AWS Academy Data Engineering Practice Test (Sample)

Study Guide



Everything you need from our exam experts!

Copyright © 2025 by Examzify - A Kaluba Technologies Inc. product.

ALL RIGHTS RESERVED.

No part of this book may be reproduced or transferred in any form or by any means, graphic, electronic, or mechanical, including photocopying, recording, web distribution, taping, or by any information storage retrieval system, without the written permission of the author.

Notice: Examzify makes every reasonable effort to obtain from reliable sources accurate, complete, and timely information about this product.



Questions



- 1. Which AWS services can be used to monitor and troubleshoot an AWS Glue job?
 - A. AWS CloudTrail
 - **B.** Amazon CloudWatch
 - C. Amazon S3
 - D. Amazon RDS
- 2. Why is unstructured data considered more flexible?
 - A. It follows a strict schema.
 - B. It can be stored in any format.
 - C. It is easier to query than structured data.
 - D. It has no fixed format and adapts to various use cases.
- 3. What format is commonly used for data exchange in AWS services?
 - A. XML (Extensible Markup Language)
 - **B. CSV (Comma-Separated Values)**
 - C. JSON (JavaScript Object Notation)
 - D. HTML (HyperText Markup Language)
- 4. What does feature extraction and selection reduce in machine learning?
 - A. Noise
 - **B.** Bias
 - C. Dimensionality
 - D. Variance
- 5. What does Amazon VPC enhance in terms of AWS resources?
 - A. Network performance
 - **B.** Network isolation
 - C. Data transfer speed
 - D. Resource cost efficiency

- 6. Which AWS service is designed to ingest data from file systems?
 - A. AWS DataSync
 - **B.** Amazon Glue
 - C. Amazon S3
 - D. Amazon Lambda
- 7. Which stages are part of every modern data pipeline?
 - A. Storage and processing
 - **B.** Analysis and visualization
 - C. Rechecking and auditing
 - D. Storage and analysis
- 8. Which option describes a best practice when cleaning data?
 - A. Use random sampling for evaluation.
 - B. Come to an agreement on what clean looks like.
 - C. Discard any data older than one year.
 - D. Only keep processed data.
- 9. How does AWS provide data ingestion from edge devices?
 - A. AWS Lambda
 - **B. AWS IoT Core**
 - C. Amazon EC2
 - **D. AWS Snowball**
- 10. Which AWS service is known for data warehousing?
 - A. Amazon RDS
 - B. Amazon DynamoDB
 - C. Amazon Redshift
 - D. Amazon S3

Answers



- 1. A 2. D 3. C 4. C 5. B 6. A 7. A 8. B 9. B 10. C



Explanations



1. Which AWS services can be used to monitor and troubleshoot an AWS Glue job?

- A. AWS CloudTrail
- **B.** Amazon CloudWatch
- C. Amazon S3
- **D. Amazon RDS**

AWS Glue jobs can be effectively monitored and troubleshot using Amazon CloudWatch. This service provides a robust set of features that allows users to track metrics, log events, and set alerts for their AWS Glue jobs. CloudWatch collects and visualizes logs and metrics related to Glue jobs, enabling users to gain insights into job performance, execution time, and potential issues. By using CloudWatch, users can also set custom alarms and notifications based on specific metrics related to AWS Glue, which helps in maintaining the reliability and efficiency of data processing workflows. This capability is essential for troubleshooting, as it provides visibility into the job's execution state and any errors that may occur. While AWS CloudTrail tracks API calls and user activity in your AWS account, it does not provide job-specific metrics or logs for troubleshooting Glue jobs directly. Amazon S3 is used for storage and does not inherently offer monitoring capabilities for Glue jobs. Amazon RDS is a managed database service and is not relevant to monitoring job execution in AWS Glue.

2. Why is unstructured data considered more flexible?

- A. It follows a strict schema.
- B. It can be stored in any format.
- C. It is easier to query than structured data.
- D. It has no fixed format and adapts to various use cases.

Unstructured data is considered more flexible primarily because it has no fixed format and can be easily adapted to a variety of use cases. This characteristic allows organizations to collect and utilize a wide range of data types, such as text, images, audio, and video, without the constraints that often come with structured data. Structured data typically conforms to a predefined schema, which can limit how the data can be used and analyzed. In contrast, unstructured data's lack of rigid structure means that it can be stored in any format, making it easier to accommodate new types of data as they arise and to integrate disparate data sources. Another aspect of unstructured data's flexibility is its potential for insights. Since it can contain valuable information that isn't confined to specific fields or formats, organizations can derive insights through processing techniques such as natural language processing, image analysis, and more. This adaptability can be instrumental in today's rapidly evolving data landscape.

- 3. What format is commonly used for data exchange in AWS services?
 - A. XML (Extensible Markup Language)
 - **B. CSV (Comma-Separated Values)**
 - C. JSON (JavaScript Object Notation)
 - D. HTML (HyperText Markup Language)

JSON (JavaScript Object Notation) is the most commonly used format for data exchange in AWS services due to several key factors. Its light-weight nature makes it easy to read and write by both humans and machines, which is essential for the efficient transfer of data across different AWS services and APIs. JSON's simplicity and flexibility allow it to represent complex data structures, including arrays and nested objects, making it suitable for a wide range of applications. In the context of AWS, many services, such as AWS Lambda, Amazon API Gateway, and Amazon SNS, utilize JSON for their input and output data formats. This standardization facilitates seamless integration and interoperability among service components, which is particularly important in cloud environments where services often communicate and share data. While other formats such as XML and CSV are also used within specific use cases or services, JSON's popularity stems from its ease of use and native compatibility with modern web technologies and programming languages, making it the preferred choice for data exchange in AWS. HTML, being primarily a markup language for web pages, is not designed for data exchange and is therefore not applicable in this context.

- 4. What does feature extraction and selection reduce in machine learning?
 - A. Noise
 - **B.** Bias
 - C. Dimensionality
 - D. Variance

Feature extraction and selection primarily focus on reducing dimensionality in machine learning. Dimensionality refers to the number of features or attributes in a dataset, and high-dimensional data can lead to challenges such as increased computational costs, the risk of overfitting, and difficulties in visualization. By extracting and selecting only the most relevant features, models can become more efficient, as they can operate on a clearer subset of data that is more informative. This process helps to simplify the model, making it easier to interpret and often leads to improved performance by ensuring that the model learns from the most pertinent information. Reducing dimensionality also helps in making the models faster and more effective at generalizing to unseen data, as less irrelevant or redundant information is included in the training process. While feature extraction and selection can potentially impact aspects such as noise and variance, their primary and most significant contribution is the reduction of dimensionality, which enhances the model's efficiency and effectiveness.

5. What does Amazon VPC enhance in terms of AWS resources?

- A. Network performance
- **B.** Network isolation
- C. Data transfer speed
- D. Resource cost efficiency

Amazon VPC (Virtual Private Cloud) primarily enhances network isolation for AWS resources. It allows users to create a private, isolated network within the AWS cloud, enabling them to control their own networking environment. This includes defining IP address ranges, creating subnets, and setting up route tables and network gateways, which ensures that resources within the VPC can operate securely and independently from other AWS customers' resources. This isolation feature is particularly important for organizations that must adhere to regulatory compliance standards or have specific security requirements. By controlling access to the VPC and using security groups and network access control lists, users can tightly regulate who or what can access their resources, thus ensuring that sensitive data remains protected. The other options relate to aspects of network performance or cost but do not specifically capture the essence of what Amazon VPC emphasizes. While network performance and resource cost efficiency can be influenced by how VPC is configured and utilized, they are not the primary focus of VPC. Hence, the strength of Amazon VPC lies in its ability to provide a secure and isolated network environment for AWS resources.

6. Which AWS service is designed to ingest data from file systems?

- A. AWS DataSync
- **B.** Amazon Glue
- C. Amazon S3
- D. Amazon Lambda

AWS DataSync is designed specifically for transferring large amounts of data between on-premises storage systems and AWS storage services, such as Amazon S3 or Amazon EFS. It automates the process of moving data, allowing users to efficiently and securely ingest data from file systems into the cloud. DataSync simplifies the task of data ingestion by handling operations like bandwidth management and data validation, making it a suitable choice for users looking to synchronize files or migrate large datasets to AWS. In contrast, while Amazon Glue is a data integration service primarily focused on preparing and transforming data for analytics, it does not directly handle the ingestion from file systems. Amazon S3 is a storage solution where data resides but does not inherently perform the ingestion process; rather, it serves as a destination for ingested data. Amazon Lambda is a compute service that allows the execution of code in response to events, but it is not specifically designed for data ingestion from file systems.

7. Which stages are part of every modern data pipeline?

- A. Storage and processing
- **B.** Analysis and visualization
- C. Rechecking and auditing
- D. Storage and analysis

In a modern data pipeline, the stages of storage and processing are fundamental components that enable the effective handling and transformation of data. Storage refers to the methods and technologies used to retain data in a way that it can be easily accessed and used later. This is essential for ensuring that data is available for processing, analysis, and eventual consumption. Various storage options are utilized in modern pipelines, including cloud storage solutions like Amazon S3 and databases, which ensure that data is both scalable and reliable. Processing involves the transformation and manipulation of the stored data. This step can include data cleaning, aggregation, and enrichment, allowing raw data to be converted into a format that is useful for analysis and other applications. Tools and frameworks, such as Apache Spark and AWS Glue, are often employed to efficiently handle and process large volumes of data. Together, these stages form the backbone of any modern data pipeline, enabling the pathway from raw data collection to meaningful insights, which are critical for informed decision-making in organizations. Other stages, while important, are not universally required in every pipeline; therefore, storage and processing stand as essential to all modern data workflows.

8. Which option describes a best practice when cleaning data?

- A. Use random sampling for evaluation.
- B. Come to an agreement on what clean looks like.
- C. Discard any data older than one year.
- D. Only keep processed data.

Coming to an agreement on what "clean" data looks like is essential in data cleaning processes, as it establishes a clear baseline and understanding among all stakeholders involved in the data management workflow. This consensus helps ensure that everyone involved in the project—data engineers, analysts, and stakeholders—has the same expectations regarding the criteria and standards for data cleanliness. This agreement might involve defining acceptable ranges for numerical values, determining how to handle missing data, or setting rules for categorizing and formatting data. By agreeing on these definitions upfront, teams can work more efficiently and avoid confusion or misinterpretation later in the project. A well-defined understanding of cleanliness can also guide the development of automated data validation and cleaning processes, ultimately improving data quality and reliability. While other options touch upon aspects of data management, they do not tackle the foundational issue of establishing a clear standard for what constitutes clean data, which is vital for effective and collaborative data cleaning efforts.

9. How does AWS provide data ingestion from edge devices?

- A. AWS Lambda
- **B. AWS IoT Core**
- C. Amazon EC2
- **D. AWS Snowball**

AWS provides data ingestion from edge devices primarily through AWS IoT Core. This service is designed specifically to facilitate the connection and management of Internet of Things (IoT) devices, allowing them to communicate securely with cloud applications and other devices. AWS IoT Core supports multiple protocols (like MQTT, HTTP, and WebSockets), which are commonly used by edge devices to send data to the cloud. Additionally, it offers features such as device management, data processing, and integration with other AWS services, making it an ideal solution for handling data generated by edge devices. The ability to process data in real-time, route messages, and trigger actions based on device status adds to its effectiveness for data ingestion purposes. Other options, while they have their respective use cases, do not specifically focus on direct data ingestion from edge devices in the same comprehensive manner that AWS IoT Core does. AWS Lambda, for instance, is more about running code in response to events rather than managing device connectivity and data ingestion itself. Amazon EC2 provides compute capacity but does not inherently focus on the unique requirements of edge devices. AWS Snowball is primarily aimed at large-scale data transfer rather than real-time data ingestion from devices.

10. Which AWS service is known for data warehousing?

- A. Amazon RDS
- **B.** Amazon DynamoDB
- C. Amazon Redshift
- D. Amazon S3

Amazon Redshift is a fully managed data warehousing service in AWS that is designed specifically for online analytical processing (OLAP) and large-scale data storage and analysis. It allows organizations to run complex queries across vast amounts of structured and semi-structured data, making it ideal for business intelligence (BI) workloads and analytics applications. One of the key features of Redshift is its ability to handle petabyte-scale data warehousing, providing fast query performance through techniques like columnar storage, data compression, and advanced query optimization. Additionally, it integrates seamlessly with various data visualization tools and other AWS services, enabling users to easily analyze their data and derive insights. In contrast, Amazon RDS is a relational database service optimized for OLTP (transaction processing) rather than data warehousing. Amazon DynamoDB is a NoSQL database designed for high-availability and low-latency data access, which is not specifically intended for data warehousing functionalities. Amazon S3 is a scalable object storage service that can be used to store data for analytics but doesn't provide the specialized features of a data warehousing solution like Amazon Redshift does.