



Which of the following is an example of big data utilized in action today?

- A. Individual, Unconnected Hospital Databases
- B. The Internet
- C. Wi-Fi Networks
- D. Social Media

D. Social Media

What reasoning was given for the following: why is the "data storage to price ratio" relevant to big data?

- A. Larger storage means easier accessibility to big data for every user because it allows users to download in bulk.
- B. It isn't, it was just an arbitrary example of big data usage.
- C. Companies can't afford to own, maintain, and spend the energy to support large data storage unless the cost is sufficiently low.
- D. Lower prices mean larger storage becomes easier to access for everyone, creating bigger amounts of data for client-facing services to work with.

D. Lower prices mean larger storage becomes easier to access for everyone, creating bigger amounts of data for client-facing services to work with.

What is the best description of personalized marketing enabled by big data?

- A. Marketing to each customer on an individual level and suiting to their needs.
- B. Being able to obtain and use customer information for groups of consumers and utilize them for marketing needs.
- C. Being able to use personalized data from every single customer for personalized marketing needs.

C. Being able to use personalized data from every single customer for personalized marketing needs.

Of the following, which are some examples of personalized marketing related to big data?

- A. Facebook revealing posts that cater towards similar interests.
- B. News outlets gathering information from the internet in order to report them to the public.
- C. A survey that asks your age and markets to you a specific brand.

A. Facebook revealing posts that cater towards similar interests.

What is the workflow for working with big data?

- A. Extrapolation -> Understanding -> Reproducing
- B. Theory -> Models -> Precise Advice
- C. Big Data -> Better Models -> Higher Precision

C. Big Data -> Better Models -> Higher Precision

Which is the most compelling reason why mobile advertising is related to big data?

- A. Mobile advertising benefits from data integration with location which requires big data.
- B. Since almost everyone owns a cell/mobile phone, the mobile advertising market is large and thus requires big data to contain all the information.
- C. Mobile advertising allows massive cellular/mobile texting to a wide audience, thus providing large amounts of data.
- D. Mobile advertising in and of itself is always associated with big data.

A. Mobile advertising benefits from data integration with location which requires big data.

What are the three types of diverse data sources?

- A. Information Networks, Map Data, and People
- B. Machine Data, Map Data, and Social Media
- C. Machine Data, Organizational Data, and People
- D. Sensor Data, Organizational Data, and Social Media

C. Machine Data, Organizational Data, and People

What is an example of machine data?

- A. Weather station sensor output.
- B. Social Media
- C. Sorted data from Amazon regarding customer info.

A. Weather station sensor output.

What is an example of organizational data?

- A. Social Media
- B. Satellite Data
- C. Disease data from Center for Disease Control.

C. Disease data from Center for Disease Control.



Of the three data sources, which is the hardest to implement and streamline into a model?

- A. Machine Data
- B. Organizational Data
- C. People

Which of the following summarizes the process of using data streams?

- A. Integration -> Personalization -> Precision
- B. Big Data -> Better Models -> Higher Precision
- C. Theory -> Models -> Precise Advice
- D. Extrapolation -> Understanding -> Reproducing

C. People

A. Integration -> Personalization -> Precision

Where does the real value of big data often come from?

- A. Using the three major data sources: Machines, People, and Organizations.
- B. Combining streams of data and analyzing them for new insights.
- C. Having data-enabled decisions and actions from the insights of new data.
- D. Size of the data.

B. Combining streams of data and analyzing them for new insights.

What does it mean for a device to be "smart"?

- A. Collect data and services autonomously.
- B. Having a specific processing speed in order to keep up with the demands of data processing.
- C. Must have a way to interact with the user.

A. Collect data and services autonomously.

What does the term "in situ" mean in the context of big data?

- A. Accelerometers.
- B. Bringing the computation to the location of the data.
- C. The sensors used in airplanes to measure altitude.
- D. In the situation

B. Bringing the computation to the location of the data.

Which of the following are reasons mentioned for why data generated by people are hard to process? Choose all that apply.

- A. Skilled people to analyze the data are hard to come by.
- B. They cannot be modeled and stored.
- C. Very unstructured data.
- D. The velocity of the data is very high.

ACD

What is the purpose of retrieval and storage; pre-processing; and analysis in order to convert multiple data sources into valuable data?

- A. Since the multi-layered process is built into the Neo4j database connection.
- B. To allow scalable analytical solutions to big data.
- C. Designed to work like the ETL process.
- D. To enable ETL methods.

B. To allow scalable analytical solutions to big data.

Which of the following are benefits of organization-generated data? Choose all that apply.

- A. Higher Sales
- B. Better Profit Margins
- C. High Velocity
- D. Improved Safety
- E. Customer Satisfaction

ABDE

What are data silos and why are they bad?

- A. Highly unstructured data. Bad because it does not provide meaningful results for organizations.
- B. A giant centralized database to house all the data production within an organization. Bad because it hinders opportunity for data generation.
- C. A giant centralized database to house all the data produced within an organization. Bad because it is hard to maintain as highly structured data.

D. Data produced from an organization that is spread out. Bad because it creates unsynchronized and invisible data.



D. Data produced from an organization that is spread out. Bad because it creates unsynchronized and invisible data.

Which of the following are benefits of data integration? Choose all that apply.

- A. Monitoring of data.
- B. Increase data collaboration.
- C. Adds value to big data.
- D. Increase data availability.
- E. Unify your data system.
- F. Reduce data complexity.

BCDEF

Amazon has been collecting review data for a particular product. They have realized that almost 90% of the reviews were mostly a 5/5 rating. However, of the 90%, they realized that 50% of them were customers who did not have proof of purchase or customers who did not post serious reviews about the product. Of the following, which is true about the review data collected in this situation?

- A. Low Valence
- B. High Valence
- C. High Volume
- D. High Veracity
- E. Low Volume
- F. Low Veracity

F. Low Veracity

As mentioned in the slides, what are the challenges to data with a high valence?

- A. Complex Data Exploration Algorithms
- B. Difficult to Integrate
- C. Reliability of Data

A. Complex Data Exploration Algorithms

Which of the following are the 6 V's in big data?

- A. Variety
- B. Valence
- C. Vision
- D. Volume
- E. Velocity
- F. Value
- G. Veracity

ABDEFG

What is the veracity of big data?

- A. The size of the data.
- B. The speed at which data is produced.
- C. The connectedness of data.
- D. The abnormality or uncertainties of data.

D. The abnormality or uncertainties of data.

What are the challenges of data with high variety?

- A. The quality of data is low.
- B. Hard to perform emergent behavior analysis.
- C. Hard to integrate.
- D. Hard in utilizing group event detection.

C. Hard to integrate.

Which of the following is the best way to describe why it is crucial to process data in real-time?

- A. More expensive to batch process.
- B. Batch processing is an older method that is not as accurate as real-time processing.
- C. Prevents missed opportunities.
- D. More accurate.

C. Prevents missed opportunities.

What are the challenges with big data that has high volume?

- A. Speed Increase in Processing
- B. Cost, Scalability, and Performance
- C. Storage and Accessibility
- D. Effectiveness and Cost

B. Cost, Scalability, and Performance



Which of the following are parts of the 5 P's of data science and what is the additional P introduced in the slides?

- A. Platforms
- B. Purpose
- C. Process
- D. Programmability
- E. Product
- F. People
- G. Perception

ABCDEF

Which of the following are part of the four main categories to acquire, access, and retrieve data?

- A. NoSQL Storage
- B. Traditional Databases
- C. Web Services
- D. Remote Data
- E. Text Files

ABDE

What are the steps required for data analysis?

- A. Regression, Evaluate, Classification
- B. Classification, Regression, Analysis
- C. Investigate, Build Model, Evaluate
- D. Select Technique, Build Model, Evaluate

D. Select Technique, Build Model, Evaluate

Of the following, which is a technique mentioned in the videos for building a model?

- A. Validation
- B. Evaluation
- C. Analysis
- D. Investigation

C. Analysis

What is the first step in finding a right problem to tackle in data science?

- A. Ask the Right Questions
- B. Define Goals
- C. Assess the Situation
- D. Define the Problem

D. Define the Problem

What is the first step in determining a big data strategy?

- A. Collect Data
- B. Build In-House Expertise
- C. Organizational Buy-In
- D. Business Objectives

D. Business Objectives

According to Ilkay, why is exploring data crucial to better modeling?

Data exploration... <complete the sentence>

- A. enables a description of data which allows visualization.
- B. enables understanding of general trends, correlations, and outliers.
- C. enables histograms and others graphs as data visualization.
- D. leads to data understanding which allows an informed analysis of the data.

D. leads to data understanding which allows an informed analysis of the data.

Why is data science mainly about teamwork?

- A. Exhibition of curiosity is required.
- B. Data science requires a variety of expertise in different fields.
- C. Engineering solutions are preferred.
- D. Analytic solutions are required.

B. Data science requires a variety of expertise in different fields.

What are the ways to address data quality issues?

- A. Remove outliers.
- B. Data Wrangling
- C. Remove data with missing values.
- D. Generate best estimates for invalid values.
- E. Merge duplicate records.

ACDE



What is done to the data in the preparation stage?	A. Cleaning, Integrating, and Packaging B. Retrieve Data C. Identify Data Sets and Query Data D. Select Analytical Techniques E. Build Models	A. Cleaning, Integrating, and Packaging
Which of the following is the best description of why it is important to learn about the foundations for big data?	A. Foundations is all that is required to show a mastery of big data concepts. B. Foundations stand the test of time. C. Foundations allow for the understanding of practical concepts in Hadoop. D. Foundations help you revisit calculus concepts required in the understanding of big data.	C. Foundations allow for the understanding of practical concepts in Hadoop.
What is the benefit of a commodity cluster?	A. Prevents individual component failures B. Prevents network connection failure C. Much faster than a traditional super computer D. Enables fault tolerance	D. Enables fault tolerance
What is a way to enable fault tolerance?	A. Distributed Computing B. Better LAN Connection C. Data-Parallel Job Restart D. System Wide Restart	C. Data-Parallel Job Restart
What are the specific benefit(s) to a distributed file system?	A. High Concurrency B. Data Scalability C. Large Storage D. High Fault Tolerance	ABD
Which of the following are general requirements for a programming language in order to support big data models?	A. Optimization of Specific Data Types B. Utilize Map Reduction Methods C. Enable Adding of More Racks D. Support Big Data Operations E. Handle Fault Tolerance	ACDE
What does IaaS provide?	A. Software On-Demand B. Computing Environment C. Hardware Only	C. Hardware Only
What does PaaS provide?	A. Hardware Only B. Software On-Demand C. Computing Environment	C. Computing Environment
What does SaaS provide?	A. Hardware Only B. Computing Environment C. Software On-Demand	C. Software On-Demand
What are the two key components of HDFS and what are they used for?	A. FASTA for genome sequence and Rasters for geospatial data. B. NameNode for block storage and Data Node for metadata. C. NameNode for metadata and DataNode for block storage.	C. NameNode for metadata and DataNode for block storage.
What is the job of the NameNode?	A. Coordinate operations and assigns tasks to Data Nodes B. Listens from DataNode for block creation, deletion, and replication. C. For gene sequencing calculations.	A. Coordinate operations and assigns tasks to Data Nodes



What is the order of the three steps to Map Reduce?

- A. Map -> Reduce -> Shuffle and Sort
- B. Shuffle and Sort -> Reduce -> Map
- C. Map -> Shuffle and Sort -> Reduce
- D. Shuffle and Sort -> Map -> Reduce

C. Map -> Shuffle and Sort -> Reduce

What is a benefit of using pre-built Hadoop images?

- A. Quick prototyping, deploying, and validating of projects.
- B. Less software choices to choose from.
- C. Guaranteed hardware support.
- D. Quick prototyping, deploying, and guaranteed bug free.

A. Quick prototyping, deploying, and validating of projects.

What are some examples of open-source tools built for Hadoop and what does it do?

- A. Giraph, for SQL-like queries.
- B. Zookeeper, analyze social graphs.
- C. Pig, for real-time and in-memory processing of big data.
- D. Zookeeper, management system for animal named related components.

D. Zookeeper, management system for animal named related components.

What is the difference between low level interfaces and high level interfaces?

- A. Low level deals with storage and scheduling while high level deals with interactivity.
- B. Low level deals with interactivity while high level deals with storage and scheduling.

A. Low level deals with storage and scheduling while high level deals with interactivity.

Which of the following are problems to look out for when integrating your project with Hadoop?

- A. Infrastructure Replacement
- B. Random Data Access
- C. Task Level Parallelism
- D. Advanced Algorithms
- E. Data Level Parallelism

ABCD

As covered in the slides, which of the following are the major goals of Hadoop?

- A. Handle Fault Tolerance
- B. Facilitate a Shared Environment
- C. Enable Scalability
- D. Optimized for a Variety of Data Types
- E. Provide Value for Data
- F. Latency Sensitive Tasks

ABCDE

What is the purpose of YARN?

- A. Allows various applications to run on the same Hadoop cluster.
- B. Enables large scale data across clusters.
- C. Implementation of Map Reduce.

A. Allows various applications to run on the same Hadoop cluster.

What are the two main components for a data computation framework that were described in the slides?

- A. Resource Manager and Node Manager
- B. Resource Manager and Container
- C. Node Manager and Applications Master
- D. Node Manager and Container
- E. Applications Master and Container

A. Resource Manager and Node Manager

(Questions 1-3 pertain to the video lecture "Exploring the Relational Data Model of CSV")

What is the approximate population of La Paz county in the state of Arizona for the CENSUS2010POP (column H)? (Choose the best answer.)

- A. 25000
- B. 10000
- C. 15000
- D. 20000

D. 20000



What county in the state of Wyoming has the smallest estimated population?

- A. Uinta
- B. Niobrara
- C. Sweetwater
- D. Platte

B. Niobrara

At 2:45 of the video, the Instructor creates a filter for all of the counties in California with a population greater than 1,000,000. However, included in the results is the entire state of California. This anomalous value might skew our analysis if, for example, we wanted to compute the average population of these results. What additional filter might work to resolve this problem?

- A. Add a filter to detect and remove results which do not include the word "County" in column G.
- B/ Add a filter which finds all counties with population greater than 100,000 AND less than 10,000,000 for column H (CENSUS2010POP).
- C. Add a filter where the value in column E is greater than 1,000,000.
- D. None of the above

A. Add a filter to detect and remove results which do not include the word "County" in column G.

(Questions 4 and 5 pertain to the video "Exploring Sensor Data")

How often (in seconds) do the R5 measurements occur?

- A. 60
- B. 50
- C. 40
- D. 30

A. 60

What is the field for rain accumulation?

- A. Sm
- B. Rc
- C. Dx
- D. Dn

B. Rc

(Questions 6 and 7 pertain to the video lecture "Exploring the Array Data Model of an Image")

What is the (Red, Green, Blue) pixel value for location 500, 2000?

- A. (163, 118, 79)
- B. (50, 156, 182)
- C. (134, 145, 46)
- D. (100, 123, 149)

A. (163, 118, 79)

Is this value likely to be land or ocean?

- A. Land
- B. Ocean

A. Land

(Questions 8 and 9 pertain to the video lecture "Exploring the Semistructured Data Model of JSON")

Given a tweet, what path would you most likely enter to obtain a count of the number of followers for a user?

- A. user/followers_count
- B. user/statuses_count
- C. user/listed_count
- D. None of the above

A. user/followers_count

Which of the following fields are nested within the 'entities' field (select all that apply)?

- A. symbols
- B. urls
- C. views
- D. tweets
- E. user_mentions
- F. events

ABE

What is a possible pitfall of utilizing Excel as a way to manipulate small databases?



- A. Excel does not allow algorithms for data manipulation.
B. Excel is a user program and thus cannot run on a server.
C. Excel does not enforce many principles of relational data models.

What does the term "atomic" mean in the context of relational databases?

- A. A column or row of data. Depends on the context.
B. A tuple that cannot be reduced.
C. One unit of information that cannot be decomposed.
D. Fixed schema of a particular database.

What is the Pareto-Optimality problem?

- A. Find the shortest path from source node to target node.
B. Find the best possible path given two or more optimization criteria where neither constraint can be fully optimized simultaneously.
C. Find the optimal path that requires going through specific nodes given by the user.

What constitutes a community within a graph?

- A. A neighborhood defined by an integer constant K around a specific node. All K+1 nodes belong in another community.
B. High density of nodes at a certain location.
C. A dense amount of edge connections between nodes in a community and a few connections across communities.
D. Many anomalous neighborhoods within the same vicinity

Why are trees useful for semi-structured data such as XML and JSON?

- A. They are only useful for XML data as tree-like structure is apparent with tags. While JSON does not contain a tree-like structure as it contains arrays.
B. Computers can easily visualize the data with a tree structure.
C. Trees take advantage of the parent-child relationship of the data for easy navigation.
D. It is not always the case that XML and JSON can be represented as trees.

What is the general purpose of modeling data as vectors?

- A. Results can be ordered by similarity using vector projection.
B. The ability to normalize vectors allowing probability distributions.
C. Enables image searching.
D. Enables weighting of the query.

For the following questions 7, 8, and 9, suppose a registration website creates data with the following fields for each person registered (note: if the user does not input a value, NULL is stored instead): Name, Date, Address, and Account Number.

Suppose we collect data month by month. Each month, we would have a batch of data containing the fields listed above. At the end of the year, we want to summarize our registrant activities for the entire year, so we would remove redundancies in our data by removing any records with duplicate account numbers from month to month. What type of operation do we use in this scenario?

- A. Join
B. Union
C. Subsetting
D. Not an Operation

From the information given in question 7, what are the constraints, if any, which we have placed on the Account Number field for the end of year collection?

- A. There are no constraints.
B. Account should have at most n digits.
C. If we had n duplicate Account Numbers then we will remove n-1

- C. Excel does not enforce many principles of relational data models.

- C. One unit of information that cannot be decomposed.

- B. Find the best possible path given two or more optimization criteria where neither constraint can be fully optimized simultaneously.

- C. A dense amount of edge connections between nodes in a community and a few connections across communities.

- C. Trees take advantage of the parent-child relationship of the data for easy navigation.

- A. Results can be ordered by similarity using vector projection.

- B. Union

- D. Account Number should be unique.



- duplicate fields.
- D. Account Number should be unique.

Suppose 100 people signup for our system and of the 100 people, 60 of them did not input an address. The system lists the values as NULL for these empty entries in the address field. Would this situation still have structure for our data?

- A. No because the majority of data do not have a specific field filled, thus our originally defined structure is lost.
- B. Yes the data has structure because we have placed a structural constraint on the data, thus the data will always have the originally defined structure.

What is true between data modeling and the formatting of the data?

- A. The data does not necessarily need to be formatted in a way that represents the data model. Just so long as it can be extrapolated.
- B. There is a one to one correspondence between formatting data and data modeling. For every model of data, there is only one way to store the data.
- C. There is always one specific schema for storing model data that is the best and preferred method for the specific data representation.

What is streaming?

- A. Using sensors to manipulate the system, such as a smart car being able to drive by itself using sensors to detect road hazards.
- B. Utilizing real time data to compute and change the state of an application continuously.
- C. Calculating results using real time data otherwise known as streaming data.
- D. Using static data stored from a real time source in order to process and guide the application.

Of the following, what best describes the properties of working with streaming data?

- A. Data is always utilized for streaming the application.
- B. Small time windows for working with data.
- C. Data manipulation is near real time.
- D. Independent computations that do not rely on previous or future data.
- E. Does not ping the source interactively for a response upon receiving the data.
- F. Always unbounded in sequence, in other words, data is not guaranteed to be in order.

BCDE

What is a characteristic of streaming data?

- A. Data is unbounded in size but requires only finite time and space to process it.
- B. The data is unbounded in size and the size determines the time and space of processing the data.
- C. The data is finite and requires only finite time and space to process the data.
- D. Data is finite in size and size determines the time and space of processing the data.

- A. Data is unbounded in size but requires only finite time and space to process it.

What type of algorithm is required for analyzing streaming data?

- A. Fast and Complex
- B. Accurate and Memory Efficient
- C. Fast and Simple
- D. Accurate and Consistent

C. Fast and Simple

What is lambda architecture?

- A. A specific method for processing streaming data using special real time processes.
- B. A specific hardware architecture for a server made specifically



for processing real time data.

- C. A method to process streaming data by utilizing batch processing and real time processing.

Of the following, which best represents the challenge regarding the size and frequency of data?

- A. The size and frequency of the streaming data may be sporadic.
B. The size and frequency of the streaming data may be too small.
C. There may not be data to produce the notion of size and frequency.

What is the difference between data lakes and data warehouses?

- A. Data lakes contain only files while data warehouses contain only databases.
B. Data lakes house raw data while data warehouses contain pre-formatted data.
C. Data lakes utilize hierarchical systems while data warehouses use object storage.

What is schema-on-read?

- A. The process where formatted data is given structure when read.
B. Another name for data lakes.
C. The process where data is pre-formatted prior to being read but the schema is loaded on read.
D. Data is stored as raw data until it is read by an application where the application assigns structure.

The desired characteristics of a BDMS include (select all that apply):

- A. Narrow range of query sizes
B. A full query language
C. Support for common "Big Data" data types
D. Support for ACID
E. Continuous data ingestion
F. A flexible semi-structured data model

BCEF

Fill in the blank with the best answer: CAP theorem states that _____ all at once within a distributed computer system?

- A. it is necessary to have consistency, accuracy, and partial tolerance
B. it is impossible to have consistency, availability, and partition tolerance
C. it is necessary to have consistency, availability, and partition tolerance
D. it is impossible to have consistency, accuracy, and partial tolerance

- C. A method to process streaming data by utilizing batch processing and real time processing.

- A. The size and frequency of the streaming data may be sporadic.

- B. Data lakes house raw data while data warehouses contain pre-formatted data.

- D. Data is stored as raw data until it is read by an application where the application assigns structure.

What is the purpose of the acronym BASE?

- A. To overcome CAP theorem.
B. The same as ACID.
C. Enables stricter enforcement of ACID type design.
D. To impose properties on a BDMS in order to guarantee certain results.

- D. To impose properties on a BDMS in order to guarantee certain results.

What are ziplists in Redis?

- A. A compressed list that is stored within the value of the database.
B. A special type of data type that can store up to 512 mb of image data.
C. A look up table that is stored as a value in the database. Look up table points to actual values in memory.
D. A special type of data type that can store hashes that point to multiple attributes.

- A. A compressed list that is stored within the value of the database.

What is one of the main features of Aerospike?

- A. Enables real time data streaming from external sources.
B. Better equipped for string based search applications.

- C. Support for geospatial data storage and geospatial queries.



- C. Support for geospatial data storage and geospatial queries.
D. Images as values within the database.

What database would be best suited for the following scenario:
An app development company is trying to implement a cloud based storage system for their new map-based app. The cloud will manage the longitude and latitude of the data in order to track user location.

C. Aerospike

- A. Vertica
B. Redis
C. Aerospike
D. Solr

What database would be best suited for the following scenario: A big wholesale company is trying to implement a search engine for their products.

B. Solr

- A. Redis
B. Solr
C. Vertica
D. Aerospike

Which of the following data types are supported by Redis? (select all that apply)

- A. Hashes
B. Streaming Video
C. Lists
D. Images
E. Sorted Sets
F. Strings

ACEF