

**Exam** : **MLS-C01**

**Title** : AWS Certified Machine  
Learning - Specialty

**Vendor** : Amazon

**Version** : V12.95

**NO.1** A company's Machine Learning Specialist needs to improve the training speed of a time-series forecasting model using TensorFlow. The training is currently implemented on a single-GPU machine and takes approximately 23 hours to complete. The training needs to be run daily.

The model accuracy is acceptable, but the company anticipates a continuous increase in the size of the training data and a need to update the model on an hourly, rather than a daily, basis. The company also wants to minimize coding effort and infrastructure changes. What should the Machine Learning Specialist do to the training solution to allow it to scale for future demand?

- A.** Do not change the TensorFlow code. Change the machine to one with a more powerful GPU to speed up the training.
- B.** Change the TensorFlow code to implement a Horovod distributed framework supported by Amazon SageMaker. Parallelize the training to as many machines as needed to achieve the business goals.
- C.** Switch to using a built-in AWS SageMaker DeepAR model. Parallelize the training to as many machines as needed to achieve the business goals.
- D.** Move the training to Amazon EMR and distribute the workload to as many machines as needed to achieve the business goals.

**Answer:** B

**NO.2** A Machine Learning Specialist observes several performance problems with the training portion of a machine learning solution on Amazon SageMaker. The solution uses a large training dataset 2 TB in size and is using the SageMaker k-means algorithm. The observed issues include the unacceptable length of time it takes before the training job launches and poor I/O throughput while training the model. What should the Specialist do to address the performance issues with the current solution?

- A.** Use the SageMaker batch transform feature.
- B.** Compress the training data into Apache Parquet format.
- C.** Ensure that the input mode for the training job is set to Pipe.
- D.** Copy the training dataset to an Amazon EFS volume mounted on the SageMaker instance.

**Answer:** B

**NO.3** An online reseller has a large, multi-column dataset with one column missing 30% of its data. A Machine Learning Specialist believes that certain columns in the dataset could be used to reconstruct the missing data. Which reconstruction approach should the Specialist use to preserve the integrity of the dataset?

- A.** Listwise deletion
- B.** Last observation carried forward
- C.** Multiple imputation
- D.** Mean substitution

**Answer:** C

**NO.4** A Machine Learning Specialist kicks off a hyperparameter tuning job for a tree-based ensemble model using Amazon SageMaker with Area Under the ROC Curve (AUC) as the objective metric. This workflow will eventually be deployed in a pipeline that retrains and tunes hyperparameters each night to model click-through on data that goes stale every 24 hours. With the goal of decreasing the

amount of time it takes to train these models, and ultimately to decrease costs, the Specialist wants to reconfigure the input hyperparameter range(s) Which visualization will accomplish this?

- A.** A histogram showing whether the most important input feature is Gaussian.
- B.** A scatter plot with points colored by target variable that uses (-Distributed Stochastic Neighbor Embedding (t-SNE) to visualize the large number of input variables in an easier-to-read dimension.
- C.** A scatter plot showing the performance of the objective metric over each training iteration
- D.** A scatter plot showing the correlation between maximum tree depth and the objective metric.

**Answer:** B

**NO.5** A company is running a machine learning prediction service that generates 100 TB of predictions every day A Machine Learning Specialist must generate a visualization of the daily precision-recall curve from the predictions, and forward a read-only version to the Business team. Which solution requires the LEAST coding effort?

- A.** Run a daily Amazon EMR workflow to generate precision-recall data, and save the results in Amazon S3 Give the Business team read-only access to S3
- B.** Generate daily precision-recall data in Amazon QuickSight, and publish the results in a dashboard shared with the Business team
- C.** Run a daily Amazon EMR workflow to generate precision-recall data, and save the results in Amazon S3 Visualize the arrays in Amazon QuickSight, and publish them in a dashboard shared with the Business team
- D.** Generate daily precision-recall data in Amazon ES, and publish the results in a dashboard shared with the Business team.

**Answer:** C

**NO.6** A Machine Learning Specialist is working with a large cybersecurity company that manages security events in real time for companies around the world The cybersecurity company wants to design a solution that will allow it to use machine learning to score malicious events as anomalies on the data as it is being ingested The company also wants be able to save the results in its data lake for later processing and analysis What is the MOST efficient way to accomplish these tasks'?

- A.** Ingest the data using Amazon Kinesis Data Firehose, and use Amazon Kinesis Data Analytics Random Cut Forest (RCF) for anomaly detection Then use Kinesis Data Firehose to stream the results to Amazon S3
- B.** Ingest the data into Apache Spark Streaming using Amazon EMR. and use Spark MLlib with k-means to perform anomaly detection Then store the results in an Apache Hadoop Distributed File System (HDFS) using Amazon EMR with a replication factor of three as the data lake
- C.** Ingest the data and store it in Amazon S3 Use AWS Batch along with the AWS Deep Learning AMIs to train a k-means model using TensorFlow on the data in Amazon S3.
- D.** Ingest the data and store it in Amazon S3. Have an AWS Glue job that is triggered on demand transform the new data Then use the built-in Random Cut Forest (RCF) model within Amazon SageMaker to detect anomalies in the data

**Answer:** B

**NO.7** A company is running an Amazon SageMaker training job that will access data stored in its Amazon S3 bucket A compliance policy requires that the data never be transmitted across the

internet How should the company set up the job?

- A.** Launch the notebook instances in a public subnet and access the data through the public S3 endpoint
- B.** Launch the notebook instances in a private subnet and access the data through a NAT gateway
- C.** Launch the notebook instances in a public subnet and access the data through a NAT gateway
- D.** Launch the notebook instances in a private subnet and access the data through an S3 VPC endpoint.

**Answer:** D

**NO.8** A Machine Learning Specialist is building a convolutional neural network (CNN) that will classify 10 types of animals. The Specialist has built a series of layers in a neural network that will take an input image of an animal, pass it through a series of convolutional and pooling layers, and then finally pass it through a dense and fully connected layer with 10 nodes The Specialist would like to get an output from the neural network that is a probability distribution of how likely it is that the input image belongs to each of the 10 classes Which function will produce the desired output?

- A.** Dropout
- B.** Smooth L1 loss
- C.** Softmax
- D.** Rectified linear units (ReLU)

**Answer:** D

**NO.9** A Machine Learning Specialist is preparing data for training on Amazon SageMaker The Specialist is transformed into a numpy .array, which appears to be negatively affecting the speed of the training What should the Specialist do to optimize the data for training on SageMaker'?

- A.** Use the SageMaker batch transform feature to transform the training data into a DataFrame
- B.** Use AWS Glue to compress the data into the Apache Parquet format
- C.** Transform the dataset into the Recordio protobuf format
- D.** Use the SageMaker hyperparameter optimization feature to automatically optimize the data

**Answer:** C

**NO.10** A company is using Amazon Polly to translate plaintext documents to speech for automated company announcements However company acronyms are being mispronounced in the current documents How should a Machine Learning Specialist address this issue for future documents'?

- A.** Convert current documents to SSML with pronunciation tags
- B.** Create an appropriate pronunciation lexicon.
- C.** Output speech marks to guide in pronunciation
- D.** Use Amazon Lex to preprocess the text files for pronunciation

**Answer:** A

**NO.11** A Machine Learning Specialist is required to build a supervised image-recognition model to identify a cat. The ML Specialist performs some tests and records the following results for a neural network-based image classifier:

Total number of images available = 1,000 Test set images = 100 (constant test set) The ML Specialist notices that, in over 75% of the misclassified images, the cats were held upside down by their

owners.

Which techniques can be used by the ML Specialist to improve this specific test error?

- A.** Increase the training data by adding variation in rotation for training images.
- B.** Increase the number of epochs for model training.
- C.** Increase the number of layers for the neural network.
- D.** Increase the dropout rate for the second-to-last layer.

**Answer:** B

**NO.12** When submitting Amazon SageMaker training jobs using one of the built-in algorithms, which common parameters **MUST** be specified? (Select THREE.)

- A.** The training channel identifying the location of training data on an Amazon S3 bucket.
- B.** The validation channel identifying the location of validation data on an Amazon S3 bucket.
- C.** The IAM role that Amazon SageMaker can assume to perform tasks on behalf of the users.
- D.** Hyperparameters in a JSON array as documented for the algorithm used.
- E.** The Amazon EC2 instance class specifying whether training will be run using CPU or GPU.
- F.** The output path specifying where on an Amazon S3 bucket the trained model will persist.

**Answer:** A E F

**NO.13** A Machine Learning Specialist receives customer data for an online shopping website. The data includes demographics, past visits, and locality information. The Specialist must develop a machine learning approach to identify the customer shopping patterns, preferences and trends to enhance the website for better service and smart recommendations.

Which solution should the Specialist recommend?

- A.** Latent Dirichlet Allocation (LDA) for the given collection of discrete data to identify patterns in the customer database.
- B.** A neural network with a minimum of three layers and random initial weights to identify patterns in the customer database
- C.** Collaborative filtering based on user interactions and correlations to identify patterns in the customer database
- D.** Random Cut Forest (RCF) over random subsamples to identify patterns in the customer database

**Answer:** C

**NO.14** An Machine Learning Specialist discover the following statistics while experimenting on a model.

Experiment 1  
Baseline model:  
Train error = 5%  
Test error = 16%

Experiment 2  
The Specialist added more layers and neurons to the model and received the following results:  
Train error = 5.2%  
Test error = 15.7%

Experiment 3  
The Specialist reverted back to the original number of neurons from Experiment 1 and implemented regularization in the neural network, which yielded the following results:  
Train error = 4.7%  
Test error = 9.5%

What can the Specialist from the experiments?

- A.** The model In Experiment 1 had a high variance error lthat was reduced in Experiment 3 by regularization Experiment 2 shows that there is minimal bias error in Experiment 1

- B.** The model in Experiment 1 had a high bias error that was reduced in Experiment 3 by regularization Experiment 2 shows that there is minimal variance error in Experiment 1
- C.** The model in Experiment 1 had a high bias error and a high variance error that were reduced in Experiment 3 by regularization Experiment 2 shows that high bias cannot be reduced by increasing layers and neurons in the model
- D.** The model in Experiment 1 had a high random noise error that was reduced in Experiment 3 by regularization Experiment 2 shows that random noise cannot be reduced by increasing layers and neurons in the model

**Answer:** C

**NO.15** A Machine Learning Specialist is working with a media company to perform classification on popular articles from the company's website. The company is using random forests to classify how popular an article will be before it is published A sample of the data being used is below. Given the dataset, the Specialist wants to convert the Day-Of\_Week column to binary values. What technique should be used to convert this column to binary values.

Article_Title	Author	Top_Keywords	Day_Of_Week	URL_of_Article	Page_Views
Building a Big Data Platform	Jane Doe	Big Data, Spark, Hadoop	Tuesday	http://examplecorp.com/data_platform.html	1300456
Getting Started with Deep Learning	John Doe	Deep Learning, Machine Learning, Spark	Tuesday	http://examplecorp.com/started_deep_learning.html	1230661
MXNet ML Guide	Jane Doe	Machine Learning, MXNet, Logistic Regression	Thursday	http://examplecorp.com/mxnet_guide.html	937291
Intro to NoSQL Databases	Mary Major	NoSQL, Operations, Database	Monday	http://examplecorp.com/nosql_intro_guide.html	407812

- A.** Binarization
- B.** One-hot encoding
- C.** Tokenization
- D.** Normalization transformation

**Answer:** B

**NO.16** A Machine Learning Specialist is developing a daily ETL workflow containing multiple ETL jobs The workflow consists of the following processes

- \* Start the workflow as soon as data is uploaded to Amazon S3
- \* When all the datasets are available in Amazon S3, start an ETL job to join the uploaded datasets with multiple terabyte-sized datasets already stored in Amazon S3
- \* Store the results of joining datasets in Amazon S3
- \* If one of the jobs fails, send a notification to the Administrator

Which configuration will meet these requirements?

- A.** Use AWS Lambda to trigger an AWS Step Functions workflow to wait for dataset uploads to complete in Amazon S3. Use AWS Glue to join the datasets Use an Amazon CloudWatch alarm to send an SNS notification to the Administrator in the case of a failure

- B.** Develop the ETL workflow using AWS Lambda to start an Amazon SageMaker notebook instance Use a lifecycle configuration script to join the datasets and persist the results in Amazon S3 Use an Amazon CloudWatch alarm to send an SNS notification to the Administrator in the case of a failure
- C.** Develop the ETL workflow using AWS Batch to trigger the start of ETL jobs when data is uploaded to Amazon S3 Use AWS Glue to join the datasets in Amazon S3 Use an Amazon CloudWatch alarm to send an SNS notification to the Administrator in the case of a failure
- D.** Use AWS Lambda to chain other Lambda functions to read and join the datasets in Amazon S3 as soon as the data is uploaded to Amazon S3 Use an Amazon CloudWatch alarm to send an SNS notification to the Administrator in the case of a failure

**Answer:** A

**NO.17** A web-based company wants to improve its conversion rate on its landing page Using a large historical dataset of customer visits, the company has repeatedly trained a multi-class deep learning network algorithm on Amazon SageMaker However there is an overfitting problem training data shows 90% accuracy in predictions, while test data shows 70% accuracy only The company needs to boost the generalization of its model before deploying it into production to maximize conversions of visits to purchases Which action is recommended to provide the HIGHEST accuracy model for the company's test and validation data?

- A.** Increase the randomization of training data in the mini-batches used in training.
- B.** Allocate a higher proportion of the overall data to the training dataset
- C.** Apply L1 or L2 regularization and dropouts to the training.
- D.** Reduce the number of layers and units (or neurons) from the deep learning network.

**Answer:** A

**NO.18** A Machine Learning Specialist is building a supervised model that will evaluate customers' satisfaction with their mobile phone service based on recent usage The model's output should infer whether or not a customer is likely to switch to a competitor in the next 30 days Which of the following modeling techniques should the Specialist use1?

- A.** Time-series prediction
- B.** Anomaly detection
- C.** Binary classification
- D.** Regression

**Answer:** D

**NO.19** A large mobile network operating company is building a machine learning model to predict customers who are likely to unsubscribe from the service. The company plans to offer an incentive for these customers as the cost of churn is far greater than the cost of the incentive. The model produces the following confusion matrix after evaluating on a test dataset of 100 customers:

Based on the model evaluation results, why is this a viable model for production?

n = 100	PREDICTED CHURN	
	Yes	No
ACTUAL Churn Yes	10	4
Actual No	10	76

- A. The model is 86% accurate and the cost incurred by the company as a result of false negatives is less than the false positives.
- B. The precision of the model is 86%, which is less than the accuracy of the model.
- C. The model is 86% accurate and the cost incurred by the company as a result of false positives is less than the false negatives.
- D. The precision of the model is 86%, which is greater than the accuracy of the model.

**Answer:** B

**NO.20** An insurance company is developing a new device for vehicles that uses a camera to observe drivers' behavior and alert them when they appear distracted. The company created approximately 10,000 training images in a controlled environment that a Machine Learning Specialist will use to train and evaluate machine learning models. During the model evaluation, the Specialist notices that the training error rate diminishes faster as the number of epochs increases and the model is not accurately inferring on the unseen test images. Which of the following should be used to resolve this issue? (Select TWO)

- A. Add vanishing gradient to the model
- B. Perform data augmentation on the training data
- C. Make the neural network architecture complex.
- D. Use gradient checking in the model
- E. Add L2 regularization to the model

**Answer:** B D

**NO.21** A Machine Learning Specialist is developing a custom video recommendation model for an application. The dataset used to train this model is very large with millions of data points and is hosted in an Amazon S3 bucket. The Specialist wants to avoid loading all of this data onto an Amazon SageMaker notebook instance because it would take hours to move and will exceed the attached 5 GB Amazon EBS volume on the notebook instance.

Which approach allows the Specialist to use all the data to train the model?

- A. Load a smaller subset of the data into the SageMaker notebook and train locally. Confirm that the training code is executing and the model parameters seem reasonable. Initiate a SageMaker training job using the full dataset from the S3 bucket using Pipe input mode.
- B. Launch an Amazon EC2 instance with an AWS Deep Learning AMI and attach the S3 bucket to the instance. Train on a small amount of the data to verify the training code and hyperparameters. Go back to Amazon SageMaker and train using the full dataset.
- C. Use AWS Glue to train a model using a small subset of the data to confirm that the data will be compatible with Amazon SageMaker. Initiate a SageMaker training job using the full dataset from the S3 bucket using Pipe input mode.
- D. Load a smaller subset of the data into the SageMaker notebook and train locally. Confirm that the training code is executing and the model parameters seem reasonable. Launch an Amazon EC2



instance with an AWS Deep Learning AMI and attach the S3 bucket to train the full dataset.

**Answer:** A

**NO.22** The Chief Editor for a product catalog wants the Research and Development team to build a machine learning system that can be used to detect whether or not individuals in a collection of images are wearing the company's retail brand. The team has a set of training data. Which machine learning algorithm should the researchers use that BEST meets their requirements?

- A. Latent Dirichlet Allocation (LDA)
- B. Recurrent neural network (RNN)
- C. K-means
- D. Convolutional neural network (CNN)

**Answer:** C

**NO.23** A Machine Learning Specialist is using Amazon SageMaker to host a model for a highly available customer-facing application.

The Specialist has trained a new version of the model, validated it with historical data, and now wants to deploy it to production. To limit any risk of a negative customer experience, the Specialist wants to be able to monitor the model and roll it back, if needed. What is the SIMPLEST approach with the LEAST risk to deploy the model and roll it back, if needed?

- A. Create a SageMaker endpoint and configuration for the new model version. Redirect production traffic to the new endpoint by updating the client configuration. Revert traffic to the last version if the model does not perform as expected.
- B. Create a SageMaker endpoint and configuration for the new model version. Redirect production traffic to the new endpoint by using a load balancer. Revert traffic to the last version if the model does not perform as expected.
- C. Update the existing SageMaker endpoint to use a new configuration that is weighted to send 5% of the traffic to the new variant. Revert traffic to the last version by resetting the weights if the model does not perform as expected.
- D. Update the existing SageMaker endpoint to use a new configuration that is weighted to send 100% of the traffic to the new variant. Revert traffic to the last version by resetting the weights if the model does not perform as expected.

**Answer:** A

**NO.24** A Marketing Manager at a pet insurance company plans to launch a targeted marketing campaign on social media to acquire new customers. Currently, the company has the following data in Amazon Aurora:

- \* Profiles for all past and existing customers
- \* Profiles for all past and existing insured pets
- \* Policy-level information
- \* Premiums received
- \* Claims paid

What steps should be taken to implement a machine learning model to identify potential new customers on social media?

- A. Use regression on customer profile data to understand key characteristics of consumer segments. Find similar profiles on social media.

- B.** Use clustering on customer profile data to understand key characteristics of consumer segments. Find similar profiles on social media.
- C.** Use a recommendation engine on customer profile data to understand key characteristics of consumer segments. Find similar profiles on social media
- D.** Use a decision tree classifier engine on customer profile data to understand key characteristics of consumer segments. Find similar profiles on social media

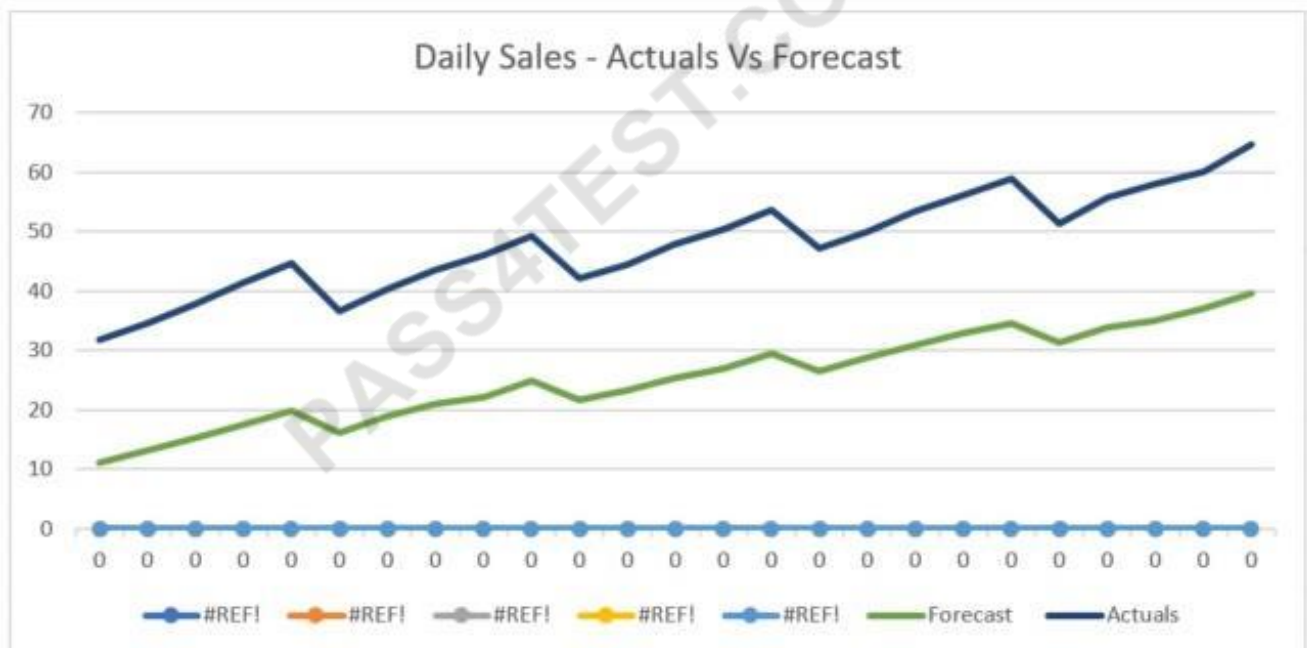
**Answer:** C

**NO.25** A Machine Learning Specialist is training a model to identify the make and model of vehicles in images. The Specialist wants to use transfer learning and an existing model trained on images of general objects. The Specialist collated a large custom dataset of pictures containing different vehicle makes and models.

- A.** Initialize the model with random weights in all layers including the last fully connected layer
- B.** Initialize the model with pre-trained weights in all layers and replace the last fully connected layer.
- C.** Initialize the model with random weights in all layers and replace the last fully connected layer
- D.** Initialize the model with pre-trained weights in all layers including the last fully connected layer

**Answer:** B

**NO.26** The displayed graph is from a forecasting model for testing a time series.



Considering the graph only, which conclusion should a Machine Learning Specialist make about the behavior of the model?

- A.** The model predicts both the trend and the seasonality well.
- B.** The model predicts the trend well, but not the seasonality.
- C.** The model predicts the seasonality well, but not the trend.
- D.** The model does not predict the trend or the seasonality well.

**Answer:** D

**NO.27** An agency collects census information within a country to determine healthcare and social

program needs by province and city. The census form collects responses for approximately 500 questions from each citizen Which combination of algorithms would provide the appropriate insights? (Select TWO )

- A. The factorization machines (FM) algorithm
- B. The Latent Dirichlet Allocation (LDA) algorithm
- C. The principal component analysis (PCA) algorithm
- D. The k-means algorithm
- E. The Random Cut Forest (RCF) algorithm

**Answer:** C D

Explanation

The PCA and K-means algorithms are useful in collection of data using census form.

**NO.28** A gaming company has launched an online game where people can start playing for free but they need to pay if they choose to use certain features The company needs to build an automated system to predict whether or not a new user will become a paid user within 1 year The company has gathered a labeled dataset from 1 million users The training dataset consists of 1.000 positive samples (from users who ended up paying within 1 year) and 999.000 negative samples (from users who did not use any paid features) Each data sample consists of 200 features including user age, device, location, and play patterns Using this dataset for training, the Data Science team trained a random forest model that converged with over 99% accuracy on the training set However, the prediction results on a test dataset were not satisfactory.

Which of the following approaches should the Data Science team take to mitigate this issue? (Select TWO.)

- A. Add more deep trees to the random forest to enable the model to learn more features.
- B. indicate a copy of the samples in the test database in the training dataset
- C. Generate more positive samples by duplicating the positive samples and adding a small amount of noise to the duplicated data.
- D. Change the cost function so that false negatives have a higher impact on the cost value than false positives
- E. Change the cost function so that false positives have a higher impact on the cost value than false negatives

**Answer:** B D

**NO.29** A Machine Learning Specialist is working for a credit card processing company and receives an unbalanced dataset containing credit card transactions. It contains 99,000 valid transactions and 1,000 fraudulent transactions The Specialist is asked to score a model that was run against the dataset The Specialist has been advised that identifying valid transactions is equally as important as identifying fraudulent transactions What metric is BEST suited to score the model?

- A. Precision
- B. Recall
- C. Area Under the ROC Curve (AUC)
- D. Root Mean Square Error (RMSE)

**Answer:** A

**NO.30** A Machine Learning team uses Amazon SageMaker to train an Apache MXNet handwritten digit classifier model using a research dataset. The team wants to receive a notification when the model is overfitting.

Auditors want to view the Amazon SageMaker log activity report to ensure there are no unauthorized API calls.

What should the Machine Learning team do to address the requirements with the least amount of code and fewest steps?

- A.** Implement an AWS Lambda function to log Amazon SageMaker API calls to Amazon S3. Add code to push a custom metric to Amazon CloudWatch. Create an alarm in CloudWatch with Amazon SNS to receive a notification when the model is overfitting.
- B.** Use AWS CloudTrail to log Amazon SageMaker API calls to Amazon S3. Add code to push a custom metric to Amazon CloudWatch. Create an alarm in CloudWatch with Amazon SNS to receive a notification when the model is overfitting.
- C.** Implement an AWS Lambda function to log Amazon SageMaker API calls to AWS CloudTrail. Add code to push a custom metric to Amazon CloudWatch. Create an alarm in CloudWatch with Amazon SNS to receive a notification when the model is overfitting.
- D.** Use AWS CloudTrail to log Amazon SageMaker API calls to Amazon S3. Set up Amazon SNS to receive a notification when the model is overfitting.

**Answer:** C

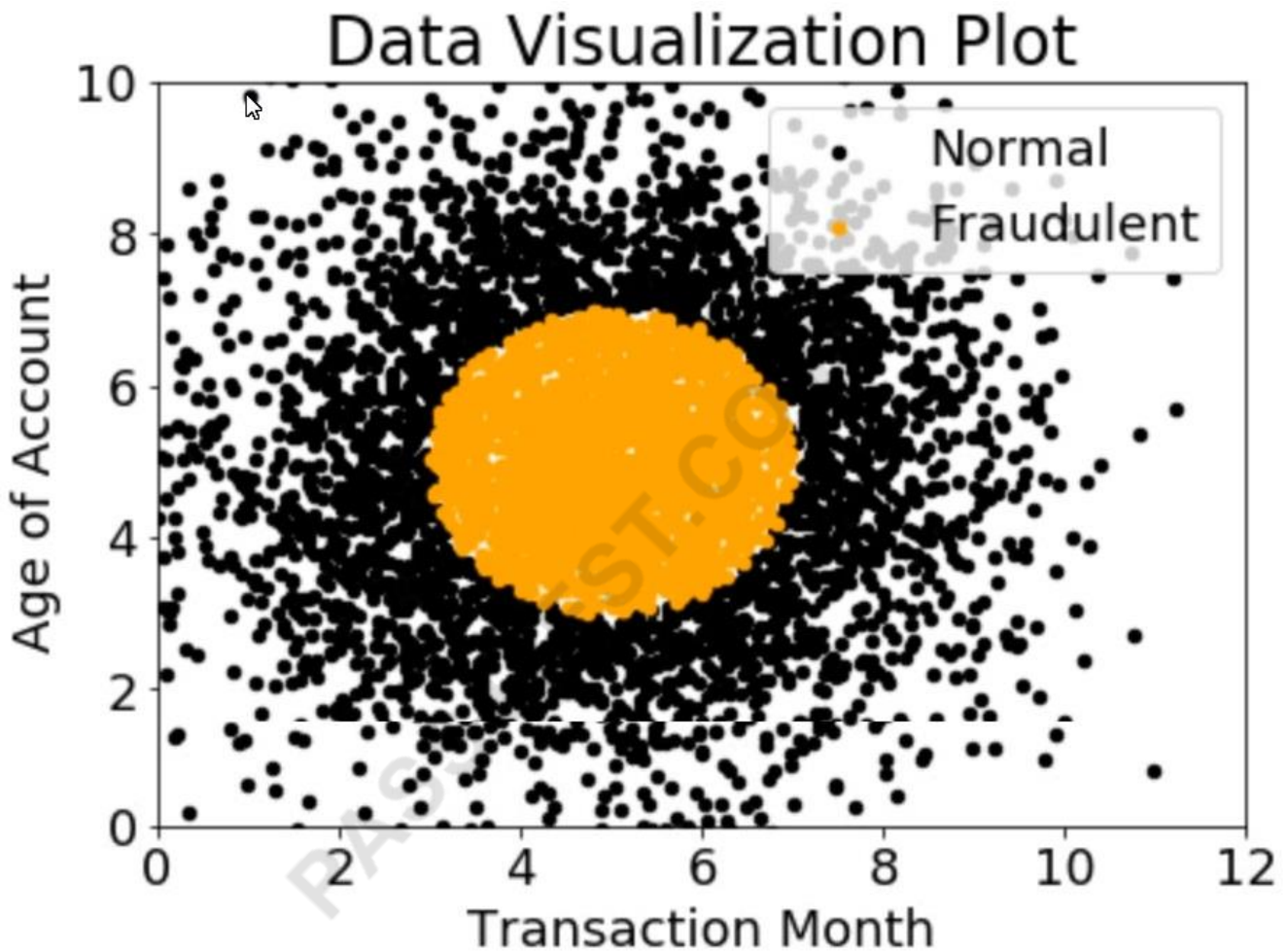
**NO.31** A Machine Learning Specialist trained a regression model, but the first iteration needs optimizing. The Specialist needs to understand whether the model is more frequently overestimating or underestimating the target.

What option can the Specialist use to determine whether it is overestimating or underestimating the target value?

- A.** Root Mean Square Error (RMSE)
- B.** Residual plots
- C.** Area under the curve
- D.** Confusion matrix

**Answer:** C

**NO.32** A company wants to classify user behavior as either fraudulent or normal. Based on internal research, a Machine Learning Specialist would like to build a binary classifier based on two features: age of account and transaction month. The class distribution for these features is illustrated in the figure provided.



Based on this information, which model would have the HIGHEST recall with respect to the fraudulent class?

- A. Decision tree
- B. Linear support vector machine (SVM)
- C. Naive Bayesian classifier
- D. Single Perceptron with sigmoidal activation function

**Answer:** C

**NO.33** Which of the following metrics should a Machine Learning Specialist generally use to compare/evaluate machine learning classification models against each other?

- A. Recall
- B. Misclassification rate
- C. Mean absolute percentage error (MAPE)
- D. Area Under the ROC Curve (AUC)

**Answer:** A

**NO.34** A Machine Learning Specialist is working with a large company to leverage machine learning within its products. The company wants to group its customers into categories based on which customers will and will not churn within the next 6 months. The company has labeled the data available to the Specialist.

Which machine learning model type should the Specialist use to accomplish this task?

- A. Linear regression
- B. Classification
- C. Clustering
- D. Reinforcement learning

**Answer:** B

Explanation

The goal of classification is to determine to which class or category a data point (customer in our case) belongs to. For classification problems, data scientists would use historical data with predefined target variables AKA labels (churner/non-churner) - answers that need to be predicted - to train a n algorithm. With classification, businesses can answer the following questions:

- \* Will this customer churn or not?
- \* Will a customer renew their subscription?
- \* Will a user downgrade a pricing plan?
- \* Are there any signs of unusual customer behavior?

**NO.35** A Machine Learning Specialist was given a dataset consisting of unlabeled data The Specialist must create a model that can help the team classify the data into different buckets What model should be used to complete this work?

- A. K-means clustering
- B. Random Cut Forest (RCF)
- C. XGBoost
- D. BlazingText

**Answer:** A

**NO.36** A Machine Learning Specialist has built a model using Amazon SageMaker built-in algorithms and is not getting expected accurate results The Specialist wants to use hyperparameter optimization to increase the model's accuracy Which method is the MOST repeatable and requires the LEAST amount of effort to achieve this?

- A. Launch multiple training jobs in parallel with different hyperparameters
- B. Create an AWS Step Functions workflow that monitors the accuracy in Amazon CloudWatch Logs and relaunches the training job with a defined list of hyperparameters
- C. Create a hyperparameter tuning job and set the accuracy as an objective metric.
- D. Create a random walk in the parameter space to iterate through a range of values that should be used for each individual hyperparameter

**Answer:** B

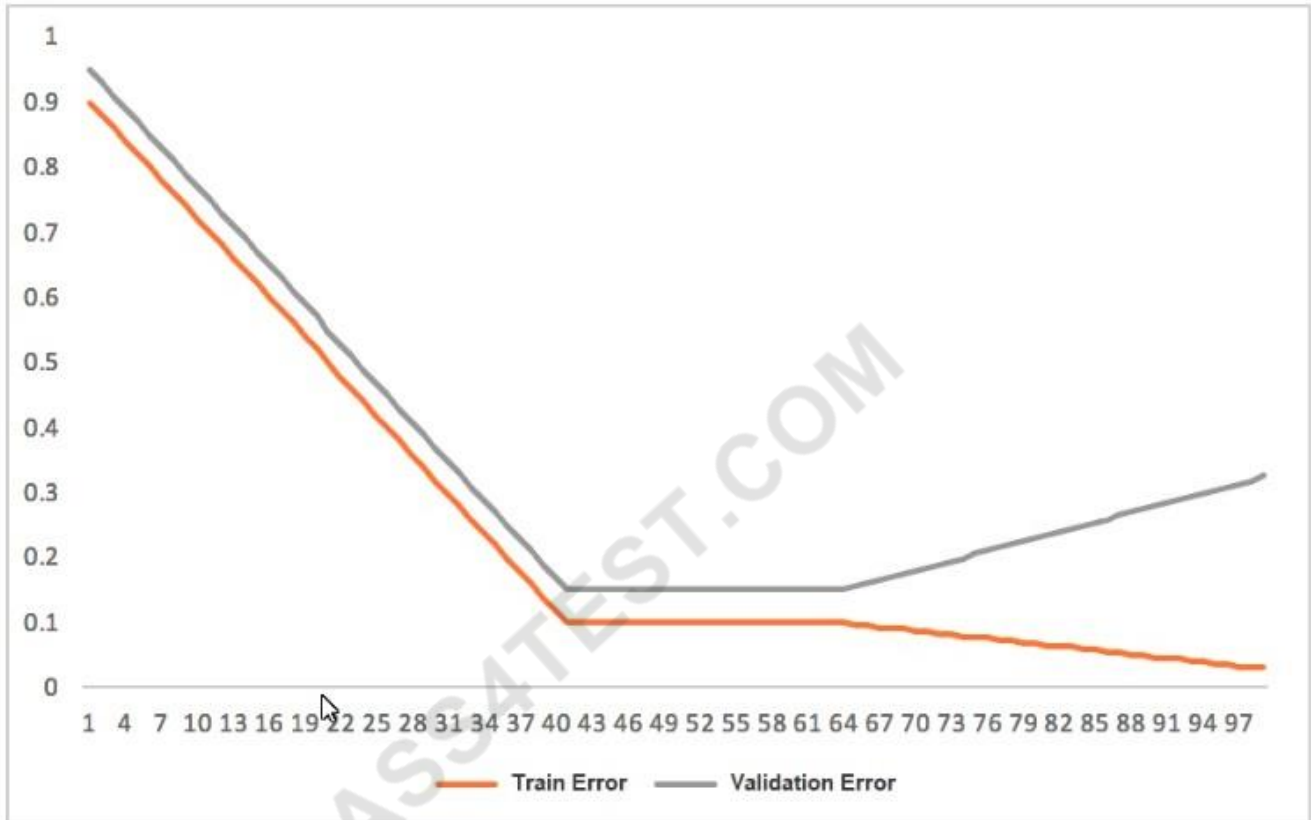
**NO.37** A bank's Machine Learning team is developing an approach for credit card fraud detection The company has a large dataset of historical data labeled as fraudulent The goal is to build a model to take the information from new transactions and predict whether each transaction is fraudulent or not Which built-in Amazon SageMaker machine learning algorithm should be used for modeling this problem?

- A. Seq2seq
- B. XGBoost
- C. K-means

**D. Random Cut Forest (RCF)****Answer:** C

**NO.38** This graph shows the training and validation loss against the epochs for a neural network. The network being trained is as follows:

- \* Two dense layers one output neuron
- \* 100 neurons in each layer
- \* 100 epochs
- \* Random initialization of weights

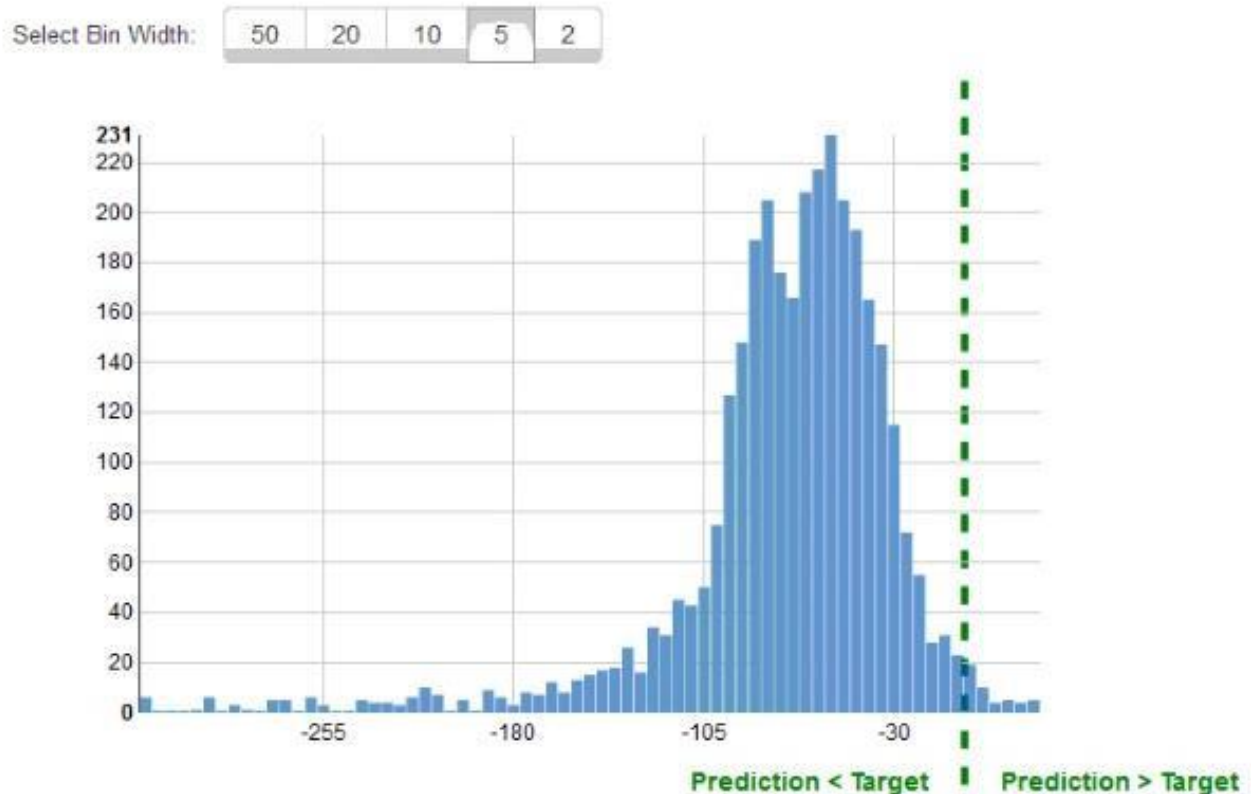


Which technique can be used to improve model performance in terms of accuracy in the validation set?

- A.** Early stopping
- B.** Random initialization of weights with appropriate seed
- C.** Increasing the number of epochs
- D.** Adding another layer with the 100 neurons

**Answer:** D

**NO.39** While reviewing the histogram for residuals on regression evaluation data a Machine Learning Specialist notices that the residuals do not form a zero-centered bell shape as shown. What does this mean?



- A. The model might have prediction errors over a range of target values.
- B. The dataset cannot be accurately represented using the regression model
- C. There are too many variables in the model
- D. The model is predicting its target values perfectly.

**Answer:** D

**NO.40** A Machine Learning Specialist is building a model that will perform time series forecasting using Amazon SageMaker. The Specialist has finished training the model and is now planning to perform load testing on the endpoint so they can configure Auto Scaling for the model variant. Which approach will allow the Specialist to review the latency, memory utilization, and CPU utilization during the load test?"

- A. Review SageMaker logs that have been written to Amazon S3 by leveraging Amazon Athena and Amazon QuickSight to visualize logs as they are being produced
- B. Generate an Amazon CloudWatch dashboard to create a single view for the latency, memory utilization, and CPU utilization metrics that are outputted by Amazon SageMaker
- C. Build custom Amazon CloudWatch Logs and then leverage Amazon ES and Kibana to query and visualize the data as it is generated by Amazon SageMaker
- D. Send Amazon CloudWatch Logs that were generated by Amazon SageMaker to Amazon ES and use Kibana to query and visualize the log data.

**Answer:** B

**NO.41** A Machine Learning Specialist is using Apache Spark for pre-processing training data. As part of the Spark pipeline, the Specialist wants to use Amazon SageMaker for training a model and hosting it. Which of the following would the Specialist do to integrate the Spark application with SageMaker? (Select THREE )



- A. Download the AWS SDK for the Spark environment
- B. Install the SageMaker Spark library in the Spark environment.
- C. Use the appropriate estimator from the SageMaker Spark Library to train a model.
- D. Compress the training data into a ZIP file and upload it to a pre-defined Amazon S3 bucket.
- E. Use the `sageMakerModel.transform` method to get inferences from the model hosted in SageMaker
- F. Convert the DataFrame object to a CSV file, and use the CSV file as input for obtaining inferences from SageMaker.

**Answer:** D E F

**NO.42** A Machine Learning Specialist is building a prediction model for a large number of features using linear models, such as linear regression and logistic regression. During exploratory data analysis, the Specialist observes that many features are highly correlated with each other. This may make the model unstable. What should be done to reduce the impact of having such a large number of features?

- A. Perform one-hot encoding on highly correlated features
- B. Use matrix multiplication on highly correlated features.
- C. Create a new feature space using principal component analysis (PCA)
- D. Apply the Pearson correlation coefficient

**Answer:** C

**NO.43** A company is setting up an Amazon SageMaker environment. The corporate data security policy does not allow communication over the internet. How can the company enable the Amazon SageMaker service without enabling direct internet access to Amazon SageMaker notebook instances?

- A. Create a NAT gateway within the corporate VPC.
- B. Route Amazon SageMaker traffic through an on-premises network.
- C. Create Amazon SageMaker VPC interface endpoints within the corporate VPC.
- D. Create VPC peering with Amazon VPC hosting Amazon SageMaker.

**Answer:** A

**NO.44** During mini-batch training of a neural network for a classification problem, a Data Scientist notices that training accuracy oscillates. What is the MOST likely cause of this issue?

- A. The class distribution in the dataset is imbalanced
- B. Dataset shuffling is disabled
- C. The batch size is too big
- D. The learning rate is very high

**Answer:** D

**NO.45** A Machine Learning Specialist is designing a system for improving sales for a company. The objective is to use the large amount of information the company has on users' behavior and product preferences to predict which products users would like based on the users' similarity to other users. What should the Specialist do to meet this objective?

- A.** Build a content-based filtering recommendation engine with Apache Spark ML on Amazon EMR.
- B.** Build a collaborative filtering recommendation engine with Apache Spark ML on Amazon EMR.
- C.** Build a model-based filtering recommendation engine with Apache Spark ML on Amazon EMR.
- D.** Build a combinative filtering recommendation engine with Apache Spark ML on Amazon EMR.

**Answer:** B

Explanation

Many developers want to implement the famous Amazon model that was used to power the "People who bought this also bought these items" feature on Amazon.com. This model is based on a method called Collaborative Filtering. It takes items such as movies, books, and products that were rated highly by a set of users and recommending them to other users who also gave them high ratings. This method works well in domains where explicit ratings or implicit user actions can be gathered and analyzed.

**NO.46** A Machine Learning Specialist must build out a process to query a dataset on Amazon S3 using Amazon Athena. The dataset contains more than 800,000 records stored as plaintext CSV files. Each record contains 200 columns and is approximately 1.5 MB in size. Most queries will span 5 to 10 columns only. How should the Machine Learning Specialist transform the dataset to minimize query runtime?

- A.** Convert the records to Apache Parquet format
- B.** Convert the records to JSON format
- C.** Convert the records to GZIP CSV format
- D.** Convert the records to XML format

**Answer:** A

**NO.47** A company is observing low accuracy while training on the default built-in image classification algorithm in Amazon SageMaker. The Data Science team wants to use an Inception neural network architecture instead of a ResNet architecture.

Which of the following will accomplish this? (Select TWO.)

- A.** Customize the built-in image classification algorithm to use Inception and use this for model training.
- B.** Create a support case with the SageMaker team to change the default image classification algorithm to Inception.
- C.** Bundle a Docker container with TensorFlow Estimator loaded with an Inception network and use this for model training.
- D.** Use custom code in Amazon SageMaker with TensorFlow Estimator to load the model with an Inception network and use this for model training.
- E.** Download and apt-get install the inception network code into an Amazon EC2 instance and use this instance as a Jupyter notebook in Amazon SageMaker.

**Answer:** A D

**NO.48** A Machine Learning Specialist built an image classification deep learning model. However, the Specialist ran into an overfitting problem in which the training and testing accuracies were 99% and 75% respectively.

How should the Specialist address this issue and what is the reason behind it?

- A.** The learning rate should be increased because the optimization process was trapped at a local minimum.
- B.** The dropout rate at the flatten layer should be increased because the model is not generalized enough.
- C.** The dimensionality of dense layer next to the flatten layer should be increased because the model is not complex enough.
- D.** The epoch number should be increased because the optimization process was terminated before it reached the global minimum.

**Answer:** D

**NO.49** A Machine Learning Specialist is assigned a TensorFlow project using Amazon SageMaker for training, and needs to continue working for an extended period with no Wi-Fi access.

Which approach should the Specialist use to continue working?

- A.** Install Python 3 and boto3 on their laptop and continue the code development using that environment.
- B.** Download the TensorFlow Docker container used in Amazon SageMaker from GitHub to their local environment, and use the Amazon SageMaker Python SDK to test the code.
- C.** Download TensorFlow from tensorflow.org to emulate the TensorFlow kernel in the SageMaker environment.
- D.** Download the SageMaker notebook to their local environment then install Jupyter Notebooks on their laptop and continue the development in a local notebook.

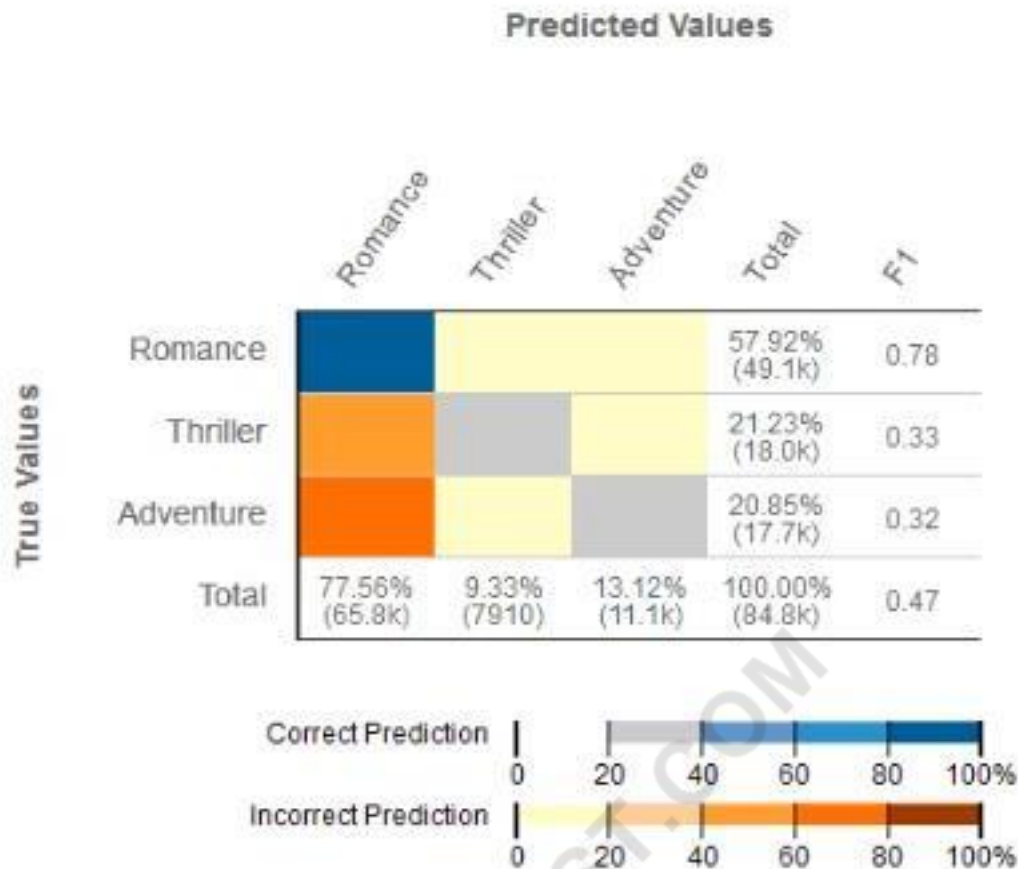
**Answer:** A

**NO.50** A Machine Learning Specialist has completed a proof of concept for a company using a small data sample and now the Specialist is ready to implement an end-to-end solution in AWS using Amazon SageMaker. The historical training data is stored in Amazon RDS. Which approach should the Specialist use for training a model using that data?

- A.** Write a direct connection to the SQL database within the notebook and pull data in
- B.** Push the data from Microsoft SQL Server to Amazon S3 using an AWS Data Pipeline and provide the S3 location within the notebook.
- C.** Move the data to Amazon DynamoDB and set up a connection to DynamoDB within the notebook to pull data in
- D.** Move the data to Amazon ElastiCache using AWS DMS and set up a connection within the notebook to pull data in for fast access.

**Answer:** B

**NO.51** Given the following confusion matrix for a movie classification model, what is the true class frequency for Romance and the predicted class frequency for Adventure?



- A.** The true class frequency for Romance is 77.56% and the predicted class frequency for Adventure is 20.85%
- B.** The true class frequency for Romance is 57.92% and the predicted class frequency for Adventure is 13.12%
- C.** The true class frequency for Romance is 0.78 and the predicted class frequency for Adventure is (0.47 - 0.32).
- D.** The true class frequency for Romance is  $77.56\% \times 0.78$  and the predicted class frequency for Adventure is  $20.85\% \times 0.32$

**Answer:** A

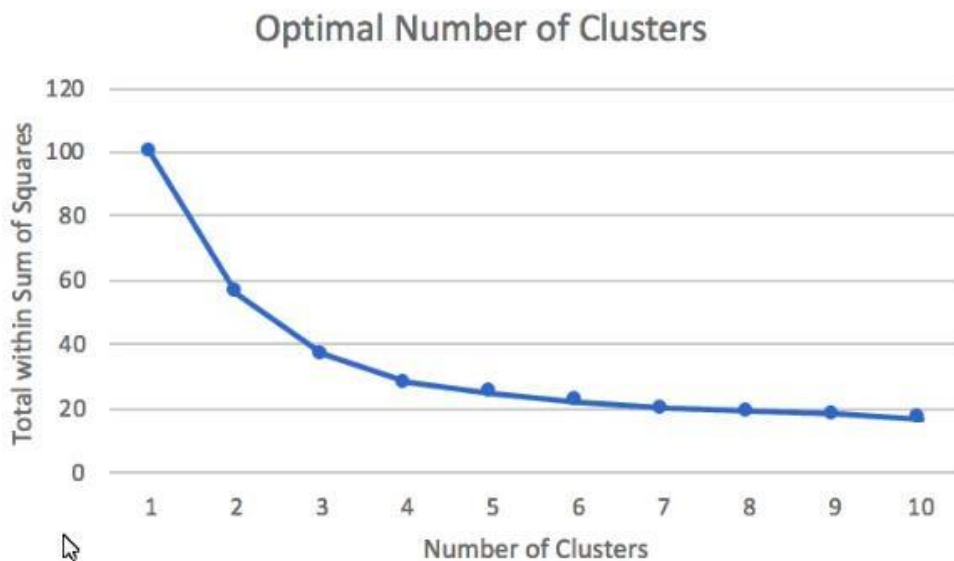
**NO.52** A Machine Learning Specialist is implementing a full Bayesian network on a dataset that describes public transit in New York City. One of the random variables is discrete, and represents the number of minutes New Yorkers wait for a bus given that the buses cycle every 10 minutes, with a mean of 3 minutes.

Which prior probability distribution should the ML Specialist use for this variable?

- A.** Poisson distribution ,
- B.** Uniform distribution
- C.** Normal distribution
- D.** Binomial distribution

**Answer:** D

**NO.53** A Machine Learning Specialist prepared the following graph displaying the results of k-means for  $k = [1:10]$



Considering the graph, what is a reasonable selection for the optimal choice of  $k$ ?

- A. 1
- B. 4
- C. 7
- D. 10

**Answer:** C

**NO.54** A Machine Learning Specialist is configuring automatic model tuning in Amazon SageMaker. When using the hyperparameter optimization feature, which of the following guidelines should be followed to improve optimization?

Choose the maximum number of hyperparameters supported by

- A. Amazon SageMaker to search the largest number of combinations possible
- B. Specify a very large hyperparameter range to allow Amazon SageMaker to cover every possible value.
- C. Use log-scaled hyperparameters to allow the hyperparameter space to be searched as quickly as possible
- D. Execute only one hyperparameter tuning job at a time and improve tuning through successive rounds of experiments

**Answer:** C

**NO.55** For the given confusion matrix, what is the recall and precision of the model?

		Actual	
		Yes	No
Predicted	Yes	12	3
	No	1	9

- A. Recall = 0.92 Precision = 0.84
- B. Recall = 0.84 Precision = 0.8
- C. Recall = 0.92 Precision = 0.8
- D. Recall = 0.8 Precision = 0.92

**Answer:** A

**NO.56** A manufacturing company asks its Machine Learning Specialist to develop a model that classifies defective parts into one of eight defect types. The company has provided roughly 100000 images per defect type for training. During the initial training of the image classification model the Specialist notices that the validation accuracy is 80%, while the training accuracy is 90%. It is known that human-level performance for this type of image classification is around 90%. What should the Specialist consider to fix this issue?

- A. A longer training time
- B. Making the network larger
- C. Using a different optimizer
- D. Using some form of regularization

**Answer:** D

**NO.57** A Machine Learning Specialist is packaging a custom ResNet model into a Docker container so the company can leverage Amazon SageMaker for training. The Specialist is using Amazon EC2 P3 instances to train the model and needs to properly configure the Docker container to leverage the NVIDIA GPUs.

What does the Specialist need to do?

- A. Bundle the NVIDIA drivers with the Docker image.
- B. Build the Docker container to be NVIDIA-Docker compatible.
- C. Organize the Docker container's file structure to execute on GPU instances.
- D. Set the GPU flag in the Amazon SageMaker CreateTrainingJob request body

**Answer:** A

**NO.58** A monitoring service generates 1 TB of scale metrics record data every minute. A Research team performs queries on this data using Amazon Athena. The queries run slowly due to the large volume of data, and the team requires better performance. How should the records be stored in Amazon S3 to improve query performance?

- A. CSV files

- B. Parquet files
- C. Compressed JSON
- D. RecordIO

**Answer:** B

**NO.59** Amazon Connect has recently been tolled out across a company as a contact call center. The solution has been configured to store voice call recordings on Amazon S3. The content of the voice calls are being analyzed for the incidents being discussed by the call operators. Amazon Transcribe is being used to convert the audio to text, and the output is stored on Amazon S3. Which approach will provide the information required for further analysis?

- A. Use Amazon Comprehend with the transcribed files to build the key topics
- B. Use Amazon Translate with the transcribed files to train and build a model for the key topics
- C. Use the AWS Deep Learning AMI with Gluon Semantic Segmentation on the transcribed files to train and build a model for the key topics
- D. Use the Amazon SageMaker k-Nearest-Neighbors (kNN) algorithm on the transcribed files to generate a word embeddings dictionary for the key topics

**Answer:** B

**NO.60** A Machine Learning Specialist needs to create a data repository to hold a large amount of time-based training data for a new model. In the source system, new files are added every hour. Throughout a single 24-hour period, the volume of hourly updates will change significantly. The Specialist always wants to train on the last 24 hours of the data.

Which type of data repository is the MOST cost-effective solution?

- A. An Amazon EBS-backed Amazon EC2 instance with hourly directories
- B. An Amazon RDS database with hourly table partitions
- C. An Amazon S3 data lake with hourly object prefixes
- D. An Amazon EMR cluster with hourly hive partitions on Amazon EBS volumes

**Answer:** C

**NO.61** A company has raw user and transaction data stored in Amazon S3, a MySQL database, and Amazon Redshift. A Data Scientist needs to perform an analysis by joining the three datasets from Amazon S3, MySQL, and Amazon Redshift, and then calculating the average of a few selected columns from the joined data. Which AWS service should the Data Scientist use?

- A. Amazon Athena
- B. Amazon Redshift Spectrum
- C. AWS Glue
- D. Amazon QuickSight

**Answer:** A

**NO.62** A manufacturing company has structured and unstructured data stored in an Amazon S3 bucket. A Machine Learning Specialist wants to use SQL to run queries on this data. Which solution requires the LEAST effort to be able to query this data?

- A. Use AWS Data Pipeline to transform the data and Amazon RDS to run queries.

- B. Use AWS Glue to catalogue the data and Amazon Athena to run queries
- C. Use AWS Batch to run ETL on the data and Amazon Aurora to run the queries
- D. Use AWS Lambda to transform the data and Amazon Kinesis Data Analytics to run queries

**Answer:** D

**NO.63** A city wants to monitor its air quality to address the consequences of air pollution. A Machine Learning Specialist needs to forecast the air quality in parts per million of contaminants for the next 2 days in the city. As this is a prototype, only daily data from the last year is available. Which model is MOST likely to provide the best results in Amazon SageMaker?

- A. Use the Amazon SageMaker k-Nearest-Neighbors (kNN) algorithm on the single time series consisting of the full year of data with a predictor\_type of regressor.
- B. Use Amazon SageMaker Random Cut Forest (RCF) on the single time series consisting of the full year of data.
- C. Use the Amazon SageMaker Linear Learner algorithm on the single time series consisting of the full year of data with a predictor\_type of regressor.
- D. Use the Amazon SageMaker Linear Learner algorithm on the single time series consisting of the full year of data with a predictor\_type of classifier.

**Answer:** C

**NO.64** A Machine Learning Specialist is configuring Amazon SageMaker so multiple Data Scientists can access notebooks, train models, and deploy endpoints. To ensure the best operational performance, the Specialist needs to be able to track how often the Scientists are deploying models, GPU and CPU utilization on the deployed SageMaker endpoints, and all errors that are generated when an endpoint is invoked.

Which services are integrated with Amazon SageMaker to track this information? (Select TWO.)

- A. AWS CloudTrail
- B. AWS Health
- C. AWS Trusted Advisor
- D. Amazon CloudWatch
- E. AWS Config

**Answer:** A D

**NO.65** A Machine Learning Specialist at a company sensitive to security is preparing a dataset for model training.

The dataset is stored in Amazon S3 and contains Personally Identifiable Information (PII). The dataset:

- \* Must be accessible from a VPC only.
- \* Must not traverse the public internet.

How can these requirements be satisfied?

- A. Create a VPC endpoint and apply a bucket access policy that restricts access to the given VPC endpoint and the VPC.
- B. Create a VPC endpoint and apply a bucket access policy that allows access from the given VPC endpoint and an Amazon EC2 instance.
- C. Create a VPC endpoint and use Network Access Control Lists (NACLs) to allow traffic between only the given VPC endpoint and an Amazon EC2 instance.



**D.** Create a VPC endpoint and use security groups to restrict access to the given VPC endpoint and an Amazon EC2 instance.

**Answer:** B

**NO.66** A Data Scientist is developing a machine learning model to predict future patient outcomes based on information collected about each patient and their treatment plans. The model should output a continuous value as its prediction. The data available includes labeled outcomes for a set of 4,000 patients. The study was conducted on a group of individuals over the age of 65 who have a particular disease that is known to worsen with age.

Initial models have performed poorly. While reviewing the underlying data, the Data Scientist notices that, out of 4,000 patient observations, there are 450 where the patient age has been input as 0. The other features for these observations appear normal compared to the rest of the sample population. How should the Data Scientist correct this issue?

**A.** Drop all records from the dataset where age has been set to 0.

**B.** Replace the age field value for records with a value of 0 with the mean or median value from the dataset.

**C.** Drop the age feature from the dataset and train the model using the rest of the features.

**D.** Use k-means clustering to handle missing features.

**Answer:** A

**NO.67** A Machine Learning Specialist is building a logistic regression model that will predict whether or not a person will order a pizza. The Specialist is trying to build the optimal model with an ideal classification threshold.

What model evaluation technique should the Specialist use to understand how different classification thresholds will impact the model's performance?

**A.** Receiver operating characteristic (ROC) curve

**B.** Misclassification rate

**C.** Root Mean Square Error (RMSE)

**D.** L1 norm

**Answer:** A

**NO.68** An office security agency conducted a successful pilot using 100 cameras installed at key locations within the main office. Images from the cameras were uploaded to Amazon S3 and tagged using Amazon Rekognition, and the results were stored in Amazon ES. The agency is now looking to expand the pilot into a full production system using thousands of video cameras in its office locations globally. The goal is to identify activities performed by non-employees in real time.

Which solution should the agency consider?

**A.** Use a proxy server at each local office and for each camera, and stream the RTSP feed to a unique Amazon Kinesis Video Streams video stream. On each stream, use Amazon Rekognition Video and create a stream processor to detect faces from a collection of known employees, and alert when non-employees are detected.

**B.** Use a proxy server at each local office and for each camera, and stream the RTSP feed to a unique Amazon Kinesis Video Streams video stream. On each stream, use Amazon Rekognition Image to detect faces from a collection of known employees and alert when non-employees are detected.

**C.** Install AWS DeepLens cameras and use the DeepLens\_Kinesis\_Video module to stream video to Amazon Kinesis Video Streams for each camera. On each stream, use Amazon Rekognition Video and create a stream processor to detect faces from a collection on each stream, and alert when nonemployees are detected.

**D.** Install AWS DeepLens cameras and use the DeepLens\_Kinesis\_Video module to stream video to Amazon Kinesis Video Streams for each camera. On each stream, run an AWS Lambda function to capture image fragments and then call Amazon Rekognition Image to detect faces from a collection of known employees, and alert when non-employees are detected.

**Answer:** D

**NO.69** A Mobile Network Operator is building an analytics platform to analyze and optimize a company's operations using Amazon Athena and Amazon S3. The source systems send data in CSV format in real time. The Data Engineering team wants to transform the data to the Apache Parquet format before storing it on Amazon S3. Which solution takes the LEAST effort to implement?

**A.** Ingest .CSV data using Apache Kafka Streams on Amazon EC2 instances and use Kafka Connect S3 to serialize data as Parquet.

**B.** Ingest .CSV data from Amazon Kinesis Data Streams and use Amazon Glue to convert data into Parquet.

**C.** Ingest .CSV data using Apache Spark Structured Streaming in an Amazon EMR cluster and use Apache Spark to convert data into Parquet.

**D.** Ingest .CSV data from Amazon Kinesis Data Streams and use Amazon Kinesis Data Firehose to convert data into Parquet.

**Answer:** C

**NO.70** A Data Engineer needs to build a model using a dataset containing customer credit card information.

How can the Data Engineer ensure the data remains encrypted and the credit card information is secure?

**A.** Use a custom encryption algorithm to encrypt the data and store the data on an Amazon SageMaker instance in a VPC. Use the SageMaker DeepAR algorithm to randomize the credit card numbers.

**B.** Use an IAM policy to encrypt the data on the Amazon S3 bucket and Amazon Kinesis to automatically discard credit card numbers and insert fake credit card numbers.

**C.** Use an Amazon SageMaker launch configuration to encrypt the data once it is copied to the SageMaker instance in a VPC. Use the SageMaker principal component analysis (PCA) algorithm to reduce the length of the credit card numbers.

**D.** Use AWS KMS to encrypt the data on Amazon S3.

**Answer:** C

**NO.71** A Machine Learning Specialist is creating a new natural language processing application that processes a dataset comprised of 1 million sentences. The aim is to then run Word2Vec to generate embeddings of the sentences and enable different types of predictions - Here is an example from the dataset

"The quck BROWN FOX jumps over the lazy dog "

Which of the following are the operations the Specialist needs to perform to correctly sanitize and prepare the data in a repeatable manner? (Select THREE)

- A. Perform part-of-speech tagging and keep the action verb and the nouns only
- B. Normalize all words by making the sentence lowercase
- C. Remove stop words using an English stopword dictionary.
- D. Correct the typography on "quck" to "quick."
- E. One-hot encode all words in the sentence
- F. Tokenize the sentence into words.

**Answer:** A B D

**NO.72** A Machine Learning Specialist is working with multiple data sources containing billions of records that need to be joined. What feature engineering and model development approach should the Specialist take with a dataset this large?

- A. Use an Amazon SageMaker notebook for both feature engineering and model development
- B. Use an Amazon SageMaker notebook for feature engineering and Amazon ML for model development
- C. Use Amazon EMR for feature engineering and Amazon SageMaker SDK for model development
- D. Use Amazon ML for both feature engineering and model development.

**Answer:** B

**NO.73** A Data Scientist wants to gain real-time insights into a data stream of GZIP files. Which solution would allow the use of SQL to query the stream with the LEAST latency?

- A. Amazon Kinesis Data Analytics with an AWS Lambda function to transform the data.
- B. AWS Glue with a custom ETL script to transform the data.
- C. An Amazon Kinesis Client Library to transform the data and save it to an Amazon ES cluster.
- D. Amazon Kinesis Data Firehose to transform the data and put it into an Amazon S3 bucket.

**Answer:** A

**NO.74** A Machine Learning Specialist has created a deep learning neural network model that performs well on the training data but performs poorly on the test data.

Which of the following methods should the Specialist consider using to correct this? (Select THREE.)

- A. Decrease regularization.
- B. Increase regularization.
- C. Increase dropout.
- D. Decrease dropout.
- E. Increase feature combinations.
- F. Decrease feature combinations.

**Answer:** B D E

**NO.75** A Machine Learning Specialist needs to move and transform data in preparation for training. Some of the data needs to be processed in near-real time and other data can be moved hourly. There are existing Amazon EMR MapReduce jobs to clean and feature engineering to perform on the data. Which of the following services can feed data to the MapReduce jobs? (Select TWO)

- A. AWS DMS
- B. Amazon Kinesis
- C. AWS Data Pipeline
- D. Amazon Athena
- E. Amazon ES

**Answer:** B D

**NO.76** Example Corp has an annual sale event from October to December. The company has sequential sales data from the past 15 years and wants to use Amazon ML to predict the sales for this year's upcoming event. Which method should Example Corp use to split the data into a training dataset and evaluation dataset?

- A. Pre-split the data before uploading to Amazon S3
- B. Have Amazon ML split the data randomly.
- C. Have Amazon ML split the data sequentially.
- D. Perform custom cross-validation on the data

**Answer:** C

**NO.77** While working on a neural network project, a Machine Learning Specialist discovers that some features in the data have very high magnitude resulting in this data being weighted more in the cost function. What should the Specialist do to ensure better convergence during backpropagation?

- A. Dimensionality reduction
- B. Data normalization
- C. Model regularization
- D. Data augmentation for the minority class

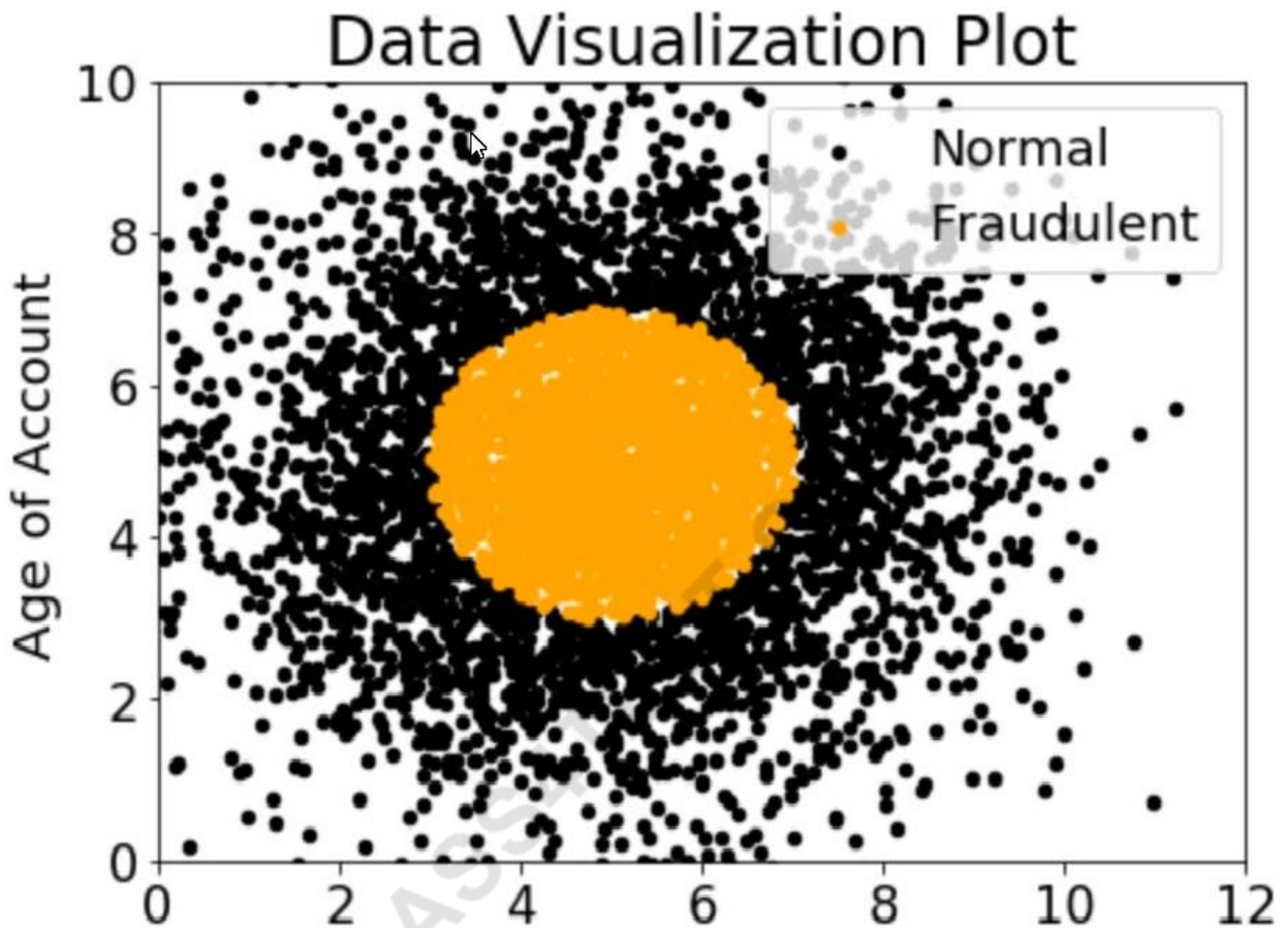
**Answer:** D

**NO.78** A Data Scientist is working on an application that performs sentiment analysis. The validation accuracy is poor and the Data Scientist thinks that the cause may be a rich vocabulary and a low average frequency of words in the dataset. Which tool should be used to improve the validation accuracy?

- A. Amazon Comprehend syntax analysis and entity detection
- B. Amazon SageMaker BlazingText allow mode
- C. Natural Language Toolkit (NLTK) stemming and stop word removal
- D. Scikit-learn term frequency-inverse document frequency (TF-IDF) vectorizers

**Answer:** D

**NO.79** A company wants to classify user behavior as either fraudulent or normal. Based on internal research, a Machine Learning Specialist would like to build a binary classifier based on two features: age of account and transaction month. The class distribution for these features is illustrated in the figure provided.



Based on this information which model would have the HIGHEST accuracy?

- A. Long short-term memory (LSTM) model with scaled exponential linear unit (SELL)
- B. Logistic regression
- C. Support vector machine (SVM) with non-linear kernel
- D. Single perceptron with tanh activation function

**Answer:** B

**NO.80** IT leadership wants to transition a company's existing machine learning data storage environment to AWS as a temporary ad hoc solution. The company currently uses a custom software process that heavily leverages SQL as a query language and exclusively stores generated CSV documents for machine learning. The ideal state for the company would be a solution that allows it to continue to use the current workforce of SQL experts. The solution must also support the storage of CSV and JSON files, and be able to query over semi-structured data. The following are high priorities for the company:

- \* Solution simplicity
- \* Fast development time
- \* Low cost
- \* High flexibility

What technologies meet the company's requirements?

- A. Amazon S3 and Amazon Athena
- B. Amazon Redshift and AWS Glue
- C. Amazon DynamoDB and DynamoDB Accelerator (DAX)

**D. Amazon RDS and Amazon ES****Answer:** B

**NO.81** A financial services company is building a robust serverless data lake on Amazon S3. The data lake should be flexible and meet the following requirements:

- \* Support querying old and new data on Amazon S3 through Amazon Athena and Amazon Redshift Spectrum.
- \* Support event-driven ETL pipelines.
- \* Provide a quick and easy way to understand metadata.

Which approach meets these requirements?

- A.** Use an AWS Glue crawler to crawl S3 data, an AWS Lambda function to trigger an AWS Glue ETL job, and an AWS Glue Data catalog to search and discover metadata.
- B.** Use an AWS Glue crawler to crawl S3 data, an AWS Lambda function to trigger an AWS Batch job, and an external Apache Hive metastore to search and discover metadata.
- C.** Use an AWS Glue crawler to crawl S3 data, an Amazon CloudWatch alarm to trigger an AWS Batch job, and an AWS Glue Data Catalog to search and discover metadata.
- D.** Use an AWS Glue crawler to crawl S3 data, an Amazon CloudWatch alarm to trigger an AWS Glue ETL job, and an external Apache Hive metastore to search and discover metadata.

**Answer:** B

**NO.82** An e-commerce company needs a customized training model to classify images of its shirts and pants products. The company needs a proof of concept in 2 to 3 days with good accuracy. Which compute choice should the Machine Learning Specialist select to train and achieve good accuracy on the model quickly?

- A.** m5.4xlarge (general purpose)
- B.** r5.2xlarge (memory optimized)
- C.** p3.2xlarge (GPU accelerated computing)
- D.** p3.8xlarge (GPU accelerated computing)

**Answer:** C

**NO.83** A Machine Learning Specialist needs to be able to ingest streaming data and store it in Apache Parquet files for exploration and analysis. Which of the following services would both ingest and store this data in the correct format?

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

**Answer:** C

**NO.84** A Machine Learning Specialist working for an online fashion company wants to build a data ingestion solution for the company's Amazon S3-based data lake.

The Specialist wants to create a set of ingestion mechanisms that will enable future capabilities comprised of:

- \* Real-time analytics

- \* Interactive analytics of historical data
- \* Clickstream analytics
- \* Product recommendations

Which services should the Specialist use?

- A.** AWS Glue as the data catalog; Amazon Kinesis Data Streams and Amazon Kinesis Data Analytics for real-time data insights; Amazon Kinesis Data Firehose for delivery to Amazon ES for clickstream analytics; Amazon EMR to generate personalized product recommendations
- B.** Amazon Athena as the data catalog; Amazon Kinesis Data Streams and Amazon Kinesis Data Analytics for near-realtime data insights; Amazon Kinesis Data Firehose for clickstream analytics; AWS Glue to generate personalized product recommendations
- C.** AWS Glue as the data catalog; Amazon Kinesis Data Streams and Amazon Kinesis Data Analytics for historical data insights; Amazon Kinesis Data Firehose for delivery to Amazon ES for clickstream analytics; Amazon EMR to generate personalized product recommendations
- D.** Amazon Athena as the data catalog; Amazon Kinesis Data Streams and Amazon Kinesis Data Analytics for historical data insights; Amazon DynamoDB streams for clickstream analytics; AWS Glue to generate personalized product recommendations

**Answer:** A

**NO.85** An Amazon SageMaker notebook instance is launched into Amazon VPC. The SageMaker notebook references data contained in an Amazon S3 bucket in another account. The bucket is encrypted using SSE-KMS. The instance returns an access denied error when trying to access data in Amazon S3.

Which of the following are required to access the bucket and avoid the access denied error? (Select THREE )

- A.** An AWS KMS key policy that allows access to the customer master key (CMK)
- B.** A SageMaker notebook security group that allows access to Amazon S3
- C.** An IAM role that allows access to the specific S3 bucket
- D.** A permissive S3 bucket policy
- E.** An S3 bucket owner that matches the notebook owner
- F.** A SageMaker notebook subnet ACL that allow traffic to Amazon S3.

**Answer:** A C F

**NO.86** A manufacturing company has structured and unstructured data stored in an Amazon S3 bucket. A Machine Learning Specialist wants to use SQL to run queries on this data.

Which solution requires the LEAST effort to be able to query this data?

- A.** Use AWS Data Pipeline to transform the data and Amazon RDS to run queries.
- B.** Use AWS Glue to catalogue the data and Amazon Athena to run queries.
- C.** Use AWS Batch to run ETL on the data and Amazon Aurora to run the queries.
- D.** Use AWS Lambda to transform the data and Amazon Kinesis Data Analytics to run queries.

**Answer:** B

**NO.87** A manufacturer of car engines collects data from cars as they are being driven. The data collected includes timestamp, engine temperature, rotations per minute (RPM), and other sensor readings. The company wants to predict when an engine is going to have a problem so it can notify

drivers in advance to get engine maintenance The engine data is loaded into a data lake for training Which is the MOST suitable predictive model that can be deployed into production'?

- A.** Add labels over time to indicate which engine faults occur at what time in the future to turn this into a supervised learning problem Use a recurrent neural network (RNN) to train the model to recognize when an engine might need maintenance for a certain fault.
- B.** This data requires an unsupervised learning algorithm Use Amazon SageMaker k-means to cluster the data
- C.** Add labels over time to indicate which engine faults occur at what time in the future to turn this into a supervised learning problem Use a convolutional neural network (CNN) to train the model to recognize when an engine might need maintenance for a certain fault.
- D.** This data is already formulated as a time series Use Amazon SageMaker seq2seq to model the time series.

**Answer:** B

**NO.88** A Machine Learning Specialist is using an Amazon SageMaker notebook instance in a private subnet of a corporate VPC. The ML Specialist has important data stored on the Amazon SageMaker notebook instance's Amazon EBS volume, and needs to take a snapshot of that EBS volume. However the ML Specialist cannot find the Amazon SageMaker notebook instance's EBS volume or Amazon EC2 instance within the VPC.

Why is the ML Specialist not seeing the instance visible in the VPC?

- A.** Amazon SageMaker notebook instances are based on the EC2 instances within the customer account, but they run outside of VPCs.
- B.** Amazon SageMaker notebook instances are based on the Amazon ECS service within customer accounts.
- C.** Amazon SageMaker notebook instances are based on EC2 instances running within AWS service accounts.
- D.** Amazon SageMaker notebook instances are based on AWS ECS instances running within AWS service accounts.

**Answer:** C

**NO.89** A Data Science team is designing a dataset repository where it will store a large amount of training data commonly used in its machine learning models. As Data Scientists may create an arbitrary number of new datasets every day the solution has to scale automatically and be cost-effective. Also, it must be possible to explore the data using SQL.

Which storage scheme is MOST adapted to this scenario?

- A.** Store datasets as files in Amazon S3.
- B.** Store datasets as files in an Amazon EBS volume attached to an Amazon EC2 instance.
- C.** Store datasets as tables in a multi-node Amazon Redshift cluster.
- D.** Store datasets as global tables in Amazon DynamoDB.

**Answer:** A

**NO.90** A retail chain has been ingesting purchasing records from its network of 20,000 stores to Amazon S3 using Amazon Kinesis Data Firehose To support training an improved machine learning model, training records will require new but simple transformations, and some attributes will be



combined The model needs to be retrained daily Given the large number of stores and the legacy data ingestion, which change will require the LEAST amount of development effort?

- A.** Require that the stores to switch to capturing their data locally on AWS Storage Gateway for loading into Amazon S3 then use AWS Glue to do the transformation
- B.** Deploy an Amazon EMR cluster running Apache Spark with the transformation logic, and have the cluster run each day on the accumulating records in Amazon S3, outputting new/transformed records to Amazon S3
- C.** Spin up a fleet of Amazon EC2 instances with the transformation logic, have them transform the data records accumulating on Amazon S3, and output the transformed records to Amazon S3.
- D.** Insert an Amazon Kinesis Data Analytics stream downstream of the Kinesis Data Firehouse stream that transforms raw record attributes into simple transformed values using SQL.

**Answer:** D

**NO.91** A Data Scientist is developing a machine learning model to classify whether a financial transaction is fraudulent. The labeled data available for training consists of 100,000 non-fraudulent observations and 1,000 fraudulent observations.

The Data Scientist applies the XGBoost algorithm to the data, resulting in the following confusion matrix when the trained model is applied to a previously unseen validation dataset. The accuracy of the model is

99.1%, but the Data Scientist has been asked to reduce the number of false negatives.

Predicted	0	1
Actual	0 99,966   34	
	1 877   123	

Which combination of steps should the Data Scientist take to reduce the number of false positive predictions by the model? (Select TWO.)

- A.** Change the XGBoost eval\_metric parameter to optimize based on rmse instead of error.
- B.** Increase the XGBoost scale\_pos\_weight parameter to adjust the balance of positive and negative weights.
- C.** Increase the XGBoost max\_depth parameter because the model is currently underfitting the data.
- D.** Change the XGBoost eval\_metric parameter to optimize based on AUC instead of error.
- E.** Decrease the XGBoost max\_depth parameter because the model is currently overfitting the data.

**Answer:** D E

**NO.92** A large JSON dataset for a project has been uploaded to a private Amazon S3 bucket The Machine Learning Specialist wants to securely access and explore the data from an Amazon SageMaker notebook instance A new VPC was created and assigned to the Specialist How can the privacy and integrity of the data stored in Amazon S3 be maintained while granting access to the Specialist for analysis?

- A.** Launch the SageMaker notebook instance within the VPC with SageMaker-provided internet access enabled Use an S3 ACL to open read privileges to the everyone group
- B.** Launch the SageMaker notebook instance within the VPC and create an S3 VPC endpoint for the

notebook to access the data Copy the JSON dataset from Amazon S3 into the ML storage volume on the SageMaker notebook instance and work against the local dataset

**C.** Launch the SageMaker notebook instance within the VPC and create an S3 VPC endpoint for the notebook to access the data Define a custom S3 bucket policy to only allow requests from your VPC to access the S3 bucket

**D.** Launch the SageMaker notebook instance within the VPC with SageMaker-provided internet access enabled. Generate an S3 pre-signed URL for access to data in the bucket

**Answer:** B

**NO.93** A Machine Learning Specialist is building a model to predict future employment rates based on a wide range of economic factors While exploring the data, the Specialist notices that the magnitude of the input features vary greatly The Specialist does not want variables with a larger magnitude to dominate the model What should the Specialist do to prepare the data for model training'?

**A.** Apply quantile binning to group the data into categorical bins to keep any relationships in the data by replacing the magnitude with distribution

**B.** Apply the Cartesian product transformation to create new combinations of fields that are independent of the magnitude

**C.** Apply normalization to ensure each field will have a mean of 0 and a variance of 1 to remove any significant magnitude

**D.** Apply the orthogonal sparse Diagram (OSB) transformation to apply a fixed-size sliding window to generate new features of a similar magnitude.

**Answer:** C

**NO.94** A Machine Learning Specialist wants to determine the appropriate SageMakerVariant Invocations Per Instance setting for an endpoint automatic scaling configuration. The Specialist has performed a load test on a single instance and determined that peak requests per second (RPS) without service degradation is about 20 RPS As this is the first deployment, the Specialist intends to set the invocation safety factor to 0.5 Based on the stated parameters and given that the invocations per instance setting is measured on a per-minute basis, what should the Specialist set as the sageMakervariantinvocationsPerinstance setting?

**A.** 10

**B.** 30

**C.** 600

**D.** 2,400

**Answer:** C

**NO.95** A Machine Learning Specialist deployed a model that provides product recommendations on a company's website Initially, the model was performing very well and resulted in customers buying more products on average However within the past few months the Specialist has noticed that the effect of product recommendations has diminished and customers are starting to return to their original habits of spending less The Specialist is unsure of what happened, as the model has not changed from its initial deployment over a year ago Which method should the Specialist try to improve model performance?

- A.** The model needs to be completely re-engineered because it is unable to handle product inventory changes
- B.** The model's hyperparameters should be periodically updated to prevent drift
- C.** The model should be periodically retrained from scratch using the original data while adding a regularization term to handle product inventory changes
- D.** The model should be periodically retrained using the original training data plus new data as product inventory changes

**Answer:** D

**NO.96** A Machine Learning Specialist is packaging a custom ResNet model into a Docker container so the company can leverage Amazon SageMaker for training. The Specialist is using Amazon EC2 P3 instances to train the model and needs to properly configure the Docker container to leverage the NVIDIA GPUs. What does the Specialist need to do?

- A.** Bundle the NVIDIA drivers with the Docker image
- B.** Build the Docker container to be NVIDIA-Docker compatible
- C.** Organize the Docker container's file structure to execute on GPU instances.
- D.** Set the GPU flag in the Amazon SageMaker Create TrainingJob request body

**Answer:** A

**NO.97** An employee found a video clip with audio on a company's social media feed. The language used in the video is Spanish. English is the employee's first language, and they do not understand Spanish. The employee wants to do a sentiment analysis.

What combination of services is the MOST efficient to accomplish the task?

- A.** Amazon Transcribe, Amazon Translate, and Amazon Comprehend
- B.** Amazon Transcribe, Amazon Comprehend, and Amazon SageMaker seq2seq
- C.** Amazon Transcribe, Amazon Translate, and Amazon SageMaker Neural Topic Model (NTM)
- D.** Amazon Transcribe, Amazon Translate, and Amazon SageMaker BlazingText

**Answer:** C

**NO.98** A large consumer goods manufacturer has the following products on sale

- \* 34 different toothpaste variants
- \* 48 different toothbrush variants
- \* 43 different mouthwash variants

The entire sales history of all these products is available in Amazon S3. Currently, the company is using custom-built autoregressive integrated moving average (ARIMA) models to forecast demand for these products. The company wants to predict the demand for a new product that will soon be launched. Which solution should a Machine Learning Specialist apply?

- A.** Train a custom ARIMA model to forecast demand for the new product.
- B.** Train an Amazon SageMaker DeepAR algorithm to forecast demand for the new product
- C.** Train an Amazon SageMaker k-means clustering algorithm to forecast demand for the new product.
- D.** Train a custom XGBoost model to forecast demand for the new product

**Answer:** B

Explanation

The Amazon SageMaker DeepAR forecasting algorithm is a supervised learning algorithm for forecasting scalar (one-dimensional) time series using recurrent neural networks (RNN). Classical forecasting methods, such as autoregressive integrated moving average (ARIMA) or exponential smoothing (ETS), fit a single model to each individual time series. They then use that model to extrapolate the time series into the future.

**NO.99** A Data Scientist needs to create a serverless ingestion and analytics solution for high-velocity, real-time streaming data.

The ingestion process must buffer and convert incoming records from JSON to a query-optimized, columnar format without data loss. The output datastore must be highly available, and Analysts must be able to run SQL queries against the data and connect to existing business intelligence dashboards. Which solution should the Data Scientist build to satisfy the requirements?

- A.** Create a schema in the AWS Glue Data Catalog of the incoming data format. Use an Amazon Kinesis Data Firehose delivery stream to stream the data and transform the data to Apache Parquet or ORC format using the AWS Glue Data Catalog before delivering to Amazon S3. Have the Analysts query the data directly from Amazon S3 using Amazon Athena, and connect to BI tools using the Athena Java Database Connectivity (JDBC) connector.
- B.** Write each JSON record to a staging location in Amazon S3. Use the S3 Put event to trigger an AWS Lambda function that transforms the data into Apache Parquet or ORC format and writes the data to a processed data location in Amazon S3. Have the Analysts query the data directly from Amazon S3 using Amazon Athena, and connect to BI tools using the Athena Java Database Connectivity (JDBC) connector.
- C.** Write each JSON record to a staging location in Amazon S3. Use the S3 Put event to trigger an AWS Lambda function that transforms the data into Apache Parquet or ORC format and inserts it into an Amazon RDS PostgreSQL database. Have the Analysts query and run dashboards from the RDS database.
- D.** Use Amazon Kinesis Data Analytics to ingest the streaming data and perform real-time SQL queries to convert the records to Apache Parquet before delivering to Amazon S3. Have the Analysts query the data directly from Amazon S3 using Amazon Athena and connect to BI tools using the Athena Java Database Connectivity (JDBC) connector.

**Answer:** A

**NO.100** A manufacturing company has a large set of labeled historical sales data. The manufacturer would like to predict how many units of a particular part should be produced each quarter. Which machine learning approach should be used to solve this problem?

- A.** Logistic regression
- B.** Random Cut Forest (RCF)
- C.** Principal component analysis (PCA)
- D.** Linear regression

**Answer:** B

**NO.101** An interactive online dictionary wants to add a widget that displays words used in similar contexts. A Machine Learning Specialist is asked to provide word features for the downstream nearest neighbor model powering the widget.

What should the Specialist do to meet these requirements?

- A.** Create one-hot word encoding vectors.
- B.** Produce a set of synonyms for every word using Amazon Mechanical Turk.
- C.** Create word embedding factors that store edit distance with every other word.
- D.** Download word embedding's pre-trained on a large corpus.

**Answer:** A

**NO.102** A Machine Learning Specialist works for a credit card processing company and needs to predict which transactions may be fraudulent in near-real time. Specifically, the Specialist must train a model that returns the probability that a given transaction may be fraudulent How should the Specialist frame this business problem'?

- A.** Streaming classification
- B.** Binary classification
- C.** Multi-category classification
- D.** Regression classification

**Answer:** A

**NO.103** A Machine Learning Specialist is developing recommendation engine for a photography blog Given a picture, the recommendation engine should show a picture that captures similar objects The Specialist would like to create a numerical representation feature to perform nearest-neighbor searches What actions would allow the Specialist to get relevant numerical representations?

- A.** Reduce image resolution and use reduced resolution pixel values as features
- B.** Use Amazon Mechanical Turk to label image content and create a one-hot representation indicating the presence of specific labels
- C.** Run images through a neural network pre-trained on ImageNet, and collect the feature vectors from the penultimate layer
- D.** Average colors by channel to obtain three-dimensional representations of images.

**Answer:** A

**NO.104** A retail company intends to use machine learning to categorize new products A labeled dataset of current products was provided to the Data Science team The dataset includes 1 200 products The labeled dataset has

15 features for each product such as title dimensions, weight, and price Each product is labeled as belonging to one of six categories such as books, games, electronics, and movies.

Which model should be used for categorizing new products using the provided dataset for training?

- A.** An XGBoost model where the objective parameter is set to multi: softmax
- B.** A deep convolutional neural network (CNN) with a softmax activation function for the last layer
- C.** A regression forest where the number of trees is set equal to the number of product categories
- D.** A DeepAR forecasting model based on a recurrent neural network (RNN)

**Answer:** B

**NO.105** A company is using Amazon Polly to translate plaintext documents to speech for automated company announcements However company acronyms are being mispronounced in the current documents How should a Machine Learning Specialist address this issue for future documents'?

- A.** Output speech marks to guide in pronunciation

- B.** Create an appropriate pronunciation lexicon.
- C.** Use Amazon Lex to preprocess the text files for pronunciation
- D.** Convert current documents to SSML with pronunciation tags

**Answer:** D

**NO.106** A Data Science team within a large company uses Amazon SageMaker notebooks to access data stored in Amazon S3 buckets. The IT Security team is concerned that internet-enabled notebook instances create a security vulnerability where malicious code running on the instances could compromise data privacy. The company mandates that all instances stay within a secured VPC with no internet access, and data communication traffic must stay within the AWS network.

How should the Data Science team configure the notebook instance placement to meet these requirements?

- A.** Associate the Amazon SageMaker notebook with a private subnet in a VPC. Place the Amazon SageMaker endpoint and S3 buckets within the same VPC.
- B.** Associate the Amazon SageMaker notebook with a private subnet in a VPC. Use IAM policies to grant access to Amazon S3 and Amazon SageMaker.
- C.** Associate the Amazon SageMaker notebook with a private subnet in a VPC. Ensure the VPC has S3 VPC endpoints and Amazon SageMaker VPC endpoints attached to it.
- D.** Associate the Amazon SageMaker notebook with a private subnet in a VPC. Ensure the VPC has a NAT gateway and an associated security group allowing only outbound connections to Amazon S3 and Amazon SageMaker

**Answer:** D