choice for real-time data processing.
- Limited Support for Real-time Processing: Hadoop's batch-oriented nature makes it less suited for real-time streaming or interactive data processing use cases.
- Limited Support for Structured Data: Hadoop is designed to work with unstructured and semi-structured data, it is not well-suited for structured data processing
- Data Security: Hadoop does not provide built-in security features such as data encryption or user authentication, which can make it difficult to secure sensitive data.
- Limited Support for Ad-hoc Queries: Hadoop's MapReduce programming model is not well-suited for ad-hoc queries, making it difficult to perform exploratory data analysis.
- Limited Support for Graph and Machine Learning: Hadoop's core component HDFS and MapReduce are not well-suited for graph and machine learning workloads, specialized components like Apache Graph and Mahout are available but have some limitations.
- Cost: Hadoop can be expensive to set up and maintain, especially for organizations with large amounts of data.
- Data Loss: In the event of a hardware failure, the data stored in a single node may be lost permanently.
- Data Governance: Data Governance is a critical aspect of data management, Hadoop does not provide a built-in feature to

BDA

manage data lineage, data quality, data cataloging, data lineage, and data audit.

Q. Hadoop 2 vs Hadoop 3

S.No.	Feature	Hadoop 2.x	Hadoop 3.x
1	License	Apache 2.0 is used for licensing which is open-source.	Apache 2.0 is used for licensing which is open-source.
2	Minimum supported Java version	JAVA 7 is the minimum compatible version.	JAVA 8 is the minimum compatible version.
3	Fault Tolerance	Replication is the only way to handle fault tolerance which is not space optimized.	Erasur coding is used for handling fault tolerance.
4	Data Balancing	HDFS balancer is used for Data Balancing.	Intra-data node balancer is used which is called via HDFS disk-balancer command-line interface.
5	Storage Scheme	3x Replication Scheme is used.	uses eraser encoding in HDFS.
6	Storage Overhead	200% of HDFS is consumed in Hadoop 2.x	50% used in Hadoop 3.x means we have more space to work.
7	YARN Timeline Service	Uses timeline service with scalability issue.	Improve the time line service along with improving scalability and reliability of this service.

BDA

S.No.	Feature	Hadoop 2.x	Hadoop 3.x
8	Scalability	Limited Scalability, can have upto 10000 nodes in a cluster.	Scalability is improved, can have more than 10000 nodes in a cluster.

Hadoop 1 vs Hadoop 2 vs Hadoop 3

Hadoop 1.X	Hadoop 2.X	Hadoop 3.X
Hadoop 1.x was released in 2011	Hadoop 2.x released in 2012	Hadoop 3.x released in 2017
It introduced MapReduce and HDFS. That is to say, the MapReduce framework is used as data processing and for resource management also.	YARN (Yet another resource negotiator) added for better resource management. As a result, it enabled multi-tenancy. Therefore, the same cluster can be used by MapReduce as well as by some other processes using YARN.	In Hadoop 3.x, the YARN resource model is generalized to support user-defined resource types beyond CPU and memory. For example, the administrator can define resources like GPUs, software licenses, or locally-attached storage. YARN tasks can then be scheduled based on the availability of these resources.
Supports single tenancy only	Supports multiple tenants using YARN	Multiple tenants are supported here.
Hadoop 1.x uses Master-Slave architecture that consists of a single master and multiple slaves. So, in case the master node gets failed then the entire clusters become unavailable.	Hadoop 2.x is also a Master-Slave architecture. However, this consists of multiple masters that includes active namenode and standby namenode. So, in this case if master node get failed then the standby master node will take over it. As a result, hadoop 2.x fixes the problem of a single point of failure.	It added supports for multiple active namenodes
Hadoop 1.x is limited to 4000 nodes per cluster.	It supports up to 10000 nodes in a cluster.	The scalability is improved in Hadoop 3.x and it can have more than 10000 nodes in one cluster.

BDA

	Manual intervention is needed for namenode recovery.	We don't need manual intervention for namenode recovery.
	Java 7 is the minimum supported version	Java 8 is the minimum supported version.
	It supports HDFS(default), FTP, Amazon S3 and Windows Azure Storage Blobs (WASB) file systems.	All file systems including Microsoft Azure Data Lake filesystem is compatible with Hadoop 3.x.
	It uses 3x replication scheme that results in 200% storage overhead.	Hadoop 3 uses eraser encoding in HDFS that helps to reduce the storage overhead. It has 50% storage overhead only.
		It added support for GPU hardware that can be used to execute deep learning algorithms on a Hadoop cluster.

Q. What is Apache Pig.

Apache Pig is a high-level data flow platform for executing MapReduce programs of Hadoop. The language used for Pig is Pig Latin.

Pig tutorial provides basic and advanced concepts of Pig. Our Pig tutorial is designed for beginners and professionals.

Pig is a high-level data flow platform for executing Map Reduce programs of Hadoop. It was developed by Yahoo. The language for Pig is pig Latin.

Our Pig tutorial includes all topics of Apache Pig with Pig usage, Pig Installation, Pig Run Modes, Pig Latin concepts, Pig Data Types, Pig example, Pig user defined functions etc.

Features of Apache Pig

Let's see the various uses of Pig technology.

1) Ease of programming

BDA

Writing complex java programs for map reduce is quite tough for non-programmers. Pig makes this process easy. In the Pig, the queries are converted to MapReduce internally.

2) Optimization opportunities

It is how tasks are encoded permits the system to optimize their execution automatically, allowing the user to focus on semantics rather than efficiency.

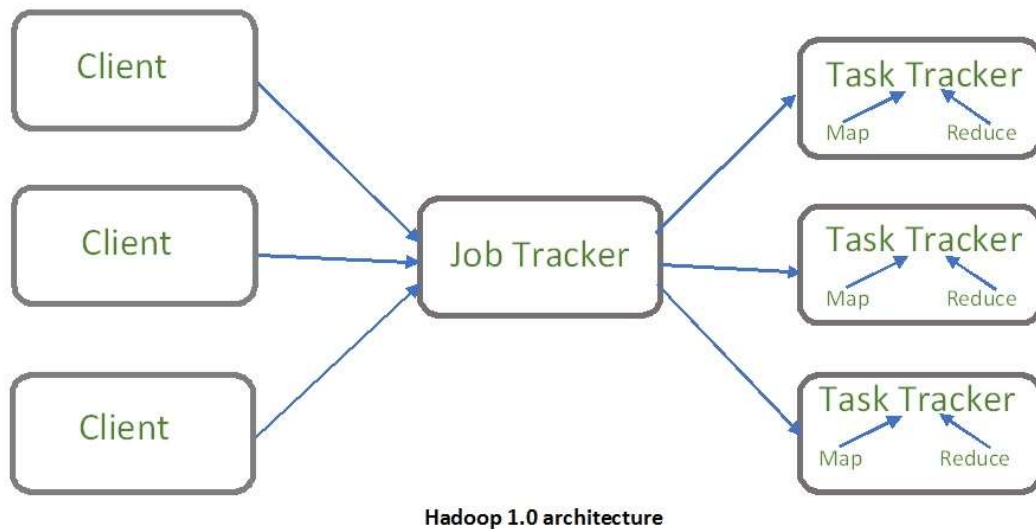
3) Extensibility

A user-defined function is written in which the user can write their logic to execute over the data set.

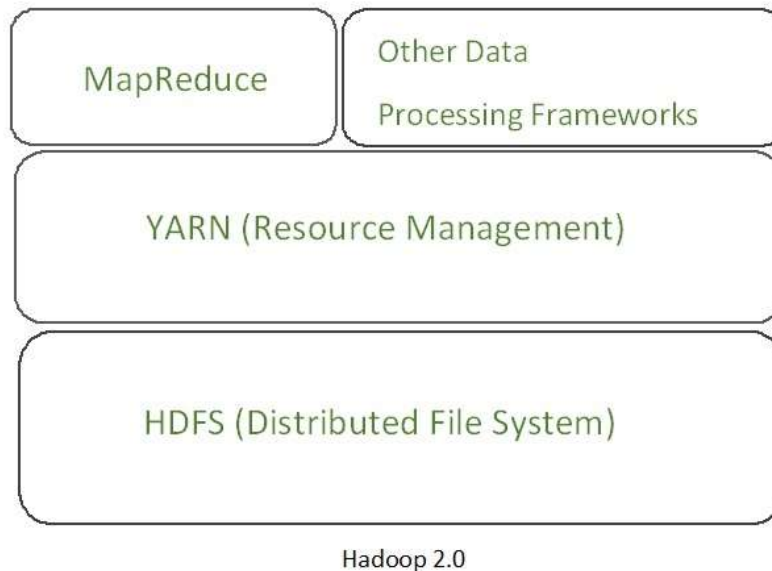
Unit – 4

Q. YARN

YARN stands for “**Yet Another Resource Negotiator**”. It was introduced in Hadoop 2.0 to remove the bottleneck on Job Tracker which was present in Hadoop 1.0. YARN was described as a “*Redesigned Resource Manager*” at the time of its launching, but it has now evolved to be known as large-scale distributed operating system used for Big Data processing.



YARN architecture basically separates resource management layer from the processing layer. In Hadoop 1.0 version, the responsibility of Job tracker is split between the resource manager and application manager.

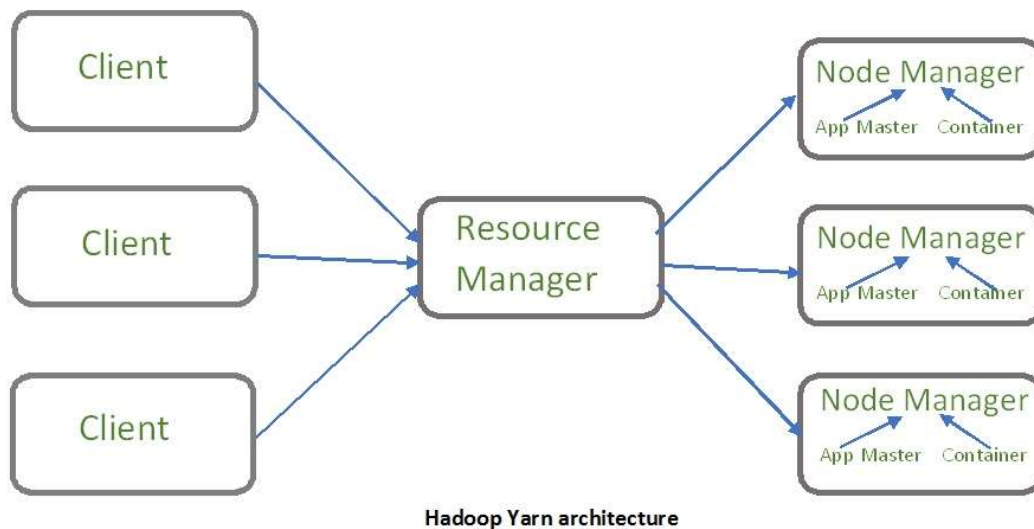


YARN also allows different data processing engines like graph processing, interactive processing, stream processing as well as batch processing to run and process data stored in HDFS (Hadoop Distributed File System) thus making the system much more efficient. Through its various components, it can dynamically allocate various resources and schedule the application processing. For large volume data processing, it is quite necessary to manage the available resources properly so that every application can leverage them.

YARN Features: YARN gained popularity because of the following features-

- **Scalability:** The scheduler in Resource manager of YARN architecture allows Hadoop to extend and manage thousands of nodes and clusters.
- **Compatibility:** YARN supports the existing map-reduce applications without disruptions thus making it compatible with Hadoop 1.0 as well.
- **Cluster Utilization:** Since YARN supports Dynamic utilization of cluster in Hadoop, which enables optimized Cluster Utilization.
- **Multi-tenancy:** It allows multiple engine access thus giving organizations a benefit of multi-tenancy.

Hadoop YARN Architecture



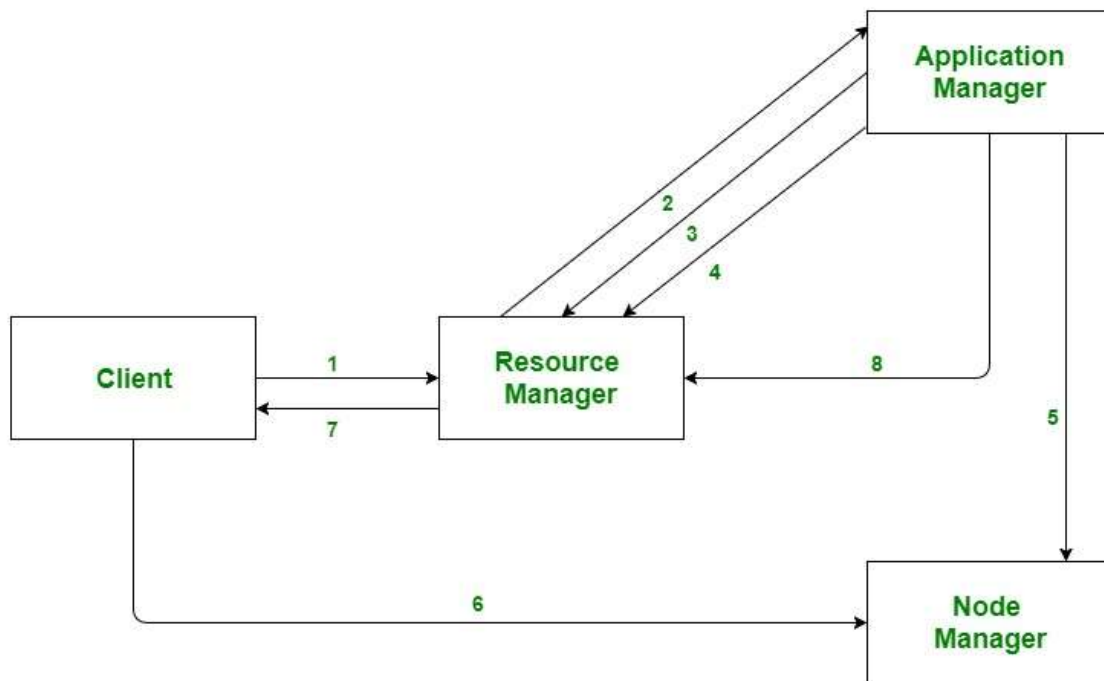
The main components of YARN architecture include:

- **Client:** It submits map-reduce jobs.
- **Resource Manager:** It is the master daemon of YARN and is responsible for resource assignment and management among all the applications. Whenever it receives a processing request, it forwards it to the corresponding node manager and allocates resources for the completion of the request accordingly. It has two major components:
 - **Scheduler:** It performs scheduling based on the allocated application and available resources. It is a pure scheduler, means it does not perform other tasks such as monitoring or tracking and does not guarantee a restart if a task fails. The YARN scheduler supports plugins such as Capacity

Scheduler and Fair Scheduler to partition the cluster resources.

- **Application manager:** It is responsible for accepting the application and negotiating the first container from the resource manager. It also restarts the Application Master container if a task fails.
- **Node Manager:** It take care of individual node on Hadoop cluster and manages application and workflow and that particular node. Its primary job is to keep-up with the Resource Manager. It registers with the Resource Manager and sends heartbeats with the health status of the node. It monitors resource usage, performs log management and also kills a container based on directions from the resource manager. It is also responsible for creating the container process and start it on the request of Application master.
- **Application Master:** An application is a single job submitted to a framework. The application master is responsible for negotiating resources with the resource manager, tracking the status and monitoring progress of a single application. The application master requests the container from the node manager by sending a Container Launch Context(CLC) which includes everything an application needs to run. Once the application is started, it sends the health report to the resource manager from time-to-time.
- **Container:** It is a collection of physical resources such as RAM, CPU cores and disk on a single node. The containers are invoked by Container Launch Context(CLC) which is a record that contains information such as environment variables, security tokens, dependencies etc.

Application workflow in Hadoop YARN:



1. Client submits an application
2. The Resource Manager allocates a container to start the Application Manager
3. The Application Manager registers itself with the Resource Manager
4. The Application Manager negotiates containers from the Resource Manager
5. The Application Manager notifies the Node Manager to launch containers
6. Application code is executed in the container
7. Client contacts Resource Manager/Application Manager to monitor application's status
8. Once the processing is complete, the Application Manager unregisters with the Resource Manager

Advantages :

- **Flexibility:** YARN offers flexibility to run various types of distributed processing systems such as Apache Spark, Apache Flink, Apache Storm, and others. It allows multiple processing engines to run simultaneously on a single Hadoop cluster.

BDA

- **Resource Management:** YARN provides an efficient way of managing resources in the Hadoop cluster. It allows administrators to allocate and monitor the resources required by each application in a cluster, such as CPU, memory, and disk space.
- **Scalability:** YARN is designed to be highly scalable and can handle thousands of nodes in a cluster. It can scale up or down based on the requirements of the applications running on the cluster.
- **Improved Performance:** YARN offers better performance by providing a centralized resource management system. It ensures that the resources are optimally utilized, and applications are efficiently scheduled on the available resources.
- **Security:** YARN provides robust security features such as Kerberos authentication, Secure Shell (SSH) access, and secure data transmission. It ensures that the data stored and processed on the Hadoop cluster is secure.

Disadvantages :

- **Complexity:** YARN adds complexity to the Hadoop ecosystem. It requires additional configurations and settings, which can be difficult for users who are not familiar with YARN.
- **Overhead:** YARN introduces additional overhead, which can slow down the performance of the Hadoop cluster. This overhead is required for managing resources and scheduling applications.
- **Latency:** YARN introduces additional latency in the Hadoop ecosystem. This latency can be caused by resource allocation, application scheduling, and communication between components.
- **Single Point of Failure:** YARN can be a single point of failure in the Hadoop cluster. If YARN fails, it can cause the entire cluster to go down. To avoid this, administrators need to set up a backup YARN instance for high availability.
- **Limited Support:** YARN has limited support for non-Java programming languages. Although it supports multiple processing engines, some engines have limited language support, which can limit the usability of YARN in certain environments.

Q. MapReduce

Certainly! Here's a simplified version of the notes on MapReduce in Hadoop:

BDA

MapReduce is a programming model for processing and generating large datasets. It's widely used in Hadoop, a framework for distributed storage and processing of big data.

Key Concepts:

1. **Map Function:**

- Takes input data and transforms it into key-value pairs.
- Each input data item is processed independently.

2. **Shuffling and Sorting:**

- The framework sorts and groups the key-value pairs from the Map phase.
- Groups data by keys, preparing it for the Reduce phase.

3. **Reduce Function:**

- Accepts key-value pairs and performs a specified operation.
- Aggregates or processes data based on the keys.

Workflow:

1. **Input Split:**

- Large dataset is divided into smaller chunks called input splits.

2. **Map Phase:**

- Map tasks run parallelly on input splits, producing intermediate key-value pairs.

3. **Shuffling and Sorting:**

- Framework organizes and groups intermediate data by keys.

4. **Reduce Phase:**

- Reduce tasks process the grouped data, producing the final output.

Example Scenario:

Let's say we have a large dataset of words and want to count the frequency of each word.

BDA

1. **Map Function:**

- Map tasks process portions of the data, emitting key-value pairs like (word, 1) for each word.

2. **Shuffling and Sorting:**

- Framework groups data by word, making it ready for the next phase.

3. **Reduce Function:**

- Reduce tasks sum up the counts for each word, giving the final word frequency.

Benefits:

- **Parallel Processing:**

- MapReduce enables parallel processing of data across a distributed cluster of computers.

- **Fault Tolerance:**

- Hadoop's MapReduce is designed to handle failures by rerunning failed tasks on other nodes.

- **Scalability:**

- Scales easily with the addition of more hardware to the cluster.

In summary, MapReduce simplifies the processing of large datasets by dividing the task into smaller, manageable parts that can be processed in parallel across a distributed environment. This makes it efficient for big data processing in systems like Hadoop.

Q. Serilization of BigData

It seems there might be a slight error in your question. I assume you meant "serialization" instead of "sterilization." In the context of big data, serialization refers to the process of converting data structures or objects into a format that can be easily stored, transmitted, or reconstructed later.

Here are some common serialization formats used in the context of big data:

1. **JSON (JavaScript Object Notation):**

- A lightweight data interchange format.
- Human-readable and easy for both machines and humans to parse.

BDA

- Widely used for configuration files and web APIs.
2. **XML (eXtensible Markup Language):**
 - A markup language that defines rules for encoding documents.
 - More verbose compared to JSON but supports a hierarchical structure.
 - Commonly used in web services and configuration files.
 3. **Avro:**
 - A binary serialization format developed within the Apache Hadoop project.
 - Compact and fast, designed for big data processing.
 - Supports schema evolution, allowing data to evolve over time.
 4. **Protocol Buffers (protobuf):**
 - Developed by Google, a language-agnostic binary serialization format.
 - Compact and efficient, suitable for high-performance applications.
 - Requires a predefined schema.
 5. **Parquet:**
 - A columnar storage file format optimized for big data processing.
 - Suitable for use with frameworks like Apache Spark and Apache Hive.
 - Supports nested data structures and compression.
 6. **ORC (Optimized Row Columnar):**
 - Another columnar storage file format designed for Hadoop workloads.
 - Provides lightweight compression and efficient read performance.
 - Suitable for analytics and querying large datasets.
 7. **MessagePack:**
 - Binary format that is more efficient in terms of size and speed compared to JSON.
 - Supports a wide range of data types.
 - Used for data exchange between languages.
 8. **Thrift:**
 - Developed by Apache, a framework for scalable cross-language services.
 - Supports data serialization in various programming languages.
 - Requires a predefined schema.

BDA

Choosing the right serialization format depends on factors such as data size, processing speed, and compatibility with the tools and frameworks in your big data ecosystem. Each format has its strengths and weaknesses, so the choice often depends on specific use cases and requirements.

Q. big Data Serialization Format

Serialization formats play a crucial role in efficiently storing, transmitting, and processing big data. Here are some popular serialization formats used in the context of big data:

1. **JSON (JavaScript Object Notation):**

- **Pros:**
 - Human-readable and easy to understand.
 - Widely supported across various programming languages.
- **Cons:**
 - Relatively verbose, which can impact storage and transmission efficiency.

2. **XML (eXtensible Markup Language):**

- **Pros:**
 - Hierarchical structure supports complex data.
 - Self-descriptive.
- **Cons:**
 - Verbosity can lead to increased storage and processing overhead.

3. **Avro:**

- **Pros:**
 - Compact binary format designed for efficient serialization.
 - Supports schema evolution.
 - Fast serialization and deserialization.
- **Cons:**
 - Binary format is not human-readable.

4. **Protocol Buffers (protobuf):**

- **Pros:**
 - Efficient binary format with smaller size compared to JSON and XML.
 - Supports schema evolution.
 - Fast serialization and deserialization.

BDA

- **Cons:**
 - Binary format is not human-readable.

5. **Parquet:**

- **Pros:**
 - Columnar storage format optimized for analytics.
 - Efficient compression and encoding.
 - Suitable for big data processing frameworks like Apache Spark and

Apache Hive.

- **Cons:**
 - Not human-readable due to its binary format.

6. **ORC (Optimized Row Columnar):**

- **Pros:**
 - Columnar storage with lightweight compression.
 - Optimized for Hadoop workloads.
- **Cons:**
 - Binary format not suitable for human readability.

7. **MessagePack:**

- **Pros:**
 - Binary format with a compact size.
 - Efficient serialization and deserialization.
 - Supports a wide range of data types.
- **Cons:**
 - Not as widely adopted as JSON or XML.

8. **Thrift:**

- **Pros:**
 - Cross-language support for data serialization.
 - Compact binary format.
 - Efficient for communication between different systems.
- **Cons:**
 - Requires a predefined schema.

Choosing the right serialization format depends on your specific use case, considering factors such as data size, processing speed, and compatibility with your big data tools and frameworks. Each format has its advantages and trade-offs, so the choice often involves balancing these considerations based on your particular requirements.
