

# AWS Certified Machine Learning - Specialty Practice Questions

**Requirement:** Share 65 ML Specialty practice questions.

**Important Note:** The practice questions should appropriately belong to ML Specialty in terms of exam objectives & difficulty level.

**Delivery Timeline:** < Date >

## Question Response Types

There are two types of questions:

- Multiple Choice Single Response – **1** correct answer **3** incorrect responses (distractors).
- Multiple Choice Multiple Response – **2** or more correct answers out of **5** or more options.

## Important Note

- Do write Question Number for quick identification. Q# 1, Q# 2 .... & so on.
- Please mention Domain (based on ML Specialty exam blueprint), Topic & Sub-Topic (If Applicable) with every question.
- Note that we're expecting standard scenario based questions & NOT straight-forward definition kind of questions.
- The options should not have any obviously incorrect option. We need to word the options so that all of them should appear correct for the students, but a subtle point should mark the correct answer without any ambiguity. So, one answer should be the best choice without any doubt.
- The answer / explanation section should contain explanations on why the answer is correct and others are incorrect. It should also contain the relevant resource link (for details) preferably from AWS documentation.
  - Example
    - Option A is incorrect because..
    - Option B is CORRECT because...
    - Option C is incorrect because..
    - Option D is incorrect because..
- Try to balance the domains based on weightage % defined in the exam blueprint.
- Any AWS service or feature must be approximately 6 months old to figure out in Practice Tests. Put a note in the explanation for any latest service or feature that might be on the borderline of appearing in the real exam.

- **Plagiarism** in any form - Question or in Explanation will be **rejected**. Questions & Explanations should reflect your own professional experience & AWS skills. Author's who indulge in plagiarism will be **blacklisted** & dropped from our author's list.
- The ownership of the questions once approved & published on Whizlabs LMS platform, lies solely with Whizlabs Software Pvt. Ltd. You can't share or publish it elsewhere in any circumstances.

## Sample Format of Questions

---

**Question : #**

**Main Topic : < >**

**Sub Topic : [optional]**

**Domain: < >**

**Question text:**

<Scenario based. Should be clear in terms of requirements. No ambiguity. No duplicate options. In case of multiple answers, at the end, you should include number of expected answers. e.g. **Choose 2 answers**, choose 3 answers etc. For single answers this is NOT required>

**A)** Option A...

**B)** Option B...

**C)** Option C...

**D)** Option D...

**Answer:** A and C

**Explanation:**

**Option A is CORRECT because...**

**Option B is incorrect because...**

**Option C is CORRECT because...**

**Option D is incorrect because...**

[Insert the explanation in clear and lucid language here.]

**Diagram:** [Optional] [Insert the architectural or conceptual diagram here.]

**Reference:** [Insert the references here - which may include links to AWS Documentation, Blog, re:Invent video, Authority YouTube video].

---

### ML Specialty has 4 Domains

S. No.	Name of the Domain	Weight	Estimated No. of Questions (out of 65 As per weightage %)
1	Data Engineering	20%	13
2	Exploratory Data Analysis	24%	16
3	Modeling	36%	23
4	ML Implementation and Operations	20%	13

-----Question Section Starts-----

Question: 1

**Main Topic :** Machine Learning

**Sub Topic :** Create data repositories for machine learning

**Domain:** Data Engineering

**Question text:**

You are a machine learning expert working for a marketing firm. You are supporting a team of data scientists and marketing managers who are running a marketing campaign. Your data scientists and marketing managers need to answer the question “Will this user subscribe to my campaign?” You have been given a dataset in the form of a CSV file which is formatted as such:

UserId, jobId, jobDescription, educationLevel, campaign, duration, willRespondToCampaign

When you build your schema for this dataset, which of the following data descriptors would you use to define the willRespondToCampaign attribute? Please choose 2 answers.

- A) CATEGORICAL
- B) targetAttributeName
- C) TEXT
- D) BINARY
- E) Numeric
- F) rowId

**Answer:** B and D

**Explanation:**

Option A is incorrect because you choose the CATEGORICAL data type for an attribute that holds a limited set of unique strings. For example, a user name, the region, and a product code are categorical values. The willRespondToCampaign attribute takes on either ‘yes’ or ‘no’ values, which are binary in nature.

Option B is correct because for each user observation you are trying to discern “Will this user subscribe to my campaign?” You assign the targetAttributeName field value to the name of the attribute that you are trying to predict. You must assign a targetAttributeName when you create or evaluate your model.

Option C is incorrect because you choose the TEXT data type for an attribute that is a string, or a set of words. Amazon ML converts text attributes into tokens and uses white space as a delimiter. For example, document title becomes document and title, and document-title here becomes document-title and here.

Option D is correct because you choose the BINARY data type for an attribute that only has two possible values, such as yes or no, or true or false. The attribute willRespondToCampaign has only two possible answers: yes or no.

Option E is incorrect because you choose the NUMERIC data type for an attribute that holds a quantity as a number. For example, count, height, and acceleration rate are numeric values.

Option F is incorrect because you choose the rowId field value as an optional flag associated with an attribute in the input data. If you specify an attribute as the rowId, it is included in the prediction output. This attribute allows you to associate each prediction with its observation. The willRespondToCampaign attribute would make a poor identifier for each observation since it only takes two values: yes, or no.

**Reference:**

Please see the AWS Developer Guide titled **Creating a Data Schema for Amazon ML** (<https://docs.aws.amazon.com/machine-learning/latest/dg/creating-a-data-schema-for-amazon-ml.html#assigning-data-types>) for a complete description of the schema attributes.

Question: 2

**Main Topic :** Machine Learning

**Sub Topic :** Identify and implement a data-ingestion solution

**Domain:** Data Engineering

**Question text:**

You work for an energy company that buys and sells energy to customers. To get the best prices for their energy customers, your company trades financial energy derivative futures contracts. The trading of these futures contracts requires accurate forecasting of energy prices. You need to build a model that compares spot prices (current commodity price) to future commodity prices (price that a commodity can be bought or sold in the future). Your model needs to assist your company's futures traders in hedging against future energy price changes based on current price predictions. To source the model with appropriate data you need to gather and process the energy price data automatically.

The data pipeline requires two sources of data:

- 1) Historic energy spot prices
- 2) Energy consumption and production rates

Based on the company analysts' requirements, you have decided you need multiple years of historical data. You also realize you'll need to update the data feed daily as the market prices

change. You can gather the required data through APIs from data provider vendor systems. Your company's traders require a forecast from your model multiple times per day to help them form their trading strategy. So your pipeline needs to call the data provider APIs multiple times per day. Your data-ingestion pipeline needs to take the data from the API calls, perform preprocessing, and then store the data in an S3 data lake from which your forecasting model will access the data.

Your data-ingestion pipeline has three main steps:

- 1) Data ingestion
- 2) Data storage
- 3) Inference generation

Assuming you have written a lambda function that interacts with the data provider APIs and stores the data in CSV format, which of the following python libraries are the best option to perform the data preprocessing to transform the data by changing raw feature vectors into a format best suited for a SageMaker batch transform job to generate your forecast?

- A) matplotlib and plotly
- B) boto3 and moto
- C) pandas and scikit-learn
- D) NLTK and scrapy

**Answer: C**

**Explanation:**

Option A is incorrect because matplotlib and plotly are data visualization python libraries which contain no data transformation functions (see <https://matplotlib.org> and <https://plot.ly/python/>).

Option B is incorrect because boto3 is a python library that is used to interface with AWS services such as S3, DynamoDB, SQS, etc. Boto3 has no data transformation functions (see <https://aws.amazon.com/sdk-for-python/>). Moto is a python library used to mock interfaces to AWS services such as S3, DynamoDB, SQS, etc. The moto library also contains no data transformation functions (see <https://pypi.org/project/moto/>).

Option C is correct because pandas is the best choice for data wrangling and manipulation of tabular data such as CSV formatted data (see <https://pypi.org/project/pandas/>). Scikit-learn is the best python package to transform raw feature vectors into a format suited to downstream estimators (see <https://scikit-learn.org/stable/modules/preprocessing.html>).

Option D is incorrect because Natural Language Toolkit (NLTK) is best suited to text tagging, classification, and tokenizing, not manipulation of tabular data (see <https://www.nltk.org>). Scrapy is best suited to crawling functionality used to gather structured data from websites, not manipulation of tabular data (see <https://scrapy.org>).

**Reference:**

Please see the scikit-learn preprocessing data documentation:

<https://scikit-learn.org/stable/modules/preprocessing.html>, and a detailed pandas example: <https://towardsdatascience.com/why-and-how-to-use-pandas-with-large-data-9594dda2ea4c>

Question: 3

**Main Topic :** Machine Learning

**Sub Topic :** Perform hyperparameter optimization

**Domain:** Modeling

**Question text:**

You work for a retail firm that wishes to conduct a direct mail campaign to attract new customers. Your marketing manager wishes to get answers to questions that can be put into discrete categories, such as “using historical customer email campaign responses, should this customer receive an email from our current campaign?” You decide to use the SageMaker Linear Learner algorithm to build your model. Which hyperparameter setting would you use to get the algorithm to produce discrete results?

- A) set the objective hyperparameter to reg:logistic.
- B) set the predictor\_type hyperparameter to binary\_classifier.
- C) set the predictor\_type hyperparameter to regressor.
- D) set the objective hyperparameter to reg:linear.

**Answer:** B

**Explanation:**

Option A is incorrect because the objective hyperparameter is set to reg:logistic when you are using the XGBoost algorithm (See the AWS SageMaker developer documentation: [https://docs.aws.amazon.com/sagemaker/latest/dg/xgboost\\_hyperparameters.html](https://docs.aws.amazon.com/sagemaker/latest/dg/xgboost_hyperparameters.html)).

Option B is correct because the AWS SageMaker documentation states that for this type of discrete classification problem, when using the Linear Learner algorithm, you set the predictor\_type hyperparameter to binary\_classifier (See the AWS SageMaker documentation: [https://sagemaker.readthedocs.io/en/stable/linear\\_learner.html](https://sagemaker.readthedocs.io/en/stable/linear_learner.html)).

Option C is incorrect because the predictor\_type hyperparameter is set to regressor when you are using the Linear Learner algorithm for answers that are quantitative, not discrete (See the AWS SageMaker documentation: [https://sagemaker.readthedocs.io/en/stable/linear\\_learner.html](https://sagemaker.readthedocs.io/en/stable/linear_learner.html)).

Option D is incorrect because the objective hyperparameter is set to reg:linear when you are using the XGBoost algorithm for answers that are quantitative in nature (See the AWS SageMaker developer documentation:

[https://docs.aws.amazon.com/sagemaker/latest/dg/xgboost\\_hyperparameters.html](https://docs.aws.amazon.com/sagemaker/latest/dg/xgboost_hyperparameters.html)).

**Reference:**

Please see the AWS SageMaker developer guide titled **Using Amazon SageMaker Built-in Algorithms**: <https://docs.aws.amazon.com/sagemaker/latest/dg/algos.html>) for a complete description of the SageMaker hyperparameter settings.

Question: 4

**Main Topic :** Machine Learning

**Sub Topic :** Select the appropriate model(s) for a given machine learning problem

**Domain:** Modeling

**Question text:**

You work for the information security department of a major corporation. You have been asked to build a solution that detects web application log anomalies to protect your organization from fraudulent activity. The system needs to have near-real-time updates to the model where log entry data points dynamically change the underlying model as the log files are updated. Which AWS service component do you use to implement the best algorithm based on these requirements?

- A) SageMaker Random Cut Forest
- B) Kinesis Data Streams Naive Bayes Classifier
- C) Kinesis Data Analytics Random Cut Forest
- D) Kinesis Data Analytics Nearest Neighbor

**Answer:** C

**Explanation:**

Option A is incorrect because SageMaker Random Cut Forest is best used for large batch data sets where you don't need to update the model frequently (See AWS Kinesis Data Analytics documentation:

<https://docs.aws.amazon.com/kinesisanalytics/latest/sqlref/sqlrf-random-cut-forest.html>).

Answer B is incorrect because the Naive Bayes Classifier is used to find independent data points. The Kinesis Data Streams service does not have machine learning algorithm capabilities (See the AWS Kinesis Streams developer documentation:

<https://docs.aws.amazon.com/streams/latest/dev/introduction.html>).



Option C is correct. The Kinesis Data Analytics Random Cut Forest algorithm works really well for near-real-time updates to your model (See the AWS Kinesis Data Analytics documentation: <https://docs.aws.amazon.com/kinesisanalytics/latest/sqlref/sqlrf-random-cut-forest.html>).

Option D is incorrect because Kinesis Data Analytics provides a hotspots function that detects higher than normal activity using the distance between a hotspot and its nearest neighbor, but it does not provide ML model update capabilities (See AWS Kinesis Data Analytics documentation: <https://docs.aws.amazon.com/kinesisanalytics/latest/sqlref/sqlrf-hotspots.html>).

**Reference:**

For an example, please see the AWS Big Data blog titled **Perform Near Real-time Analytics on Streaming Data with Amazon Kinesis and Amazon Elasticsearch Service**:

<https://aws.amazon.com/blogs/big-data/perform-near-real-time-analytics-on-streaming-data-with-amazon-kinesis-and-amazon-elasticsearch-service/>) for a complete description of the use of Kinesis Data Analytics and the random cut forest algorithm.

Question: 5

**Main Topic :** Machine Learning

**Sub Topic :** Sanitize and prepare data for modeling

**Domain:** Exploratory Data Analysis

**Question text:**

You work in the data analytics department of a ride sharing software company. You need to use the K-means machine learning algorithm to separate your company's optimized ride data into clusters based on ride coordinates. How would you best use AWS Glue to build the data tables needed to classify the ride data?

- A) Use Glue crawlers together with a K-means classifier to classify the ride data based on coordinates
- B) Use Glue FindMatches to find and remove duplicate records in you data
- C) Use Glue to automatically generate code to classify the ride data based on coordinates
- D) Use Glue to transform and flatten your data so you can classify the ride data based on coordinates

**Answer:** A

**Explanation:**

Option A is correct. The best way to classify your optimized data is to use a Glue crawler that applies the K-means algorithm. See the AWS Machine Learning documentation (See the AWS

SageMaker <https://docs.aws.amazon.com/sagemaker/latest/dg/k-means.html> and AWS Glue crawler <https://docs.aws.amazon.com/glue/latest/dg/add-crawler.html> documentation).

Answer B is incorrect because there is no stated need to remove duplicates from the data.

Option C is incorrect because you don't need to automatically generate code since Glue will classify your data based on a prioritized list of classifiers without custom code (See the AWS Glue developers guide: (<https://docs.aws.amazon.com/glue/latest/dg/add-classifier.html>)).

Option D is incorrect because there is no stated requirement to flatten the ride data.

**Reference:**

For an example, please see the AWS Machine Learning blog titled **Serverless unsupervised machine learning with AWS Glue and Amazon Athena**:

<https://aws.amazon.com/blogs/machine-learning/serverless-unsupervised-machine-learning-with-aws-glue-and-amazon-athena/>).