# Databricks Data Engineering Professional Practice Exam (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. What is the significance of parallel execution in Spark?
  - A. It allows for single-threaded processing
  - B. It enables faster processing by distributing tasks across nodes
  - C. It minimizes the need for data partitioning
  - D. It automates resource allocation
- 2. What kind of performance benefit does Delta Lake offer compared to traditional formats when handling large datasets?
  - A. It allows versioning for historical queries
  - B. It provides better schema enforcement
  - C. It allows for ACID transactions on data operations
  - D. It enables faster read operations through efficient storage format
- 3. What challenges may arise from cross-region reads and writes in Databricks?
  - A. Enhanced security measures may be required
  - B. Cost and latency may significantly increase
  - C. Data format compatibility issues may occur
  - D. Only certain regions can support extensive data processing
- 4. Which command should be removed from a notebook before scheduling it as a job?
  - A. Cmd 2
  - B. Cmd 3
  - C. Cmd 4
  - D. Cmd 5
- 5. What is the purpose of "auto optimization" in Databricks?
  - A. Automatically generates machine learning models
  - B. Manages and optimizes performance-related configurations
  - C. Enhances data encryption processes
  - D. Facilitates data visualization automatically

- 6. In a Lakehouse architecture, what is the primary function of gold tables?
  - A. To act as the raw data layer
  - B. To support machine learning operations
  - C. To serve business intelligence and reporting needs
  - D. To store intermediate processing results
- 7. Why is Delta Lake an important component in Databricks?
  - A. It serves as a relational database
  - B. It provides ACID transactions and scalable metadata handling
  - C. It is primarily used for data visualization
  - D. It limits data storage to on-premise servers
- 8. How does querying a registered view of recent orders process data from Delta Lake tables?
  - A. All logic is executed and stored when the view is defined, returning cached results.
  - B. Results are computed and cached, incrementally updating with new records.
  - C. All logic executes at query time, returning current valid versions of source tables.
  - D. All logic executes at query time based on versions of source tables at the moment the query started.
- 9. How long is the job run history retained in Databricks?
  - A. Until you export or delete job run logs
  - B. For 30 days, during which you can deliver job logs to DBFS or S3
  - C. For 60 days, during which you can export notebook run results to HTML
  - D. For 90 days or until the run-id is reused

- 10. What does the run\_id field represent when a job run request is submitted in Databricks?
  - A. The job\_id and number of times the job has been run are concatenated and returned.
  - B. The total number of jobs that have been run in the workspace.
  - C. The number of times the job definition has been run in this workspace.
  - D. The globally unique ID of the newly triggered run.



### **Answers**



- 1. B 2. C
- 3. B

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



### **Explanations**



- 1. What is the significance of parallel execution in Spark?
  - A. It allows for single-threaded processing
  - B. It enables faster processing by distributing tasks across nodes
  - C. It minimizes the need for data partitioning
  - D. It automates resource allocation

The significance of parallel execution in Spark is captured in the concept that it enables faster processing by distributing tasks across nodes. This parallelism is a core feature of Spark's design, allowing it to leverage a distributed computing environment effectively. In Spark, data is divided into partitions that can be processed independently across various nodes in a cluster. By executing tasks in parallel, Spark can significantly reduce the time required to perform operations on large data sets, as multiple computations are carried out simultaneously rather than sequentially. This architecture is particularly beneficial for workloads that involve large-scale data processing, enabling scalability and efficiency. It allows Spark to take full advantage of the available computing resources, leading to optimized performance, especially for iterative algorithms and large-scale transformations. Consequently, applications that process huge volumes of data can see notable improvements in both speed and efficiency due to this parallel execution strategy.

- 2. What kind of performance benefit does Delta Lake offer compared to traditional formats when handling large datasets?
  - A. It allows versioning for historical queries
  - B. It provides better schema enforcement
  - C. It allows for ACID transactions on data operations
  - D. It enables faster read operations through efficient storage format

Delta Lake offers significant performance benefits for handling large datasets, particularly through its support of ACID transactions. ACID (Atomicity, Consistency, Isolation, Durability) transactions ensure that all operations on the data are completed successfully and reliably. This means that when large numbers of transactions or data updates occur, Delta Lake can maintain the integrity of the dataset, allowing safe concurrent access without conflicts. When managing large datasets, the ability to perform transactional operations leads to improved data consistency and accuracy. This is especially crucial in environments where multiple processes may be reading from or writing to the dataset at the same time. Using ACID transactions helps avoid issues such as lost updates and dirty reads, which could occur in traditional data formats that do not support such robust transactional guarantees. While versioning for historical queries, improved schema enforcement, and enhanced read operations are also valuable features of Delta Lake, the specific aspect of ACID transactions is particularly impactful for ensuring reliable and consistent performance in data operations at scale, making it a standout benefit.

- 3. What challenges may arise from cross-region reads and writes in Databricks?
  - A. Enhanced security measures may be required
  - B. Cost and latency may significantly increase
  - C. Data format compatibility issues may occur
  - D. Only certain regions can support extensive data processing

The challenges of cross-region reads and writes in Databricks primarily relate to cost and latency. When data is accessed from one region to another, the network latency can introduce delays in data retrieval and processing times, which is particularly crucial for real-time or time-sensitive applications. Furthermore, from a cost perspective, transferring data across geographic regions can lead to increased costs due to data egress fees and the potential need for higher throughput that can further escalate expenses. This context emphasizes why managing data across regions must be carefully considered, especially for applications that demand performance efficiency and cost-effectiveness. The other options, while they may be relevant in different contexts, do not directly address the core challenges associated with cross-region data operations in Databricks as clearly as the increase in cost and latency.

- 4. Which command should be removed from a notebook before scheduling it as a job?
  - A. Cmd 2
  - B. Cmd 3
  - C. Cmd 4
  - **D.** Cmd 5

When scheduling a notebook as a job in Databricks, it is essential to ensure that the job runs without interruptions or manual input. Typically, commands that are interactive or rely on user prompts should not be included in a scheduled job. These commands can cause the execution to halt, awaiting user interaction, which is not possible in a fully automated job environment. In this context, the correct answer indicates that the specific command (Cmd 5) is likely an interactive or non-executable command in the context of a job. It could be a display command or a command that requires user confirmation, causing it to fail when executed in a non-interactive job environment. By removing it from the notebook before scheduling, you can ensure that the job runs smoothly and achieves its intended results without manual intervention. In contrast, other commands (like Cmd 2, Cmd 3, and Cmd 4) are typically straightforward, non-interactive commands that can be executed automatically as part of the scheduled job process. These might include data manipulations, queries, or transformations that are suitable for automated execution without requiring user inputs or interactions. Thus, they can remain intact in the notebook when converting it into a scheduled job.

#### 5. What is the purpose of "auto optimization" in Databricks?

- A. Automatically generates machine learning models
- B. Manages and optimizes performance-related configurations
- C. Enhances data encryption processes
- D. Facilitates data visualization automatically

The purpose of "auto optimization" in Databricks primarily revolves around enhancing the performance of data processing workflows by managing various configurations automatically. This feature is designed to help users optimize their data processing jobs without requiring manual intervention, thereby streamlining operations and improving efficiency. Auto optimization can include adjustments to parameters such as caching strategies, data partitioning, and resource allocation, which are critical for enhancing the performance of Spark queries and reducing execution times. By leveraging these automatic adjustments, users can experience faster job execution and better resource utilization, ultimately leading to cost savings and improved performance in large-scale data engineering tasks. The other choices relate to functionalities that are not the focus of auto optimization in Databricks. For instance, generating machine learning models, enhancing encryption processes, or facilitating data visualization does not fall under the umbrella of performance optimizations but rather addresses different aspects of data processing and analysis.

## 6. In a Lakehouse architecture, what is the primary function of gold tables?

- A. To act as the raw data layer
- B. To support machine learning operations
- C. To serve business intelligence and reporting needs
- D. To store intermediate processing results

In a Lakehouse architecture, gold tables play a crucial role in serving business intelligence and reporting needs. These tables are typically the most refined level of data within the architecture and are designed to provide clean, high-quality datasets that can be easily accessed for analysis and reporting purposes. By aggregating and transforming raw and intermediate data from silver and bronze tables, gold tables enable organizations to derive valuable insights and make data-driven decisions. The focus on usability and accessibility in gold tables makes them essential for business analysts and decision-makers who require reliable and well-organized data to generate reports and visualizations. The optimization of gold tables for speed and efficiency enhances the performance of business intelligence tools that leverage this data, further supporting analysis and reporting functions. In contrast, the other options refer to different stages or functionalities within the data lifecycle. For example, raw data typically resides in the bronze layer, while intermediate processing results are handled in the silver layer. Machine learning operations, although they may utilize data from gold tables, are not the primary focus of these tables, highlighting the distinct purpose of each type of table within the Lakehouse framework.

#### 7. Why is Delta Lake an important component in Databricks?

- A. It serves as a relational database
- B. It provides ACID transactions and scalable metadata handling
- C. It is primarily used for data visualization
- D. It limits data storage to on-premise servers

Delta Lake is a crucial component of Databricks because it provides ACID transactions and scalable metadata handling, which are essential for ensuring data integrity, reliability, and performance in big data environments. ACID stands for Atomicity, Consistency, Isolation, and Durability, which are principles that guarantee reliable processing of database transactions. This is particularly important in environments where concurrent writes and reads occur, as it prevents data corruption and ensures that users see consistent views of data. The ability to manage metadata effectively is also significant; Delta Lake allows for efficient handling of large datasets and provides features such as schema enforcement and evolution. This means that as data evolves, Delta Lake can adapt without requiring extensive manual intervention. The combination of ACID transactions and robust metadata management makes Delta Lake a powerful solution for handling enterprise data workflows that demand high levels of reliability and performance. The other options do not capture the key functionalities of Delta Lake. While a relational database serves a different purpose, data visualization is typically addressed through separate tools and technologies within the Databricks ecosystem, and limiting data storage to on-premise servers contradicts the cloud-based and flexible nature of Databricks' offerings. Thus, the specific capabilities of Delta Lake in providing transactional integrity and efficient data management highlight its

- 8. How does querying a registered view of recent orders process data from Delta Lake tables?
  - A. All logic is executed and stored when the view is defined, returning cached results.
  - B. Results are computed and cached, incrementally updating with new records.
  - C. All logic executes at query time, returning current valid versions of source tables.
  - D. All logic executes at query time based on versions of source tables at the moment the query started.

Querying a registered view of recent orders processes data from Delta Lake tables by executing all logic at query time based on the versions of the source tables at the moment the query started. This means that when a query is made against the view, the underlying data from Delta Lake is accessed in real-time to ensure that the most up-to-date and valid versions of the table are returned. Delta Lake's transaction log enables this capability, allowing the query to reflect any changes made to the underlying data since the view was created or last queried. This approach provides consistency and correctness, as it dynamically pulls the latest data rather than relying on potentially stale or cached information. The real-time execution of logic ensures that users have access to the most accurate data, reflecting any additions, deletions, or modifications to the underlying Delta Lake tables that have occurred up to the time the query is executed. This is essential, especially in environments where data changes frequently, such as in the case of recent orders, which may be populated continually. In contrast, the other options suggest different mechanisms for handling the data, such as pre-computation or caching results, which would not adequately address the need for real-time accuracy inherent in querying views built on top of Delta Lake tables.

- 9. How long is the job run history retained in Databricks?
  - A. Until you export or delete job run logs
  - B. For 30 days, during which you can deliver job logs to DBFS or S3
  - C. For 60 days, during which you can export notebook run results to HTML
  - D. For 90 days or until the run-id is reused

The retention period for job run history in Databricks is designed to provide users with timely access to historical job details while managing storage effectively. Specifically, job run history is retained for 60 days, allowing users to reference and analyze past jobs easily. During this period, exporting the notebook run results to HTML is a functional benefit, giving teams the option to archive significant job outputs for reporting or auditing purposes. The ability to retain logs for 60 days ensures that users can troubleshoot and monitor job performance without losing important historical data. This retention timeframe aligns well with common data processing needs, where recent job runs are typically of the most interest. In contrast, other options may suggest either different retention periods or functionalities that don't precisely align with how Databricks manages job run history. This highlights the importance of understanding both the retention policy and the features available during that retention window.

- 10. What does the run\_id field represent when a job run request is submitted in Databricks?
  - A. The job\_id and number of times the job has been run are concatenated and returned.
  - B. The total number of jobs that have been run in the workspace.
  - C. The number of times the job definition has been run in this workspace.
  - D. The globally unique ID of the newly triggered run.

The run\_id field represents the globally unique ID of the newly triggered run. When a job is submitted in Databricks, the run\_id is generated to uniquely identify that specific run instance of the job. This unique identifier is essential for monitoring and managing job runs, as it allows users to track the status, log details, and results of that particular execution without confusion with other runs. Each job execution may have the same job\_id if the job is triggered multiple times, but the run\_id will always be different, thus providing a clear reference for each individual run. This functionality is critical for maintaining operational clarity and organizational efficiency in managing jobs within a Databricks workspace. Understanding the run\_id allows developers and data engineers to debug and analyze job performance over time, ensuring effective monitoring and troubleshooting of their data pipelines.