# AWS Certified Machine Learning Specialty (MLS-C01) 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. What does overfitting refer to in machine learning?
  - A. A model that performs poorly on training data
  - B. A model that generalizes well to new, unseen data
  - C. A model that performs exceptionally well on training data but poorly on unseen data
  - D. A model that is too simple to capture the underlying data patterns
- 2. What method is used to fill missing values between the item start and item end date of a data set?
  - A. Middle filling
  - B. Back filling
  - C. Future filling
  - D. Average filling
- 3. What is the developer kit that helps learn machine learning concepts with computer vision?
  - A. AWS RoboMaker
  - **B. AWS DeepLens**
  - C. Amazon Rekognition
  - D. Amazon SageMaker
- 4. Which service is best suited for analyzing streaming data in real-time using Apache Flink?
  - A. Amazon Kinesis Data Streams
  - **B.** Amazon Kinesis Data Analytics
  - C. Amazon Kinesis Firehose
  - D. AWS Glue
- 5. Which of the following services would you use to implement real-time data stream processing?
  - A. Amazon S3
  - **B.** Amazon Kinesis Data Analytics
  - C. AWS Lambda
  - D. Amazon RDS

- 6. What metric evaluates the fraction of true positive instances among all the positive predictions?
  - A. Accuracy
  - **B. Precision**
  - C. Recall
  - **D. Specificity**
- 7. What is the main function of multiclass classification in machine learning?
  - A. To label target data into multiple categories
  - B. To detect anomalies in data
  - C. To visualize high-dimensional data
  - D. To interpolate missing values
- 8. Which algorithm would you utilize in Amazon SageMaker for topic modeling?
  - A. Amazon SageMaker k-Nearest Neighbors
  - **B. Amazon SageMaker Support Vector Machine**
  - C. Amazon SageMaker Latent Dirichlet Allocation (LDA)
  - D. Amazon SageMaker Random Forest
- 9. Which service facilitates real-time collection and analysis of video and data streams?
  - A. Amazon Kinesis Data Streams
  - B. Amazon QuickSight
  - C. Amazon Kinesis
  - D. AWS Batch
- 10. How is the performance of regression models typically assessed?
  - A. By looking at the accuracy percentage alone.
  - B. Using metrics such as MAE, MSE, and R<sup>2</sup>.
  - C. By visual inspection of prediction errors.
  - D. Only through cross-validation techniques.

#### **Answers**



- 1. C 2. A 3. B

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



#### **Explanations**



- 1. What does overfitting refer to in machine learning?
  - A. A model that performs poorly on training data
  - B. A model that generalizes well to new, unseen data
  - C. A model that performs exceptionally well on training data but poorly on unseen data
  - D. A model that is too simple to capture the underlying data patterns

Overfitting occurs when a machine learning model learns the details and noise in the training data to the extent that it negatively impacts its performance on new, unseen data. When a model is overfitted, it has essentially memorized the training data rather than learning the underlying patterns. This typically results in high accuracy on the training dataset but significantly reduced accuracy on validation or test datasets. In contrast, a model that generalizes well will perform reasonably on both the training and unseen data, indicating that it has effectively captured the essential patterns without being too tailored to the training data. Models that perform poorly on training data would indicate underfitting, where the model is unable to learn enough from the data, and a model that is too simple would also suffer from underfitting, not capturing the complexity required for accurate predictions.

- 2. What method is used to fill missing values between the item start and item end date of a data set?
  - A. Middle filling
  - **B.** Back filling
  - C. Future filling
  - D. Average filling

The correct method for filling missing values between item start and item end dates is commonly referred to as "middle filling." This technique typically involves inferring or estimating values that fall between two known points, effectively creating a smoother transition across gaps in the data. In scenarios where you have start and end dates, middle filling would aim to populate the values with those that logically fit the sequence of time-based data. For example, if you have a date range with missing entries, middle filling ensures that the gaps are filled using information from surrounding points, enabling a coherent timeline. The other methods mentioned have different specific applications. Back filling and future filling are used to propagate known values backward or forward in time, respectively. Average filling generally means replacing missing values with the average of available data points, a method not specifically suited for time series data where the sequence is crucial. Thus, "middle filling" aligns best with the objective of accurately completing the dataset based on the temporal context.

## 3. What is the developer kit that helps learn machine learning concepts with computer vision?

- A. AWS RoboMaker
- **B. AWS DeepLens**
- C. Amazon Rekognition
- D. Amazon SageMaker

The developer kit that helps learn machine learning concepts with a focus on computer vision is AWS DeepLens. This device is equipped with a built-in camera and runs deep learning models locally, allowing users to experiment with real-time image and video analysis directly on the device. DeepLens supports frameworks such as TensorFlow and MXNet, enabling users to deploy pre-trained models or create their own to perform various computer vision tasks, such as object recognition, image classification, and facial recognition. Engaging with AWS DeepLens provides a practical, hands-on approach to understanding machine learning concepts by integrating hardware and software, making it an ideal choice for learning. The ability to process video feeds in real time at the edge also enhances the learning experience, allowing developers to see tangible results immediately. Other options, while relevant to machine learning, do not specifically address the learning of computer vision concepts in the same way. AWS RoboMaker is focused on robotics and simulation; Amazon Rekognition is a service for image and video analysis but does not provide a hands-on learning kit, and Amazon SageMaker is a comprehensive service for building, training, and deploying machine learning models but is not specifically tailored to computer vision learning in a practical sense.

## 4. Which service is best suited for analyzing streaming data in real-time using Apache Flink?

- A. Amazon Kinesis Data Streams
- **B.** Amazon Kinesis Data Analytics
- C. Amazon Kinesis Firehose
- D. AWS Glue

Amazon Kinesis Data Analytics is the most suitable service for analyzing streaming data in real-time using Apache Flink. This service is specifically designed to process and analyze large streams of data in real-time, leveraging Apache Flink as its underlying engine. It allows users to write SQL queries on live data streams, facilitating immediate insights and enabling prompt decision-making. With Kinesis Data Analytics, you can easily integrate data from the Kinesis Data Streams or other sources, apply complex processing and transformations with Flink, and then either store the results in databases or stream them to other services for further analysis or action. This makes it highly effective for scenarios where real-time data processing is required. Other services, while related to streaming data, serve different primary functions. For example, while Amazon Kinesis Data Streams is ideal for collecting and processing large streams of data, it does not provide the built-in capability for sophisticated data analysis that Kinesis Data Analytics does. Amazon Kinesis Firehose is focused on delivering data streams to destinations like Amazon S3 or Redshift, rather than the analysis of the data itself. AWS Glue is primarily a data integration service that prepares data for analytics but does not specialize in real-time analysis like Kinesis Data Analytics. Thus, Kinesis Data Analytics

## 5. Which of the following services would you use to implement real-time data stream processing?

- A. Amazon S3
- **B.** Amazon Kinesis Data Analytics
- C. AWS Lambda
- D. Amazon RDS

Amazon Kinesis Data Analytics is designed specifically for real-time data stream processing. This service allows users to analyze streaming data using standard SQL queries. It can take in large volumes of data from various sources, providing the ability to continuously process and analyze this data for insights in real-time. Kinesis Data Analytics enables real-time data analysis, making it ideal for scenarios such as monitoring applications, financial transactions, or social media feeds. By integrating with other AWS services, it can easily fetch data from sources like Kinesis Data Streams or Amazon Kinesis Data Firehose, process it on the fly, and deliver structured outputs to dashboards or other AWS services. While AWS Lambda is also a tool that supports event-driven architectures and can process data in real-time, it is not fundamentally designed for stream processing in the way that Kinesis Data Analytics is. Amazon S3 is primarily a storage solution and does not provide real-time processing capabilities on its own. Amazon RDS, a managed relational database service, is used for structured data storage and retrieval but does not support real-time data stream processing. Therefore, for the requirement of implementing real-time data stream processing, Amazon Kinesis Data Analytics is the most suitable choice.

## 6. What metric evaluates the fraction of true positive instances among all the positive predictions?

- A. Accuracy
- **B. Precision**
- C. Recall
- D. Specificity

The correct answer is precision. Precision measures the ratio of true positive instances to the total number of instances that were predicted as positive. In other words, it evaluates how many of the positively identified instances are actually correct. Precision is particularly important in scenarios where the cost of false positives is high, as it helps gauge the reliability of the positive predictions made by the model. It focuses on the quality of the positive predictions, ensuring that when a prediction indicates a positive outcome, there is a high likelihood that it is indeed a true positive. In contrast, accuracy measures the overall correctness of the model across both positive and negative classes, which can be misleading, especially in cases with imbalanced datasets. Recall, on the other hand, focuses on the ability of the model to identify all relevant instances (true positives) among the actual positives, rather than evaluating the predicted positives. Specificity is concerned with the model's ability to correctly identify all actual negatives and does not relate to positive predictions.

#### 7. What is the main function of multiclass classification in machine learning?

- A. To label target data into multiple categories
- B. To detect anomalies in data
- C. To visualize high-dimensional data
- D. To interpolate missing values

The primary function of multiclass classification in machine learning is to label target data into multiple categories. In many real-world scenarios, the target variable can take on more than two values, and multiclass classification is specifically designed to handle such cases. Unlike binary classification, which deals with two classes, multiclass classification allows an algorithm to predict which category a given input belongs to among a set of multiple categories. For example, if you were classifying types of fruits, your classes might be apples, oranges, and bananas. The model learns from the training data, distinguishing the various characteristics of each of the classes, and then accurately classifies new instances based on that learned information. In this context, the other options represent different machine learning tasks that do not directly relate to multiclass classification. While anomaly detection focuses on identifying rare or unexpected observations in data, visualizing high-dimensional data is more about understanding complex relationships rather than categorization. Similarly, interpolating missing values pertains to estimating gaps in data rather than classifying data into distinct groups.

## 8. Which algorithm would you utilize in Amazon SageMaker for topic modeling?

- A. Amazon SageMaker k-Nearest Neighbors
- B. Amazon SageMaker Support Vector Machine
- C. Amazon SageMaker Latent Dirichlet Allocation (LDA)
- D. Amazon SageMaker Random Forest

Latent Dirichlet Allocation (LDA) is a generative probabilistic model commonly used for topic modeling in large collections of documents. It works by assuming that documents are mixtures of topics and that each topic is characterized by a distribution of words. This means LDA can effectively identify and extract the latent topics that are present in a set of text data. In the context of Amazon SageMaker, choosing LDA for topic modeling is particularly suitable because it is designed specifically for handling the unsupervised nature of topic classification across documents. This allows data scientists and machine learning practitioners to uncover hidden thematic structures in their text data without needing labeled examples. Other algorithms mentioned, such as k-Nearest Neighbors, Support Vector Machine, and Random Forest, are primarily designed for classification, regression, or supervised learning tasks. They do not inherently fit the requirements for topic modeling, which relies on the ability to decipher underlying themes or topics within unlabeled text data. Therefore, LDA stands out as the appropriate choice for effectively conducting topic modeling in Amazon SageMaker.

#### 9. Which service facilitates real-time collection and analysis of video and data streams?

- A. Amazon Kinesis Data Streams
- B. Amazon QuickSight
- C. Amazon Kinesis
- D. AWS Batch

Amazon Kinesis is designed to facilitate the real-time collection, processing, and analysis of streaming data, including video and other types of data streams. It allows users to build applications that can continuously ingest large amounts of data in real time, making it ideal for scenarios involving video feeds or other data streams that require immediate processing and analytics. Kinesis offers capabilities such as Kinesis Data Streams, which specifically allows for real-time ingestion and processing, making it suitable for analyzing video data as it arrives. This flexibility enables users to create various applications like live dashboards, data analytics, and real-time alerts based on incoming data. In contrast, while Amazon Kinesis Data Streams is a component of Kinesis focused on streaming data, it does not encompass the entirety of the Kinesis offerings. Amazon QuickSight is a business intelligence service for dashboards and visualizations but does not specialize in real-time collection or processing of data streams. AWS Batch is used for running batch computing workloads and is not suitable for real-time data handling or analysis.

#### 10. How is the performance of regression models typically assessed?

- A. By looking at the accuracy percentage alone.
- B. Using metrics such as MAE, MSE, and R<sup>2</sup>.
- C. By visual inspection of prediction errors.
- D. Only through cross-validation techniques.

The performance of regression models is assessed using various statistical metrics that provide insights into how well the model predicts the target variable. Specifically, metrics like Mean Absolute Error (MAE), Mean Squared Error (MSE), and the coefficient of determination (R2) are commonly utilized for this purpose. MAE measures the average magnitude of the errors in a set of predictions, without considering their direction, which allows it to provide a clear indication of the average prediction error. MSE, on the other hand, squares the errors before averaging them, giving more weight to larger errors and thus providing a useful measure to understand the variance in predictions. Lastly, R<sup>2</sup> indicates the proportion of variance in the dependent variable that can be explained by the independent variables in the model, providing a sense of how well the model fits the data. These metrics collectively give a comprehensive understanding of a regression model's performance, far beyond what a simple accuracy percentage could offer, especially since accuracy is more relevant for classification tasks rather than for regression scenarios. Visual inspections can be helpful but are often subjective and do not provide the quantitative analysis that metrics like these can offer. Additionally, while cross-validation is an important technique to evaluate models, it is not the singular method for performance assessment; rather