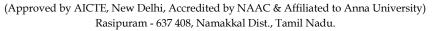


## **MUTHAYAMMAL ENGINEERING COLLEGE**

(An Autonomous Institution)





## **MUST KNOW CONCEPTS**

**MKC** 

MCA 2021-2022

Course Code & Course Name : 21CAB13 & Big Data Analytics

Year/Sem/Sec : I/II

I cairs	em/sec		. 1/	11		
S.No.	T	erm	Notation (Symbol)		Concept/Definition/Meaning/ Units/Equation/Expression	Units
			Unit-I : Int	rodu	ction to Big Data	
1.	Big data		7	of s data so	data is defined as the voluminous amount structured, unstructured or semi-structured a that has huge potential for mining but is large that it cannot be processed using itional database systems.	I
2.	Big data a	nalytics	3	of d	data analytics examines large amounts ata to uncover hidden patterns, correlations other insights.	I
3.	Types of I	Big Data	8	4	es of Big Data  1. Structured  2. Unstructured  3. Semi-structured	I
4.	Characteri Data	istics of Big		Cha	<ul> <li>variety</li> <li>Velocity</li> <li>Variability</li> </ul>	I
5.	Volume	533	G Wi N	to a a ve data		I
6.	Variety		Estd	the	iety refers to heterogeneous sources and nature of data, both structured and tructured.	Ι
7.	Velocity			gen gen	term 'velocity' refers to the speed of eration of data. How fast the data is erated and processed to meet the demands, ermines real potential in the data.	I
8.	Variability	ý		whi han	iability – This refers to the inconsistency ch can be shown by the data at times, thus apering the process of being able to handle manage the data effectively.	Ι
9.	Big data p	latform		con	data platform is a type of IT solution that abines the features and capabilities of eral big data application and utilities within agle solution	Ι

Intelligent Data   Analysis	_		T		
between analysts and collections of aggregated data that may have been reformulated into alternate representational forms as a means for improved analytical performance.  Business analytics tools are types of application software that retrieve data from one or more business systems and combine it in a repository, such as a data warehouse, to be reviewed and analyzed.  Reporting   Reporting is the process of organizing data into informational summaries in order to monitor how different areas of a business are performing.  Analysis is the process of exploring data and reports in order to extract meaningful, actionable insights, which can be used to better understand and improve business performance.  R is the leading analytics tool in the industry and widely used for statistics and data modeling. It can easily manipulate your data and present in different ways.  Tableau Public is a free software that connects any data source be it corporate Data Warehouse.  Python is an object-oriented scripting language which is easy to read, write, maintain and is a free open source tool. It was developed by Guido van Rossum in late 1980's which supports both functional and structured programming methods.  Sas is a programming environment and language for data manipulation and a leader in analytics, developed by the SAS Institute in 1966 and further developed in 1980's and 1990's. SAS is easily accessible, managable and can analyze data from any sources.  Apache Spark — Hadoop clusters 100 times faster in memory and 10 times faster on disk.  Excel is a basic, popular and widely used analytical tool almost in all industries. Whether you are an expert in Sas, Ro I	10.	_		and information. Intelligent data analysis reveals implicit, previously unknown and potentially valuable information or	I
application software that retrieve data from one or more business systems and combine it in a repository, such as a data warehouse, to be reviewed and analyzed.  Reporting is the process of organizing data into informational summaries in order to monitor how different areas of a business are performing.  Analysis is the process of exploring data and reports in order to extract meaningful, actionable insights, which can be used to better understand and improve business performance.  R is the leading analytics tool in the industry and widely used for statistics and data modeling. It can easily manipulate your data and present in different ways.  Tableau Public is a free software that connects any data source be it corporate Data Warehouse.  Python is an object-oriented scripting language which is easy to read, write, maintain and is a free open source tool. It was developed by Guido van Rossum in late 1980's which supports both functional and structured programming methods.  Sas is a programming environment and language for data manipulation and a leader in analytics, developed by the SAS Institute in 1966 and further developed in 1980's and 1990's. SAS is easily accessible, managable and can analyze data from any sources.  Apache Spark is a fast large-scale data processing engine and executes applications in Hadoop clusters 100 times faster in memory and 10 times faster on disk.  Excel is a basic, popular and widely used analytical tool almost in all industries. Whether you are an expert in Sas, R or Tableau, you will still need to use Excel.	11.	Analytical processing		between analysts and collections of aggregated data that may have been reformulated into alternate representational forms as a means for	I
13. Reporting into informational summaries in order to monitor how different areas of a business are performing.  Analysis is the process of exploring data and reports in order to extract meaningful, actionable insights, which can be used to better understand and improve business performance.  R is the leading analytics tool in the industry and widely used for statistics and data modeling. It can easily manipulate your data and present in different ways.  Tableau Public is a free software that connects any data source be it corporate Data Warehouse.  Python is an object-oriented scripting language which is easy to read, write, maintain and is a free open source tool. It was developed by Guido van Rossum in late 1980's which supports both functional and structured programming methods.  Sas is a programming environment and language for data manipulation and a leader in analytics, developed by the SAS Institute in 1966 and further developed in 1980's and 1990's. SAS is easily accessible, managable and can analyze data from any sources.  Apache Spark   Apache Spark is a fast large-scale data processing engine and executes applications in Hadoop clusters 100 times faster in memory and 10 times faster on disk.  Excel is a basic, popular and widely used analytical tool almost in all industries. Whether you are an expert in Sas, R or Tableau, you will still need to use Excel.	12.			application software that retrieve data from one or more business systems and combine it in a repository, such as a data warehouse, to be reviewed and analyzed.	I
reports in order to extract meaningful, actionable insights, which can be used to better understand and improve business performance.  R is the leading analytics tool in the industry and widely used for statistics and data modeling. It can easily manipulate your data and present in different ways.  Tableau Public is a free software that connects any data source be it corporate Data  Warehouse.  Python is an object-oriented scripting language which is easy to read, write, maintain and is a free open source tool. It was developed by Guido van Rossum in late 1980's which supports both functional and structured programming methods.  Sas is a programming environment and language for data manipulation and a leader in analytics, developed by the SAS Institute in 1966 and further developed in 1980's and 1990's. SAS is easily accessible, managable and can analyze data from any sources.  Apache Spark is a fast large-scale data processing engine and executes applications in Hadoop clusters 100 times faster in memory and 10 times faster on disk.  Excel is a basic, popular and widely used analytical tool almost in all industries. Whether you are an expert in Sas, R or Tableau, you will still need to use Excel.	13.	Reporting		into informational summaries in order to monitor how different areas of a business are	I
and widely used for statistics and data modeling. It can easily manipulate your data and present in different ways.  Tableau Public is a free software that connects any data source be it corporate Data Warehouse.  Python is an object-oriented scripting language which is easy to read, write, maintain and is a free open source tool. It was developed by Guido van Rossum in late 1980's which supports both functional and structured programming methods.  Sas is a programming environment and language for data manipulation and a leader in analytics, developed by the SAS Institute in 1966 and further developed in 1980's and 1990's. SAS is easily accessible, managable and can analyze data from any sources.  Apache Spark   Apache Spark is a fast large-scale data processing engine and executes applications in Hadoop clusters 100 times faster in memory and 10 times faster on disk.  Excel is a basic, popular and widely used analytical tool almost in all industries. Whether you are an expert in Sas, R or Tableau, you will still need to use Excel.	14.	Analysis	20	reports in order to extract meaningful, actionable insights, which can be used to better understand and improve business performance.	I
16. Tableau	15.	R Language	39	and widely used for statistics and data modeling. It can easily manipulate your data	I
which is easy to read, write, maintain and is a free open source tool. It was developed by Guido van Rossum in late 1980's which supports both functional and structured programming methods.  Sas is a programming environment and language for data manipulation and a leader in analytics, developed by the SAS Institute in 1966 and further developed in 1980's and 1990's. SAS is easily accessible, managable and can analyze data from any sources.  Apache Spark is a fast large-scale data processing engine and executes applications in Hadoop clusters 100 times faster in memory and 10 times faster on disk.  Excel is a basic, popular and widely used analytical tool almost in all industries. Whether you are an expert in Sas, R or Tableau, you will still need to use Excel.  A sampling distribution is	16.	Tableau		any data source be it corporate Data	I
language for data manipulation and a leader in analytics, developed by the SAS Institute in 1966 and further developed in 1980's and 1990's. SAS is easily accessible, managable and can analyze data from any sources.  Apache Spark is a fast large-scale data processing engine and executes applications in Hadoop clusters 100 times faster in memory and 10 times faster on disk.  Excel is a basic, popular and widely used analytical tool almost in all industries. Whether you are an expert in Sas, R or Tableau, you will still need to use Excel.  A sampling distribution is a leader in analytical in analytical tool almost in all industries.	17.	Python	SMIN	which is easy to read, write, maintain and is a free open source tool. It was developed by Guido van Rossum in late 1980's which supports both functional and structured	I
Apache Spark  Apache Spark is a fast large-scale data processing engine and executes applications in Hadoop clusters 100 times faster in memory and 10 times faster on disk.  Excel is a basic, popular and widely used analytical tool almost in all industries. Whether you are an expert in Sas, R or Tableau, you will still need to use Excel.  A sampling distribution is a L	18.	Sas	istd	language for data manipulation and a leader in analytics, developed by the SAS Institute in 1966 and further developed in 1980's and 1990's. SAS is easily accessible, managable	I
20. Excel  analytical tool almost in all industries. Whether you are an expert in Sas, R or Tableau, you will still need to use Excel.  A sampling distribution is a L	19.	Apache Spark		Apache Spark is a fast large-scale data processing engine and executes applications in Hadoop clusters 100 times faster in memory	I
21 Sampling distribution A sampling distribution is a I	20.	Excel		analytical tool almost in all industries. Whether you are an expert in Sas, R or	I
21.   Sampling distribution   3-	21.	Sampling distribution		A sampling distribution is a	I

			probability distribution of a statistic obtained			
			from a larger number of samples drawn from a			
			specific population.			
			Thus given fortuna of a compline			
			Three primary factors of a sampling distribution:			
	Three primary feators		distribution.			
22.	Three primary factors of a sampling		The number observed in a population	I		
22.	distribution		<ul> <li>The number observed in a population</li> <li>The number observed in the sample</li> </ul>			
	distribution		The method of choosing the sample			
			The meaner of choosing the sample			
			Resampling is the method that consists of			
			drawing repeated samples from the			
23.	Resampling		original data samples. The method	I		
			of Resampling is a nonparametric method of			
			statistical inference.			
			Statistical inference is the process of			
			using data analysis to deduce properties of an			
24	Ctatisticalinform		underlying distribution of probability. It is	I		
24.	Statistical inference		assumed that the observed <i>data</i> set is sampled			
			from a larger population.			
		10 Albert	nom a rarger population.			
			A prediction error is the failure of some			
25.	Prediction error		expected event to occur. Applying that type of	I		
25.	1 rediction error	The same of	knowledge can inform decisions and improve			
			the quality of future predictions.			
		Unit-II : Mining Data Streams				
		77.0	Streaming Applications			
	h	70				
		7.9	Streaming Applications			
		75	Streaming Applications Sensor networks  - Monitor habitat and environmental parameters			
26	Streaming		Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend			
26.	Streaming Applications		Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis	II		
26.	<u>o</u>		Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies	II		
26.	<u>o</u>		Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies  - Monitor power grid, customer usage	Ш		
26.	<u>o</u>		Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies  - Monitor power grid, customer usage patterns etc.	II		
26.	<u>o</u>		Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies  - Monitor power grid, customer usage patterns etc.  - Alerts and rapid response in case of	II		
26.	<u>o</u>		Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies  - Monitor power grid, customer usage patterns etc.  - Alerts and rapid response in case of problems	II		
26.	<u>o</u>	Sto	Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies  - Monitor power grid, customer usage patterns etc.  - Alerts and rapid response in case of	II		
26. 27.	<u>o</u>	Sto	Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies  - Monitor power grid, customer usage patterns etc.  - Alerts and rapid response in case of problems  Streaming data is data that is continuously	II		
	Applications	Surv	Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies  - Monitor power grid, customer usage patterns etc.  - Alerts and rapid response in case of problems  Streaming data is data that is continuously generated by different sources.			
	Applications	SNIN	Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies  - Monitor power grid, customer usage patterns etc.  - Alerts and rapid response in case of problems  Streaming data is data that is continuously generated by different sources. Such data should be processed incrementally			
	Applications	Stid	Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies  - Monitor power grid, customer usage patterns etc.  - Alerts and rapid response in case of problems  Streaming data is data that is continuously generated by different sources. Such data should be processed incrementally using Stream Processing techniques without			
27.	Applications  Streaming data  Benefits of streaming	Sto	Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies  - Monitor power grid, customer usage patterns etc.  - Alerts and rapid response in case of problems  Streaming data is data that is continuously generated by different sources. Such data should be processed incrementally using Stream Processing techniques without having access to all of the data.  The top benefits of streaming analytics are:  • Improve operational efficiencies.	II		
	Applications  Streaming data	Std	Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies  - Monitor power grid, customer usage patterns etc.  - Alerts and rapid response in case of problems  Streaming data is data that is continuously generated by different sources. Such data should be processed incrementally using Stream Processing techniques without having access to all of the data.  The top benefits of streaming analytics are:			
27.	Applications  Streaming data  Benefits of streaming	Std	Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies  - Monitor power grid, customer usage patterns etc.  - Alerts and rapid response in case of problems  Streaming data is data that is continuously generated by different sources. Such data should be processed incrementally using Stream Processing techniques without having access to all of the data.  The top benefits of streaming analytics are:  • Improve operational efficiencies.	II		
27.	Applications  Streaming data  Benefits of streaming	Sto	Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies  - Monitor power grid, customer usage patterns etc.  - Alerts and rapid response in case of problems  Streaming data is data that is continuously generated by different sources. Such data should be processed incrementally using Stream Processing techniques without having access to all of the data.  The top benefits of streaming analytics are:  • Improve operational efficiencies.  • Reduce infrastructure cost.  • Provide faster insights and actions.  Streaming refers to any media content – live or	II		
27.	Applications  Streaming data  Benefits of streaming analytics	Surv	Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies  - Monitor power grid, customer usage patterns etc.  - Alerts and rapid response in case of problems  Streaming data is data that is continuously generated by different sources.  Such data should be processed incrementally using Stream Processing techniques without having access to all of the data.  The top benefits of streaming analytics are:  • Improve operational efficiencies.  • Reduce infrastructure cost.  • Provide faster insights and actions.  Streaming refers to any media content – live or recorded – delivered to computers and mobile	II		
27.	Applications  Streaming data  Benefits of streaming	Sto	Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies  - Monitor power grid, customer usage patterns etc.  - Alerts and rapid response in case of problems  Streaming data is data that is continuously generated by different sources.  Such data should be processed incrementally using Stream Processing techniques without having access to all of the data.  The top benefits of streaming analytics are:  • Improve operational efficiencies.  • Reduce infrastructure cost.  • Provide faster insights and actions.  Streaming refers to any media content – live or recorded – delivered to computers and mobile devices via the internet and played back in real	II		
27. 28.	Applications  Streaming data  Benefits of streaming analytics  Streaming	Surv	Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies  - Monitor power grid, customer usage patterns etc.  - Alerts and rapid response in case of problems  Streaming data is data that is continuously generated by different sources.  Such data should be processed incrementally using Stream Processing techniques without having access to all of the data.  The top benefits of streaming analytics are:  • Improve operational efficiencies.  • Reduce infrastructure cost.  • Provide faster insights and actions.  Streaming refers to any media content – live or recorded – delivered to computers and mobile devices via the internet and played back in real time.	II		
27.	Applications  Streaming data  Benefits of streaming analytics	Sto	Streaming Applications Sensor networks  - Monitor habitat and environmental parameters  - Track many objects, intrusions, trend analysis Utility Companies  - Monitor power grid, customer usage patterns etc.  - Alerts and rapid response in case of problems  Streaming data is data that is continuously generated by different sources.  Such data should be processed incrementally using Stream Processing techniques without having access to all of the data.  The top benefits of streaming analytics are:  • Improve operational efficiencies.  • Reduce infrastructure cost.  • Provide faster insights and actions.  Streaming refers to any media content – live or recorded – delivered to computers and mobile devices via the internet and played back in real	II		

the data and streaming it back out as a single flow.  Stream sampling is the process of collecting a representative sample of the elements of a data stream.  Stream sampling is the process of collecting a representative sample of the elements of a data stream.  Stream sampling is the process of collecting a representative sample of the elements of a data stream.  Four main types of probability sample  Four main types of probability sample  Simple random sampling  Stratified sampling  In computer science, the count-distinct elements in a data stream with repeated elements.  In computer science, the count-distinct elements in a data stream with repeated elements.  Different streaming data types  Permutations, Graph Data, Geometric Data (Location Streams)  Different streaming processing models  Stratified sampling  Stratified sampling  In the stream it is independent of other items of the same stream item is independent of other items of data stream.  In computer science, the count-distinct elements in a data stream.  In computer science, the count-distinct elements in a data stream.  In computer science, the count-distinct elements of finding the number of distinct pro					
Stream sampling   Stream sampling is the process of collecting a representative sample of the elements of a data stream.				to mean pulling in streams of data; processing	
31. Stream sampling Stream sampling is the process of collecting a representative sample of the elements of a data stream.  32. Stream sampling Stream sampling is the process of collecting a representative sample of the elements of a data stream.  Stream sampling is the process of collecting a representative sample of the elements of a data stream.  Four main types of probability sample Simple random sampling Systematic sampling Stratified sampling Stratified sampling Cluster sampling Stratified sampling Cluster sampling Stratified sampling Cluster sampling Simple random sampling Systematic sampling Systematic sampling Systematic sampling Systematic sampling Systematic sampling Simple random sampling Simple random sampling Systematic sampling Simple random sampling Systematic sampling Simple random sampling Simple random sampling Systematic sampling Simple random				_	
31. Stream sampling representative sample of the elements of a data stream.  32. Stream sampling representative sample of the elements of a data stream.  Stream sampling is the process of collecting a representative sample of the elements of a data stream.  Four main types of probability sample  Simple random sampling Streatfied sampling Stream is independent of other items of the same stream or any other data stream.  In computer science, the count-distinct problem is the problem of finding the number of distinct elements in a data stream with repeated elements.  Different streaming data types - Permutations, Graph Data, Geometric Data (Location Streams)  Different streaming processing models - Sliding Windows, Exponential and other decay, Duplicate sensitivity, Random order streams, Skewed streams  Different streaming scenarios - Distributed computations, sensor network computations  Pattern finding: finding common patterns or features  - Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues - Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring Learning and Predicting					
Stream sampling   Stream sampling is the process of collecting a representative sample of the elements of a data stream.   Four main types of probability sample   Simple random sampling   Systematic sampling   Systematic sampling   Systematic sampling   Systematic sampling   Stratified sampling   Systematic sampling   Systemat					
32. Stream sampling Stream sampling is the process of collecting a representative sample of the elements of a data stream.  Four main types of probability sample  Four main types of probability sample  Simple random sampling Systematic sampling Stratified samplin	31.	Stream sampling		representative sample of the elements of a data	II
32. Stream sampling representative sample of the elements of a data stream.  Four main types of probability sample  Simple random sampling Systematic sampling Systematic sampling Cluster sampling Cluster sampling Filtering condition of a stream item is independent of other items of the same stream or any other data stream.  In computer science, the count-distinct problem  Different streaming 36. data types  Different streaming processing models  Siding Windows, Exponential and other decay, Duplicate sensitivity, Random order streams, Skewed streams  Different streaming 38. scenarios  Different streaming constraints computations  Different streaming processing models  Sliding Windows, Exponential and other decay, Duplicate sensitivity, Random order streams, Skewed streams  Different streaming computations  Pattern finding: finding common patterns or features  Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues  Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring  Learning and Predicting				stream.	
Four main types of probability sample  Four main types of probability sample  Simple random sampling Systematic sampling Systematic sampling Cluster sampling Cluster sampling Filtering condition of a stream item is independent of other items of the same stream or any other data stream.  Filtering stream  Filtering condition of a stream item is independent of other items of the same stream or any other data stream.  In computer science, the count-distinct problem of finding the number of distinct elements in a data stream with repeated elements.  Different streaming data types  Permutations, Graph Data, Geometric Data (Location Streams)  Different streaming processing models  Sliding Windows, Exponential and other decay, Duplicate sensitivity, Random order streams, Skewed streams  Different streaming scenarios  Different streaming conditions, sensor network computations  Pattern finding: finding common patterns or features  Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues  Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring  Learning and Predicting				Stream sampling is the process of collecting a	
Four main types of probability sample  Four main types of probability sample  Simple random sampling Systematic sampling Systematic sampling Cluster sampling Cluster sampling Cluster sampling Filtering condition of a stream item is independent of other items of the same stream or any other data stream.  In computer science, the count-distinct problem of finding the number of distinct elements in a data stream with repeated elements.  Different streaming processing models  Different streaming processing models Similar problem  Different streaming processing models Siding Windows, Exponential and other decay, Duplicate sensitivity, Random order streams, Skewed streams  Different streaming scenarios Different streaming scenarios Different streaming scenarios Pattern finding: finding common patterns or features Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations Data Quality Issues  Learning and  Filtering condition of a stream item is independent of other items of the same stream. III	32.	Stream sampling		representative sample of the elements of a data	II
Four main types of probability sample				stream.	
Four main types of probability sample  Simple random sampling Systematic sampling Stratified sampling Cluster sampling Filtering condition of a stream item is independent of other items of the same stream or any other data stream.  In computer science, the count-distinct problem is the problem of finding the number of distinct elements in a data stream with repeated elements.  Different streaming data types  Different streaming processing models  Different streaming processing models  Siding Windows, Exponential and other decay, Duplicate sensitivity, Random order streams, Skewed streams  Different streaming scenarios Different streaming scenarios Different finding: finding common patterns or features  Pattern finding: finding common patterns or features  Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues Learning and  Earning and Predicting				Four main types of probability sample	
33. probability sample		Four main types of			
Systematic sampling  Cluster sampling  Cluster sampling  Cluster sampling  Cluster sampling  Filtering condition of a stream item is independent of other items of the same stream or any other data stream.  In computer science, the count-distinct problem is the problem of finding the number of distinct elements in a data stream with repeated elements.  Different streaming data types  Permutations, Graph Data, Geometric Data (Location Streams)  Different streaming processing models  Sliding Windows, Exponential and other decay, Duplicate sensitivity, Random order streams, Skewed streams  Different streaming scenarios  Different streaming scenarios  Different streaming scenarios  Different streaming scenarios  Pattern finding: finding common patterns or features  Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues  Data Quality Issues  Learning and  Learning and Predicting	22	· ·			TT
Filtering stream  Filtering condition of a stream item is independent of other items of the same stream or any other data stream.  In computer science, the count-distinct problem is the problem of finding the number of distinct elements in a data stream with repeated elements.  Different streaming data types  Permutations, Graph Data, Geometric Data (Location Streams)  Different streaming processing models  Siding Windows, Exponential and other decay, Duplicate sensitivity, Random order streams, Skewed streams  Different streaming scenarios  Different streaming scenarios  Different streaming scenarios  Pattern finding: finding common patterns or features  Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues  Learning and  Learning and Predicting	<i>55.</i>	productine y sample			11
34. Filtering stream					
34. Filtering stream			-	Cluster sampling	
Same stream or any other data stream.   In computer science, the count-distinct problem of finding the number of distinct elements in a data stream with repeated elements.   If number of distinct elements in a data stream with repeated elements.				Filtering condition of a stream item is	
In computer science, the count-distinct problem of finding the number of distinct elements in a data stream with repeated elements.    Different streaming data types   Different streaming data types   Permutations, Graph Data, Geometric Data (Location Streams)	34.	Filtering stream		independent of other items of the	II
distinct problem is the problem of finding the number of distinct elements in a data stream with repeated elements.  Different streaming data types - Permutations, Graph Data, Geometric Data (Location Streams)  Different streaming processing models - Sliding Windows, Exponential and other decay, Duplicate sensitivity, Random order streams, Skewed streams  Different streaming scenarios - Distributed computations, sensor network computations Pattern finding: finding common patterns or features - Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues - Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring Learning and Predicting			- Table	same stream or any other data stream.	
distinct problem    Different streaming   Different streaming   Clocation Streams   Different streaming   Different streaming   Different streaming   Different streaming   Different streaming   Different streaming processing models   Different streaming   Different streaming processing models   Sliding Windows, Exponential and other decay, Duplicate sensitivity, Random order streams, Skewed streams   Different streaming scenarios   Different streaming scenarios   Distributed computations, sensor network computations   Pattern finding: finding common patterns or features   Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations   Data Quality Issues   Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring   Learning and Predicting   Data Quality Issued   Data Quality I				In computer science, the count-	
Different streaming data types  36. data types  Different streaming data types  Permutations, Graph Data, Geometric Data (Location Streams)  Different streaming processing models  Sliding Windows, Exponential and other decay, Duplicate sensitivity, Random order streams, Skewed streams  Different streaming scenarios  Different streaming processing models  III	25	Count-		distinct problem is the problem of finding the	TT
Different streaming data types  - Permutations, Graph Data, Geometric Data (Location Streams)  Different streaming processing models - Sliding Windows, Exponential and other decay, Duplicate sensitivity, Random order streams, Skewed streams  Different streaming Different streaming scenarios - Distributed computations, sensor network computations  Pattern finding: finding common patterns or features - Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues - Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring  Learning and  Learning and Predicting	33.	distinct problem	-	number of distinct elements in a data	Ш
Different streaming data types  - Permutations, Graph Data, Geometric Data (Location Streams)  Different streaming processing models - Sliding Windows, Exponential and other decay, Duplicate sensitivity, Random order streams, Skewed streams  Different streaming Different streaming scenarios - Distributed computations, sensor network computations  Pattern finding: finding common patterns or features - Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues - Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring  Learning and  Learning and Predicting				stream with repeated elements.	
36. data types Permutations, Graph Data, Geometric Data (Location Streams)  Different streaming processing models Sliding Windows, Exponential and other decay, Duplicate sensitivity, Random order streams, Skewed streams  Different streaming scenarios Distributed computations, sensor network computations  Pattern finding: finding common patterns or features Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring  Learning and Learning and Predicting		Different streaming	- A		
Different streaming processing models   Different streaming processing models   Sliding Windows, Exponential and other decay, Duplicate sensitivity, Random order streams, Skewed streams	36.	_	27		II
Different streaming processing models  - Sliding Windows, Exponential and other decay, Duplicate sensitivity, Random order streams, Skewed streams  Different streaming scenarios  - Distributed computations, sensor network computations  Pattern finding: finding common patterns or features  - Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues  - Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring  Learning and  Learning and Predicting					
37. Different streaming processing models Sliding Windows, Exponential and other decay, Duplicate sensitivity, Random order streams, Skewed streams  Different streaming scenarios Distributed computations, sensor network computations  Pattern finding: finding common patterns or features Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring  Learning and Learning and Predicting		D.CC			
decay, Duplicate sensitivity, Random order streams, Skewed streams  Different streaming scenarios  Pattern finding: finding common patterns or features  Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues  Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring  Learning and  Learning and Predicting	0.77		b. 400		TT
Streams, Skewed streams  Different streaming scenarios - Distributed computations, sensor network computations  Pattern finding: finding common patterns or features - Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues - Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring  Learning and  Learning and Learning and Predicting	37.	processing models	Callet To		Ш
Different streaming scenarios  - Distributed computations, sensor network computations  Pattern finding: finding common patterns or features  - Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues  - Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring  Learning and  Learning and Predicting			100	· · ·	
38. scenarios  - Distributed computations, sensor network computations  Pattern finding: finding common patterns or features  - Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues  - Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring  Learning and  Learning and Predicting		Different streaming			
29. Pattern finding Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring  Learning and Learning and Predicting	38.	U	Part W		II
Pattern finding: finding common patterns or features  Association rule mining, Clustering, Histograms, Wavelet & Fourier Representations  Data Quality Issues Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring  Learning and Learning and Predicting		1 0		<u>-</u>	
features  - Association rule mining, Clustering, II Histograms, Wavelet & Fourier Representations  Data Quality Issues  - Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring  Learning and Learning and Predicting				1	
Histograms, Wavelet & Fourier Representations  Data Quality Issues Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring Learning and Learning and Predicting					
Histograms, Wavelet & Fourier Representations  Data Quality Issues Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring Learning and Learning and Predicting	39.	Pattern finding		- Association rule mining, Clustering,	II
Au.  Data Quality Issues  Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring  Learning and Learning and Predicting					
40. Data Quality Issues Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring Learning and Learning and Predicting		1.073	COMPANI	Representations	
40. Data Quanty Issues  Change Detection, Data Cleaning, Anomaly detection, Continuous Distributed Monitoring  Learning and Learning and Predicting		Data Ovalita I			
detection, Continuous Distributed Monitoring Learning and Learning and Predicting	40.	Data Quality Issues			II
Learning and Predicting			LOT M		
		Learning and			
41.   Predicting     Building Decision Trees, Regression,   II	41.	Predicting		- Building Decision Trees, Regression,	II
Supervised Learning					
Six rules to represent a stream by buckets				Six rules to represent a stream by buckets	
The right end of a bucket is always a					
position with a 1.				position with a 1.	
Every position with a 1 is in some				<u> </u>	
hucket		Civ mulas to marragent			
42. Six rules to represent  No position is in more than one bucket. II	42.	-		No position is in more than one bucket.	II
1 a stream by blickets		a stream by buckets		There are one or two buckets of any	
				given size, up to some maximum size.	
				All sizes must be a power of 2.	
given size, up to some maximum size.				Buckets cannot decrease in size as we	
given size, up to some maximum size.  • All sizes must be a power of 2.				move to the left (back in time).	

43.	Decaying window		In a decaying window, you assign a score or weight to every element of the incoming data stream. Further, you need to calculate the aggregate sum for each distinct element by	П
			adding all the weights assigned to that element. The element with the highest total score is listed as trending or the most popular.	
			Real-time analytics • Refers to finding meaningful patterns in data	
44.	Real-time analytics		at the actual time of receiving  • Real-Time Analytics Platform (RTAP) analyses the data, correlates, and predicts the	II
		7-1	outcomes in the real time.	
45.	Benefits of RTAP		Benefits of RTAP  • Manages and processes data and helps timely decision-making  • Helps to develop dynamic analysis applications	П
			• Leads to evolution of business intelligence	
46.	Widely used RTAPs		<ul> <li>Widely used RTAPs</li> <li>Apache Spark Streaming—a Big Data platform for data stream analytics in real time.</li> <li>Cisco Connected Streaming Analytics (CSA)—a platform that delivers insights from high-velocity streams of live data from multiple sources and enables immediate action.</li> </ul>	П
47.	IBM Stream Computing	SMA	IBM Stream Computing  —a data streaming tool that analyzes a broad range of streaming data  — unstructured text, video, audio, geospatial, sensor  — helping organizations spot the opportunities and risks and make decisions in real time	П
48.	Sentiment Analysis other names	std	Sentiment Analysis other names  Opinion extraction Opinion mining Sentiment mining Subjectivity analysis	П
49.	Why Sentiment analysis?		Why Sentiment analysis?  Movie: Is this review positive or negative?  Products: What do people think about the new iPhone?  Public sentiment: How is consumer confidence? Is despair increasing?  Politics: What do people think about this candidate or issue?  Prediction: Predict election outcomes or market trends from sentiment	II

50.	Sentiment Analysis		Sentiment Analysis is the process of determining whether a piece of writing is positive, negative or neutral. Sentiment analysis helps data analysts within large enterprises gauge public opinion, conduct nuanced market research, monitor brand and product reputation, and understand customer experiences.	П
		Unit-III : F	Hadoop Environment	
51.	Hadoop features		Hadoop features:  Open Source Highly Scalable Runs on Commodity Hardware Has a good ecosystem	III
52.	YARN components		YARN components: Resource Manager: Runs on a master daemon and manages the resource allocation in the cluster. Node Manager: They run on the slave daemons and are responsible for the execution of a task on every single Data Node.	III
53.	YARN application components	39	YARN application components: Client ApplicationMaster(AM) Container	III
54.	Hosts View		Hosts View The host name, IP address, number of cores, memory, disk usage, current load average, and Hadoop components are listed in this window in tabular form.	III
55.	HDFS in Safe Mode - command	4	HDFS in Safe Mode - command: To Enter hdfs dfsadmin -safemode enter To Leave hdfs dfsadmin -safemode leave	III
56.	fsck	istd	fsck stands for File System Check. It is a command used by HDFS. This command is used to check inconsistencies and if there is any problem in the file. For example, if there are any missing blocks for a file, HDFS gets notified through this command.	III
57.	Components of HDFS		NameNode – This is the master node for processing metadata information for data blocks within the HDFS  DataNode/Slave node – This is the node which acts as slave node to store the data, for processing and use by the NameNode	III
58.	NameNode		NameNode – This is the master node for processing metadata information for data blocks within the HDFS	III

59.	DataNode/Slave node		DataNode/Slave node – This is the node which acts as slave node to store the data, for processing and use by the NameNode	III
60.	BackupNode		BackupNode- It is a read-only NameNode which contains file system metadata information excluding the block locations	III
61.	What happens when two users try to access the same file in the HDFS		HDFS NameNode supports exclusive write only. Hence, only the first user will receive the grant for file access and the second user will be rejected.	III
62.	Rack Awareness		It is an algorithm applied to the NameNode to decide how blocks and its replicas are placed. Depending on rack definitions network traffic is minimized between DataNodes within the same rack.	III
63.	HDFS Block Vs Input Split		The HDFS divides the input data physically into blocks for processing which is known as HDFS Block.  Input Split is a logical division of data by mapper for mapping operation	III
64.	Common input formats in Hadoop	20	Text Input Format Sequence File Input Format Key Value Input	III
65.	Pseudo-Distributed Mode		Pseudo-Distributed Mode – In the pseudo-distributed mode, Hadoop runs on a single node just like the Standalone mode. In this mode, each daemon runs in a separate Java process. As all the daemons run on a single node, there is the same node for both the Master and Slave nodes.	III
66.	Standalone (Local) Mode	stid	Standalone (Local) Mode – By default, Hadoop runs in a local mode i.e. on a non-distibuted, single node. This mode uses the local file system to perform input and output operation.	III
67.	Fully – Distributed Mode		Fully – Distributed Mode – In the fully-distributed mode, all the daemons run on separate individual nodes and thus forms a multi-node cluster. There are different nodes for Master and Slave nodes.	III
68.	Hadoop default block size		Hadoop default block size The default block size in Hadoop 1 is: 64 MB The default block size in Hadoop 2 is: 128 MB	III
69.	Distributed Cache		Distributed Cache is a feature of Hadoop MapReduce framework to cache files for applications. Hadoop framework makes cached files available for every map/reduce tasks running on the data nodes.	III

70.	core-site.xml		core-site.xml – This configuration file contains Hadoop core configuration settings, for example, I/O settings, very common for MapReduce and HDFS. It uses hostname a port.	III
71.	mapred-site.xml		mapred-site.xml – This configuration file specifies a framework name for MapReduce by setting mapreduce.framework.name	III
72.	hdfs-site.xml		hdfs-site.xml – This configuration file contains HDFS daemons configuration settings. It also specifies default block permission and replication checking on HDFS.	III
73.	yarn-site.xml		yarn-site.xml – This configuration file pecifies configuration settings for ResourceManager and NodeManager	III
74.	MapReduce	H	MapReduce is a programming model in Hadoop for processing large data sets over a cluster of computers, commonly known as HDFS. It is a parallel programming model.	III
75.	Two phases of MapReduce operation		Map phase – In this phase, the input data is split by map tasks. The map tasks run in parallel. These split data is used for analysis purpose.  Reduce phase- In this phase, the similar split data is aggregated from the entire collection and shows the result.	III
	Unit-IV	: Data Analy	ysis Systems and Visualization	
76.	Link Analysis	7.9	Link Analysis deals with mining useful information from linked structures like graphs. Graphs have vertices representing objects and	IV
		4.0	links among those vertices representing relationships among those objects.	1,
77.	Link mining	S NI A	relationships among those objects.  Link mining works with graph structures that have nodes with defined set of properties.  These nodes may be of the same type (homogeneous) or different (heterogeneous).	IV
77.	Link mining  Hyperlink	s vi v	relationships among those objects.  Link mining works with graph structures that have nodes with defined set of properties.  These nodes may be of the same type	
	10.83	std	relationships among those objects.  Link mining works with graph structures that have nodes with defined set of properties. These nodes may be of the same type (homogeneous) or different (heterogeneous).  The most common interpretation of the word link today is hyperlink—a means of connecting two web documents wherein activating a special element embedded in one	IV
78.	Hyperlink	std	relationships among those objects.  Link mining works with graph structures that have nodes with defined set of properties. These nodes may be of the same type (homogeneous) or different (heterogeneous).  The most common interpretation of the word link today is hyperlink—a means of connecting two web documents wherein activating a special element embedded in one document takes you to the other.  A link represents a relationship and connects two objects that are related to each other in	IV IV

82.	Link analysis		Link analysis is a data-analysis technique used to evaluate relationships (connections) between nodes. Relationships may be identified among various types of nodes (objects), including organizations, people and transactions.	IV
83.	LOC		LOC (Link-based Object Classification) is a technique used to assign class labels to nodes according to their link characteristics.	IV
84.	PageRank		PageRank is an algorithm that addresses the Link-based Object Ranking (LOR) problem. The objective is to assign a numerical rank or priority to each web page by exploiting the "link" structure of the web.	IV
85.	importance of a web page rating		The importance of a web page can be rated based on the number of backlinks to that page and the importance of the web pages that provide these backlinks, i.e., a web page referred to by important and reliable web pages, is important and reliable.	IV
86.	Backlink		A backlink of a page Pu is a citation to Pu from another page	IV
87.	In-degree, out-degree	30	deg (P) – The number of links coming into a page P (in-degree of P) deg (P) + The number of links going out of a page P (outdegree of P)	IV
88.	HITS		The Hyperlink-Induced Topic Search (HITS) algorithm was originally proposed by Kleinberg (1999) as a method of filtering results from web page search engines in order to identify results most relevant to a user query.	IV
89.	Recommender system	5414	Recommender system – The objective is to develop a system that recommends choices based on user behavior. Netflix is the characteristic example of this data product, where based on the ratings of users, other movies are recommended	IV
90.	Dashboard	istd	Dashboard – Business normally needs tools to visualize aggregated data. A dashboard is a graphical mechanism to make this data accessible.	IV
91.	content based recommender		A content based recommender works with data that the user provides, either explicitly (rating) or implicitly (clicking on a link). Based on that data, a user profile is generated, which is then used to make suggestions to the user.	IV
92.	Core components of recommender system		Data collection and processing Recommender model Recommendation post-processing Online modules User interface	IV
93.	Collaborative filtering		Collaborative filtering is a technique that can filter out items that a user might like on	IV

		T		
			the basis of reactions by similar users. It works	
			by searching a large group of people and	
			finding a smaller set of users with tastes	
94.	Dimensionality reduction in recommender systems		similar to a particular user  There are two ways of using dimensionality reduction in recommender systems: The first is creating latent factor models which reduce the dimensions of both users and items simultaneously, and produce a dense matrix, which can generate rating predictions.	IV
95.	Data visualization		Data visualization is the graphical representation of information and data. In the world of Big Data, data visualization tools and technologies are essential to analyze massive amounts of information and make data-driven decisions.	IV
96.	VR		Virtual reality is going to have a huge impact on the potential for data visualizations, allowing people to interact with data in the third dimension for the first time.	IV
97.	Common general types of data visualization		Common general types of data visualization:	IV
98.	Big Data visualization	X	Big Data visualization involves the presentation of data of almost any type in a graphical format that makes it easy to understand and interpret.	IV
99.	Interaction techniques	<u> </u>	Interaction techniques essentially involve data entry and manipulation, and thus place greater emphasis on input than output. Output is merely used to convey affordances and provide user feedback.	IV
100.	Four stages of Visualization	istd	Four stages of Visualization	IV
		nit-V : Fram	eworks and Applications	
			HBase is a distributed column-oriented	
101.	Hbase		database built on top of the Hadoop file system.	V
102.	Hive		Hive: It is a platform used to develop SQL type scripts to do MapReduce operations	V
103.	Features of Hive		It stores schema in a database and processed data into HDFS. It provides SQL type language for querying	V

			called HiveQL or HQL.	
			It is familiar, fast, scalable, and extensible	
			Column Types	
			• Literals	
104.	Hive - Data Types		Null Values	V
			- 1 - 1 - 1 - 1 - 1 - 1 - 1 - 1 - 1 - 1	
			Complex Types	
			Arrays: Arrays in Hive are used the same way	
4.05	Hive - Complex		they are used in Java.	• •
105.	Types		Maps: Maps in Hive are similar to Java Maps.	V
	71 ***		Structs: Structs in Hive is similar to using	
			complex data with comment	
			• Apache HBase is used to have random, real-	
			time read/write access to Big Data.	
	Where to Use HBase		• It hosts very large tables on top of clusters	
106.	THE TO USE TIDASE		of commodity hardware.	V
	, iii		<ul> <li>Apache HBase is a non-relational database</li> </ul>	
			modeled after Google's Bigtable. Bigtable acts	
			up on Google File System, likewise Apache	
			YARN components:	
			Resource Manager: Runs on a master daemon	
	VADN components		and manages the resource allocation in the	
107.	YARN components	Zo.	cluster.	V
			Node Manager: They run on the slave	
			daemons and are responsible for the execution	
			of a task on every single Data Node.	
		-	YARN application components:	
	YARN application components	h "100"	• Client	
108.			ApplicationMaster(AM)	V
			• Container	
		F-10-	Region- This component contains memory data store and Hfile.	
			Region Server-This monitors the Region.	
			HBase Master-It is responsible for monitoring	
			the region server.	
109.	Key components of		Zookeeper- It takes care of the coordination	V
	HBase	COMPANIE	between the HBase Master component and the	
			client.	
			Catalog Tables-The two important catalog	
		Estd	tables are ROOT and META.ROOT table	
			tracks where the META table is and META	
			table stores all the regions in the system.	
110	Region		Region- This component contains memory	V
110.	Region		data store and Hfile.	V
			data store and Hfile.  Zookeeper- It takes care of the coordination	
<ul><li>110.</li><li>111.</li></ul>	Region Zookeeper		data store and Hfile.  Zookeeper- It takes care of the coordination between the HBase Master component and the	V
			data store and Hfile.  Zookeeper- It takes care of the coordination	
			data store and Hfile.  Zookeeper- It takes care of the coordination between the HBase Master component and the	
			data store and Hfile.  Zookeeper- It takes care of the coordination between the HBase Master component and the client.	
	Zookeeper		data store and Hfile.  Zookeeper- It takes care of the coordination between the HBase Master component and the client.  Record Level Operational Commands in	
111.	Zookeeper  Operational		data store and Hfile.  Zookeeper- It takes care of the coordination between the HBase Master component and the client.  Record Level Operational Commands in HBase are –put, get, increment, scan and delete.	V
111.	Zookeeper  Operational		data store and Hfile.  Zookeeper- It takes care of the coordination between the HBase Master component and the client.  Record Level Operational Commands in HBase are –put, get, increment, scan and delete.  Table Level Operational Commands in HBase	V
111.	Zookeeper  Operational		data store and Hfile.  Zookeeper- It takes care of the coordination between the HBase Master component and the client.  Record Level Operational Commands in HBase are –put, get, increment, scan and delete.	V

			RDBMS does not have support for in-built	
			partitioning whereas in HBase there is	
			automated partitioning.	
			RDBMS stores normalized data whereas	
			HBase stores de-normalized data.	
	Catalog tables in		The two important catalog tables in HBase, are	
114.	HBase		ROOT and META. ROOT table tracks where	V
114.	Tibase		the META table is and META table stores all	V
			the regions in the system.	
			HBase and Hive both are completely different	
			hadoop based technologies-Hive is a data	
115.	HBase Vs Hive		warehouse infrastructure on top of Hadoop	V
115.			whereas HBase is a NoSQL key value store	•
			that runs on top of Hadoop.	
			Licence based (also Open Source)	
	Mongo DP footures		NoSQL Database	
116.	MongoDB features		NosQL Database     Document Oriented	V
		-		
			Aggregation Pipeline etc.	
			Open Source     N. GOL B	
<del>.</del> .	Cassandra features		NoSQL Database	
117.			Log-Structured Storage	V
			Includes Cassandra Structure Language	
			(CQL)	
		100	NoSQL Database is a non-relational Data	
	11 75		Management System, that does not require a	
118.	NoCOL Databasa	The second second	fixed schema. It avoids joins, and is easy to	V
110.	NoSQL Database	b"	scale. The major purpose of using a NoSQL	V
			database is for distributed data stores with	
		100	humongous data storage needs	
110	Scaleup or Vertical		Scaleup or Vertical Scaling: Increase of RAM,	<b>T</b> 7
119.	Scaling	P	CPU, and HDD	V
	Scaleout or		Scaleout or Horizontal Scaling: Increase of	
120.	Horizontal Scaling		Commodity hardware	V
	Tionzonar zeamg		Types of NoSQL Databases:	
			Key-value Pair Based	
121.	Types of NoSQL	GALA	Column-oriented Graph	V
141,	Databases		Graphs based	¥
	10.5		Document-oriented	
		المراط مراثا	Key Value Pair Based	
			Data is stored in key/value pairs. It is designed	
	W Wi D' D		in such a way to handle lots of data and heavy	
122.	Key Value Pair Based		load. Key-value pair storage databases store	V
			data as a hash table where each key is unique,	•
			and the value can be a JSON, BLOB (Binary	
			Large Objects), string, etc. eg. DynamoDB,	
			Redis, etc.	
			Column-oriented databases work on columns	
	Column boood		and are based on BigTable paper by Google.	
123.	Column-based		Every column is treated separately. Values of	V
			single column databases are stored	
			contiguously. eg.Cassandra, HBase, etc.	
			Document-Oriented NoSQL DB stores and	
124.	Documents-Oriented		retrieves data as a key value pair but the value	V
147.	2 ocamona official		part is stored as a document. The document is	•
<u> </u>			part is stored as a document. The document is	

		<u> </u>	11 1001 177 0						
			stored in JSON or XML formats. The value is						
			understood by the DB and can be queried. eg.						
			CouchDB, MongoDB,etc.						
125.	Graph-Based		A graph type database stores entities as well the relations amongst those entities. The entity is stored as a node with the relationship as edges. An edge gives a relationship between nodes. Every node and edge has a unique identifier. eg. Neo4j, OrientDB,etc.	V					
Placement Questions									
			Text mining is the art and science of						
126.	Text mining		discovering knowledge, insights, and patterns from an organized collection of textual databases.						
127.	Naïve Bayes technique		Naïve Bayes technique is a supervised machine learning technique that that uses probability theory based analysis.						
128.	Support Vector Machine	26	Support Vector Machine (SVM) is a supervised machine learning algorithm which can be used for both classification and regression challenges.						
129.	Web mining	70	Web mining is the art and science of discovering patterns and insights from the World Wide Web so as to improve it.						
130.	Business Intelligence		Business Intelligence (BI) is an umbrella term that includes a variety of IT applications that are used to analyze an organization's data and communicate the information to relevant users.						
131.	Applications of BI and data mining	7.0	Retail, Telecom, Customer Relationship Management, Healthcare and Wellness, Education, Banking, Financial Services, Insurance, Manufacturing, and Public Sector						
132.	Data warehouse	istd	A data warehouse (DW) is an organized collection of integrated, subject oriented databases designed to support decision support functions						
133.	Data mining		Data mining is the art and science of discovering knowledge, insights, and patterns in data.						
134.	Classification techniques		Classification techniques are called supervised learning as there is a way to supervise whether the model's prediction is right or wrong.						
135.	Decision tree		A decision tree is a hierarchically organized branched, structured to help make decision in an easy and logical manner.						
136.	Regression		Regression is a relatively simple and the most popular statistical data mining technique. The goal is to fit a smooth well-defined curve to						

				-
			the data. Regression analysis techniques, for example, can be used to model and predict the	
			energy consumption as a function of daily	
			temperature.	
			Artificial neural network (ANN) is a	
			sophisticated data mining technique from the	
	Artificial neural		Artificial Intelligence stream in Computer	
137.			Science. It mimics the behavior of human	
137.	network		neural structure: Neurons receive stimuli,	
			process them, and communicate their results to	
			other neurons successively, and eventually a	
			neuron outputs a decision.	
			Cluster analysis is an exploratory learning	
100	Classian and last		technique that helps in identifying a set of	
138.	Cluster analysis		similar groups in the data. It is a technique	
			used for automatic identification of natural	
			groupings of things.  Association rules are a popular data mining	
			method in business, especially where selling is	
139.	Association rules		involved. Also known as market basket	
107.	1 100001ation 1 titos		analysis, it helps in answering questions about	
			cross-selling opportunities	
			NFS (Network File System) is one of the	
		1	oldest and popular distributed file storage	
140.	NFS Vs HDFS		systems whereas HDFS (Hadoop Distributed	
	11 22		File System) is the recently used and popular	
		b_ 700	one to handle big data.	
			Data which can be stored in traditional	
			database systems in the form of rows and	
141.	Structured Data	b 100	columns, for example the online purchase	
		Alba."	transactions can be referred to as Structured	
			Data.	
		7	Data which can be stored only partially in traditional database systems, for example, data	
142.	Semi structured data.		in XML records can be referred to as semi	
			structured data.	
	n n n n	CONTRACTOR	Unorganized and raw data that cannot be	
			categorized as semi structured or structured	
140	I In atms at seed date		data is referred to as unstructured data.	
143.	Unstructured data	LOT M	Facebook updates, Tweets on Twitter,	
	_	200	Reviews, web logs, etc. are all examples of	
			unstructured data.	
	Two ways of Big		Two ways of Big Data processing	
144.	Data processing		1. Batch processing	
	p1300001115		2. Stream processing	
			Data Science Vs Big Data	
	Data Science Vs Big Data		Data science is a broad spectrum of      stirities involving analysis of Big	
145.			activities involving analysis of Big	
			Data, finding patterns, trends in data,	
			interpreting statistical terms and predicting future trends.	
			<ul> <li>Big Data is just one part of Data</li> </ul>	
			Science. Though Data Science is a	
			science. Though Data science is a	

			<ul> <li>broad term and very important in the overall Business operations, it is nothing without Big Data.</li> <li>All the activities we perform in Data Science are based on Big Data. Thus Big Data and Data Science are</li> </ul>	
			interrelated and cannot be seen in isolation.	
146.	Cloud computing		Cloud computing is internet-based computing. It relies on sharing computing resources ondemand rather than having local servers or PCS and other devices.	
147.	Rule induction		Rule induction is an area of machine learning in which formal rules are extracted from a set of observations.	
148.	Sensor networks		Sensor networks are a huge source of data occurring in streams. They are used in numerous situations that require constant monitoring of several variables, based on which important decisions are made.	
149.	Bloom Filter	20	A Bloom Filter is a space-efficient probabilistic data structure, conceived by Burton Howard Bloom in 1970, that is used to test whether an element is a member of set.	
150.	Reservoir sampling	70	Biased reservoir sampling is defined as bias function to regulate the sampling from the stream.	

Dr.M.Moorthy

**Faculty Prepared** 

Signature

Estd. 2000

HoD