

Databricks Data Engineering Professional Practice Exam (Sample)

Study Guide



Everything you need from our exam experts!

Copyright © 2026 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 accurate, complete, and timely information about this product from reliable sources.

SAMPLE

Table of Contents

Copyright	1
Table of Contents	2
Introduction	3
How to Use This Guide	4
Questions	5
Answers	8
Explanations	10
Next Steps	16

SAMPLE

Introduction

Preparing for a certification exam can feel overwhelming, but with the right tools, it becomes an opportunity to build confidence, sharpen your skills, and move one step closer to your goals. At Examzify, we believe that effective exam preparation isn't just about memorization, it's about understanding the material, identifying knowledge gaps, and building the test-taking strategies that lead to success.

This guide was designed to help you do exactly that.

Whether you're preparing for a licensing exam, professional certification, or entry-level qualification, this book offers structured practice to reinforce key concepts. You'll find a wide range of multiple-choice questions, each followed by clear explanations to help you understand not just the right answer, but why it's correct.

The content in this guide is based on real-world exam objectives and aligned with the types of questions and topics commonly found on official tests. It's ideal for learners who want to:

- Practice answering questions under realistic conditions,
- Improve accuracy and speed,
- Review explanations to strengthen weak areas, and
- Approach the exam with greater confidence.

We recommend using this book not as a stand-alone study tool, but alongside other resources like flashcards, textbooks, or hands-on training. For best results, we recommend working through each question, reflecting on the explanation provided, and revisiting the topics that challenge you most.

Remember: successful test preparation isn't about getting every question right the first time, it's about learning from your mistakes and improving over time. Stay focused, trust the process, and know that every page you turn brings you closer to success.

Let's begin.

How to Use This Guide

This guide is designed to help you study more effectively and approach your exam with confidence. Whether you're reviewing for the first time or doing a final refresh, here's how to get the most out of your Examzify study guide:

1. Start with a Diagnostic Review

Skim through the questions to get a sense of what you know and what you need to focus on. Your goal is to identify knowledge gaps early.

2. Study in Short, Focused Sessions

Break your study time into manageable blocks (e.g. 30 - 45 minutes). Review a handful of questions, reflect on the explanations.

3. Learn from the Explanations

After answering a question, always read the explanation, even if you got it right. It reinforces key points, corrects misunderstandings, and teaches subtle distinctions between similar answers.

4. Track Your Progress

Use bookmarks or notes (if reading digitally) to mark difficult questions. Revisit these regularly and track improvements over time.

5. Simulate the Real Exam

Once you're comfortable, try taking a full set of questions without pausing. Set a timer and simulate test-day conditions to build confidence and time management skills.

6. Repeat and Review

Don't just study once, repetition builds retention. Re-attempt questions after a few days and revisit explanations to reinforce learning. Pair this guide with other Examzify tools like flashcards, and digital practice tests to strengthen your preparation across formats.

There's no single right way to study, but consistent, thoughtful effort always wins. Use this guide flexibly, adapt the tips above to fit your pace and learning style. You've got this!

Questions

SAMPLE

- 1. What is the role of Databricks Runtime?**
 - A. To manage data security**
 - B. To serve as an interface for users**
 - C. To provide an optimized version of Apache Spark with performance improvements**
 - D. To perform data integration tasks**
- 2. What approach allows a developer to share code updates without overwriting the work of teammates in Databricks Repos?**
 - A. Use Repos to checkout all changes and send the git diff log to the team.**
 - B. Use Repos to create a fork of the remote repository, commit all changes, and make a pull request.**
 - C. Use Repos to pull changes from the remote Git repository, then commit and push changes to a branch.**
 - D. Use Repos to create a new branch, commit all changes, and push changes to the remote Git repository.**
- 3. Which data formats are supported by Delta Lake?**
 - A. Only text and XML formats**
 - B. Parquet, ORC, Avro, CSV, and JSON**
 - C. Proprietary formats developed by Databricks**
 - D. Only structured formats like SQL tables**
- 4. What does "data governance" involve?**
 - A. Management of application deployment processes**
 - B. Control of data quality, accessibility, and usage policies**
 - C. Creation of data visualization tools**
 - D. Deployment of data applications for analytics**
- 5. What does "unified analytics" refer to in the context of Databricks?**
 - A. Integration of data storage and processing solutions**
 - B. Combination of data engineering, data science, and analytics on a single platform**
 - C. Seamless data migration to third-party platforms**
 - D. Real-time analytics without any physical data storage**

6. What is a key characteristic of data engineering?

- A. Focus on user interface design**
- B. Design of systems for processing large datasets**
- C. Utilization of only unstructured data**
- D. Emphasis on manual data entry tasks**

7. What is the correct approach for a developer to review the current logic in a notebook after using an outdated branch?

- A. Use Repos to make a pull request use the Databricks REST API to update the current branch to dev-2.3.9**
- B. Use Repos to pull changes from the remote Git repository and select the dev-2.3.9 branch**
- C. Use Repos to checkout the dev-2.3.9 branch and auto-resolve conflicts with the current branch**
- D. Merge all changes back to the main branch in the remote Git repository and clone the repo again**

8. Does the logic used to delete records from a Delta Lake table guarantee that they are no longer accessible?

- A. Yes; Delta Lake ACID guarantees provide assurance of complete deletion.**
- B. No; the Delta cache may return records until the cluster is restarted.**
- C. Yes; Delta cache updates to reflect latest data files.**
- D. No; deleted records may still be accessible with time travel until a VACUUM command is used.**

9. Which method can be utilized for incremental data loads?

- A. Full data refreshes**
- B. Delta Lake with change tracking**
- C. In-memory data processing**
- D. Static data ingests only**

10. Which adjustment will reduce cloud storage costs for a Structured Streaming job processing less than 10 minutes?

- A. Set the trigger interval to 3 seconds**
- B. Increase the number of shuffle partitions**
- C. Set the trigger interval to 10 minutes**
- D. Reduce the number of active clusters**

Answers

SAMPLE

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

SAMPLE

Explanations

SAMPLE

1. What is the role of Databricks Runtime?

- A. To manage data security
- B. To serve as an interface for users
- C. To provide an optimized version of Apache Spark with performance improvements**
- D. To perform data integration tasks

The role of Databricks Runtime is primarily to provide an optimized version of Apache Spark that includes various performance improvements and additional capabilities tailored for big data analytics and machine learning workloads. This optimized runtime is designed to enhance the performance of Spark jobs by incorporating optimizations related to execution, memory management, and resource allocation. Databricks Runtime enables users to take advantage of specific features and efficiencies that are not present in the vanilla version of Apache Spark. These improvements result in faster execution times and better resource utilization, making it easier for data engineers and data scientists to process and analyze large datasets in a more efficient manner. In contrast, while data security and user interface management are important aspects of working with data platforms, they do not define the core purpose of Databricks Runtime. Its main focus is on enhancing the performance of Spark, which is foundational for executing data processing jobs effectively. Data integration tasks, while also essential in a data engineering context, are typically carried out using other tools and functionalities provided within the Databricks environment, rather than the runtime itself.

2. What approach allows a developer to share code updates without overwriting the work of teammates in Databricks Repos?

- A. Use Repos to checkout all changes and send the git diff log to the team.
- B. Use Repos to create a fork of the remote repository, commit all changes, and make a pull request.
- C. Use Repos to pull changes from the remote Git repository, then commit and push changes to a branch.
- D. Use Repos to create a new branch, commit all changes, and push changes to the remote Git repository.**

The approach that allows a developer to share code updates without overwriting the work of teammates in Databricks Repos involves creating a new branch, committing all changes, and then pushing these changes to the remote Git repository. By using branches, developers can work on their own features or fixes independently, without affecting the main codebase or the ongoing work of their teammates. When a developer creates a new branch, they can develop their feature or fix in isolation. This means their changes won't interfere with the main branch or other branches that teammates might be working on. Once the work is complete and tested, the developer can push the new branch to the remote repository. This facilitates code reviews and collaborative work before changes are merged back into the main branch, ensuring that the integrity of the codebase is maintained. This branching strategy encourages better collaboration and reduces the risk of conflicts between different team members' work. It is a best practice in version control systems to use branches for new features or bug fixes so that integration into the main codebase can be managed smoothly and systematically.

3. Which data formats are supported by Delta Lake?

- A. Only text and XML formats
- B. Parquet, ORC, Avro, CSV, and JSON**
- C. Proprietary formats developed by Databricks
- D. Only structured formats like SQL tables

Delta Lake is designed to work with various data formats, making it highly versatile for data engineering tasks. It primarily leverages the Parquet format, which is a columnar storage format optimized for use with massive data processing frameworks. Additionally, Delta Lake supports other formats, including ORC, Avro, CSV, and JSON. This compatibility allows users to easily read from and write to Delta tables using data stored in these formats, facilitating data integration and transformation across diverse datasets. This broad support for multiple formats is crucial because it allows organizations to utilize Delta Lake in a wide range of scenarios, whether they are ingesting streaming data, processing batch data, or integrating data from different sources. Consequently, it enhances the overall flexibility and efficiency of data storage and analytics processes. Other choices are limited in their scope, either by focusing on a narrow range of formats or specifying unsupported types.

4. What does "data governance" involve?

- A. Management of application deployment processes
- B. Control of data quality, accessibility, and usage policies**
- C. Creation of data visualization tools
- D. Deployment of data applications for analytics

Data governance primarily focuses on the policies and processes that ensure high data quality, proper accessibility, and the management of data usage throughout an organization. This involves establishing standards and guidelines that dictate how data is to be handled, who can access it, and under what conditions. Effective data governance helps organizations maintain regulatory compliance, secure sensitive data, and enhance the overall trustworthiness of their data, which is critical for informed decision-making and strategic planning. The other options pertain to different aspects of data management or application lifecycle and do not encapsulate the core concepts of data governance. For instance, managing application deployment processes, creating visualization tools, or deploying analytics applications are all valuable tasks but do not inherently involve the governance of data itself. Therefore, the focus on quality, accessibility, and usage policies in choice B accurately captures the essence of data governance.

5. What does "unified analytics" refer to in the context of Databricks?

- A. Integration of data storage and processing solutions**
- B. Combination of data engineering, data science, and analytics on a single platform**
- C. Seamless data migration to third-party platforms**
- D. Real-time analytics without any physical data storage**

In the context of Databricks, "unified analytics" refers to the combination of data engineering, data science, and analytics on a single platform. This concept underscores the primary objective of Databricks, which is to provide an integrated environment where data engineers and data scientists can collaborate effectively. It streamlines workflows by allowing teams to work on data ingestion, transformation, machine learning, and analytics together, facilitating a more cohesive approach to managing data. The unified nature of the platform is particularly beneficial as it eliminates silos between teams that traditionally worked in isolation. By bringing these domains together, organizations can enhance productivity, reduce the time to insights, and leverage the power of collaborative analytics. This approach allows users to work with large datasets efficiently while both developing and deploying machine learning models, thereby driving actionable insights faster and more reliably. In contrast, the other options do not accurately capture the essence of what unified analytics means in this context. For example, while integration of data storage and processing solutions is essential, it is not the sole focus of unified analytics. Similarly, seamless data migration or real-time analytics without data storage, while important in specific scenarios, do not encompass the broad, collaborative approach promoted by unified analytics in the Databricks environment.

6. What is a key characteristic of data engineering?

- A. Focus on user interface design**
- B. Design of systems for processing large datasets**
- C. Utilization of only unstructured data**
- D. Emphasis on manual data entry tasks**

The key characteristic of data engineering is the design of systems for processing large datasets. Data engineering involves creating and maintaining architectures, pipelines, and data systems that support the collection, storage, transformation, and analysis of large volumes of data efficiently. This process includes the use of various tools and technologies to ensure data is processed in a timely manner, allowing organizations to derive insights and make data-driven decisions. In contrast to the other choices, which do not accurately represent the core functions of data engineering, the focus on user interface design pertains more to front-end development and user experience, which is not the primary concern of data engineers. The assertion that data engineering is limited to only unstructured data is misleading, as data engineers work with structured, semi-structured, and unstructured data, depending on the requirements of the data architecture. Finally, emphasizing manual data entry tasks does not align with the goals of data engineering, which strives for automation and efficiency in data processing, minimizing manual intervention as much as possible.

7. What is the correct approach for a developer to review the current logic in a notebook after using an outdated branch?

- A. Use Repos to make a pull request use the Databricks REST API to update the current branch to dev-2.3.9
- B. Use Repos to pull changes from the remote Git repository and select the dev-2.3.9 branch**
- C. Use Repos to checkout the dev-2.3.9 branch and auto-resolve conflicts with the current branch
- D. Merge all changes back to the main branch in the remote Git repository and clone the repo again

To effectively review the current logic in a notebook after working with an outdated branch, using the Repos feature to pull changes from the remote Git repository and select the appropriate branch is the most logical approach. This method allows the developer to directly access the latest updates from the target branch, which in this case is the dev-2.3.9 branch. By pulling changes, the developer can sync their local environment with the latest state of the branch in the remote repository. This ensures that they are working with the most current version of the code, which is crucial for understanding the latest implementations, bug fixes, or feature additions. This approach facilitates a smoother review process since the developer can examine how the logic has evolved and how it may impact their ongoing work. In this context, other options do not align with optimal practices for reviewing notebook logic after an outdated branch usage. Checking out a branch and leveraging automatic conflict resolution could lead to unintended changes that obscure the original intent of both the current and target branches. Merging changes back to the main branch or cloning the repository again would be more cumbersome and could disrupt the workflow, complicating the review process instead of simplifying it.

8. Does the logic used to delete records from a Delta Lake table guarantee that they are no longer accessible?

- A. Yes; Delta Lake ACID guarantees provide assurance of complete deletion.
- B. No; the Delta cache may return records until the cluster is restarted.
- C. Yes; Delta cache updates to reflect latest data files.
- D. No; deleted records may still be accessible with time travel until a VACUUM command is used.**

The logic behind deleting records from a Delta Lake table does not guarantee that they are no longer accessible due to the nature of Delta Lake's time travel feature. When records are deleted, they are marked as such but are not immediately purged from the storage layer. Instead, Delta Lake maintains a version history of the data, allowing users to perform time travel queries to access previous states of the table, including records that may have been marked as deleted. Time travel in Delta Lake enables users to query data as it existed at a specific point in time, which means that even after a delete operation, the records may still be retrieved until a VACUUM command is executed. The VACUUM command is crucial as it physically removes the old data files that are no longer needed, ensuring that deleted records are permanently removed from the system and cannot be accessed via time travel. In summary, while the delete logic marks records for deletion, they remain potentially accessible through time travel until a VACUUM command is performed, thus providing the rationale for why the chosen answer is correct.

9. Which method can be utilized for incremental data loads?

- A. Full data refreshes
- B. Delta Lake with change tracking**
- C. In-memory data processing
- D. Static data ingests only

The best method for conducting incremental data loads is Delta Lake with change tracking. This approach enables efficient data operations by only processing the changes made since the last load, rather than reloading the entire dataset. Change tracking allows for the identification of inserts, updates, and deletions at a granular level. Delta Lake ensures data consistency with ACID (Atomicity, Consistency, Isolation, Durability) transactions, which is crucial for data reliability when performing incremental updates. Additionally, it supports features like time travel, which helps understand the state of the data at various points in time, making it easier to manage and track incremental changes. In contrast, using full data refreshes would entail reloading the entire dataset every time data is updated, which is often inefficient, especially for large datasets. In-memory data processing might optimize data access during compute operations but does not inherently facilitate incremental loading. Static data ingests refer to loading data as a one-time batch process, which also does not support the continuous and efficient updating of data that incremental loading requires.

10. Which adjustment will reduce cloud storage costs for a Structured Streaming job processing less than 10 minutes?

- A. Set the trigger interval to 3 seconds
- B. Increase the number of shuffle partitions
- C. Set the trigger interval to 10 minutes**
- D. Reduce the number of active clusters

Setting the trigger interval to 10 minutes is a strategy that can help reduce cloud storage costs for a Structured Streaming job processing data for less than 10 minutes. When the trigger interval is set to a longer duration, the job processes the data in larger batches rather than immediately for each incoming micro-batch. This means that fewer files will be created in the cloud storage, which can lead to decreased storage costs because cloud storage pricing often depends on the number of files and the frequency of writes. When the trigger interval is shorter, such as 3 seconds, the system creates and writes data to storage much more frequently, resulting in a higher number of files being stored, which can increase costs. Increasing the number of shuffle partitions typically impacts performance rather than storage costs, as it partitions the data being shuffled among different nodes but does not directly affect how data is written to cloud storage. Reducing the number of active clusters can influence compute costs but does not directly reduce storage costs associated with how data is written or batched in the storage layer. Thus, adjusting the trigger interval for a Structured Streaming job to 10 minutes results in fewer writes to cloud storage, thereby reducing the associated costs effectively.

Next Steps

Congratulations on reaching the final section of this guide. You've taken a meaningful step toward passing your certification exam and advancing your career.

As you continue preparing, remember that consistent practice, review, and self-reflection are key to success. Make time to revisit difficult topics, simulate exam conditions, and track your progress along the way.

If you need help, have suggestions, or want to share feedback, we'd love to hear from you. Reach out to our team at hello@examzify.com.

Or visit your dedicated course page for more study tools and resources:

<https://databricksdataengrpro.examzify.com>

We wish you the very best on your exam journey. You've got this!

SAMPLE