

Databricks Machine Learning (ML) Associate Practice Test (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 are Shapley Values (SHAP) used for in machine learning?**
 - A. To reduce model complexity**
 - B. To estimate importance of features to model predictions**
 - C. To enhance model training speed**
 - D. To visualize model accuracy**
- 2. What is a common output of the `model.fit()` method in Spark ML?**
 - A Transformer which can be used for predictions.**
 - A DataFrame containing the original data.**
 - A log of training statistics.**
 - An estimator that performs transformations.**
- 3. How does MLflow tracking differ from the Model Registry?**
 - A. Tracking is for experimentation; the Model Registry is for production**
 - B. Tracking stores all data, while the Model Registry only retains best models**
 - C. Tracking focuses on user access, while the Model Registry is for analytics**
 - D. Tracking is less secure than the Model Registry**
- 4. What is a primary function of a machine learning model's label?**
 - A. To classify the input features**
 - B. To determine error rates**
 - C. To define the expected output variable for prediction**
 - D. To enhance data visualization**
- 5. What is a typical evaluation metric used for regression problems?**
 - A. Accuracy**
 - B. R2 Score**
 - C. F1 Score**
 - D. Recall**

6. Which of the following is a key feature of the Databricks Runtime for Machine Learning?

- A. Integrated support for SQL queries**
- B. Pre-installed popular ML libraries such as TensorFlow, PyTorch, and Scikit-learn**
- C. Real-time data ingestion capabilities**
- D. On-demand data cleaning tools**

7. What is the goal of regularization techniques in machine learning?

- A. To prevent overfitting**
- B. To increase the training accuracy**
- C. To simplify the model**
- D. To enhance computational efficiency**

8. What distinguishes dense vectors from sparse vectors in the context of machine learning?

- A. Sparse vectors store only non-zero entries**
- B. Dense vectors are only used for classification tasks**
- C. Sparse vectors are simpler and easier to manipulate**
- D. Dense vectors cannot handle large datasets**

9. Which scenario might require the use of anomaly detection?

- A. Identifying duplicate records**
- B. Determining customer buying habits**
- C. Monitoring for fraudulent transactions**
- D. Grouping similar products together**

10. What are hyperparameters in machine learning?

- A. User-defined values that never change.**
- B. Configurations that control the learning process of the algorithm.**
- C. Output variables that report performance.**
- D. Parameters that are automatically optimized by the model.**

Answers

SAMPLE

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

SAMPLE

Explanations

SAMPLE

1. What are Shapley Values (SHAP) used for in machine learning?

- A. To reduce model complexity
- B. To estimate importance of features to model predictions**
- C. To enhance model training speed
- D. To visualize model accuracy

Shapley Values, commonly abbreviated as SHAP, are used in machine learning to provide a quantifiable measure of the contribution of each feature to a model's predictions. They stem from cooperative game theory and help to fairly distribute the "payout" (in this case, the prediction) among the "players" (the features). By attributing the prediction of a model to its input features, SHAP values allow for a deep understanding of how features influence the outcome. Each feature is analyzed to see how it impacts the prediction when considered alone, compared to when it is included with other features. This makes SHAP a powerful tool for interpretability in machine learning models, as it gives insights into why certain predictions are made based on the input data. This important distinction makes SHAP valuable in contexts such as model evaluation, troubleshooting, and ensuring fairness and transparency in automated decision-making systems. Understanding feature importance is crucial for stakeholders to trust and validate the predictions made by models, thereby playing a significant role in areas where interpretability is paramount.

2. What is a common output of the `model.fit()` method in Spark ML?

- A. A Transformer which can be used for predictions.**
- B. A DataFrame containing the original data.
- C. A log of training statistics.
- D. An estimator that performs transformations.

The `model.fit()` method in Spark ML is designed to train a machine learning model on a given dataset. The primary output of this method is a Transformer, which represents the trained model. Once the model is trained, it can be used to make predictions on new data. In the context of Spark ML, a Transformer is an abstraction that includes both the learned parameters from the training dataset and the logic needed to apply those parameters to new data for predictions, making it a practical and reusable component in the machine learning pipeline. Thus, using the output of `model.fit()` as a Transformer allows you to seamlessly integrate this trained model into your data processing flow for predicting outcomes on unseen data. The other options do not accurately reflect the output of the `model.fit()` method in this context. For example, while a DataFrame containing the original data is part of the input process, it is not a direct output of model fitting. Similarly, training statistics may be logged or tracked during training for monitoring purposes, but they are not the primary output of the fit method. Lastly, an estimator refers to an abstract class in Spark ML that includes fitting and transforming logic but is distinct from the output of the fit method itself, which specifically produces a Transformer.

3. How does MLflow tracking differ from the Model Registry?

- A. Tracking is for experimentation; the Model Registry is for production**
- B. Tracking stores all data, while the Model Registry only retains best models**
- C. Tracking focuses on user access, while the Model Registry is for analytics**
- D. Tracking is less secure than the Model Registry**

The distinction between MLflow tracking and the Model Registry is primarily based on their intended use cases and the lifecycle stages they address within machine learning workflows. Tracking is designed to facilitate experimentation by capturing and storing information about various experiments, including metrics, parameters, artifacts, and other relevant data as models are trained and evaluated. This allows data scientists and engineers to analyze different approaches, compare results, and iterate on model development in an environment focused on experimentation. On the other hand, the Model Registry serves a different purpose: it is specifically designed for production management of machine learning models. This involves tracking models that have been selected for deployment, ensuring they are versioned, and maintaining records of the associated metadata. The Model Registry streamlines the deployment process by ensuring that only vetted, high-quality models make their way into production, allowing teams to manage model lifecycle transitions from experimentation to production systematically. By highlighting these key functional differences, it becomes clear that tracking is more about the iterative experimentation process, while the Model Registry is about managing and deploying models in a structured and secure manner for production use.

4. What is a primary function of a machine learning model's label?

- A. To classify the input features**
- B. To determine error rates**
- C. To define the expected output variable for prediction**
- D. To enhance data visualization**

The primary function of a machine learning model's label is to define the expected output variable for prediction. In supervised learning, which encompasses many machine learning tasks, labels serve as the ground truth for training algorithms. This means they represent the outcomes that correspond to the input features within a dataset. During the training phase, a model learns to map the input features to these labeled outputs, enabling it to make predictions on unseen data. The correct understanding of labels is crucial, as they guide the learning process by providing the model with examples of what it needs to predict. Without accurate and well-defined labels, the model would struggle to learn relationships between the input data and the desired output, leading to poor performance during inference. Since labels are foundational to the training process, they directly impact a model's ability to analyze patterns and make accurate predictions. Understanding this can help one appreciate the significance of labels in any machine learning application and highlight the importance of having a labeled dataset for model training.

5. What is a typical evaluation metric used for regression problems?

- A. Accuracy
- B. R2 Score**
- C. F1 Score
- D. Recall

In regression problems, the primary focus is on predicting continuous outcomes, and traditional classification metrics like accuracy, F1 score, or recall are not applicable. The R2 Score, also known as the coefficient of determination, is a common evaluation metric used in regression tasks to assess the proportion of variance in the dependent variable that can be explained by the independent variables in the model. A high R2 Score indicates that a significant portion of the variability in the output can be captured by the model, while an R2 Score of zero suggests that the model does not explain any variability, and a negative score indicates that the model is performing worse than a horizontal line representing the mean of the target variable. The other metrics mentioned, such as accuracy, F1 Score, and recall, are designed for classification tasks where the output is categorical rather than continuous. Hence, they would not provide meaningful insights into the performance of a regression model. This makes the R2 Score the most appropriate evaluation metric for assessing regression models.

6. Which of the following is a key feature of the Databricks Runtime for Machine Learning?

- A. Integrated support for SQL queries
- B. Pre-installed popular ML libraries such as TensorFlow, PyTorch, and Scikit-learn**
- C. Real-time data ingestion capabilities
- D. On-demand data cleaning tools

The choice highlighting the pre-installed popular ML libraries such as TensorFlow, PyTorch, and Scikit-learn is significant because it reflects the goal of the Databricks Runtime for Machine Learning to streamline the data science workflow. This runtime environment is specifically optimized for machine learning tasks, ensuring that practitioners have immediate access to widely used tools without the overhead of individual installations and configurations. Having these libraries pre-installed allows data scientists and machine learning engineers to quickly start building, training, and deploying models. It enhances productivity as users can focus on developing their algorithms and experiments rather than worrying about setting up their machine learning environment. The convenience of a ready-to-use environment also promotes consistency across teams, ensuring that everyone is working with the same set of tools and libraries, fostering collaboration and reducing compatibility issues. While integrated support for SQL queries, real-time data ingestion capabilities, and on-demand data cleaning tools are valuable features in various contexts, they do not specifically define the core strength of the Databricks Runtime for Machine Learning, which is aimed at providing an efficient and effective foundation for machine learning development and deployment.

7. What is the goal of regularization techniques in machine learning?

- A. To prevent overfitting**
- B. To increase the training accuracy**
- C. To simplify the model**
- D. To enhance computational efficiency**

The primary goal of regularization techniques in machine learning is to prevent overfitting. Overfitting occurs when a model learns not only the underlying patterns in the training data but also the noise, leading to poor performance on unseen data. Regularization introduces additional information or constraints into the model, effectively adding a penalty for complexity. This encourages the model to focus on the most significant features and leads to simpler, more generalized solutions that perform better during validation and testing phases. While some techniques may also simplify the model by reducing the number of parameters or features considered, the central objective is to enhance the model's ability to generalize by mitigating the adverse effects of overfitting. Regularization methods, such as L1 and L2 regularization, directly target this issue by modifying the loss function used during training, thus ensuring the model remains robust against variations in new data.

8. What distinguishes dense vectors from sparse vectors in the context of machine learning?

- A. Sparse vectors store only non-zero entries**
- B. Dense vectors are only used for classification tasks**
- C. Sparse vectors are simpler and easier to manipulate**
- D. Dense vectors cannot handle large datasets**

Dense vectors and sparse vectors are two different ways of representing data, particularly in the context of machine learning and numerical computations. The distinguishing feature of sparse vectors is that they are designed to store only non-zero entries, which makes them memory efficient when dealing with high-dimensional data that contains many zero values. In practical terms, this means that when using sparse vectors, the only values stored are the entries that have significance (non-zero values), along with their corresponding indices. This helps save memory and computational resources, especially in scenarios where the dataset features a large number of dimensions but only a few of these dimensions contain meaningful information. On the other hand, dense vectors store every entry in the vector, regardless of whether it is zero or non-zero. While dense vectors can be used for any type of machine learning task, they are less efficient in terms of memory usage for high-dimensional sparse data. The incorrect options highlight misconceptions about the characteristics and applications of these vector types. For instance, saying that dense vectors are only used for classification tasks overlooks their broad application across various machine learning methods, including regression. The notion that sparse vectors are simpler and easier to manipulate isn't accurate, as their advantages come from their specific use case in handling predominantly sparse data. Lastly, the claim

9. Which scenario might require the use of anomaly detection?

- A. Identifying duplicate records**
- B. Determining customer buying habits**
- C. Monitoring for fraudulent transactions**
- D. Grouping similar products together**

Anomaly detection is a powerful technique used to identify patterns in data that do not conform to expected behavior. In the context of monitoring for fraudulent transactions, anomaly detection is particularly effective. Fraudulent activities often exhibit unusual patterns or behaviors compared to typical transaction data, such as an abrupt increase in transaction size, frequency, or changes in purchase locations that deviate from a user's normal behavior. Implementing anomaly detection models allows organizations to flag these atypical transactions for further investigation, thereby enhancing security and reducing financial losses caused by fraud. This capability to detect outliers in transaction data is critical in real-time monitoring systems used by banks and online payment systems to quickly respond to potential fraud. Other scenarios such as identifying duplicate records, determining customer buying habits, and grouping similar products do not inherently require anomaly detection techniques. Duplicate records are identified through data cleaning processes, customer buying habits are analyzed through traditional statistical analysis, and similar products are typically grouped through clustering algorithms that focus on similarity rather than identifying deviations from standard patterns.

10. What are hyperparameters in machine learning?

- A. User-defined values that never change.**
- B. Configurations that control the learning process of the algorithm.**
- C. Output variables that report performance.**
- D. Parameters that are automatically optimized by the model.**

Hyperparameters are crucial components in machine learning as they define the configurations that control the learning process of an algorithm. Unlike parameters that the model learns during training, hyperparameters are set before the training begins and dictate aspects such as the learning rate, the number of hidden layers in a neural network, the batch size, and the number of trees in a random forest, among others. Adjusting these hyperparameters can significantly influence the performance and effectiveness of the model. Understanding hyperparameters helps practitioners optimize the model's learning process and performance, as tuning them can lead to better generalization on unseen data. Hyperparameter tuning is typically performed using techniques like grid search or random search to find the best combination for a given dataset and task, indicating their foundational role in shaping how a model learns. The other options do not correctly capture the essence of hyperparameters. For instance, the idea of user-defined values that never change doesn't reflect their dynamic nature in training configurations. Output variables that report performance refer to metrics derived after the model has been trained, while parameters that are automatically optimized by the model usually pertain to the weights and biases learned during training, rather than the hyperparameters set beforehand.

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://databricksmlassociate.examzify.com>

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

SAMPLE