# CertyIQ

## Premium exam material

Get certification quickly with the CertyIQ Premium exam material.
Everything you need to prepare, learn & pass your certification exam easily. Lifetime free updates
First attempt guaranteed success.

https://www.CertyIQ.com

# About CertyIQ

We here at CertyIQ eventually got enough of the industry's greedy exam paid for. Our team of IT professionals comes with years of experience in the IT industry Prior to training CertIQ we worked in test areas where we observed the horrors of the paywall exam preparation system.

The misuse of the preparation system has left our team disillusioned. And for that reason, we decided it was time to make a difference. We had to make In this way, CertyIQ was created to provide quality materials without stealing from everyday people who are trying to make a living.

# Doubt Support

We have developed a very scalable solution using which we are able to solve 400+ doubts every single day with an average rating of 4.8 out of 5.

https://www.certyiq.com

Mail us on - certyiqofficial@gmail.com

### Lifetime Free Updates
We provide lifetime free updates to our customers. To make life easier for our valued customers and fulfill their needs

### Free Exam PDF
You are sure to pass the exam completely free of charge

### Money Back Guarantee
We Provide 100% money back guarantee to our customer in case of any failure

---

**John**

October 19, 2022

★★★★★

Thanks you so much for your help. I scored 972 in my exam today. More than 90% were from your PDFs!

**Dana**

September 04, 2022

★★★★★

Thanks a lot for this updated AZ-900 Q&A. I just passed my exam and got 974, I followed both of your Az-900 videos and the 6 PDF, the PDFs are very much valid, all answers are correct. Could you please create a similar video/PDF for DP900, your content/PDF's is really awesome. The team did a really good job. Thank You 😊.

**Ahamed Shibly**

2 months ago

★★★★★

Customer support is realy fast and helpful, I just finished my exam and this video along with the 6 PDF helped me pass! Definitely recommend getting the PDFs. Thank you!

---

October 22, 2022

★★★★★

Passed my exam today with 891 marks. Out of 52 questions, 51 were from certyiq PDFs including Contoso case study. Thank You certyiq team!

**Henry Rome**

2 months ago

★★★★★

These questions are real and 100 % valid. Thank you so much for your efforts, also your 4 PDFs are awesome, I passed the DP900 exam on 1 Sept. With 968 marks. Thanks a lot, buddy!

**Esmaria**

2 months ago

★★★★★

Simple easy to understand explanations. To anyone out there wanting to write AZ900, I highly recommend 6 PDF's.Thank you so much, appreciate all your hard work in having such great content. Passed my exam Today - 3 September with 942 score.

# Amazon

(AWS Certified AI Practitioner AIF-C01)

AWS Certified AI Practitioner AIF-C01

## Question: 1

A company makes forecasts each quarter to decide how to optimize operations to meet expected demand. The company uses ML models to make these forecasts.
An AI practitioner is writing a report about the trained ML models to provide transparency and explainability to company stakeholders.
What should the AI practitioner include in the report to meet the transparency and explainability requirements?

    A.Code for model training

    B.Partial dependence plots (PDPs)

    C.Sample data for training

    D.Model convergence tables

**Answer: B**

**Explanation:**

The correct answer is **B. Partial dependence plots (PDPs)**. Here's why:

Transparency and explainability in AI/ML refer to the ability to understand how a model makes its predictions. This is crucial for trust, accountability, and debugging. While code, sample data, and model convergence tables are valuable for model development and understanding its performance during training, they don't directly address why the model makes a particular prediction for a given input.

Partial Dependence Plots (PDPs) directly contribute to explainability. A PDP visualizes the marginal effect of one or two features on the predicted outcome of a machine learning model. It shows how the model's prediction changes as the selected feature(s) vary, while holding all other features constant (in a sense, averaging out their effects). By showing these relationships, stakeholders can gain insight into which features are most influential in driving the model's forecasts and in what direction.

Code (A) is important for reproducibility, but it doesn't directly explain why the model is making certain forecasts. Sample training data (C) is useful for verifying the data quality and distribution, but it doesn't reveal the model's internal logic. Model convergence tables (D) show how well the model learned during training, but don't illustrate feature importance or impact on predictions. Only PDPs effectively showcase the feature-outcome relationships learned by the model, enhancing understanding for stakeholders. For example, a PDP might show that as advertising spend increases, the predicted demand also increases, but the effect plateaus at a certain point, giving actionable insights for decision-making.

Therefore, PDPs are the most relevant tool for meeting the transparency and explainability requirements in the context of explaining ML model forecasts to company stakeholders. They enable stakeholders to understand the model's reasoning and build trust in its predictions.

Relevant links:

Interpretable Machine Learning by Christoph Molnar: https://christophm.github.io/interpretable-ml-book/pdp.html
scikit-learn PDP documentation: https://scikit-learn.org/stable/modules/partial_dependence.html

## Question: 2

A law firm wants to build an AI application by using large language models (LLMs). The application will read legal documents and extract key points from the documents.
Which solution meets these requirements?

    A.Build an automatic named entity recognition system.

    B.Create a recommendation engine.

C.Develop a summarization chatbot.

D.Develop a multi-language translation system.

**Answer: C**

**Explanation:**

The correct answer is **C. Develop a summarization chatbot.**

Here's why: The law firm needs a solution that can process legal documents and extract key information. A summarization chatbot is designed to ingest text and produce a condensed, coherent summary highlighting the main points. LLMs are particularly effective for text summarization due to their ability to understand context, identify crucial information, and generate human-quality summaries. A chatbot interface allows users to interact with the LLM, specify documents, and receive the extracted key points in a conversational manner.

Option A, building an automatic named entity recognition (NER) system, while helpful for identifying entities like names, organizations, and dates, doesn't inherently summarize or extract key points. It only identifies and classifies pre-defined entity types. Option B, creating a recommendation engine, is irrelevant as it focuses on suggesting items based on user preferences or historical data. Option D, developing a multi-language translation system, addresses language translation and doesn't extract or summarize content.

Therefore, a summarization chatbot leverages LLMs to best satisfy the requirement of reading legal documents and extracting key points, providing a user-friendly interface for accessing the summaries.

For more information on text summarization using LLMs:

**Amazon SageMaker JumpStart text summarization:** https://aws.amazon.com/sagemaker/jumpstart/ (search for summarization)

**Generative AI on AWS:** https://aws.amazon.com/machine-learning/generative-ai/

---

**Question: 3**                                                              **CertyIQ**

A company wants to classify human genes into 20 categories based on gene characteristics. The company needs an ML algorithm to document how the inner mechanism of the model affects the output.
Which ML algorithm meets these requirements?

A.Decision trees

B.Linear regression

C.Logistic regression

D.Neural networks

**Answer: A**

**Explanation:**

The correct answer is A, Decision Trees. Here's why:

Decision Trees are highly interpretable machine learning algorithms. Their structure mimics a flowchart, with each internal node representing a test on an attribute (gene characteristic in this case), each branch representing the outcome of the test, and each leaf node representing a class label (one of the 20 gene categories). This clear, hierarchical structure allows users to easily trace the path from input features to the final classification. You can see exactly which features were used to make a decision, and in what order.

The key requirement is understanding the model's inner workings. Linear regression and logistic regression, while interpretable to some extent (coefficients show feature importance), don't provide the same level of granular, step-by-step explanation as decision trees. Neural networks, especially deep neural networks, are

notoriously "black boxes." It's difficult to understand exactly why a neural network made a specific prediction.

Decision Trees are also suitable for multi-class classification problems like classifying genes into 20 categories. While they might not always be the most accurate model (prone to overfitting), their interpretability makes them ideal when understanding the why behind the prediction is crucial.

Here's a breakdown of why the other options are less suitable:

**B. Linear Regression:** Primarily for predicting continuous values, not suitable for multi-class classification. While coefficients indicate feature importance, it doesn't detail decision paths.
**C. Logistic Regression:** Primarily for binary classification. While multi-class extensions exist, they still lack the detailed, traceable paths of decision trees.
**D. Neural Networks:** Extremely difficult to interpret due to their complex, non-linear structure. Techniques like LIME and SHAP can help with post-hoc interpretability, but don't inherently reveal the internal decision-making process.

Decision tree's ease of visualization and inherent traceability of the path taken through feature evaluation makes it a powerful tool when you need to understand and document the reasoning behind the ML model.

**Authoritative Links:**

**Decision Trees:** (Scikit-learn documentation - a popular Python ML library) - https://scikit-learn.org/stable/modules/tree.html
**Interpretability of Machine Learning:** (Google AI Blog) - https://ai.googleblog.com/2023/03/interpretability-beyond-feature.html (While this blog is a broader topic, it provides context on the value of interpretable models.)

---

**Question: 4**                                                           **CertyIQ**

A company has built an image classification model to predict plant diseases from photos of plant leaves. The company wants to evaluate how many images the model classified correctly.
Which evaluation metric should the company use to measure the model's performance?

   A.R-squared score

   B.Accuracy

   C.Root mean squared error (RMSE)

   D.Learning rate

**Answer: B**

**Explanation:**

Here's a detailed justification for why Accuracy is the most appropriate evaluation metric in this scenario:

The company needs to determine how often its image classification model correctly identifies plant diseases in leaf photos. Image classification is a specific type of supervised machine learning problem where the model predicts a category or class label for each input image.

Accuracy is the ratio of correctly classified instances to the total number of instances. In this case, it represents the proportion of plant leaf images that the model correctly classified as having a specific disease (or not). A higher accuracy score indicates better performance.

Here's why the other options are not suitable:

R-squared score: This metric is used to evaluate the performance of regression models, which predict continuous values, not classification models.

Root Mean Squared Error (RMSE): Similar to R-squared, RMSE is a regression metric, measuring the average magnitude of errors between predicted and actual continuous values.

Learning rate: This is a hyperparameter that controls how much the model's weights are adjusted during training. It doesn't measure the model's performance after training.

Accuracy provides a straightforward and easily interpretable measure of how well the model is performing at its primary task: correctly classifying plant diseases based on images. This makes it the best choice for the company's evaluation needs.

Therefore, the correct answer is B.

Further reading:

AWS Machine Learning documentation: https://aws.amazon.com/machine-learning/
Classification Accuracy: https://developers.google.com/machine-learning/crash-course/classification/accuracy

---

**Question: 5**                                                              **CertyIQ**

A company is using a pre-trained large language model (LLM) to build a chatbot for product recommendations. The company needs the LLM outputs to be short and written in a specific language.
Which solution will align the LLM response quality with the company's expectations?

　　A.Adjust the prompt.

　　B.Choose an LLM of a different size.

　　C.Increase the temperature.

　　D.Increase the Top K value.

**Answer: A**

**Explanation:**

The correct answer is A: Adjust the prompt. Here's a detailed justification:

The primary goal is to influence the LLM's output to be shorter and in a specific language. Prompt engineering is the most direct and efficient way to achieve this. A well-crafted prompt acts as a precise instruction set, guiding the LLM to generate the desired response. By including specific instructions within the prompt (e.g., "Answer in less than 20 words, in Spanish"), the company can directly shape the output to meet its requirements.

Option B, choosing an LLM of a different size, is a less targeted approach. While different LLMs might have varying tendencies, switching models introduces significant complexity and potential disruption without directly addressing the specific needs of output length and language.

Option C, increasing the temperature, introduces more randomness and creativity into the output. This would likely make the output less predictable and harder to control in terms of length and language adherence, which is the opposite of the desired effect. Temperature controls the randomness of the token selection; a higher temperature leads to more diverse and potentially less coherent responses.

Option D, increasing the Top K value, expands the pool of potential tokens the LLM considers before making a prediction. While it can lead to more diverse outputs, it doesn't directly enforce length constraints or language specifications. Like temperature, increasing Top K primarily affects diversity, not specific style or formatting.

Therefore, prompt engineering offers the most granular and effective control over LLM output. Specific prompts can be designed to directly address the needs of producing short, language-specific responses,

making it the best solution in this scenario. Techniques such as few-shot prompting (providing examples of desired output) can further refine the LLM's responses.

For further research:

**Prompt Engineering Guide:** https://www.promptingguide.ai/
**Google AI - Prompt Design:** https://developers.google.com/machine-learning/prompt-design (This link might be outdated, but searches on "Google AI Prompt Design" will show the current similar pages.)

These resources will provide comprehensive information on crafting effective prompts for LLMs.

---

**Question: 6**                                                                                             **CertyIQ**

A company uses Amazon SageMaker for its ML pipeline in a production environment. The company has large input data sizes up to 1 GB and processing times up to 1 hour. The company needs near real-time latency.
Which SageMaker inference option meets these requirements?

    A.Real-time inference

    B.Serverless inference

    C.Asynchronous inference

    D.Batch transform

**Answer: C**

**Explanation:**

Here's a detailed justification for why Asynchronous Inference (Option C) is the most suitable SageMaker inference option for the company's requirements:

The company deals with input data sizes up to 1 GB and processing times extending to 1 hour while demanding near real-time latency. This combination presents a unique challenge.

**Real-time Inference (Option A):** While designed for low latency, real-time inference typically handles smaller payloads and shorter processing times. One-hour processing times are highly unusual and often impractical for real-time endpoints. This could lead to timeouts and degraded performance as the endpoint needs to stay active for a prolonged duration.

**Serverless Inference (Option B):** Serverless Inference automatically scales resources based on incoming request traffic. While cost-effective for sporadic traffic, the 1 GB input size and 1-hour processing time could potentially cause issues with Lambda function execution limits, if Serverless Inference relies on Lambda under the hood. Furthermore, cold starts associated with serverless deployments may introduce unacceptable latency for near real-time needs, especially for such lengthy inference times.

**Asynchronous Inference (Option C):** Asynchronous Inference is explicitly designed for handling large payloads and long processing times without blocking the client. It decouples the request from the response, allowing the client to submit a request and retrieve the results later. This makes it ideal for scenarios where immediate responses are not crucial, but timely delivery of results is. The near real-time requirement can be met if the processing time is generally within the acceptable range, and retrieval mechanism is efficient. SageMaker manages queuing the requests, executing the model, and storing the output.

**Batch Transform (Option D):** Batch Transform is intended for offline inference on large datasets. It's not suitable for near real-time latency requirements as the data is processed in batches, and there would be a significant delay.

Therefore, Asynchronous Inference is the best fit because it supports large input sizes and long processing

times, allowing the company to meet the near real-time latency requirement by decoupling the inference request and enabling a timely output retrieval. The 1-hour processing time is a key factor favoring asynchronous over real-time solutions.

**Supporting Resources:**

**AWS SageMaker Asynchronous Inference:** https://docs.aws.amazon.com/sagemaker/latest/dg/async-inference.html
**AWS SageMaker Inference Options:** https://docs.aws.amazon.com/sagemaker/latest/dg/deploy-model.html

---

## Question: 7                                                                 CertyIQ

A company is using domain-specific models. The company wants to avoid creating new models from the beginning. The company instead wants to adapt pre-trained models to create models for new, related tasks.
Which ML strategy meets these requirements?

    A.Increase the number of epochs.

    B.Use transfer learning.

    C.Decrease the number of epochs.

    D.Use unsupervised learning.

**Answer: B**

**Explanation:**

The correct answer is **B. Use transfer learning.**

Transfer learning is a machine learning technique where a model developed for a task is reused as the starting point for a model on a second task. In this scenario, the company already possesses domain-specific models. Transfer learning allows them to leverage the knowledge these pre-trained models have acquired, rather than building entirely new models from scratch. By adapting these existing models to new, related tasks, the company can significantly reduce training time, computational resources, and the amount of new labeled data required. This is because the pre-trained model already understands relevant features and patterns from its initial training. The pre-trained model's weights serve as a valuable initialization for the new task, accelerating convergence and improving performance. Techniques like fine-tuning, feature extraction, or a combination of both can be employed to adapt the pre-trained model.

Option A, increasing the number of epochs, only extends the training of the same model, not leveraging knowledge from existing models. Option C, decreasing the number of epochs, would likely result in underfitting. Option D, unsupervised learning, is generally used to discover patterns in unlabeled data, not to adapt existing models for new tasks.

Further Research:

**AWS Documentation on Transfer Learning:** https://aws.amazon.com/machine-learning/transfer-learning/
(While not specific to AWS, it is a general ML concept AWS utilizes.)
**Transfer Learning (Stanford CS231n):** http://cs231n.github.io/transfer-learning/

---

## Question: 8                                                                 CertyIQ

A company is building a solution to generate images for protective eyewear. The solution must have high accuracy and must minimize the risk of incorrect annotations.
Which solution will meet these requirements?

A.Human-in-the-loop validation by using Amazon SageMaker Ground Truth Plus

B.Data augmentation by using an Amazon Bedrock knowledge base

C.Image recognition by using Amazon Rekognition

D.Data summarization by using Amazon QuickSight Q

**Answer: A**

**Explanation:**

Here's a detailed justification for why option A is the correct answer:

The core requirement is high accuracy and minimizing incorrect annotations for image generation related to protective eyewear. This points to a need for meticulous review and correction of the generated images and/or their training data.

**Amazon SageMaker Ground Truth Plus (A):** This service provides a managed workforce to perform data labeling and validation. Its human-in-the-loop (HITL) capabilities directly address the accuracy requirement. Trained workers can review generated images, correct annotations, and ensure the generated images align with the desired outcome. This minimizes the risk of feeding flawed data into the image generation model, consequently improving the overall quality and accuracy of the generated eyewear images. Ground Truth Plus's managed workforce takes care of sourcing and managing annotators, and includes project management tools to ensure high-quality labels are delivered.

**Amazon Bedrock knowledge base (B):** Data augmentation through a knowledge base in Bedrock primarily focuses on expanding the dataset using existing knowledge. While it might improve the model's robustness, it doesn't guarantee the correction of incorrect annotations in the first place. It's better for improving a model's understanding of concepts than for ensuring initial annotation accuracy.

**Amazon Rekognition (C):** Rekognition is an image recognition service. While helpful for object detection after image generation, it doesn't directly address the initial annotation or image generation quality control needed to minimize errors during the process.

**Amazon QuickSight (D):** QuickSight is a data visualization and business intelligence tool. It provides insights from data, but it doesn't contribute to image generation accuracy or annotation correction.

Therefore, the optimal choice is SageMaker Ground Truth Plus because it uniquely offers human validation to minimize the risk of incorrect annotations during image generation and enhance accuracy.

**Further Research:**

Amazon SageMaker Ground Truth Plus: https://aws.amazon.com/sagemaker/groundtruth/
Amazon Bedrock: https://aws.amazon.com/bedrock/
Amazon Rekognition: https://aws.amazon.com/rekognition/
Amazon QuickSight: https://aws.amazon.com/quicksight/

**Question: 9**                                                        **CertyIQ**

A company wants to create a chatbot by using a foundation model (FM) on Amazon Bedrock. The FM needs to access encrypted data that is stored in an Amazon S3 bucket. The data is encrypted with Amazon S3 managed keys (SSE-S3).
The FM encounters a failure when attempting to access the S3 bucket data.
Which solution will meet these requirements?

A.Ensure that the role that Amazon Bedrock assumes has permission to decrypt data with the correct encryption key.

B.Set the access permissions for the S3 buckets to allow public access to enable access over the internet.

C.Use prompt engineering techniques to tell the model to look for information in Amazon S3.

D.Ensure that the S3 data does not contain sensitive information.

**Answer: A**

**Explanation:**

The correct answer is A. Here's a detailed justification:

Amazon Bedrock needs specific permissions to access and process data within an S3 bucket, especially when that data is encrypted. Since the data is encrypted with SSE-S3 (Amazon S3 managed keys), Bedrock must have the necessary authorization to decrypt the data before it can be used by the foundation model for chatbot functionalities. This authorization is granted through an IAM role that Bedrock assumes when interacting with other AWS services, including S3.

Option A directly addresses the root cause of the failure: lack of decryption permissions. By ensuring the IAM role assigned to Bedrock has the s3:GetObject permission on the S3 bucket and the necessary permission to decrypt the data using the SSE-S3 key, the foundation model can successfully access and utilize the encrypted data. Specifically, the role needs kms:Decrypt permissions if using KMS managed keys, but SSE-S3 keys don't require explicit KMS permissions.

Option B is incorrect and a security risk. Granting public access to the S3 bucket exposes sensitive data to the internet, violating the principle of least privilege and potentially leading to data breaches. It's a completely unacceptable solution for handling encrypted data.

Option C is also incorrect. Prompt engineering is a technique used to refine the instructions given to a foundation model to influence its output. It does not grant the model the necessary permissions to access and decrypt data in S3. The model needs authorization to access the data before prompt engineering comes into play.

Option D is related to data governance but does not solve the immediate problem of the foundation model failing to access the encrypted data. Even if the data doesn't contain sensitive information, the model still needs the required permissions to access and decrypt it. Furthermore, simply removing sensitive data doesn't address the core requirement of using encrypted data for the chatbot. The data is encrypted for security, and removing the sensitive information does not obviate the need for encryption.

In conclusion, the only way to fix the failure is to ensure Bedrock's IAM role has the appropriate s3:GetObject permissions on the S3 bucket and access to decrypt the data using the key associated with SSE-S3. This adheres to best practices for secure access to AWS resources.

References:

**IAM Roles:** https://docs.aws.amazon.com/IAM/latest/UserGuide/id_roles.html
**Amazon S3 Encryption:** https://docs.aws.amazon.com/AmazonS3/latest/userguide/serv-side-encryption.html
**AWS Bedrock:** https://aws.amazon.com/bedrock/

**Question: 10**                                                                    **CertyIQ**

A company wants to use language models to create an application for inference on edge devices. The inference must have the lowest latency possible.
Which solution will meet these requirements?

A.Deploy optimized small language models (SLMs) on edge devices.

B.Deploy optimized large language models (LLMs) on edge devices.

C.Incorporate a centralized small language model (SLM) API for asynchronous communication with edge

devices.

D.Incorporate a centralized large language model (LLM) API for asynchronous communication with edge devices.

**Answer: A**

**Explanation:**

The correct answer is A: Deploy optimized small language models (SLMs) on edge devices.

Here's why: The primary requirement is the lowest possible latency for inference. Latency refers to the delay between a request and the response. Edge devices process data locally, reducing network hops and round-trip times compared to centralized servers.

Options C and D involve centralized APIs, meaning data must travel to the server and back, adding significant latency due to network communication. Asynchronous communication, while improving throughput in some cases, does not directly address the need for lowest latency on a single inference request.

LLMs (Large Language Models), while powerful, are computationally intensive and require more resources to run. Deploying them on edge devices would likely result in higher latency due to the processing constraints of edge hardware.

SLMs (Small Language Models) are designed to be more efficient and can be optimized for edge deployment. They require less computational power and memory, enabling faster inference on edge devices. By optimizing the SLM, the company can further reduce latency. Deploying directly on the edge eliminates network latency, fulfilling the requirement of the lowest possible latency.

Therefore, deploying optimized SLMs on edge devices is the most suitable approach for minimizing inference latency in this scenario. This approach aligns with the principles of edge computing, where data processing is moved closer to the source to reduce latency and bandwidth consumption.

For further research on Edge Computing and model deployment:

**AWS Documentation on Edge Computing:** https://aws.amazon.com/edge/
**AWS SageMaker Edge Manager:** https://aws.amazon.com/sagemaker/edge-manager/
**Paperspace Blog on SLMs vs. LLMs:** (search on Google for "paperspace small language model vs large language model" - there isn't a definitive direct AWS link for this concept but many vendor sites will have helpful comparisons)

**Question: 11**                                                          **CertyIQ**

A company wants to build an ML model by using Amazon SageMaker. The company needs to share and manage variables for model development across multiple teams.
Which SageMaker feature meets these requirements?

A.Amazon SageMaker Feature Store

B.Amazon SageMaker Data Wrangler

C.Amazon SageMaker Clarify

D.Amazon SageMaker Model Cards

**Answer: A**

**Explanation:**

The correct answer is **A. Amazon SageMaker Feature Store**. Here's why:

SageMaker Feature Store is a fully managed, centralized repository for storing, managing, and sharing ML features. It enables organizations to define, store, and retrieve features in a consistent and scalable manner, making them readily available for model training, inference, and feature exploration. For a company looking to share and manage variables (features) across multiple teams for model development in SageMaker, Feature Store provides the ideal solution.

**Centralized Repository:** Feature Store provides a central place to define and store features, preventing feature duplication and ensuring consistency across teams.
**Collaboration:** Multiple teams can access and use the same features, fostering collaboration and reducing redundant feature engineering efforts.
**Versioning:** Feature Store often supports versioning, allowing teams to track changes to features and ensure reproducibility of models.
**Discoverability:** Feature Store provides a mechanism to discover and understand available features, making it easier for teams to find the right features for their models.
**Scalability:** Feature Store is designed to handle large volumes of feature data and can scale to meet the needs of enterprise-level ML projects.

Alternatives are not suited:

**B. Amazon SageMaker Data Wrangler:** primarily focuses on data preparation and feature engineering; it's not designed for sharing and managing features across teams in the same way as Feature Store.
**C. Amazon SageMaker Clarify:** is used for bias detection and explainability of ML models, not for feature sharing.
**D. Amazon SageMaker Model Cards:** is used to document and track information about ML models, but it doesn't manage or share features.

In summary, SageMaker Feature Store directly addresses the requirement of sharing and managing variables (features) for model development across multiple teams within the SageMaker environment. It facilitates collaboration, ensures consistency, and simplifies feature management, making it the most suitable choice.

Relevant links:

**Amazon SageMaker Feature Store:** https://aws.amazon.com/sagemaker/feature-store/
**SageMaker Feature Store documentation:** https://docs.aws.amazon.com/sagemaker/latest/dg/feature-store.html

---

**Question: 12**

A company wants to use generative AI to increase developer productivity and software development. The company wants to use Amazon Q Developer.
What can Amazon Q Developer do to help the company meet these requirements?

  A.Create software snippets, reference tracking, and open source license tracking.

  B.Run an application without provisioning or managing servers.

  C.Enable voice commands for coding and providing natural language search.

  D.Convert audio files to text documents by using ML models.

**Answer: A**

**Explanation:**

The most suitable answer is A because Amazon Q Developer is specifically designed to boost developer productivity by directly assisting with coding tasks. Option A accurately reflects this by listing key functionalities like creating software snippets, which accelerates code development by providing ready-to-

use code blocks. Reference tracking within Amazon Q helps developers understand the origin and context of code elements, improving maintainability and reducing errors. Open-source license tracking is crucial for compliance and avoiding legal issues when using open-source components. These features directly address the company's goal of enhancing developer productivity and software development.

Option B describes serverless computing, which is a characteristic of services like AWS Lambda, and while beneficial for certain applications, it doesn't directly contribute to enhanced coding assistance provided by Amazon Q Developer.

Option C, enabling voice commands for coding, is an interesting capability but is not the core focus or primary function advertised for Amazon Q Developer. While natural language search might be incorporated, the main thrust is on code generation and understanding.

Option D pertains to audio transcription, a function handled by services like Amazon Transcribe, and is unrelated to the core objective of increasing developer productivity through code assistance.

Therefore, given the company's need to increase developer productivity and software development using generative AI, the most relevant answer is A.For more information, research Amazon Q Developer on the AWS website: https://aws.amazon.com/q/developer/

## Question: 13 <span>CertyIQ</span>

A financial institution is using Amazon Bedrock to develop an AI application. The application is hosted in a VPC. To meet regulatory compliance standards, the VPC is not allowed access to any internet traffic.
Which AWS service or feature will meet these requirements?

    A.AWS PrivateLink
    B.Amazon Macie
    C.Amazon CloudFront
    D.Internet gateway

**Answer: A**

**Explanation:**

The correct answer is **A. AWS PrivateLink**. Here's a detailed justification:

AWS PrivateLink enables you to access AWS services and services hosted by other AWS accounts (referred to as endpoint services) in a private and secure manner, without exposing your traffic to the public internet. This is achieved by establishing private connectivity between your VPC and the service using Elastic Network Interfaces (ENIs) within your VPC.

In the scenario described, the financial institution requires a secure connection to Amazon Bedrock within a VPC that doesn't allow internet access due to regulatory compliance. AWS PrivateLink directly addresses this requirement. By creating a VPC endpoint for Amazon Bedrock powered by PrivateLink, the application within the VPC can privately access Bedrock's API without traversing the internet. This connection is isolated within the AWS network.

Let's examine why the other options are incorrect:

**B. Amazon Macie:** Macie is a data security and data privacy service that uses machine learning and pattern matching to discover and protect sensitive data in AWS. It doesn't establish private connectivity.

**C. Amazon CloudFront:** CloudFront is a content delivery network (CDN) used to distribute content with low latency and high transfer speeds. It typically involves internet access and isn't suitable for completely

isolating traffic within a VPC.

**D. Internet Gateway:** An Internet Gateway allows resources within a VPC to access the internet. This directly contradicts the requirement of no internet access.

In summary, AWS PrivateLink provides the necessary private connectivity to Amazon Bedrock from within the restricted VPC, fulfilling the regulatory compliance requirements of the financial institution.

**Authoritative Links:**

AWS PrivateLink Documentation
Amazon Bedrock Documentation

---

**Question: 14**                                                                                    **CertyIQ**

A company wants to develop an educational game where users answer questions such as the following: "A jar contains six red, four green, and three yellow marbles. What is the probability of choosing a green marble from the jar?"
Which solution meets these requirements with the LEAST operational overhead?

   A.Use supervised learning to create a regression model that will predict probability.

   B.Use reinforcement learning to train a model to return the probability.

   C.Use code that will calculate probability by using simple rules and computations.

   D.Use unsupervised learning to create a model that will estimate probability density.

**Answer: C**

**Explanation:**

The correct answer is **C. Use code that will calculate probability by using simple rules and computations.**

**Justification:**

The problem describes a scenario where probability calculation is straightforward and based on well-defined mathematical formulas. The question explicitly involves counting and applying basic probability rules (number of favorable outcomes divided by total number of outcomes). Implementing this calculation through code (e.g., Python) directly addresses the problem efficiently and accurately.

**Operational Overhead:** Using code for direct calculation introduces the least operational overhead. It requires minimal infrastructure, no model training or deployment, and is computationally inexpensive.
**Supervised Learning (Option A):** Supervised learning would require a large dataset of questions and correct probabilities to train a regression model. This introduces significant overhead for data collection, labeling, model training, and deployment. It is an over-engineered solution for a simple problem.
**Reinforcement Learning (Option B):** Reinforcement learning would involve training an agent to answer probability questions. This is highly inappropriate and inefficient for a problem with a deterministic solution. It's complex to implement and requires extensive training and reward engineering.
**Unsupervised Learning (Option D):** Unsupervised learning techniques like probability density estimation are not relevant here. The task is to compute a specific probability, not to understand the underlying distribution of probability values.

Since there is already a formula that calculates probability, utilizing it would lead to the least overhead. Using machine learning would increase complexity while providing no extra benefit.

In cloud computing, minimizing operational overhead is a key design principle. By using a simple code-based solution, the company reduces infrastructure costs, maintenance efforts, and overall complexity while still meeting the requirements of the educational game.

Here are some authoritative links for further research:

**Basic Probability:** https://www.mathsisfun.com/data/probability.html
**Calculating Probability:** https://www.khanacademy.org/math/statistics-probability/counting-permutations-and-combinations

## Question: 15

Which metric measures the runtime efficiency of operating AI models?

    A.Customer satisfaction score (CSAT)

    B.Training time for each epoch

    C.Average response time

    D.Number of training instances

**Answer: C**

**Explanation:**

Here's a detailed justification for why average response time is the best metric for measuring the runtime efficiency of operating AI models:

Average response time directly reflects how quickly an AI model provides a prediction or output in a real-world application. It's a crucial indicator of the model's operational performance from a user's perspective. A shorter average response time implies that the model is processing requests efficiently, leading to a better user experience. This is especially vital in latency-sensitive applications like real-time recommendations, fraud detection, or conversational AI, where delays can negatively impact usability and effectiveness.

While training time (option B) is important during model development, it doesn't directly measure runtime efficiency. The training phase focuses on learning patterns from data, while runtime refers to how quickly the model provides predictions after it has been deployed. Customer satisfaction (CSAT), option A, is a broad measure of user experience but can be affected by numerous factors beyond the model's runtime performance. Finally, option D, "Number of training instances," pertains to the dataset used for training and is not a metric for runtime efficiency.

The principle of minimizing latency is fundamental in cloud computing and AI deployment. Efficient models contribute to lower operational costs by consuming fewer resources and minimizing the need for scaling infrastructure. Average response time allows DevOps and MLOps engineers to monitor model performance, identify bottlenecks, and optimize resource allocation. Slow response times could indicate issues such as inefficient model code, insufficient compute resources, or network latency problems.

In cloud environments like AWS, services such as Amazon CloudWatch can be used to monitor and alarm on average response time metrics for deployed AI models. Techniques like model optimization, caching, and the selection of appropriate instance types can be employed to improve response times.

Further research on model deployment and performance monitoring on AWS can be found at the following links:

**AWS Documentation on Monitoring Machine Learning Models:**
https://docs.aws.amazon.com/sagemaker/latest/dg/model-monitor.html
**AWS Documentation on CloudWatch:** https://docs.aws.amazon.com/cloudwatch/

## Question: 16

A company is building a contact center application and wants to gain insights from customer conversations. The company wants to analyze and extract key information from the audio of the customer calls.
Which solution meets these requirements?

A.Build a conversational chatbot by using Amazon Lex.

B.Transcribe call recordings by using Amazon Transcribe.

C.Extract information from call recordings by using Amazon SageMaker Model Monitor.

D.Create classification labels by using Amazon Comprehend.

**Answer: B**

**Explanation:**

The correct answer is **B. Transcribe call recordings by using Amazon Transcribe.**

Here's a detailed justification:

The core requirement is to gain insights and extract key information from customer conversation audio. Amazon Transcribe directly addresses this by converting audio into text. This transcription process allows for subsequent analysis to identify key phrases, sentiment, topics discussed, and other relevant data points.

Option A (Amazon Lex) is more suitable for building conversational interfaces, not directly analyzing existing audio recordings. While Lex could be integrated later to process the transcribed text, it's not the initial step to extract information from the audio itself.

Option C (Amazon SageMaker Model Monitor) is for monitoring the performance of machine learning models, not for directly transcribing or analyzing audio. It's a downstream tool that might be relevant after analysis, but not for the primary task of extracting information from the calls.

Option D (Amazon Comprehend) performs natural language processing (NLP) tasks like sentiment analysis and entity recognition. Comprehend is useful, but it requires text as input. It cannot directly process audio; it needs transcribed text from a service like Transcribe.

Therefore, Amazon Transcribe is the most logical first step as it bridges the gap between audio data and text-based analysis tools like Comprehend. By transcribing the call recordings, the company will then be able to leverage NLP tools or even manual analysis to glean the desired insights from the conversations.

Essentially, Transcribe provides the textual data necessary for subsequent analysis to uncover valuable information from the calls.

Authoritative Links:

Amazon Transcribe: https://aws.amazon.com/transcribe/
Amazon Comprehend: https://aws.amazon.com/comprehend/

---

## Question: 17                                                    CertyIQ

A company has petabytes of unlabeled customer data to use for an advertisement campaign. The company wants to classify its customers into tiers to advertise and promote the company's products.
Which methodology should the company use to meet these requirements?

A.Supervised learning

B.Unsupervised learning

C.Reinforcement learning

D.Reinforcement learning from human feedback (RLHF)

**Answer: B**

**Explanation:**

The company needs to classify its customers into tiers using petabytes of unlabeled data. This is a classic clustering problem, where the goal is to group similar data points together without any prior knowledge of the correct labels.

Supervised learning (Option A) requires labeled data to train a model. Since the data is unlabeled, supervised learning is not applicable. Supervised learning algorithms learn a mapping function from input features to output labels using labeled training data.

Unsupervised learning (Option B) is ideal for this scenario. Unsupervised learning algorithms discover patterns and structure in unlabeled data. Clustering algorithms, a type of unsupervised learning, can automatically group customers into tiers based on their inherent similarities, without requiring pre-defined labels. Common unsupervised learning methods include k-means clustering, hierarchical clustering, and DBSCAN. The company can use algorithms like k-means to automatically segment the customers based on their characteristics derived from the data.

Reinforcement learning (Option C) involves training an agent to make decisions in an environment to maximize a reward. It's not relevant for classifying data into tiers. Reinforcement learning is typically applied when an agent needs to learn optimal actions through trial and error in an environment.

Reinforcement learning from human feedback (RLHF) (Option D) is an advanced technique within reinforcement learning where human feedback is used to guide the agent's learning process. It's irrelevant because the task does not involve an agent and rewards.

Therefore, unsupervised learning (Option B) is the correct methodology because it allows the company to discover inherent groupings within their unlabeled customer data and classify them into tiers, thereby meeting the requirements for their advertisement campaign.

Further Reading:

AWS AI Services: https://aws.amazon.com/ai/
Unsupervised Learning on AWS: https://aws.amazon.com/blogs/machine-learning/performing-a-k-means-clustering-analysis-with-amazon-athena/

---

**Question: 18**                                                       **CertyIQ**

An AI practitioner wants to use a foundation model (FM) to design a search application. The search application must handle queries that have text and images.
Which type of FM should the AI practitioner use to power the search application?

   A.Multi-modal embedding model

   B.Text embedding model

   C.Multi-modal generation model

   D.Image generation model

**Answer: A**

**Explanation:**

The AI practitioner needs a model that can understand and compare both text and images to power the search application. A multi-modal embedding model is the most appropriate choice.

Here's why:

**Multi-modal:** The application requires handling two different modalities: text and images. A multi-modal model is designed to process and relate information from different modalities.

**Embedding:** Embedding models create vector representations of the input data (text and images in this case). These vector representations capture the semantic meaning of the data.

**Search Application:** The core of a search application is comparing queries to the indexed content. Embedding models allow for efficient similarity search in a vector space. The text and image queries can be converted into embeddings and compared against embeddings of the indexed data.

**Comparison:** By embedding both text and images into a common vector space, the model allows for the comparison of text queries to images and vice versa. The system can then retrieve the most relevant content based on the similarity of their embeddings.

**Alternatives are unsuitable:** Text embedding models only handle text, and image generation models create new images rather than find relevant ones. Multi-modal generation models create new content from different inputs, whereas, the search application requires comparing existing content.

In contrast, a text embedding model would only handle text queries, an image generation model would create new images rather than find relevant ones, and a multi-modal generation model would create new text or images based on the inputs rather than identify similarity and relevance.

**Further Reading:**

**Amazon Bedrock Multi-Modal Embeddings:** https://aws.amazon.com/bedrock/multi-modal-embeddings/
(This link talks about multi-modal capabilities in Amazon Bedrock, illustrating the concept).
**Multi-Modal Learning:** https://www.cs.cmu.edu/~mmv/ (Carnegie Mellon University's research on multi-modal learning.)

---

**Question: 19**      **CertyIQ**

A company uses a foundation model (FM) from Amazon Bedrock for an AI search tool. The company wants to fine-tune the model to be more accurate by using the company's data.
Which strategy will successfully fine-tune the model?

    A.Provide labeled data with the prompt field and the completion field.

    B.Prepare the training dataset by creating a .txt file that contains multiple lines in .csv format.

    C.Purchase Provisioned Throughput for Amazon Bedrock.

    D.Train the model on journals and textbooks.

**Answer: A**

**Explanation:**

The correct answer is **A. Provide labeled data with the prompt field and the completion field.**

Fine-tuning a foundation model (FM) in Amazon Bedrock involves adapting the pre-trained model to perform better on a specific task using your own data. This is achieved by providing the model with examples of inputs (prompts) and the desired outputs (completions). This teaches the model to generate more accurate and relevant responses for your use case.

Option A directly addresses this fine-tuning process. By providing labeled data in a format of prompt-completion pairs, you are guiding the model to learn the desired relationships and improve its accuracy on your specific data. The prompt acts as the input to the model, and the completion represents the ideal output for that prompt. This is a common method of fine-tuning large language models.

Option B is incorrect. While a .txt file with .csv format might be used for data storage, it doesn't inherently define the prompt-completion structure required for fine-tuning. The data must be specifically formatted to indicate which part is the input and which is the desired output.

Option C is irrelevant. Provisioned Throughput in Amazon Bedrock relates to ensuring dedicated capacity and performance for inference, not to the fine-tuning process itself. It addresses the speed and availability of predictions from a model but does not contribute to its accuracy.

Option D is too generic. While training on journals and textbooks could improve a model's general knowledge, it doesn't guarantee improved performance on the company's specific AI search task. Fine-tuning with relevant, labeled data from the company is far more effective for achieving the desired accuracy.

In summary, the prompt-completion format enables the model to learn the desired relationships within your data, resulting in a fine-tuned model that is more accurate for the AI search tool.

Relevant links for further research:

Amazon Bedrock documentation
Fine-tuning Large Language Models

---

**Question: 20** <span style="color:orange">Certy</span>**IQ**

A company wants to use AI to protect its application from threats. The AI solution needs to check if an IP address is from a suspicious source.
Which solution meets these requirements?

 A.Build a speech recognition system.

 B.Create a natural language processing (NLP) named entity recognition system.

 C.Develop an anomaly detection system.

 D.Create a fraud forecasting system.

**Answer: C**

**Explanation:**

The correct answer is C, developing an anomaly detection system. Here's why:

The problem describes a scenario requiring the identification of suspicious IP addresses accessing an application. Anomaly detection, as a branch of AI, is specifically designed to identify deviations from normal behavior patterns. In this context, "normal" would represent typical IP address access patterns, and "anomalous" would be unusual IPs indicating potential threats. AWS offers services like Amazon GuardDuty and Amazon CloudWatch Anomaly Detection that can be leveraged for this purpose. GuardDuty, for example, analyzes VPC Flow Logs, DNS logs, and CloudTrail logs to identify malicious or unauthorized behavior.

Options A, B, and D are incorrect because they address different AI problems. Speech recognition (A) focuses on converting audio into text, irrelevant to network security. NLP named entity recognition (B) identifies and categorizes entities (e.g., people, organizations) in text, also unrelated to IP address analysis. Fraud forecasting (D) predicts future fraudulent activities, whereas the requirement is to detect current suspicious activity, making anomaly detection the more appropriate solution.

Anomaly detection systems can be trained on historical network traffic data to establish a baseline of normal activity. When new IP addresses access the application, the system compares their behavior against this baseline. If an IP address exhibits unusual access patterns (e.g., accessing resources at unusual times, generating an abnormally high number of requests), the system flags it as suspicious. The detected anomalies can trigger alerts or automated responses, such as blocking the IP address.

**Supporting Links:**

  1. **Amazon GuardDuty:** https://aws.amazon.com/guardduty/
  2. **Amazon CloudWatch Anomaly Detection:** https://aws.amazon.com/cloudwatch/features/anomaly-

detection/

3. **Anomaly detection:** https://en.wikipedia.org/wiki/Anomaly_detection

## Question: 21                                                                                          Certy**IQ**

Which feature of Amazon OpenSearch Service gives companies the ability to build vector database applications?

A.Integration with Amazon S3 for object storage

B.Support for geospatial indexing and queries

C.Scalable index management and nearest neighbor search capability

D.Ability to perform real-time analysis on streaming data

**Answer: C**

**Explanation:**

The correct answer is C: Scalable index management and nearest neighbor search capability. This is because vector databases excel at storing and querying high-dimensional vector embeddings, which represent data points in a semantic space. Amazon OpenSearch Service's ability to manage indexes at scale is crucial for handling the large datasets typically associated with vector embeddings. More importantly, the "nearest neighbor search capability" allows for efficient similarity searches, finding vectors closest to a query vector, which is the fundamental operation in vector database applications like semantic search, recommendation systems, and image recognition. Options A, B, and D, while valid features of OpenSearch Service, are not specific to the needs of vector database applications. A relates to general object storage, B to location-based data, and D to real-time data analysis. While real-time analysis can leverage vector embeddings, the core function is provided by the nearest neighbor search and scalable index management. Feature stores leverage vector embeddings; therefore, the ability to efficiently search and scale is critical. A vector database built upon OpenSearch needs those capabilities directly.

For further research, explore the following resources:

**Amazon OpenSearch Service Documentation:** https://aws.amazon.com/opensearch-service/
**Amazon OpenSearch Service k-NN:** https://opensearch.org/docs/latest/search-plugins/knn/index/
**AWS AI and Machine Learning:** https://aws.amazon.com/ai/

## Question: 22                                                                                          Certy**IQ**

Which option is a use case for generative AI models?

A.Improving network security by using intrusion detection systems

B.Creating photorealistic images from text descriptions for digital marketing

C.Enhancing database performance by using optimized indexing

D.Analyzing financial data to forecast stock market trends

**Answer: B**

**Explanation:**

The correct answer is B, creating photorealistic images from text descriptions for digital marketing. Generative AI models are designed to generate new, original content. Option B directly aligns with this core function. Text-to-image models, a specific type of generative AI, excel at producing visual content based on textual prompts. This makes them highly valuable for digital marketing, where visually appealing and unique

images can enhance campaigns and attract customers.

Option A, improving network security with intrusion detection systems, relates to AI in cybersecurity, particularly anomaly detection. While AI is used, this doesn't necessarily involve generating content. Option C, enhancing database performance with optimized indexing, employs AI for optimization, a different application area. Option D, analyzing financial data for stock market trends, falls under predictive analytics, again, an area separate from generative AI's content creation focus.

Generative AI models like DALL-E 2, Stable Diffusion, and Midjourney showcase the ability to produce photorealistic images, artwork, and variations based on text prompts. These models are transforming digital marketing by allowing users to create customized visuals without needing traditional photography or design skills.

For more information on Generative AI and its applications:

**AWS Documentation on Generative AI:** (Search AWS Documentation for "Generative AI" or "Amazon Bedrock" for their offerings)
**Google AI Blog on Generative Models:** https://ai.googleblog.com/ (Search for "Generative Models")
**OpenAI's DALL-E 2:** https://openai.com/dall-e-2/

---

## Question: 23 <span>CertyIQ</span>

A company wants to build a generative AI application by using Amazon Bedrock and needs to choose a foundation model (FM). The company wants to know how much information can fit into one prompt.
Which consideration will inform the company's decision?

A.Temperature

B.Context window

C.Batch size

D.Model size

**Answer: B**

**Explanation:**

The correct answer is **B. Context window**.

The context window of a foundation model (FM) in Amazon Bedrock dictates the maximum amount of text that can be included within a single prompt and its associated response. This is a critical consideration for a generative AI application because it directly impacts the amount of information the model can effectively process and utilize to generate relevant and coherent outputs.

A larger context window allows the model to consider more data, enabling it to handle longer documents, maintain better context in conversations, and generate more detailed and nuanced responses. Conversely, a smaller context window limits the model's ability to leverage extensive information, potentially leading to less accurate or complete results.

**Why the other options are incorrect:**

**A. Temperature:** Temperature controls the randomness of the model's output. While important for controlling creativity and predictability, it doesn't affect how much data can be fed into the model.
**C. Batch size:** Batch size refers to the number of prompts processed simultaneously. This is relevant for throughput and efficiency, but not the size of a single prompt.
**D. Model size:** Model size, referring to the number of parameters in the FM, indicates the model's complexity and potential performance. However, it does not directly define the amount of text the model can process

within a prompt.

Therefore, when selecting an FM for a generative AI application in Amazon Bedrock, the company must consider the context window of each model to ensure it can accommodate the expected input size for their use case.

**Supporting Links:**

Amazon Bedrock Documentation: provides comprehensive details regarding the models and their capabilities.
Understanding LLM Context Windows: External resources for understanding context windows.

---

**Question: 24**                                                                 **CertyIQ**

A company wants to make a chatbot to help customers. The chatbot will help solve technical problems without human intervention.
The company chose a foundation model (FM) for the chatbot. The chatbot needs to produce responses that adhere to company tone.
Which solution meets these requirements?

   A.Set a low limit on the number of tokens the FM can produce.

   B.Use batch inferencing to process detailed responses.

   C.Experiment and refine the prompt until the FM produces the desired responses.

   D.Define a higher number for the temperature parameter.

**Answer: C**

**Explanation:**

The correct answer is C: Experiment and refine the prompt until the FM produces the desired responses. This is because prompt engineering is a key technique for tailoring the output of a foundation model to specific requirements and stylistic guidelines.

Here's a detailed justification:

Foundation models are pre-trained on vast datasets and, while powerful, may not inherently understand or adhere to a specific company tone or desired response style. Fine-tuning a pre-trained FM can be a heavy lift in time and resources. Therefore, prompt engineering becomes a crucial, more cost-effective method for influencing the model's output.

By carefully crafting prompts, which are the instructions given to the model, the company can guide the chatbot to generate responses that align with their brand voice and target audience. Experimenting with different prompt structures, keywords, and examples helps to discover the prompts that consistently produce the desired responses. This iterative process of experimentation and refinement allows the company to fine-tune the chatbot's behavior without directly retraining the FM.

Option A, setting a low token limit, primarily controls the length of the response, not the tone or style. While a shorter response might indirectly influence the tone, it doesn't guarantee adherence to company guidelines. Option B, using batch inferencing, is more related to processing large volumes of requests offline rather than influencing the response style. Option D, increasing the temperature parameter, introduces more randomness and creativity into the responses. While this might sometimes be desirable, it's generally not suitable for a chatbot that needs to adhere to a specific, controlled tone and provide reliable information. A high temperature can lead to unpredictable and inconsistent outputs, potentially deviating from the desired company tone.

In conclusion, prompt engineering provides a direct and effective mechanism to guide the FM towards producing responses that meet the specific requirements of the company's chatbot, making it the optimal

solution in this scenario.

Relevant links:

**Prompt Engineering Guide:** https://www.promptingguide.ai/
**Amazon Bedrock Documentation on Prompt Engineering:** (Check AWS Documentation directly as links expire)

---

**Question: 25**

A company wants to use a large language model (LLM) on Amazon Bedrock for sentiment analysis. The company wants to classify the sentiment of text passages as positive or negative.

Which prompt engineering strategy meets these requirements?

A.Provide examples of text passages with corresponding positive or negative labels in the prompt followed by the new text passage to be classified.

B.Provide a detailed explanation of sentiment analysis and how LLMs work in the prompt.

C.Provide the new text passage to be classified without any additional context or examples.

D.Provide the new text passage with a few examples of unrelated tasks, such as text summarization or question answering.

**Answer: A**

**Explanation:**

The correct answer is A, providing examples of text passages with corresponding sentiment labels in the prompt. This approach leverages a prompt engineering strategy known as "few-shot learning." Large Language Models (LLMs) like those available through Amazon Bedrock excel when given examples of the desired task. By including examples of text and their corresponding sentiment (positive or negative), you're effectively demonstrating to the LLM how you want it to classify the new, unlabeled text.

This method guides the LLM to understand the nuances of positive and negative sentiment within the specific context of the company's data. It helps the LLM generalize from these examples to the new input, improving the accuracy and relevance of the sentiment classification. Options B, C, and D are less effective because they don't provide the LLM with the specific examples needed for the targeted sentiment analysis task.

Providing a detailed explanation of sentiment analysis (Option B) is unnecessary, as LLMs are pre-trained on vast amounts of text data and already possess a general understanding of the concept. Simply providing the text passage without context (Option C) relies entirely on the LLM's inherent knowledge and might not be sufficient for accurate sentiment classification in a specific domain or style. Presenting unrelated tasks (Option D) would confuse the LLM and hinder its ability to focus on the sentiment analysis goal. Few-shot learning gives the model a tangible and tailored guide to perform the specific task, leading to better results. Few-shot learning optimizes accuracy by presenting the model with relevant examples, which enables the model to adapt its responses to match the user's specific needs.https://aws.amazon.com/bedrock/https://towardsdatascience.com/prompt-engineering-guide-for-developers-a49d27609e88

---

**Question: 26**

A security company is using Amazon Bedrock to run foundation models (FMs). The company wants to ensure that only authorized users invoke the models. The company needs to identify any unauthorized access attempts to set appropriate AWS Identity and Access Management (IAM) policies and roles for future iterations of the FMs. Which AWS service should the company use to identify unauthorized users that are trying to access Amazon

Bedrock?

    A.AWS Audit Manager

    B.AWS CloudTrail

    C.Amazon Fraud Detector

    D.AWS Trusted Advisor

**Answer: B**

**Explanation:**

The correct answer is **B. AWS CloudTrail.**

AWS CloudTrail is a service that enables governance, compliance, operational auditing, and risk auditing of your AWS account. CloudTrail logs API calls made to AWS services, including Amazon Bedrock. This means that every attempt to invoke a Bedrock model, whether successful or unsuccessful, is recorded in CloudTrail logs. These logs capture information about the identity of the caller (the IAM user or role used to make the request), the time of the request, the source IP address, the specific API call made (e.g., InvokeModel), and whether the request was authorized.

By analyzing these CloudTrail logs, the security company can identify unauthorized access attempts. They can filter the logs for events where the IAM user or role attempting to access Amazon Bedrock does not have the necessary permissions. This information is crucial for identifying users who are trying to access the models without proper authorization. They can then refine their IAM policies and roles to prevent future unauthorized access, ensuring only authorized users can invoke the models.

AWS Audit Manager (A) helps you continuously audit your AWS usage to simplify how you assess risk and compliance with regulations and industry standards. While useful for compliance, it doesn't directly identify specific unauthorized access attempts like CloudTrail. Amazon Fraud Detector (C) is used to detect fraudulent activities and is not designed for tracking API access or unauthorized access attempts to AWS services like Bedrock. AWS Trusted Advisor (D) provides recommendations for cost optimization, performance, security, fault tolerance, and service limits. It does not track API calls or provide detailed information on unauthorized access attempts.

In summary, CloudTrail is the appropriate service because it specifically logs API calls, including unauthorized attempts, allowing the company to identify the users and the actions they were trying to perform, directly addressing the need to identify unauthorized access attempts.

Here are some authoritative links for further research:

**AWS CloudTrail:** https://aws.amazon.com/cloudtrail/
**Logging Amazon Bedrock API calls with AWS CloudTrail:**
https://docs.aws.amazon.com/bedrock/latest/userguide/security-logging-using-cloudtrail.html

**Question: 27**         **CertyIQ**

A company has developed an ML model for image classification. The company wants to deploy the model to production so that a web application can use the model.
The company needs to implement a solution to host the model and serve predictions without managing any of the underlying infrastructure.
Which solution will meet these requirements?

    A.Use Amazon SageMaker Serverless Inference to deploy the model.

    B.Use Amazon CloudFront to deploy the model.

    C.Use Amazon API Gateway to host the model and serve predictions.

D.Use AWS Batch to host the model and serve predictions.

**Answer: A**

**Explanation:**

The correct answer is **A. Use Amazon SageMaker Serverless Inference to deploy the model.**

Here's why:

Amazon SageMaker Serverless Inference is specifically designed to deploy machine learning models for inference without managing underlying infrastructure. It automatically provisions and scales compute resources based on the request volume, eliminating the need for manual capacity planning and management. This aligns perfectly with the company's requirement to host the model and serve predictions without infrastructure overhead. The model is invoked through an endpoint, making it suitable for integration with a web application.

Option B, Amazon CloudFront, is a content delivery network (CDN) used for caching and distributing static and dynamic web content. It's not designed for hosting and serving ML model predictions.

Option C, Amazon API Gateway, is used to create, publish, maintain, monitor, and secure APIs. While it can be used as a front-end for accessing an ML model, it doesn't host or manage the model itself. You would still need a compute resource behind API Gateway to serve the predictions.

Option D, AWS Batch, is a batch computing service that allows you to run batch computing workloads at any scale. It's designed for running discrete jobs and isn't suitable for real-time, low-latency inference required by a web application.

SageMaker Serverless Inference is the best option because it provides a fully managed environment for hosting and serving ML models, automatically scaling to meet demand, and abstracting away the complexity of infrastructure management.

**Relevant links:**

Amazon SageMaker Serverless Inference

---

**Question: 28**

**CertyIQ**

An AI company periodically evaluates its systems and processes with the help of independent software vendors (ISVs). The company needs to receive email message notifications when an ISV's compliance reports become available.
Which AWS service can the company use to meet this requirement?

A.AWS Audit Manager

B.AWS Artifact

C.AWS Trusted Advisor

D.AWS Data Exchange

**Answer: B**

**Explanation:**

The correct answer is B (AWS Artifact). Here's why:

AWS Artifact is a service that provides on-demand access to AWS's compliance reports and agreements. Critically, it also provides a mechanism for customers to download compliance reports from third-party

vendors who have chosen to share them through Artifact. AWS Artifact allows you to subscribe to notifications when new reports are available. This directly addresses the requirement of receiving email notifications when an ISV's compliance reports become available.

AWS Audit Manager automates the process of auditing your AWS usage and provides evidence to support audits. While important for compliance, it doesn't directly manage or distribute third-party compliance reports, nor does it provide notifications about their availability.

AWS Trusted Advisor provides recommendations to optimize your AWS infrastructure for cost, performance, security, and fault tolerance. It does not provide access to or notifications about third-party compliance reports.

AWS Data Exchange is a service for finding, subscribing to, and using third-party data in the cloud. While it deals with third-party data, it's focused on data sets for analysis and use, not compliance reports.

Therefore, AWS Artifact is the only service that specifically facilitates accessing and being notified about the availability of third-party compliance reports, fulfilling the exam question's requirement.

For further information, refer to the official AWS Artifact documentation: https://aws.amazon.com/artifact/ and the AWS Audit Manager documentation https://aws.amazon.com/audit-manager/. Also refer to AWS Trusted Advisor https://aws.amazon.com/premiumsupport/technology/trusted-advisor/ and AWS Data Exchange documentation https://aws.amazon.com/data-exchange/.

---

**Question: 29**  **CertyIQ**

A company wants to use a large language model (LLM) to develop a conversational agent. The company needs to prevent the LLM from being manipulated with common prompt engineering techniques to perform undesirable actions or expose sensitive information.
Which action will reduce these risks?

A.Create a prompt template that teaches the LLM to detect attack patterns.

B.Increase the temperature parameter on invocation requests to the LLM.

C.Avoid using LLMs that are not listed in Amazon SageMaker.

D.Decrease the number of input tokens on invocations of the LLM.

**Answer: A**

**Explanation:**

The correct answer is A: Create a prompt template that teaches the LLM to detect attack patterns.

Here's a detailed justification:

Large Language Models (LLMs) are susceptible to prompt injection attacks, where malicious users manipulate the input to bypass intended security measures or elicit undesirable responses. A robust defense strategy involves proactively training the LLM to recognize and neutralize these attack patterns. This is achieved by crafting prompt templates that explicitly instruct the model on how to identify, categorize, and respond to potential attacks. Such templates can include examples of common attack vectors like prompt leaking, denial-of-service, and jailbreaking attempts.

By teaching the LLM to detect and respond to such patterns, you create a defensive layer that mitigates the risk of the model being manipulated. This proactive approach is more effective than simply relying on input sanitization or rate limiting. Prompt templates allow you to define boundaries for the LLM's behavior and instruct it to reject or flag prompts that deviate from intended use cases.

Option B is incorrect because increasing the temperature parameter increases the randomness and creativity of the LLM's responses, potentially making it more vulnerable to manipulation and generating unpredictable or harmful outputs.

Option C is incorrect because limiting the LLMs to only those on Amazon SageMaker does not inherently protect from prompt injection. The security of the LLM ultimately depends on how it is configured and secured, not solely on its availability within a particular platform.

Option D is incorrect because decreasing the number of input tokens may limit some complex attacks, but it also restricts the LLM's ability to understand the user's intent and perform its intended function. This approach is also insufficient to guard against all attack types and would negatively affect the usability of the conversational agent.

Therefore, creating a prompt template that teaches the LLM to identify and neutralize attack patterns offers the most effective means of mitigating the risks of prompt injection and ensuring that the conversational agent operates within defined security boundaries.

Supporting Links:

**Prompt Injection:** https://owasp.org/www-project-top-ten-for-llm-applications/
**Prompt Engineering:** https://www.promptingguide.ai/

---

**Question: 30**                                                              **Certy**IQ

A company is using the Generative AI Security Scoping Matrix to assess security responsibilities for its solutions. The company has identified four different solution scopes based on the matrix.
Which solution scope gives the company the MOST ownership of security responsibilities?

   A.Using a third-party enterprise application that has embedded generative AI features.

   B.Building an application by using an existing third-party generative AI foundation model (FM).

   C.Refining an existing third-party generative AI foundation model (FM) by fine-tuning the model by using data specific to the business.

   D.Building and training a generative AI model from scratch by using specific data that a customer owns.

**Answer: D**

**Explanation:**

The answer is D because it represents the scenario where the company has the most control and therefore the most responsibility for security. The Generative AI Security Scoping Matrix, which aims to clarify security responsibilities across different engagement levels, would classify building and training a model from scratch as requiring the most ownership.

Option A involves using a third-party application where the vendor largely manages the security of the embedded AI features. Option B, building an application using an existing FM, shifts more security responsibility to the company but still relies on the FM provider's security measures. Option C, fine-tuning an FM, increases the company's security responsibility over the data used for fine-tuning, but the underlying FM's security remains primarily the vendor's concern.

However, in Option D, the company controls every aspect of the model's lifecycle, including data ingestion, model training, infrastructure security, and monitoring. This complete control directly translates to complete accountability for ensuring the model's security and preventing vulnerabilities like data poisoning, model evasion, or unintended bias. The entire security burden falls squarely on the company, as they're building the system from the ground up. This includes not only the model's code and architecture, but also the security of the data used for training and inference.

Therefore, the solution scope that necessitates the most ownership of security responsibilities is building and training a generative AI model from scratch using specific data that the customer owns. This "from scratch" approach means no reliance on external vendors for the core model's security, increasing the company's responsibilities significantly.

For further research, consider these resources:

**OWASP (Open Web Application Security Project):** Provides guidance on AI security risks and mitigations.
https://owasp.org/
**NIST AI Risk Management Framework:** Offers a framework for managing risks associated with AI systems.
https://www.nist.gov/itl/ai-risk-management-framework
**ENISA (European Union Agency for Cybersecurity) - AI Cybersecurity:** provides analysis of AI specific cybersecurity threats and vulnerabilities https://www.enisa.europa.eu/topics/emerging-technologies/artificial-intelligence

---

**Question: 31**                                                                                    **CertyIQ**

An AI practitioner has a database of animal photos. The AI practitioner wants to automatically identify and categorize the animals in the photos without manual human effort.
Which strategy meets these requirements?

A.Object detection

B.Anomaly detection

C.Named entity recognition

D.Inpainting

---

**Answer: A**

**Explanation:**

The correct answer is **A. Object detection** because it directly addresses the requirement of automatically identifying and categorizing objects (animals, in this case) within images.

Here's why the other options are less suitable:

**B. Anomaly detection:** This technique is used to identify data points that deviate significantly from the norm. While it could potentially flag unusual animals, it wouldn't categorize common ones or identify the specific animal type. It primarily focuses on identifying outliers, not classification.

**C. Named entity recognition (NER):** NER extracts named entities (like people, organizations, or locations) from text. Since the input is images, not text, NER is not applicable.

**D. Inpainting:** Inpainting fills in missing parts of an image. It's a useful tool for image restoration but doesn't classify objects within an image.

Object detection algorithms excel at identifying and localizing multiple objects within an image and assigning a class label to each. In the context of animal photos, an object detection model could be trained to recognize various animal species (e.g., dog, cat, bird, lion) and draw bounding boxes around each animal in the image, along with its predicted label.

Furthermore, cloud services like Amazon Rekognition provide pre-trained object detection models that can be readily used or fine-tuned for custom use cases. This aligns with the "AI Practitioner" context, suggesting a practical approach. Therefore, object detection provides the necessary functionality for automatic identification and categorization of animals in photos with minimal manual intervention.

**Supporting Links:**

## Question: 32

**Certy**IQ

A company wants to create an application by using Amazon Bedrock. The company has a limited budget and prefers flexibility without long-term commitment.
Which Amazon Bedrock pricing model meets these requirements?

A.On-Demand

B.Model customization

C.Provisioned Throughput

D.Spot Instance

**Answer: A**

**Explanation:**

The correct answer is A. On-Demand pricing for Amazon Bedrock is the most suitable option for the company's needs. On-Demand pricing allows users to pay only for what they use, offering flexibility and avoiding long-term commitments. This aligns perfectly with the company's desire to avoid large upfront costs and maintain flexibility.

Option B, Model Customization, involves costs related to fine-tuning foundation models, which may not be immediately necessary or aligned with a limited budget. Option C, Provisioned Throughput, involves committing to a specific throughput level for a sustained period, which contradicts the company's preference for flexibility and aversion to long-term commitment. Finally, Spot Instances (Option D) are not directly applicable to Amazon Bedrock. Spot Instances are typically associated with EC2 and involve bidding for unused compute capacity, which doesn't apply to the consumption model of Bedrock's AI services. The On-Demand option offers the best balance of cost-effectiveness and flexibility.

Authoritative links:

Amazon Bedrock Pricing: https://aws.amazon.com/bedrock/pricing/

## Question: 33

**Certy**IQ

Which AWS service or feature can help an AI development team quickly deploy and consume a foundation model (FM) within the team's VPC?

A.Amazon Personalize

B.Amazon SageMaker JumpStart

C.PartyRock, an Amazon Bedrock Playground

D.Amazon SageMaker endpoints

**Answer: B**

**Explanation:**

The correct answer is **B. Amazon SageMaker JumpStart**. Here's why:

Amazon SageMaker JumpStart provides pre-trained models, pre-built solutions, and example notebooks for a variety of machine learning tasks, including those utilizing foundation models (FMs). It allows AI development

teams to quickly get started with FMs without needing to build everything from scratch. Crucially, SageMaker JumpStart allows you to deploy these FMs directly within your Virtual Private Cloud (VPC), ensuring that data and model access remain secure and compliant with your organizational policies. This is essential for many enterprises that need to keep their AI/ML workloads within a defined network boundary.

Amazon Personalize (A) focuses on recommendation systems, which is a specific AI use case, not a general FM deployment tool. PartyRock (C) is a playground environment for experimenting with Amazon Bedrock and FMs, but it's not designed for production deployment within a VPC. Amazon SageMaker endpoints (D) can be used to deploy FMs, but SageMaker JumpStart simplifies the initial deployment process by providing pre-configured models and deployment templates, making it much faster and easier for a team to get started. The pre-built aspect of JumpStart significantly accelerates the time-to-market for deploying and consuming FMs in your VPC compared to manually configuring SageMaker endpoints. In essence, JumpStart provides a curated and simplified onboarding experience for FMs within the SageMaker ecosystem.

Further Reading:

Amazon SageMaker JumpStart Documentation
Deploying Models from SageMaker JumpStart

## Question: 34 CertyIQ

How can companies use large language models (LLMs) securely on Amazon Bedrock?

A.Design clear and specific prompts. Configure AWS Identity and Access Management (IAM) roles and policies by using least privilege access.

B.Enable AWS Audit Manager for automatic model evaluation jobs.

C.Enable Amazon Bedrock automatic model evaluation jobs.

D.Use Amazon CloudWatch Logs to make models explainable and to monitor for bias.

**Answer: A**

**Explanation:**

The correct answer is A because security when using LLMs on Amazon Bedrock hinges on two critical aspects: prompt engineering and access control. Designing clear and specific prompts minimizes the chances of unintended model behavior or malicious manipulation through prompt injection. Ambiguous or overly broad prompts increase the attack surface.

More importantly, access control via IAM roles and policies implemented using the principle of least privilege is paramount. Least privilege means granting users and services only the permissions they absolutely need to perform their tasks, preventing unauthorized access and data breaches. This restricts the blast radius if a vulnerability is exploited. It confines each user and process to access only the bedrock models required for its dedicated purpose.

While model evaluation (options B and C) is crucial for assessing performance and bias, it doesn't directly address security vulnerabilities related to data access and model manipulation. AWS Audit Manager's focus is on compliance auditing, not real-time model security. Similarly, while Amazon CloudWatch Logs (option D) helps with monitoring and detecting anomalies, it doesn't prevent unauthorized access or prompt injection attacks proactively. Explainability and bias monitoring are important aspects of responsible AI, but a strong security posture relies on controlled access and secure prompting practices. The best approach is to reduce the risk of misuse through prompt design combined with least privilege access.

Therefore, configuring IAM with least privilege access and designing clear prompts is the foundational security step when leveraging LLMs on Bedrock, making option A the most pertinent response.

Further research:

**AWS IAM best practices:** https://docs.aws.amazon.com/IAM/latest/UserGuide/best-practices.html
**Principle of Least Privilege:** https://en.wikipedia.org/wiki/Principle_of_least_privilege
**Prompt Engineering Guide:** https://www.promptingguide.ai/

## Question: 35

A company has terabytes of data in a database that the company can use for business analysis. The company wants to build an AI-based application that can build a SQL query from input text that employees provide. The employees have minimal experience with technology.
Which solution meets these requirements?

    A.Generative pre-trained transformers (GPT)

    B.Residual neural network

    C.Support vector machine

    D.WaveNet

**Answer: A**

**Explanation:**

The correct answer is A, Generative Pre-trained Transformers (GPT). Here's a detailed justification:

The problem requires a solution that can translate natural language (employee input text) into SQL queries. GPT excels at natural language processing (NLP) tasks, specifically text-to-text generation. This is because GPT models are trained on massive datasets of text and code, enabling them to understand the relationship between human language and structured languages like SQL. They can learn to map the intent expressed in the input text to the corresponding SQL syntax.

The other options are not well-suited for this task. Residual neural networks (ResNets) are primarily used for image recognition and other computer vision tasks, while Support Vector Machines (SVMs) are used for classification and regression. Neither directly addresses the need for translating natural language to SQL. WaveNet is designed for generating raw audio waveforms and isn't applicable to the task.

GPT's ability to perform zero-shot or few-shot learning makes it particularly attractive. It means the model can generate SQL queries based on input prompts with minimal or no specific training on the company's data. A fine-tuned GPT model, further trained on the company's specific database schema and business terminology, can produce even more accurate and reliable SQL queries.

Therefore, GPT provides a ready-to-use approach for converting natural language questions into SQL, which is essential for empowering employees with limited technical experience to access and analyze the data effectively. Amazon offers services like Amazon Bedrock that provide access to powerful pre-trained models including those that can perform text-to-SQL conversion.

For further research:

**GPT:** https://openai.com/research/gpt-3
**Amazon Bedrock:** https://aws.amazon.com/bedrock/
**Text-to-SQL:** https://arxiv.org/abs/1709.00103

## Question: 36

A company built a deep learning model for object detection and deployed the model to production.

Which AI process occurs when the model analyzes a new image to identify objects?

    A.Training

    B.Inference

    C.Model deployment

    D.Bias correction

**Answer: B**

**Explanation:**

The correct answer is B, Inference. Here's why:

Inference is the process of using a trained machine learning model to make predictions on new, unseen data. In this scenario, the company has already built and deployed the object detection model. When the model analyzes a new image to identify objects, it is applying its learned knowledge to classify and locate objects within the image. This act of applying the trained model to new data to obtain predictions is precisely what inference entails.

Training, on the other hand, is the process of creating the model itself, where it learns patterns from a dataset of labeled examples. Model deployment is the process of making the trained model available for use. Bias correction focuses on identifying and mitigating unfair or discriminatory outcomes from the model's predictions or underlying data.

Since the model is already built and deployed, and the question describes it analyzing new images to identify objects, this specifically refers to inference. The model isn't being trained, deployed, or having its bias corrected; it's simply being used to make predictions.

Therefore, the correct answer is B.

Further reading:

**Amazon SageMaker Inference:** https://aws.amazon.com/sagemaker/inference/
**Machine Learning Inference:** https://docs.aws.amazon.com/machine-learning/latest/dg/machinelearning-process-inference.html

**Question: 37**                     **CertyIQ**

An AI practitioner is building a model to generate images of humans in various professions. The AI practitioner discovered that the input data is biased and that specific attributes affect the image generation and create bias in the model.
Which technique will solve the problem?

    A.Data augmentation for imbalanced classes

    B.Model monitoring for class distribution

    C.Retrieval Augmented Generation (RAG)

    D.Watermark detection for images

**Answer: A**

**Explanation:**

The correct answer is A, Data augmentation for imbalanced classes. Here's why:

The core issue is bias stemming from imbalanced data regarding attributes within the dataset used to train the image generation model. The model reflects the biases present in the training data.

Data augmentation addresses this directly. If, for example, the training data has fewer images of female doctors compared to male doctors, data augmentation can artificially increase the number of images of female doctors. This is achieved by applying transformations to existing images of female doctors (e.g., rotations, zooms, slight color variations) to create new, synthetic images. By balancing the representation of different attributes (gender, profession, race, etc.), the model is less likely to generate biased outputs. It's a proactive approach to data imbalance.

Option B, Model monitoring for class distribution, only detects the presence of bias after the model is deployed. It doesn't solve the underlying problem of biased training data. It can flag issues, but not prevent them from occurring in the first place.

Option C, Retrieval Augmented Generation (RAG), is irrelevant here. RAG is typically used to improve the factual accuracy of language models by grounding them in external knowledge sources. It's not designed to address data imbalance or bias in image generation.

Option D, Watermark detection for images, is focused on identifying if an image has been watermarked, possibly for copyright or authenticity purposes. It's not relevant to mitigating bias in image generation model.

Therefore, data augmentation for imbalanced classes directly tackles the root cause of the problem: the skewed representation of attributes in the training data. By creating a more balanced dataset, the model learns to generate images without unfairly favoring certain groups or attributes. This is a crucial step in building fair and unbiased AI systems.

Relevant Resources:

**Data Augmentation:** https://www.tensorflow.org/tutorials/images/data_augmentation (TensorFlow documentation explaining data augmentation techniques)
**Fairness in Machine Learning:** https://developers.google.com/machine-learning/fairness-ai (Google's guide to understanding and mitigating bias in machine learning)

## Question: 38 <span style="float:right">CertyIQ</span>

A company is implementing the Amazon Titan foundation model (FM) by using Amazon Bedrock. The company needs to supplement the model by using relevant data from the company's private data sources.
Which solution will meet this requirement?

A.Use a different FM.

B.Choose a lower temperature value.

C.Create an Amazon Bedrock knowledge base.

D.Enable model invocation logging.

**Answer: C**

**Explanation:**

The correct answer is **C. Create an Amazon Bedrock knowledge base.** Here's why:

Amazon Bedrock knowledge bases directly address the need to augment foundation models with data from private data sources. A knowledge base allows you to securely connect to your data repositories (like S3, databases, etc.) and ingest the data. This ingested data is then used during inference time to ground the foundation model's responses, making them more relevant and accurate in the context of your company's specific information.

The Titan FM, while powerful, doesn't inherently "know" about your private data. Bedrock knowledge bases fill this gap by providing the model with access to that data. This is a Retrieval Augmented Generation (RAG)

approach, where relevant information is retrieved from the knowledge base and provided to the model alongside the user's prompt.

Option A (Using a different FM) is incorrect. Switching FMs doesn't solve the problem of needing to incorporate private data. While other FMs might have advantages, they still require a mechanism to access and utilize company-specific knowledge.

Option B (Choosing a lower temperature value) controls the randomness of the model's output but doesn't influence its ability to use external data sources. Lower temperature makes the output more deterministic and less creative, but irrelevant to the data incorporation problem.

Option D (Enable model invocation logging) is for auditing and monitoring purposes. It doesn't help the model learn from or utilize the company's private data. It records the inputs and outputs of the model, but it doesn't enrich the model with external knowledge.

In summary, Amazon Bedrock knowledge bases are specifically designed to integrate private data with foundation models like Amazon Titan, enabling the model to provide more informed and relevant responses. The other options do not address this core requirement.

**Further Research:**

**Amazon Bedrock Knowledge Bases:** https://aws.amazon.com/bedrock/knowledge-bases/
**Retrieval Augmented Generation (RAG):** (Search for "Retrieval Augmented Generation" to find numerous articles and research papers on the concept). While no specific AWS documentation directly defines RAG in the context of Bedrock, knowledge bases are a prime implementation of this pattern.

---

## Question: 39                                                      CertyIQ

A medical company is customizing a foundation model (FM) for diagnostic purposes. The company needs the model to be transparent and explainable to meet regulatory requirements.
Which solution will meet these requirements?

    A.Configure the security and compliance by using Amazon Inspector.

    B.Generate simple metrics, reports, and examples by using Amazon SageMaker Clarify.

    C.Encrypt and secure training data by using Amazon Macie.

    D.Gather more data. Use Amazon Rekognition to add custom labels to the data.

**Answer: B**

**Explanation:**

The correct answer is B: Generate simple metrics, reports, and examples by using Amazon SageMaker Clarify.

Here's a detailed justification:

The question emphasizes the need for transparency and explainability of the customized foundation model (FM). Regulatory requirements often necessitate understanding how a model arrives at its predictions, particularly in sensitive domains like medical diagnostics.

**Amazon SageMaker Clarify** directly addresses this requirement. It provides tools to detect potential bias in machine learning models and helps explain their predictions. It generates reports with feature importance scores and examples, making the model's decision-making process more understandable. This is crucial for demonstrating compliance with regulations that demand explainable AI.

Let's examine why the other options are not suitable:

**A. Amazon Inspector:** Primarily focuses on security vulnerabilities in your infrastructure and applications. It does not provide insights into the internal workings or explainability of an AI model.

**C. Amazon Macie:** Is a data security and privacy service that discovers, classifies, and protects sensitive data stored in Amazon S3. While important for data governance, it doesn't contribute to model explainability.

**D. Amazon Rekognition:** Is an image and video analysis service. While potentially useful for adding labels to data, it does not explain the decision-making process of the customized FM. More data alone doesn't guarantee explainability; the model itself needs to be interpretable.

In summary, SageMaker Clarify is specifically designed to help understand and explain model behavior, fulfilling the regulatory requirements for transparency and explainability, making it the most appropriate solution in this scenario.

**Supporting Links:**

**Amazon SageMaker Clarify:** https://aws.amazon.com/sagemaker/clarify/
**Explainable AI (XAI):** https://aws.amazon.com/machine-learning/explainable-ai/

---

## Question: 40                                              <span>Certy**IQ**</span>

A company wants to deploy a conversational chatbot to answer customer questions. The chatbot is based on a fine-tuned Amazon SageMaker JumpStart model. The application must comply with multiple regulatory frameworks.
Which capabilities can the company show compliance for? (Choose two.)

    A.Auto scaling inference endpoints

    B.Threat detection

    C.Data protection

    D.Cost optimization

    E.Loosely coupled microservices

---

**Answer: BC**

**Explanation:**

The correct answer is **B. Threat detection** and **C. Data protection**. Here's a detailed justification:

**B. Threat Detection:** Deploying a chatbot application, especially one dealing with potentially sensitive customer data, necessitates robust threat detection mechanisms. Compliance frameworks often mandate security measures to identify and mitigate potential threats, such as unauthorized access, data breaches, and malicious attacks. Utilizing services like Amazon GuardDuty or AWS Security Hub, which can monitor API calls, network activity, and identify potential vulnerabilities, demonstrably supports compliance regarding threat detection. SageMaker JumpStart models themselves don't inherently provide threat detection; rather, the deployment environment must be secured.

**C. Data Protection:** Data protection is a cornerstone of most regulatory frameworks. A conversational chatbot inevitably processes and potentially stores customer data. Compliance mandates implementing appropriate data protection measures, including encryption (both in transit and at rest), access control, and adherence to data residency requirements. Fine-tuning a SageMaker JumpStart model does not automatically guarantee data protection. Instead, data protection relies on how the company handles the input data to the chatbot, the data processed during conversation, and the model's outputs. Using AWS KMS (Key Management Service) for encryption, IAM (Identity and Access Management) for granular access control, and adhering to data residency requirements using appropriate AWS regions are crucial for demonstrating data protection compliance.

**Why the other options are less appropriate:**

**A. Auto scaling inference endpoints:** While auto-scaling can improve performance and availability, it's primarily related to operational efficiency, not direct regulatory compliance. Though high availability can be indirectly linked to certain business continuity compliance requirements, it's not a primary compliance aspect for chatbots concerning regulated data.

**D. Cost optimization:** Cost optimization is an important aspect of cloud management, but it's not a core requirement for demonstrating compliance with regulatory frameworks pertaining to data handling and security.

**E. Loosely coupled microservices:** Microservices architecture offers benefits such as scalability and maintainability, but it doesn't directly address the core requirements of regulatory compliance in the context of a chatbot application. The compliance focus remains on data protection and threat detection, regardless of the application's architecture.

**Supporting Links:**

**AWS Security:** https://aws.amazon.com/security/
**Amazon GuardDuty:** https://aws.amazon.com/guardduty/
**AWS IAM:** https://aws.amazon.com/iam/
**AWS KMS:** https://aws.amazon.com/kms/
**AWS Compliance:** https://aws.amazon.com/compliance/

**Question: 41**                                                    **CertyIQ**

A company is training a foundation model (FM). The company wants to increase the accuracy of the model up to a specific acceptance level.
Which solution will meet these requirements?

    A.Decrease the batch size.

    B.Increase the epochs.

    C.Decrease the epochs.

    D.Increase the temperature parameter.

**Answer: B**

**Explanation:**

The correct answer is **B. Increase the epochs.**

Here's a detailed justification:

The goal is to improve the accuracy of a foundation model (FM) until it reaches an acceptable level. Epochs refer to the number of complete passes through the entire training dataset during the training process. Increasing the number of epochs allows the model to see the data more times, refining its internal parameters and learning the underlying patterns in the data more effectively.

Think of it like studying for an exam. The more times you review the material (the more epochs), the better your understanding and recall become, leading to higher accuracy on the exam (the FM's accuracy). Each epoch exposes the model to the data, allowing it to adjust its weights and biases to minimize the difference between its predictions and the actual values (loss). By increasing epochs, the model has more opportunities to optimize its parameters and converge towards a more accurate representation of the data.

Option A, decreasing the batch size, can impact training speed and potentially generalization, but it doesn't directly address the fundamental need for the model to learn more from the data. A smaller batch size might help with escaping local minima, but its primary impact isn't increasing overall accuracy.

Option C, decreasing the epochs, would be counterproductive. Fewer epochs would mean less exposure to the training data and, therefore, a less accurate model.

Option D, increasing the temperature parameter, is relevant to generating text or outputs from a trained model, not to the training process itself. The temperature parameter controls the randomness of the model's output. A higher temperature increases randomness, which is generally undesirable during training as it will reduce the model's certainty on the correct parameters.

Therefore, increasing the epochs is the most direct and effective way to improve the accuracy of a foundation model up to a specified level of acceptance by allowing the model to learn more thoroughly from the training data.

For further research, consider exploring these resources:

**AWS documentation on machine learning:** https://aws.amazon.com/machine-learning/
**Understanding Epochs in Machine Learning:** https://towardsdatascience.com/epoch-vs-iterations-training-deep-neural-networks-4a603c35f1b0 (This is a general resource on epochs)

---

A company is building a large language model (LLM) question answering chatbot. The company wants to decrease the number of actions call center employees need to take to respond to customer questions.
Which business objective should the company use to evaluate the effect of the LLM chatbot?

    A.Website engagement rate

    B.Average call duration

    C.Corporate social responsibility

    D.Regulatory compliance

**Answer: B**

**Explanation:**

The correct answer is **B. Average call duration**. Here's a detailed justification:

The primary goal of the LLM-powered chatbot is to reduce the workload of call center employees by automating responses to customer inquiries. This directly translates to decreasing the time call center employees spend on each call. Therefore, the business objective should be a metric that reflects this efficiency gain. Average call duration is a key performance indicator (KPI) that directly measures the length of customer service interactions.

An LLM chatbot that successfully answers customer questions will resolve issues faster, leading to shorter call durations. If the chatbot reduces the number of actions employees need to take (as stated in the question), the time spent on each call must go down if the chatbot works as intended. This is because agents will not have to search for answers or perform as many steps to assist the customer.

The other options are less relevant:

**A. Website engagement rate:** While a chatbot might influence website engagement, it's not the primary focus. The core objective is call center efficiency.
**C. Corporate social responsibility:** While CSR is important, it's not a direct metric for evaluating the chatbot's effectiveness in reducing call center workload.
**D. Regulatory compliance:** The chatbot must be compliant, but compliance is not the measuring objective.

Therefore, tracking the average call duration provides a clear and quantifiable way to assess the success of the LLM chatbot in achieving its intended purpose: to reduce call center workload. A significant decrease in

average call duration indicates that the chatbot is effectively resolving customer issues and freeing up call center employees to handle more complex or urgent matters.

**Authoritative Links for Further Research:**

**Key Performance Indicators (KPIs) in Customer Service:** Search on reputable business and technology websites (e.g., Forbes, McKinsey) for articles on KPIs in customer service.
**Contact Center Metrics:** Research contact center industry standards for relevant metrics.
**LLM Chatbots in Customer Service:** Search for case studies and articles on how LLMs are being used to improve customer service efficiency.

---

## Question: 43

Which functionality does Amazon SageMaker Clarify provide?

A.Integrates a Retrieval Augmented Generation (RAG) workflow

B.Monitors the quality of ML models in production

C.Documents critical details about ML models

D.Identifies potential bias during data preparation

**Answer: D**

**Explanation:**

Amazon SageMaker Clarify is specifically designed to detect and mitigate potential biases in machine learning models and data. Its primary function is to identify biases before model training, during the data preparation stage, and after model training, by evaluating model predictions. It achieves this by calculating a variety of fairness metrics and providing insights into the characteristics of the data and model that might contribute to biased outcomes. This allows data scientists to address biases proactively and build fairer, more equitable models.

Option A is incorrect because SageMaker Clarify doesn't focus on RAG workflows. RAG is related to LLMs and how they search and retrieve information before generating their response, which is a separate aspect. Option B is partially true, as SageMaker Model Monitor handles ongoing model quality monitoring in production. While Clarify can play a role in the initial model evaluation prior to deployment, the ongoing monitoring is not its core purpose. Option C describes model card documentation which while useful, is not Clarify's primary function. Clarify's primary focus is on bias detection and explanation, not general documentation.

Therefore, identifying potential bias during data preparation, as highlighted in option D, aligns perfectly with SageMaker Clarify's core functionality.

Amazon SageMaker Clarify DocumentationAWS Machine Learning Blog on Fairness and Explainability

---

## Question: 44

A company is developing a new model to predict the prices of specific items. The model performed well on the training dataset. When the company deployed the model to production, the model's performance decreased significantly.
What should the company do to mitigate this problem?

A.Reduce the volume of data that is used in training.

B.Add hyperparameters to the model.

C.Increase the volume of data that is used in training.

D.Increase the model training time.

**Question: 45**                                                                    **CertyIQ**

An ecommerce company wants to build a solution to determine customer sentiments based on written customer reviews of products.
Which AWS services meet these requirements? (Choose two.)

  A.Amazon Lex
  B.Amazon Comprehend
  C.Amazon Polly
  D.Amazon Bedrock
  E.Amazon Rekognition

valuable insights from text. Its key capability, sentiment analysis, directly addresses the requirement of determining customer sentiment from written reviews. Comprehend can identify whether a review expresses positive, negative, or neutral sentiment, along with a confidence score, allowing the e-commerce company to understand customer opinions about their products at scale. It requires no machine learning expertise and is a managed service, streamlining the process. (https://aws.amazon.com/comprehend/)

Amazon Bedrock allows you to access a wide range of foundation models (FMs) from different providers using a single API. The service provides a single endpoint that can be used to access a diverse set of foundation models (FMs) from AI21 Labs, Anthropic, Cohere, Meta, Stability AI, and Amazon. Foundation models can be trained on large text datasets and can be used to solve specific business cases for generating unique and realistic content. (https://aws.amazon.com/bedrock/) Amazon Bedrock can be used to perform sentiment analysis using one of its foundation models.

Amazon Lex (A) is a service for building conversational interfaces using voice and text. It is used to build chatbots but does not directly address the core requirement of analyzing existing text reviews for sentiment.

Amazon Polly (C) is a text-to-speech service that converts text into lifelike speech. While useful for other applications, it doesn't perform sentiment analysis.

Amazon Rekognition (E) is an image and video analysis service, primarily used for object and facial recognition, and has no direct application in analyzing text-based customer reviews for sentiment.

## Question: 46                                     CertyIQ

A company wants to use large language models (LLMs) with Amazon Bedrock to develop a chat interface for the company's product manuals. The manuals are stored as PDF files.
Which solution meets these requirements MOST cost-effectively?

    A.Use prompt engineering to add one PDF file as context to the user prompt when the prompt is submitted to Amazon Bedrock.

    B.Use prompt engineering to add all the PDF files as context to the user prompt when the prompt is submitted to Amazon Bedrock.

    C.Use all the PDF documents to fine-tune a model with Amazon Bedrock. Use the fine-tuned model to process user prompts.

    D.Upload PDF documents to an Amazon Bedrock knowledge base. Use the knowledge base to provide context when users submit prompts to Amazon Bedrock.

**Answer: D**

**Explanation:**

The most cost-effective solution is option D, utilizing an Amazon Bedrock knowledge base. This approach avoids the limitations and costs associated with the other options.

Option A is impractical because adding only one PDF file at a time limits the scope of the chat interface. Users would not have access to information spread across multiple manuals, reducing the solution's usefulness.

Option B, embedding all PDF files in the user prompt, is also inefficient and expensive. Large prompts consume significant resources from the LLM, leading to higher processing costs. Additionally, LLMs have context window limitations, meaning that very long prompts could be truncated or mishandled, degrading response quality.

Option C, fine-tuning a model using all PDF documents, is the most expensive option. Fine-tuning involves significant compute resources and time to train a new model based on the specific dataset. This is not only costly upfront, but also adds complexity to model management and deployment. Additionally, fine-tuning

might not be the most efficient way to retrieve information, as the LLM would need to generate the answer based on the learned data, rather than directly retrieve relevant sections.

Option D, leveraging a knowledge base, provides a cost-effective and efficient solution. A knowledge base allows the LLM to access relevant information from the PDF documents on demand, instead of relying on the full document to be included in the prompt. This reduces prompt size, lowers processing costs, and overcomes context window limitations. Furthermore, Bedrock knowledge bases use techniques like embeddings and vector search to quickly identify the most relevant passages within the documents, enabling more accurate and faster responses. This approach provides a balance between cost, performance, and accuracy, making it the optimal choice for this scenario.

Amazon Bedrock Knowledge BasesAmazon Bedrock pricing

---

**Question: 47**                                     **Certy**IQ

A social media company wants to use a large language model (LLM) for content moderation. The company wants to evaluate the LLM outputs for bias and potential discrimination against specific groups or individuals.
Which data source should the company use to evaluate the LLM outputs with the LEAST administrative effort?

    A.User-generated content

    B.Moderation logs

    C.Content moderation guidelines

    D.Benchmark datasets

**Answer: D**

**Explanation:**

Here's a detailed justification for why benchmark datasets are the best choice for evaluating LLM outputs for bias and discrimination in content moderation, with minimal administrative effort:

Benchmark datasets, specifically designed for bias detection, offer a pre-existing, standardized, and readily available resource. These datasets are meticulously curated to include diverse examples that highlight potential biases related to gender, race, religion, and other protected attributes. Compared to user-generated content (A) which is unstructured, potentially offensive, and requires extensive pre-processing and annotation, benchmark datasets offer a controlled and ethical environment. Moderation logs (B), while useful for understanding past moderation decisions, don't inherently provide a systematic way to identify new or subtle biases in LLM output. Content moderation guidelines (C) define the expected behavior but don't provide the data required to test if the LLM adheres to those guidelines.

Using benchmark datasets allows the company to directly compare the LLM's output against known problematic scenarios and quantify the presence of bias using established metrics. This approach reduces the administrative overhead of data collection, annotation, and bias identification because these tasks are already completed. Furthermore, utilizing established benchmarks aligns with best practices in responsible AI development. These datasets are often open-source or publicly available, minimizing cost and administrative burden. The pre-defined format of benchmark datasets also allows for automated evaluation pipelines, further reducing manual effort. By leveraging these datasets, the social media company can proactively address bias concerns before deploying the LLM into production, ensuring fairer and more equitable content moderation.

Here are some links for further research:

**AI Fairness 360:** https://aif360.mybluemix.net/ (Provides datasets and metrics for bias detection.)
**Papers with Code - Bias in NLP:** https://paperswithcode.com/task/bias-in-nlp (Lists relevant papers and

resources.)
**TensorFlow Data Validation (TFDV):** https://www.tensorflow.org/tfx/data_validation/get_started (Helps identify data anomalies and biases.)

---

**Question: 48**

A company wants to use a pre-trained generative AI model to generate content for its marketing campaigns. The company needs to ensure that the generated content aligns with the company's brand voice and messaging requirements.
Which solution meets these requirements?

   A.Optimize the model's architecture and hyperparameters to improve the model's overall performance.

   B.Increase the model's complexity by adding more layers to the model's architecture.

   C.Create effective prompts that provide clear instructions and context to guide the model's generation.

   D.Select a large, diverse dataset to pre-train a new generative model.

**Answer: C**

**Explanation:**

The correct answer is C: Create effective prompts that provide clear instructions and context to guide the model's generation.

Here's why:

Pre-trained generative AI models are powerful tools, but their output isn't inherently aligned with a specific brand's voice or marketing strategy. The most direct and efficient way to control the generated content's alignment with specific requirements is through prompt engineering. Prompts serve as instructions, providing context and guiding the model to produce outputs that match the desired brand voice, style, and messaging. By crafting precise and detailed prompts, the company can influence the model to generate content consistent with its brand guidelines.

Option A (optimizing architecture and hyperparameters) focuses on improving the model's overall performance but doesn't directly address the need for brand alignment. It is more concerned with aspects such as reducing inference time or improving accuracy in general generation tasks, not specific styling.

Option B (increasing model complexity) could potentially improve the model's ability to learn, but it's a computationally expensive and time-consuming approach that doesn't guarantee alignment with the brand voice. The increase in complexity does not equate to the generation of the desired outputs.

Option D (pre-training a new model) is the most resource-intensive option. Training a new model from scratch requires a massive dataset, significant computational resources, and a dedicated team of experts. This approach is unnecessary when the goal is to leverage an existing model for a specific purpose. It's also indirect; even with a carefully curated dataset, prompts are still needed to control the output.

Therefore, crafting effective prompts is the most practical and targeted solution for ensuring the generated content adheres to the company's brand voice and messaging requirements when using a pre-trained generative AI model.

For more information on prompt engineering, refer to these resources:

**Prompt Engineering Guide:** https://www.promptingguide.ai/
**AWS Documentation on AI/ML:** https://aws.amazon.com/machine-learning/

**Question: 49**

A loan company is building a generative AI-based solution to offer new applicants discounts based on specific business criteria. The company wants to build and use an AI model responsibly to minimize bias that could negatively affect some customers.
Which actions should the company take to meet these requirements? (Choose two.)

    A.Detect imbalances or disparities in the data.

    B.Ensure that the model runs frequently.

    C.Evaluate the model's behavior so that the company can provide transparency to stakeholders.

    D.Use the Recall-Oriented Understudy for Gisting Evaluation (ROUGE) technique to ensure that the model is 100% accurate.

    E.Ensure that the model's inference time is within the accepted limits.

---

**Answer: AC**

**Explanation:**

The correct answer is A and C.

To minimize bias and use AI responsibly, the loan company must first **detect imbalances or disparities in the data (A)**. This involves analyzing the training dataset for skewed representation of certain demographics or features that might lead to discriminatory outcomes. If the model is trained on biased data, it will likely perpetuate and amplify those biases in its predictions. Identifying and addressing these imbalances is a crucial step in building a fair AI model.

Further, the company should **evaluate the model's behavior to provide transparency to stakeholders (C)**. Model evaluation, including fairness metrics and explainability techniques, is essential to understand how the model makes decisions and to identify potential biases in its predictions. Transparency involves documenting the model's performance on various subgroups and communicating this information to stakeholders. This fosters trust and allows for ongoing monitoring and refinement of the model to mitigate any negative impacts on specific customer groups. The stakeholders will be better informed of the fairness and potential biases of the model which is in line with responsible AI.

Option B is incorrect because the frequency of model execution has little bearing on bias or fairness. Option D is incorrect because ROUGE is a metric used for evaluating text summarization models, and it is not relevant to bias detection or fairness in a loan application scenario. Even if the model achieves high ROUGE scores, it says nothing about fairness considerations. Option E is about model performance, which is about speed and resource consumption, but not about ethical considerations.

Here are some authoritative links for further research:

**Fairness in Machine Learning:** https://fairlearn.org/
**Amazon AI Fairness Checklist:** https://aws.amazon.com/machine-learning/fairness/
**Responsible AI Microsoft:** https://www.microsoft.com/en-us/ai/responsible-ai

---

**Question: 50**

A company is using an Amazon Bedrock base model to summarize documents for an internal use case. The company trained a custom model to improve the summarization quality.
Which action must the company take to use the custom model through Amazon Bedrock?

    A.Purchase Provisioned Throughput for the custom model.

    B.Deploy the custom model in an Amazon SageMaker endpoint for real-time inference.

    C.Register the model with the Amazon SageMaker Model Registry.

D.Grant access to the custom model in Amazon Bedrock.

**Answer: D**

**Explanation:**

The correct answer is D, granting access to the custom model in Amazon Bedrock. Here's why:

Amazon Bedrock is a fully managed service that allows you to access foundation models (FMs) from leading AI companies, including custom models you've trained. To use a custom model with Bedrock, you need to explicitly grant Bedrock access to it. This involves configuring the necessary permissions to allow Bedrock to invoke your model.

Option A is incorrect because Provisioned Throughput is a mechanism used for guaranteeing inference capacity for certain Bedrock models, not specifically a requirement for using a custom model.

Option B is incorrect because deploying the custom model to a SageMaker endpoint is an alternative method for serving the model but not a direct requirement for using it through Bedrock. Bedrock's strength lies in abstracting away the direct endpoint management.

Option C is incorrect because while SageMaker Model Registry is a valuable tool for managing and versioning models, it's not a mandatory step for integrating a custom model with Amazon Bedrock. Bedrock uses a different mechanism (resource access) for granting usage permissions. The specific implementation details for granting access depends on where the custom model is stored, such as S3.

Fundamentally, Bedrock needs the correct permissions to use the model, regardless of the location or registration state. This is accomplished by granting access to the model within Bedrock's configuration settings. This aligns with the principles of least privilege and controlled access to resources within AWS.

To summarize, utilizing custom models within Amazon Bedrock relies on proper access controls, which option D most accurately addresses. Options A, B, and C describe functionalities that are related to model deployment but not direct requirements for Bedrock access.

Here are some resources for more information:

**Amazon Bedrock documentation:** https://aws.amazon.com/bedrock/
**Controlling access to Amazon Bedrock resources:**
https://docs.aws.amazon.com/bedrock/latest/userguide/security-iam.html

**Question: 51**                                                          **CertyIQ**

A company needs to choose a model from Amazon Bedrock to use internally. The company must identify a model that generates responses in a style that the company's employees prefer.
What should the company do to meet these requirements?

A.Evaluate the models by using built-in prompt datasets.

B.Evaluate the models by using a human workforce and custom prompt datasets.

C.Use public model leaderboards to identify the model.

D.Use the model InvocationLatency runtime metrics in Amazon CloudWatch when trying models.

**Answer: B**

**Explanation:**

The correct answer is B because the company needs to identify a model that generates responses in a style preferred by its employees, which is a subjective requirement. Evaluating models using a human workforce

and custom prompt datasets allows the company to directly assess the stylistic preferences of the employees. Custom prompt datasets can be tailored to reflect the types of interactions employees typically have, ensuring the evaluation is relevant. A human workforce can then judge the generated responses based on style, tone, and overall suitability for internal use.

Option A, using built-in prompt datasets, may not capture the specific stylistic nuances the company requires. Public model leaderboards (option C) often focus on objective metrics like accuracy and speed and might not reflect stylistic preferences. Option D, using InvocationLatency metrics in CloudWatch, focuses solely on the model's response time and doesn't address the stylistic aspect at all. A human evaluation component using custom prompts directly addresses the subjective nature of the requirement.

By employing a human workforce and custom prompt datasets, the company can gather valuable feedback on how well each model aligns with the desired style. This approach provides a more nuanced and accurate assessment than relying solely on automated metrics or publicly available leaderboards. This ensures the chosen model not only performs well but also resonates with the company's employees.

For further research on evaluating language models and using human-in-the-loop approaches, you can refer to resources on Amazon Bedrock and general best practices for evaluating AI models:

**Amazon Bedrock Documentation:** https://aws.amazon.com/bedrock/ (search for evaluation methods and best practices)
**Human-in-the-Loop for Machine Learning:** Research papers and articles on incorporating human feedback into the evaluation and training of AI models.

---

**Question: 52** **Certy**IQ

A student at a university is copying content from generative AI to write essays.
Which challenge of responsible generative AI does this scenario represent?

A.Toxicity
B.Hallucinations
C.Plagiarism
D.Privacy

**Answer: C**

**Explanation:**

The correct answer is **C. Plagiarism**.

Plagiarism, in the context of generative AI, refers to the act of using AI-generated content without proper attribution or citation, presenting it as one's own original work. In the given scenario, the student is directly copying content from a generative AI model and incorporating it into essays without acknowledging the source. This directly violates academic integrity principles and constitutes plagiarism.

While other options might have tangential relevance, they don't directly address the core issue.

**Toxicity** refers to the generation of offensive, biased, or harmful content by the AI model. While possible with generative AI, the scenario doesn't explicitly mention the AI producing toxic outputs.
**Hallucinations** occur when the AI model generates information that is factually incorrect or nonsensical. Although AI models can hallucinate, the primary issue here is the student's unethical use of existing content, not the AI's factual accuracy.
**Privacy** concerns the protection of sensitive or personal information. The scenario doesn't involve the AI model's ability to expose private user data.

Therefore, plagiarism is the most direct and relevant challenge of responsible generative AI illustrated by the student's actions. It emphasizes the critical need for education and ethical guidelines surrounding AI usage in academic settings. Students need to be taught to properly cite AI-generated content and to use AI as a tool to aid their work, not to replace original thought and composition.

Further reading:

UNESCO AI Ethics Recommendations: https://www.unesco.org/en/articles/recommendation-ethics-artificial-intelligence
OECD Principles on AI: https://www.oecd.org/science/oecd-principles-on-ai-5ed8f95a-en.htm

## Question: 53

A company needs to build its own large language model (LLM) based on only the company's private data. The company is concerned about the environmental effect of the training process.
Which Amazon EC2 instance type has the LEAST environmental effect when training LLMs?

   A.Amazon EC2 C series

   B.Amazon EC2 G series

   C.Amazon EC2 P series

   D.Amazon EC2 Trn series

**Answer: D**

**Explanation:**

The correct answer is **D. Amazon EC2 Trn series.**

Here's why:

Amazon EC2 Trn series instances are specifically designed for high-performance deep learning training and inference, with a focus on efficiency and minimizing the environmental impact. They are powered by AWS Trainium chips, custom-built by AWS for deep learning workloads. Trainium is optimized for training deep learning models and excels in performance per watt.

**Energy Efficiency:** The Trainium chips in Trn instances are designed for energy efficiency. This means that for the same amount of computation, they consume less power compared to other GPU-based or CPU-based instances. Less power consumption directly translates to a lower environmental footprint, assuming the energy source has a carbon footprint.
**Purpose-Built for Deep Learning:** The Trainium architecture is tailored to deep learning workloads like LLM training. This specialization allows it to perform the required computations more efficiently than general-purpose instances or even GPU-based instances not specifically optimized for deep learning model training.
**Reduced Carbon Footprint:** By using less power and performing computations efficiently, Trn instances help reduce the carbon footprint of the training process. Training large language models can be exceptionally energy-intensive, so selecting an energy-efficient instance type is crucial for minimizing environmental impact.

Now let's consider why the other options are less suitable:

**Amazon EC2 C series:** These are compute-optimized instances that are designed for general-purpose compute-intensive workloads. They do not necessarily have the same level of specialized energy efficiency as the Trn series, particularly for deep learning training.
**Amazon EC2 G series:** These instances are designed for graphics-intensive workloads and machine learning inference. While they can be used for training, they are generally more focused on graphics and inference than deep learning training optimization.

**Amazon EC2 P series:** These are GPU-based instances that can be used for machine learning and other compute-intensive tasks. While useful for model training, the GPUs used are often not as energy efficient or optimized for deep learning as the Trainium chips found in Trn instances. The power consumption of GPUs used in P series instances can be relatively high compared to specialized accelerators like Trainium.

In summary, the Amazon EC2 Trn series instances, with their AWS Trainium chips, are the best option for a company aiming to minimize the environmental impact of training large language models due to their superior energy efficiency and specialization for deep learning training.

Supporting documentation:

AWS Trainium: https://aws.amazon.com/machine-learning/trainium/
Amazon EC2 Trn1 instances: https://aws.amazon.com/ec2/instance-types/trn1/

---

**Question: 54**                                                                                     **CertyIQ**

A company wants to build an interactive application for children that generates new stories based on classic stories. The company wants to use Amazon Bedrock and needs to ensure that the results and topics are appropriate for children.
Which AWS service or feature will meet these requirements?

   A.Amazon Rekognition

   B.Amazon Bedrock playgrounds

   C.Guardrails for Amazon Bedrock

   D.Agents for Amazon Bedrock

**Answer: C**

**Explanation:**

The correct answer is **C. Guardrails for Amazon Bedrock**.

Guardrails for Amazon Bedrock are specifically designed to implement safeguards and content moderation policies on large language models (LLMs) accessed through Amazon Bedrock. This is crucial for applications targeted at children, where content must be appropriate and safe. Guardrails enable you to set up deny lists (prohibiting specific keywords or topics), toxicity filters, and hallucination detection. By defining these safety parameters, you can ensure the stories generated by the application align with the desired content guidelines for children.

Amazon Rekognition (A) focuses on image and video analysis and wouldn't be suitable for moderating text generation. Amazon Bedrock Playgrounds (B) are for experimentation and model exploration but don't inherently provide content moderation features. Agents for Amazon Bedrock (D) allow LLMs to connect to external data sources and perform actions but don't address content safety needs directly.

Therefore, Guardrails for Amazon Bedrock is the most appropriate choice because it provides the necessary tools to control the output of the LLM, ensuring that the generated stories are safe and appropriate for children. It offers granular control over content, allowing you to proactively prevent the generation of unsuitable material.

For further research on Guardrails for Amazon Bedrock, refer to the official AWS documentation: https://aws.amazon.com/bedrock/guardrails/

---

**Question: 55**                                                                                     **CertyIQ**

A company is building an application that needs to generate synthetic data that is based on existing data. Which type of model can the company use to meet this requirement?

    A.Generative adversarial network (GAN)

    B.XGBoost

    C.Residual neural network

    D.WaveNet

**Answer: A**

**Explanation:**

The correct answer is A, Generative Adversarial Network (GAN). GANs are specifically designed to generate new, synthetic data that resembles an existing dataset. This makes them ideal for creating synthetic data based on real data, which is precisely what the company needs.

Here's why: GANs consist of two neural networks, a generator and a discriminator, that compete against each other. The generator creates synthetic data, and the discriminator tries to distinguish between real and synthetic data. Through this adversarial process, the generator learns to produce increasingly realistic synthetic data that fools the discriminator. This allows the creation of datasets suitable for use in scenarios where using the original data might not be permissible due to privacy or other reasons. XGBoost is a gradient boosting algorithm generally used for classification and regression, not synthetic data generation. Residual neural networks are typically used for image classification or similar tasks, not synthetic data. WaveNet is designed for generating audio, not generic synthetic data based on any given input dataset.

The company's requirement is to generate synthetic data based on existing data. GANs excel at this task because they learn the underlying distribution of the real data and generate new samples from that distribution. XGBoost, Residual neural networks, and WaveNet do not generate synthetic data in this manner. Therefore, GANs are the most suitable choice for the company.

For more information on GANs:

**Generative Adversarial Networks:** https://developers.google.com/machine-learning/gan
**GANs - Stanford CS230:** https://cs230.stanford.edu/section/generative-adversarial-networks-gans/

**Question: 56**                                                                 Certy**IQ**

A digital devices company wants to predict customer demand for memory hardware. The company does not have coding experience or knowledge of ML algorithms and needs to develop a data-driven predictive model. The company needs to perform analysis on internal data and external data.
Which solution will meet these requirements?

    A.Store the data in Amazon S3. Create ML models and demand forecast predictions by using Amazon SageMaker built-in algorithms that use the data from Amazon S3.

    B.Import the data into Amazon SageMaker Data Wrangler. Create ML models and demand forecast predictions by using SageMaker built-in algorithms.

    C.Import the data into Amazon SageMaker Data Wrangler. Build ML models and demand forecast predictions by using an Amazon Personalize Trending-Now recipe.

    D.Import the data into Amazon SageMaker Canvas. Build ML models and demand forecast predictions by selecting the values in the data from SageMaker Canvas.

**Answer: D**

**Explanation:**

The correct answer is D, utilizing Amazon SageMaker Canvas. Here's why:

SageMaker Canvas is specifically designed for business users without coding experience to build machine learning models. It offers a visual, point-and-click interface allowing users to directly select data and features without writing code, directly addressing the company's lack of coding expertise. Option A, while involving SageMaker, requires knowledge of algorithms and S3 data handling, which the company lacks. Options B and C use SageMaker Data Wrangler, which focuses on data preparation and also involve ML algorithm selections beyond the company's abilities. Personalize (Option C) is geared toward personalization/recommendation and might not be optimal for general demand forecasting.

Canvas abstracts away the complexities of ML algorithms by automatically suggesting appropriate models based on the data provided. Users can explore different features and observe their impact on model performance through a visual interface. The company can directly import both internal and external data into Canvas and perform the necessary analysis for demand forecasting using its intuitive features, such as selecting columns of data relevant to the forecast. Canvas simplifies the entire ML process, from data import to model deployment and prediction, making it the most suitable option for the given scenario, directly adhering to the requirement for no coding experience and the need to develop predictive models based on the data provided. It allows for the creation of models by simply selecting values in the data, aligning with the user's technical skill level.

Amazon SageMaker Canvas Documentation

---

## Question: 57 <span>CertyIQ</span>

A company has installed a security camera. The company uses an ML model to evaluate the security camera footage for potential thefts. The company has discovered that the model disproportionately flags people who are members of a specific ethnic group.
Which type of bias is affecting the model output?

    A.Measurement bias

    B.Sampling bias

    C.Observer bias

    D.Confirmation bias

**Answer: B**

**Explanation:**

Here's a breakdown of why the correct answer is sampling bias and why the other options are not suitable in the context of the question:

**Correct Answer: B. Sampling bias**

Sampling bias occurs when the data used to train a machine learning model does not accurately represent the overall population that the model will be applied to. In this case, the model disproportionately flags people of a specific ethnic group as potential thieves. This strongly suggests that the training dataset used to build the ML model contained an imbalance. For example, perhaps the dataset contained a significantly higher percentage of individuals from that ethnic group who were labeled as thieves (even if they weren't). Or, the data simply had a disproportionate representation of that ethnic group in various scenarios, leading the model to incorrectly associate certain features or patterns with that group and a higher likelihood of theft. Therefore, the sample of data used to train the model did not fairly represent the diversity of the actual population it's being used to monitor, leading to the biased output. If the model was trained on a dataset that showed a particular demographic committing theft more often, the model would learn that correlation, regardless of whether that correlation exists in the broader, real-world population.

**Why the other options are incorrect:**

**A. Measurement bias:** Measurement bias arises when the data itself is systematically inaccurate or flawed in a way that skews the results. For instance, if the camera quality was consistently worse when recording individuals of a certain skin tone, leading to inaccurate feature extraction. However, the problem described focuses on the data used to train the model, rather than inherent flaws in the captured footage, making sampling bias more likely.

**C. Observer bias:** Observer bias (also known as experimenter bias) happens when the researchers or individuals collecting and labeling the data unconsciously influence the data in a way that confirms their pre-existing beliefs or expectations. For example, if annotators were more likely to label actions by individuals of a certain ethnic group as "suspicious" due to their own biases. While observer bias could contribute to biased training data, the description explicitly points to the model disproportionately flagging a specific ethnic group, suggesting a deeper, more fundamental issue with the composition of the data.

**D. Confirmation bias:** Confirmation bias is the tendency to interpret new evidence as confirmation of one's existing beliefs or theories. While someone might experience confirmation bias when interpreting the model's results, it doesn't explain why the model is systematically producing biased outputs in the first place.

**Relevant Cloud Computing Concepts:**

**Data Governance:** Proper data governance practices are essential for mitigating bias in machine learning models. This includes ensuring data quality, diversity, and representativeness in training datasets.

**Model Monitoring:** Continuous monitoring of model performance is crucial for detecting and addressing bias. This involves tracking metrics for different demographic groups and investigating any discrepancies.

**Fairness Metrics:** Various fairness metrics can be used to quantify and assess bias in machine learning models. These metrics help ensure that the model is not unfairly discriminating against any particular group.

**Authoritative Links:**

**AI Fairness 360:** https://aif360.mybluemix.net/ - An open-source toolkit developed by IBM for detecting and mitigating bias in machine learning models.

**Google's Responsible AI Practices:** https://ai.google/responsibility/ - Provides guidance on building and deploying AI systems responsibly, including addressing fairness and bias.

**AWS AI & ML Responsible AI:** https://aws.amazon.com/machine-learning/responsible-ai/ - Provides tools and resources for responsible AI development on the AWS platform.

---

**Question: 58**                                                                      **CertyIQ**

A company is building a customer service chatbot. The company wants the chatbot to improve its responses by learning from past interactions and online resources.
Which AI learning strategy provides this self-improvement capability?

   A.Supervised learning with a manually curated dataset of good responses and bad responses

   B.Reinforcement learning with rewards for positive customer feedback

   C.Unsupervised learning to find clusters of similar customer inquiries

   D.Supervised learning with a continuously updated FAQ database

**Answer: B**

**Explanation:**

Here's a detailed justification for why option B (Reinforcement learning with rewards for positive customer feedback) is the most suitable AI learning strategy for the customer service chatbot scenario, along with supporting concepts and links:

The core requirement is for the chatbot to improve its responses over time by learning from past interactions. Reinforcement learning (RL) is specifically designed for this kind of iterative improvement through trial and

error. In RL, an agent (the chatbot) interacts with an environment (customer interactions), taking actions (generating responses), and receiving feedback in the form of rewards (positive customer feedback) or penalties (negative feedback or unresolved inquiries).

Option B directly leverages this mechanism. The chatbot learns a policy — a mapping from states (customer inquiries) to actions (responses) — that maximizes the cumulative reward. When a customer provides positive feedback (e.g., a high satisfaction rating, successful resolution), the chatbot's actions leading to that outcome are reinforced, making it more likely to produce similar responses in similar situations in the future.

Supervised learning (options A and D), while useful for initial training, relies on pre-labeled data. Option A's curated dataset is static and doesn't allow the chatbot to learn from live interactions. Option D, continuously updating the FAQ database, is helpful for knowledge updates but doesn't address the subtleties of conversational AI and adapting to nuanced user requests. The chatbot will not be able to adapt to new user queries.

Unsupervised learning (option C) can find clusters of similar inquiries, which can be helpful for topic modeling and routing, but it doesn't provide a mechanism for the chatbot to learn better responses. This is beneficial in discovering new topics for the chatbot to respond to, but it is not ideal for improving responses over time.

Reinforcement learning's capacity for adaptation and iterative improvement makes it ideally suited to the dynamic nature of customer service.https://aws.amazon.com/machine-learning/reinforcement-learning/https://azure.microsoft.com/en-us/resources/cloud-computing-dictionary/what-is-reinforcement-learning

## Question: 59

An AI practitioner has built a deep learning model to classify the types of materials in images. The AI practitioner now wants to measure the model performance.

Which metric will help the AI practitioner evaluate the performance of the model?

A.Confusion matrix

B.Correlation matrix

C.R2 score

D.Mean squared error (MSE)

**Answer: A**

**Explanation:**

Here's a detailed justification for why a confusion matrix is the best metric for evaluating the performance of a deep learning model classifying materials in images:

A confusion matrix is a powerful tool specifically designed for evaluating the performance of classification models. It provides a breakdown of the model's predictions, showing the counts of true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN) for each class. This allows the AI practitioner to see exactly where the model is making correct and incorrect predictions.

For a material classification problem, the confusion matrix would show how often the model correctly identified each material (e.g., how many images of "wood" were correctly classified as "wood" - TP) and how often it misclassified one material as another (e.g., how many images of "metal" were incorrectly classified as "plastic" - FP). This granular view is crucial for understanding the model's strengths and weaknesses. It allows for the calculation of other important performance metrics like precision, recall, F1-score, and accuracy for each class individually and overall.

In contrast, a correlation matrix measures the linear relationship between variables, which isn't relevant for evaluating the performance of a classification model. The R2 score (coefficient of determination) and mean squared error (MSE) are metrics primarily used for regression problems, where the goal is to predict a continuous value, not to classify discrete categories. They are not suitable for assessing the accuracy of a classification model attempting to identify different types of materials in images. Therefore, for this specific classification task, the confusion matrix provides the most insightful and actionable information about the model's performance. Analyzing the confusion matrix directly indicates what kinds of classification errors are happening most frequently, facilitating model refinement by adjusting training data, architecture, or hyperparameters.

**Authoritative Links for further research:**

**Confusion Matrix:** https://www.sciencedirect.com/topics/computer-science/confusion-matrix
**Classification Evaluation Metrics:** https://developers.google.com/machine-learning/crash-course/classification/true-false-positive-negative

---

## Question: 60                                                      CertyIQ

A company has built a chatbot that can respond to natural language questions with images. The company wants to ensure that the chatbot does not return inappropriate or unwanted images.
Which solution will meet these requirements?

    A.Implement moderation APIs.

    B.Retrain the model with a general public dataset.

    C.Perform model validation.

    D.Automate user feedback integration.

**Answer: A**

**Explanation:**

The correct answer is A, implementing moderation APIs. Here's why:

The primary goal is to prevent the chatbot from returning inappropriate or unwanted images. Moderation APIs are specifically designed to analyze images (and other content) and identify content that violates predefined policies or guidelines. These APIs use machine learning to detect explicit content, violence, hate speech, and other unwanted categories.

Option B, retraining the model with a general public dataset, is unlikely to directly address this problem and could potentially introduce new biases or unwanted content. A general public dataset might contain diverse data, but it's not guaranteed to exclude inappropriate images.

Option C, performing model validation, is a crucial step in model development, but it primarily focuses on evaluating the model's performance on its intended task, such as accuracy and efficiency. It does not directly address content moderation. Validation helps ensure the model works as expected, but it won't prevent the model from selecting inappropriate images if those images are present in the dataset or if the model isn't explicitly trained to avoid them.

Option D, automating user feedback integration, can provide valuable information about the chatbot's performance and user experience. While helpful for identifying areas for improvement, it relies on users to report inappropriate content after it has been shown, which is not a proactive solution to prevent such images from being displayed in the first place.

Implementing moderation APIs (like Amazon Rekognition Image Moderation) offers a proactive approach by filtering images before they are returned to the user. These APIs can be integrated into the chatbot's pipeline

to automatically screen images and flag those that violate specified criteria. This ensures inappropriate images are never presented to users, satisfying the company's requirements.

Therefore, using a moderation API provides the most direct and effective solution to the problem.

Relevant links for further research:

**Amazon Rekognition Image Moderation:** https://aws.amazon.com/rekognition/image-moderation/
**AWS AI Services:** https://aws.amazon.com/ai/

---

**Question: 61**

An AI practitioner is using an Amazon Bedrock base model to summarize session chats from the customer service department. The AI practitioner wants to store invocation logs to monitor model input and output data.
Which strategy should the AI practitioner use?

   A.Configure AWS CloudTrail as the logs destination for the model.

   B.Enable invocation logging in Amazon Bedrock.

   C.Configure AWS Audit Manager as the logs destination for the model.

   D.Configure model invocation logging in Amazon EventBridge.

**Answer: B**

**Explanation:**

The correct answer is **B. Enable invocation logging in Amazon Bedrock.**

Amazon Bedrock natively supports invocation logging to directly capture and store input and output data for model interactions. This functionality is built into the Bedrock service itself, providing a straightforward and efficient way to monitor model usage and performance. By enabling invocation logging, the AI practitioner can automatically record the data sent to and received from the base model, enabling auditing, debugging, and analysis of model behavior.

Option A is incorrect because while AWS CloudTrail tracks API calls made to Bedrock, it doesn't capture the detailed input and output data of the model invocations themselves. CloudTrail focuses on management events (e.g., creating or deleting models) rather than data plane operations.

Option C is incorrect because AWS Audit Manager is used for automating compliance audits and doesn't directly capture model invocation logs. It focuses on assessing and managing compliance across AWS services, not specifically on monitoring model input and output.

Option D is incorrect because Amazon EventBridge is an event bus service, primarily used for building event-driven applications. While EventBridge can be integrated with other AWS services, it is not the primary mechanism for enabling and storing invocation logs for Amazon Bedrock models. While you could potentially route invocation events to EventBridge after they are logged, it's an unnecessary complexity when Bedrock natively supports invocation logging.

Therefore, enabling invocation logging directly within Amazon Bedrock is the most appropriate and straightforward strategy for storing invocation logs and monitoring model input and output data. This allows the AI practitioner to directly access and analyze the model's interactions without relying on external services or complex configurations.

For more information, refer to the Amazon Bedrock documentation on invocation logging:https://docs.aws.amazon.com/bedrock/latest/userguide/security-data-protection.html

A company is building an ML model to analyze archived data. The company must perform inference on large datasets that are multiple GBs in size. The company does not need to access the model predictions immediately. Which Amazon SageMaker inference option will meet these requirements?

   A.Batch transform

   B.Real-time inference

   C.Serverless inference

   D.Asynchronous inference

**Answer: A**

**Explanation:**

The correct answer is **A. Batch transform**.

Here's why:

**Batch Transform** is specifically designed for performing inference on large datasets where immediate predictions are not required. It processes data in batches, making it suitable for scenarios involving archived data analysis. The inference results are stored for later retrieval.
https://docs.aws.amazon.com/sagemaker/latest/dg/batch-transform.html
**Real-time inference** is for low-latency predictions, serving requests in real time. This is unsuitable for the given scenario that does not need immediate access to the model predictions. Real-time inference involves deploying a model to a persistent endpoint.
**Serverless inference** is suitable for applications with intermittent or unpredictable traffic patterns. While serverless inference scales automatically, it is more ideal for applications needing low latency, which is not a requirement in the described scenario.
**Asynchronous inference** addresses workloads with large payload sizes, long processing times or near real-time latency needs. Although asynchronous inference provides asynchronous request processing and buffering capabilities, it might not be as efficient as batch transform when dealing with large datasets where high throughput is prioritized and lower latency is not a stringent need.

Since the company needs to analyze GBs of archived data without the need for immediate predictions, Batch Transform offers the most cost-effective and efficient solution. It avoids the overhead of maintaining a persistent endpoint (like real-time inference) while providing optimized processing for large batch datasets. Batch Transform is ideal for data which is already stored, as in archived data.

Which term describes the numerical representations of real-world objects and concepts that AI and natural language processing (NLP) models use to improve understanding of textual information?

   A.Embeddings

   B.Tokens

   C.Models

   D.Binaries

**Answer: A**

**Explanation:**

The correct answer is A, Embeddings. Here's why:

Embeddings are crucial in AI and NLP because they transform words, phrases, or even entire documents into numerical vectors. These vectors capture the semantic meaning and relationships between different elements of textual data. AI/NLP models, particularly deep learning models, inherently operate on numbers. Therefore, converting textual information into numerical representations is a prerequisite for these models to process and learn from text.

Tokens (Option B) are individual units of text, like words or subwords, often used as the initial step in preprocessing text. However, tokens themselves don't represent semantic meaning; they're merely the building blocks for creating embeddings.

Models (Option C) are the trained algorithms that use the embeddings (among other data) to perform tasks like text classification, sentiment analysis, or machine translation. They rely on embeddings to understand the input text. The model itself isn't the numerical representation; it uses the numerical representation generated by embeddings.

Binaries (Option D) refers to executable files or data represented in a binary format (0s and 1s). While machine learning models are ultimately stored as binaries, this option doesn't describe the numerical representation of real-world objects and concepts the models use. Binaries are a lower-level representation of the program code itself, not the semantic data.

Embeddings, on the other hand, leverage techniques like Word2Vec, GloVe, or Transformer-based models (like BERT) to create these numerical vectors. The position and direction of a vector in the embedding space encode semantic similarity; words with similar meanings will have vectors that are closer together. These vector representations allow the model to perform mathematical operations on the text data, enabling it to learn complex patterns and relationships. Embeddings help AI understand nuanced meanings and context, thereby improving performance in tasks involving textual information.

For further research, you can explore these resources:

**Word Embeddings (TensorFlow):** https://www.tensorflow.org/tutorials/text/word_embeddings
**Understanding Word Vectors:** https://www.analyticsvidhya.com/blog/2017/06/word-vectors-nlp-tutorial/
**GloVe: Global Vectors for Word Representation:** https://nlp.stanford.edu/projects/glove/

---

**Question: 64**                                                   **CertyIQ**

A research company implemented a chatbot by using a foundation model (FM) from Amazon Bedrock. The chatbot searches for answers to questions from a large database of research papers.
After multiple prompt engineering attempts, the company notices that the FM is performing poorly because of the complex scientific terms in the research papers.
How can the company improve the performance of the chatbot?

   A.Use few-shot prompting to define how the FM can answer the questions.

   B.Use domain adaptation fine-tuning to adapt the FM to complex scientific terms.

   C.Change the FM inference parameters.

   D.Clean the research paper data to remove complex scientific terms.

**Answer: B**

**Explanation:**

The best way to improve the chatbot's performance when struggling with complex scientific terms is **B. Use domain adaptation fine-tuning to adapt the FM to complex scientific terms.** Here's why:

The core issue is the foundation model's (FM) lack of understanding of specialized scientific vocabulary. Domain adaptation fine-tuning addresses this directly by training the FM on a dataset of research papers

containing those complex terms. This process adjusts the model's parameters to better recognize and process the specific language used in the research domain.

Here's why other options are less ideal:

**A. Use few-shot prompting to define how the FM can answer the questions:** Few-shot prompting can provide examples, but it's unlikely to fully overcome the FM's fundamental lack of understanding of the complex scientific terms. It's more suitable for guiding response style, not for fundamentally expanding the model's vocabulary or comprehension.

**C. Change the FM inference parameters:** Modifying parameters like temperature or top-p might influence response diversity or determinism, but it won't improve the model's comprehension of the scientific terms themselves. Inference parameters primarily control the generation of the response, not the understanding of the input.

**D. Clean the research paper data to remove complex scientific terms:** This defeats the purpose of the chatbot. The chatbot is supposed to answer questions about the research papers, which inherently contain complex scientific terms. Removing them would severely limit the chatbot's usefulness.

Fine-tuning, specifically domain adaptation, is designed for this exact scenario: improving a model's performance on a specific domain (in this case, scientific research) by training it on data from that domain. This allows the FM to learn the nuances and specific vocabulary of the research papers, leading to better understanding and more accurate answers. It's a targeted approach to improve comprehension, rather than merely tweaking response generation or simplifying the input data to an unacceptable degree. Fine-tuning allows the chatbot to properly utilize all information within the provided research data.

For further research on domain adaptation and fine-tuning, you can refer to these resources:

**Amazon Bedrock Documentation:** https://aws.amazon.com/bedrock/ (Look for sections related to model customization and fine-tuning. Specific documentation will depend on which FM in Bedrock is being used.)

**Research Papers on Domain Adaptation in NLP:** Search for academic papers on "domain adaptation" and "fine-tuning" in the context of Natural Language Processing (NLP). Sites like Google Scholar (https://scholar.google.com/) are great resources.

---

## Question: 65 <span style="float:right">CertyIQ</span>

A company wants to use a large language model (LLM) on Amazon Bedrock for sentiment analysis. The company needs the LLM to produce more consistent responses to the same input prompt.
Which adjustment to an inference parameter should the company make to meet these requirements?

    A.Decrease the temperature value.

    B.Increase the temperature value.

    C.Decrease the length of output tokens.

    D.Increase the maximum generation length.

**Answer: A**

**Explanation:**

The company seeks consistent LLM responses for sentiment analysis using Amazon Bedrock. The key to achieving this lies in understanding the "temperature" parameter.

The temperature parameter controls the randomness of the LLM's output. A higher temperature (closer to 1) introduces more randomness and creativity, making the model explore less probable word sequences. This

leads to more diverse but less predictable responses. Conversely, a lower temperature (closer to 0) makes the model more deterministic, favoring the most likely word sequences based on its training data. This results in more consistent and predictable outputs, ideal when seeking repeatable results for tasks like sentiment analysis.

Decreasing the temperature value forces the LLM to choose the most probable and predictable tokens. This minimizes variability and helps ensure that the model returns similar responses to the same prompt each time. This consistency is crucial for reliable sentiment scoring. Options B, C, and D do not contribute to increased consistency. Increasing the temperature introduces more randomness. Adjusting the length of output tokens impacts the length of the response, not its consistency. The maximum generation length limits the overall length of the generated text, which is unrelated to the consistency in responses to the same input.

Therefore, the best way to enforce consistent LLM responses for sentiment analysis on Amazon Bedrock is by decreasing the temperature value.

For further learning, you can consult the official Amazon Bedrock documentation regarding inference parameters: https://docs.aws.amazon.com/bedrock/latest/userguide/model-parameters.html

## Question: 66                                                                                    CertyIQ

A company wants to develop a large language model (LLM) application by using Amazon Bedrock and customer data that is uploaded to Amazon S3. The company's security policy states that each team can access data for only the team's own customers.
Which solution will meet these requirements?

A.Create an Amazon Bedrock custom service role for each team that has access to only the team's customer data.

B.Create a custom service role that has Amazon S3 access. Ask teams to specify the customer name on each Amazon Bedrock request.

C.Redact personal data in Amazon S3. Update the S3 bucket policy to allow team access to customer data.

D.Create one Amazon Bedrock role that has full Amazon S3 access. Create IAM roles for each team that have access to only each team's customer folders.

### Answer: A

### Explanation:

The most secure and compliant solution for accessing customer data in Amazon S3 through Amazon Bedrock, while adhering to team-based access restrictions, is to create separate custom service roles for each team. This approach aligns with the principle of least privilege, granting each team access only to the specific customer data they require.

Option A is superior because it implements fine-grained access control at the Amazon Bedrock service role level. By creating a dedicated role for each team and restricting its S3 access to only the team's customer data folders, the risk of unauthorized data access or modification is significantly minimized. This approach ensures data isolation and supports compliance with the company's security policy.

Option B is flawed because it relies on teams to manually specify customer names in each request. This is prone to human error and does not guarantee that teams will only access authorized data. It also introduces an auditing and monitoring challenge.

Option C, redacting personal data, might address PII concerns, but it doesn't prevent unauthorized access to other customer data within the S3 buckets. Updating the S3 bucket policy alone might not be sufficient to enforce the team-based access control efficiently and securely.

Option D, granting full S3 access to one Bedrock role and relying on IAM roles for teams, creates an overly permissive Bedrock role, violating the principle of least privilege. A compromised Bedrock role could lead to unauthorized access to all customer data.

Therefore, using separate Bedrock service roles with specific S3 access permissions for each team is the most secure and controlled way to meet the requirements. This approach is consistent with AWS best practices for security and access management. This reduces the blast radius of a potential compromise of a service role.

https://docs.aws.amazon.com/IAM/latest/UserGuide/id_roles.htmlhttps://docs.aws.amazon.com/IAM/latest/UserGu
practices.htmlhttps://aws.amazon.com/blogs/security/how-to-use-iam-to-grant-access-to-your-amazon-s3-
bucket/

---

**Question: 67**

A medical company deployed a disease detection model on Amazon Bedrock. To comply with privacy policies, the company wants to prevent the model from including personal patient information in its responses. The company also wants to receive notification when policy violations occur.
Which solution meets these requirements?

A.Use Amazon Macie to scan the model's output for sensitive data and set up alerts for potential violations.

B.Configure AWS CloudTrail to monitor the model's responses and create alerts for any detected personal information.

C.Use Guardrails for Amazon Bedrock to filter content. Set up Amazon CloudWatch alarms for notification of policy violations.

D.Implement Amazon SageMaker Model Monitor to detect data drift and receive alerts when model quality degrades.

**Answer: C**

**Explanation:**

The correct answer is C because it directly addresses the requirements of preventing the model from including personal patient information in its responses and providing notification when policy violations occur.

Here's a detailed breakdown:

**Guardrails for Amazon Bedrock:** Guardrails allow you to set up content filters based on defined policies. This ensures that the model's output adheres to the privacy requirements by filtering out any personal patient information before it is returned. This proactive content filtering prevents the model from violating privacy policies in the first place. (Source: https://aws.amazon.com/bedrock/guardrails/)

**Amazon CloudWatch Alarms:** By integrating CloudWatch with the Guardrails, you can receive notifications when a policy violation is detected. Guardrails log violation events, which CloudWatch alarms can monitor. When a violation occurs, the CloudWatch alarm will trigger a notification, alerting the company that the model has attempted to include personal information in its response. (Source: https://docs.aws.amazon.com/bedrock/latest/userguide/guardrails-cloudwatch.html)

Option A is incorrect because Amazon Macie is primarily designed for discovering and protecting sensitive data at rest in Amazon S3 and other data stores, not for real-time filtering of model responses.

Option B is incorrect because while CloudTrail monitors API activity, it doesn't provide the content filtering capabilities needed to prevent personal information from being included in model responses. It would only record that a response was generated, not whether it contained sensitive data.

Option D is incorrect because SageMaker Model Monitor focuses on detecting data drift and model quality degradation. It doesn't address the specific requirement of preventing and alerting on the inclusion of personal information in model responses.

## Question: 68

A company manually reviews all submitted resumes in PDF format. As the company grows, the company expects the volume of resumes to exceed the company's review capacity. The company needs an automated system to convert the PDF resumes into plain text format for additional processing.
Which AWS service meets this requirement?

A.Amazon Textract

B.Amazon Personalize

C.Amazon Lex

D.Amazon Transcribe

**Answer: A**

**Explanation:**

The correct answer is A, Amazon Textract. Here's why:

Amazon Textract is a fully managed AWS service that uses machine learning to automatically extract text, handwriting, and data from scanned documents. It goes beyond simple Optical Character Recognition (OCR) to identify and extract data from forms and tables as well. This makes it perfectly suited for converting PDF resumes, which often contain formatted text, tables (for work history), and even handwriting in some cases, into plain text.

The company needs to convert PDF resumes into plain text for further processing, likely for tasks such as parsing skills, experience, and contact information. Amazon Textract's ability to accurately extract text and structured data from documents makes it the ideal solution for this use case.

Options B, C, and D are not appropriate for this scenario. Amazon Personalize is a service for creating personalized recommendations. Amazon Lex is a service for building conversational interfaces using voice and text. Amazon Transcribe is a service for converting speech to text. None of these directly address the requirement of converting PDF documents to plain text.

Therefore, Amazon Textract directly addresses the need for automated PDF to text conversion, making it the correct choice.https://aws.amazon.com/textract/

## Question: 69

An education provider is building a question and answer application that uses a generative AI model to explain complex concepts. The education provider wants to automatically change the style of the model response depending on who is asking the question. The education provider will give the model the age range of the user who has asked the question.
Which solution meets these requirements with the LEAST implementation effort?

A.Fine-tune the model by using additional training data that is representative of the various age ranges that the application will support.

B.Add a role description to the prompt context that instructs the model of the age range that the response should target.

C.Use chain-of-thought reasoning to deduce the correct style and complexity for a response suitable for that

user.

D.Summarize the response text depending on the age of the user so that younger users receive shorter responses.

**Answer: B**

**Explanation:**

The correct answer is B: Add a role description to the prompt context that instructs the model of the age range that the response should target. This approach is the most efficient because it leverages the existing capabilities of generative AI models to adapt their output based on the provided context. Prompt engineering involves crafting effective prompts to guide the model's behavior. By including the age range in the prompt (e.g., "Explain this concept to a 10-year-old" or "Explain this concept as if addressing a college student"), the model can tailor its response accordingly.

Fine-tuning (option A) is a more involved process requiring significant data preparation and retraining of the model. While it can be effective, it is overkill for simply adjusting the style based on age range and carries a much higher implementation cost. Chain-of-thought reasoning (option C), while useful in certain scenarios, is not directly applicable to adapting writing style based on age; it is more suited for complex problem-solving. Summarizing the response (option D) after generation adds another layer of complexity and might lose important nuances present in the original output. It's also a post-processing step, adding to the overall latency.

Prompt engineering, specifically utilizing contextual information, allows for immediate and flexible control over the model's output style. It avoids the need for model retraining or complex post-processing, making it the solution with the least implementation effort and allowing for rapid iteration and experimentation with different prompt styles. This method leverages the models pre-trained understanding of language styles related to different audiences.

For further research on prompt engineering, refer to resources from leading AI providers such as OpenAI and Google AI, which often publish guides and best practices for effectively using their models. Look for resources on techniques like "zero-shot prompting" and "few-shot prompting," which are relevant to understanding how to influence model behavior through prompt design.

**Question: 70**

Which strategy evaluates the accuracy of a foundation model (FM) that is used in image classification tasks?

A.Calculate the total cost of resources used by the model.

B.Measure the model's accuracy against a predefined benchmark dataset.

C.Count the number of layers in the neural network.

D.Assess the color accuracy of images processed by the model.

**Answer: B**

**Explanation:**

The correct strategy for evaluating the accuracy of a foundation model used in image classification is to measure its performance against a predefined benchmark dataset (Option B). This approach allows for a quantifiable assessment of how well the model classifies images.

Here's why:

**Benchmark Datasets as Ground Truth:** Benchmark datasets, such as ImageNet or CIFAR, provide a known set

of images with established labels. These datasets serve as the "ground truth" for evaluating the model's predictions. By comparing the model's predicted labels with the actual labels in the benchmark dataset, we can determine its accuracy.

**Quantifiable Metrics:** Measuring against a benchmark dataset allows for the calculation of metrics like accuracy, precision, recall, and F1-score. These metrics provide a clear, quantifiable measure of the model's performance. For instance, accuracy represents the percentage of correctly classified images.

**Standardized Evaluation:** Using a common benchmark dataset enables standardized evaluation and comparison of different models. This allows researchers and practitioners to objectively compare the performance of different FMs and choose the most suitable one for their specific application.

**Reproducibility:** Evaluating models on publicly available benchmark datasets ensures reproducibility of results. Others can use the same dataset to verify the reported performance and compare it with their own models.

**Option A is incorrect** because cost, while important for deployment, doesn't directly reflect classification accuracy.

**Option C is incorrect** because the number of layers doesn't dictate accuracy; a deeper model isn't necessarily better.

**Option D is incorrect** because while color accuracy might be a factor in specific applications, the primary metric for image classification is the correctness of the identified object, not the color representation.

**Relevant Cloud Computing Concepts:**

**Model Evaluation Pipelines:** In a cloud environment like AWS, model evaluation is typically integrated into a model training and deployment pipeline. This pipeline often includes steps to automatically evaluate the model against a benchmark dataset and generate performance metrics.

**AWS Services:** AWS offers services like Amazon SageMaker, which provides tools for model evaluation and benchmarking.

**Authoritative Links:**

**Amazon SageMaker Model Monitor:** https://aws.amazon.com/sagemaker/model-monitor/
**Machine Learning Model Evaluation:** https://developers.google.com/machine-learning/crash-course/classification/check-your-understanding

---

**Question: 71**                                                     **Certy**IQ

An accounting firm wants to implement a large language model (LLM) to automate document processing. The firm must proceed responsibly to avoid potential harms.
What should the firm do when developing and deploying the LLM? (Choose two.)

   A.Include fairness metrics for model evaluation.
   B.Adjust the temperature parameter of the model.
   C.Modify the training data to mitigate bias.
   D.Avoid overfitting on the training data.
   E.Apply prompt engineering techniques.

**Answer: AC**

**Explanation:**

The prompt asks for the best two actions an accounting firm should take when developing and deploying an LLM for document processing, focusing on responsible AI practices to mitigate potential harms.

Option A, "Include fairness metrics for model evaluation," is crucial. LLMs can inherit biases from their training data, leading to discriminatory outcomes. Fairness metrics (e.g., demographic parity, equal opportunity) help quantify these biases and inform mitigation strategies. This aligns with responsible AI principles of ensuring equitable outcomes.

Option C, "Modify the training data to mitigate bias," is also essential. Training data often reflects societal biases. Actively addressing and mitigating these biases in the training dataset is a fundamental step towards building a fairer LLM. This can involve techniques like re-weighting data points, data augmentation to balance underrepresented groups, or carefully curating datasets to remove biased content. Modifying the training data is a preventative measure that directly addresses the source of potential bias in the LLM.

Option B, "Adjust the temperature parameter of the model," primarily affects the randomness of the model's output, not fairness or bias. A higher temperature leads to more creative, less predictable responses, while a lower temperature makes the responses more deterministic.

Option D, "Avoid overfitting on the training data," is a standard machine learning practice to ensure good generalization performance, but it doesn't directly address the specific concern of responsible AI and mitigating bias. While important for model performance, it's secondary to fairness considerations in this context.

Option E, "Apply prompt engineering techniques," focuses on crafting effective prompts to elicit desired responses from the LLM. While important for usability, prompt engineering doesn't inherently address bias or fairness concerns in the underlying model.

Therefore, A and C are the best choices because they directly address responsible AI principles by focusing on identifying and mitigating bias, ultimately working towards a fairer and more equitable AI system.

Further Reading:

AWS AI/ML Responsible AI: https://aws.amazon.com/machine-learning/responsible-ai/
Fairness Metrics: https://developers.google.com/machine-learning/fairness-overview

## Question: 72                                                    CertyIQ

A company is building an ML model. The company collected new data and analyzed the data by creating a correlation matrix, calculating statistics, and visualizing the data.
Which stage of the ML pipeline is the company currently in?

   A.Data pre-processing
   B.Feature engineering
   C.Exploratory data analysis
   D.Hyperparameter tuning

**Answer: C**

**Explanation:**

The answer is C: Exploratory Data Analysis (EDA).

The company's actions – creating a correlation matrix, calculating statistics, and visualizing the data – are all core components of EDA. EDA is the crucial initial stage in the machine learning pipeline focused on understanding the characteristics of the dataset. The primary goal is to uncover patterns, relationships, anomalies, and potential issues within the data before applying any modeling techniques.

**Correlation matrices:** Reveal relationships between different features. For example, identifying highly correlated features can inform feature selection or highlight potential multicollinearity issues.

**Calculating statistics:** Providing summary measures such as mean, median, standard deviation, and quartiles helps understand the distribution and central tendencies of the features. This helps uncover outliers and assess the general shape of data distributions.

**Visualizing the data:** Enables the identification of trends, clusters, and outliers. Histograms, scatter plots, and box plots are examples of visualizations used in EDA.

Data pre-processing (A) comes after EDA. EDA informs the data cleaning and transformation steps needed in pre-processing. Feature engineering (B) also follows EDA and involves creating new features from existing ones based on insights gained during EDA. Hyperparameter tuning (D) is specific to model training and optimization, and occurs much later in the pipeline, once the model is chosen and trained. Therefore, based on the actions described, the company is clearly engaged in the EDA phase of the machine learning pipeline.

Useful resources from AWS and others:

**AWS Documentation on ML Pipelines:** https://docs.aws.amazon.com/sagemaker/latest/dg/ml-pipeline.html
**Overview of EDA:** https://towardsdatascience.com/exploratory-data-analysis-8e25edb5906e
**SageMaker Jumpstart (AWS provides various examples using this tool, including EDA):** https://aws.amazon.com/sagemaker/jumpstart/

---

**Question: 73**

A company has documents that are missing some words because of a database error. The company wants to build an ML model that can suggest potential words to fill in the missing text.
Which type of model meets this requirement?

A.Topic modeling

B.Clustering models

C.Prescriptive ML models

D.BERT-based models

**Answer: D**

**Explanation:**

The correct answer is D, BERT-based models, because they are specifically designed for understanding and generating text in context, a crucial requirement for filling in missing words in documents.

Let's break down why the other options are not suitable:

**A. Topic Modeling:** Topic modeling, like Latent Dirichlet Allocation (LDA), aims to discover abstract "topics" within a collection of documents. It does not fill in missing words or understand the contextual relationships needed for accurate word prediction. It groups words and documents based on shared themes but offers no word-level prediction capability.

**B. Clustering Models:** Clustering algorithms (like K-Means) group data points based on similarity. While useful for segmenting documents based on content, clustering does not predict missing words or understand sentence structure. It focuses on grouping similar documents, not on linguistic relationships within individual sentences.

**C. Prescriptive ML Models:** Prescriptive models recommend actions to take based on predicted outcomes. These models are useful for decision-making processes but are not designed for natural language processing tasks like text completion. They focus on optimization strategies, not on understanding and generating text.

Now, let's explore why BERT is the right choice:

BERT (Bidirectional Encoder Representations from Transformers) is a transformer-based model pre-trained on a massive corpus of text data. One of BERT's pre-training tasks is "Masked Language Modeling" (MLM), where some words in the input are randomly masked, and the model is trained to predict the masked words based on the surrounding context. This pre-training makes BERT highly effective at understanding the relationships between words in a sentence and predicting missing words with high accuracy. BERT considers both the left and right context (bidirectional), leading to a more nuanced understanding compared to unidirectional language models. Finetuning a pre-trained BERT model on the specific document dataset of the company can further improve its accuracy in suggesting suitable words.

**In summary, BERT's architecture and pre-training objective (MLM) make it uniquely suited for the task of filling in missing words in text, understanding the context, and suggesting words that fit grammatically and semantically within the document.**

**Authoritative Links:**

BERT: https://arxiv.org/abs/1810.04805
Hugging Face Transformers Library: https://huggingface.co/transformers/model_doc/bert.html

## Question: 74

A company wants to display the total sales for its top-selling products across various retail locations in the past 12 months.
Which AWS solution should the company use to automate the generation of graphs?

    A.Amazon Q in Amazon EC2

    B.Amazon Q Developer

    C.Amazon Q in Amazon QuickSight

    D.Amazon Q in AWS Chatbot

**Answer: C**

**Explanation:**

The correct answer is C: Amazon Q in Amazon QuickSight. Here's why:

The problem requires automatically generating graphs of sales data across multiple locations over the past 12 months. Amazon QuickSight is a fully managed, cloud-native business intelligence (BI) service that makes it easy to visualize data and derive actionable insights. Integrating Amazon Q, an AI-powered assistant, into QuickSight enhances this capability by allowing users to ask natural language questions about their data and receive immediate visual answers in the form of graphs and other visualizations.

Option A, Amazon Q in Amazon EC2, is incorrect because Amazon EC2 provides virtual servers in the cloud but lacks built-in BI and visualization functionalities. While you could potentially build a visualization solution on EC2, it wouldn't be as efficient or cost-effective as using QuickSight. Amazon Q in this context would primarily assist with server-related tasks, not data visualization.

Option B, Amazon Q Developer, focuses on aiding software developers with code generation, debugging, and optimization. It's not designed for business intelligence or data visualization.

Option D, Amazon Q in AWS Chatbot, facilitates interaction with AWS services through a chatbot interface. While useful for operational tasks, it doesn't directly address the need for automated graph generation.

Amazon Q in QuickSight directly solves the problem by allowing business users to ask questions like "Show me total sales for top-selling products by location in the last 12 months" and have QuickSight automatically

generate the corresponding graphs using the underlying data. This empowers users without requiring them to have deep technical expertise in data analysis or visualization.

Therefore, QuickSight with Amazon Q is the optimal solution for automating the generation of sales graphs as it provides a seamless way to query data, create visualizations, and gain insights.

Further reading:

**Amazon QuickSight:** https://aws.amazon.com/quicksight/
**Amazon Q:** https://aws.amazon.com/amazon-q/

---

**Question: 75**

A company is building a chatbot to improve user experience. The company is using a large language model (LLM) from Amazon Bedrock for intent detection. The company wants to use few-shot learning to improve intent detection accuracy.
Which additional data does the company need to meet these requirements?

   A.Pairs of chatbot responses and correct user intents

   B.Pairs of user messages and correct chatbot responses

   C.Pairs of user messages and correct user intents

   D.Pairs of user intents and correct chatbot responses

**Answer: C**

**Explanation:**

The correct answer is C: Pairs of user messages and correct user intents. Here's why:

Few-shot learning involves training a machine learning model (like an LLM in Amazon Bedrock) to perform a task using only a small number of examples. In the context of intent detection for a chatbot, the goal is to accurately map user messages (the input) to the intended action or topic (the output).

To effectively utilize few-shot learning for intent detection, the model needs examples demonstrating the relationship between user inputs (messages) and their corresponding intended intents. These examples serve as the "few shots" that guide the LLM in understanding the desired mapping. For example, "I want to book a flight to Seattle" (user message) should be paired with "book_flight" (user intent).

Option A, pairs of chatbot responses and correct user intents, is incorrect. This data focuses on the chatbot's output rather than the user's input, which is crucial for intent detection.

Option B, pairs of user messages and correct chatbot responses, is also incorrect. While useful for training the chatbot's response generation, it doesn't directly teach the model to understand the intent behind the user message. It's more related to end-to-end dialogue modeling.

Option D, pairs of user intents and correct chatbot responses, is incorrect because it reverses the desired mapping. The model needs to learn to predict the intent from the user message, not the other way around.

Therefore, having pairs of user messages and their corresponding intents is the most relevant and necessary data for few-shot learning to improve intent detection accuracy using an LLM from Amazon Bedrock. This type of data allows the model to learn from a limited number of examples how to accurately classify user messages based on the underlying intent.

Further Reading:

**Amazon Bedrock Documentation:** https://aws.amazon.com/bedrock/ (for general understanding of the

service)

**Few-Shot Learning:** Although AWS documentation may not explicitly detail "few-shot learning" directly within the Bedrock context, searching for general resources on "few-shot learning for NLP" will provide a good understanding of the technique used.

---

**Question: 76**

A company is using few-shot prompting on a base model that is hosted on Amazon Bedrock. The model currently uses 10 examples in the prompt. The model is invoked once daily and is performing well. The company wants to lower the monthly cost.
Which solution will meet these requirements?

    A.Customize the model by using fine-tuning.

    B.Decrease the number of tokens in the prompt.

    C.Increase the number of tokens in the prompt.

    D.Use Provisioned Throughput.

**Answer: B**

**Explanation:**

The correct answer is B: Decrease the number of tokens in the prompt. Here's why:

The goal is to lower the monthly cost associated with using the model on Amazon Bedrock. Since the model is invoked once daily, the primary cost driver will be the amount of computation required for each invocation. The size of the prompt directly impacts the computation cost. Larger prompts, containing more tokens, require the model to process more data, thus increasing computation time and cost.

Few-shot prompting involves providing a few examples to guide the model. While 10 examples provide a reasonable starting point, the company can likely achieve acceptable performance with fewer examples, especially since the model is already performing well.

Decreasing the number of tokens in the prompt by reducing the number of examples reduces the amount of data processed during each invocation. This translates directly into lower computation costs and, therefore, lower monthly expenses.

Option A, fine-tuning, is generally more expensive upfront than simply adjusting the prompt. Fine-tuning requires additional training data and computational resources to adapt the base model. While it might lead to better performance with a smaller prompt later, it introduces significant initial costs.

Option C, increasing the number of tokens, is counterproductive as it would increase the computational burden and cost.

Option D, using Provisioned Throughput, guarantees capacity and can be cost-effective if usage is high and predictable. However, since the model is invoked only once daily, the cost benefits of Provisioned Throughput are unlikely to outweigh the increased cost associated with reserving resources. In fact, for low usage scenarios it's typically more expensive than on-demand pricing.

Therefore, minimizing the prompt size by reducing the number of examples is the most straightforward and cost-effective way to lower monthly expenses without sacrificing performance.

Further Research:

Amazon Bedrock Pricing: For understanding the cost components of using Bedrock.
Prompt Engineering Techniques: Learn more about optimizing prompts for LLMs.

An AI practitioner is using a large language model (LLM) to create content for marketing campaigns. The generated content sounds plausible and factual but is incorrect.
Which problem is the LLM having?

    A.Data leakage

    B.Hallucination

    C.Overfitting

    D.Underfitting

**Answer: B**

**Explanation:**

The correct answer is B, hallucination. Here's a detailed justification:

Large Language Models (LLMs) are trained on massive datasets and learn to generate text by predicting the next word in a sequence. While they are powerful at mimicking human-like writing styles and can create seemingly coherent content, they do not possess true understanding or access to a verified knowledge base.

Hallucination, in the context of LLMs, refers to the phenomenon where the model generates content that is factually incorrect, nonsensical, or unrelated to the input prompt, while presenting it as if it were true. The generated content might sound convincing, making it difficult to discern its validity.

In this scenario, the LLM produces "plausible and factual" content that is, in reality, incorrect. This directly aligns with the definition of hallucination. The model is essentially making things up or providing information that is not grounded in reality, even though it presents it convincingly.

Here's why the other options are less likely:

**A. Data Leakage:** Data leakage occurs when information from the training dataset inadvertently appears in the output of the model. While this could manifest as incorrect information, it's not the primary characteristic of the described problem. Data leakage is more about revealing specific training data, not necessarily fabricating new information.
**C. Overfitting:** Overfitting happens when a model learns the training data too well, including the noise and outliers. This results in poor generalization to new, unseen data. While overfitting can lead to inaccuracies, it typically manifests as poor performance on novel data points similar to training data, rather than the fabrication of entirely new, incorrect "facts."
**D. Underfitting:** Underfitting occurs when a model is too simple to capture the underlying patterns in the data. This leads to poor performance on both the training and testing data. Underfitting would result in generalized poor quality of generated content, but not specifically in the invention of plausible, yet incorrect, information.

**Authoritative Links:**

**Towards Data Science - Hallucination in Large Language Models:**
https://towardsdatascience.com/hallucination-in-large-language-models-explained-a98c6429d291
**VentureBeat - AI Hallucinations:** https://venturebeat.com/ai/why-ai-chatbots-lie-the-problem-of-ai-hallucinations/

Therefore, hallucination is the most appropriate explanation for the LLM generating plausible but incorrect content for marketing campaigns.

An AI practitioner trained a custom model on Amazon Bedrock by using a training dataset that contains confidential data. The AI practitioner wants to ensure that the custom model does not generate inference responses based on confidential data.
How should the AI practitioner prevent responses based on confidential data?

    A.Delete the custom model. Remove the confidential data from the training dataset. Retrain the custom model.

    B.Mask the confidential data in the inference responses by using dynamic data masking.

    C.Encrypt the confidential data in the inference responses by using Amazon SageMaker.

    D.Encrypt the confidential data in the custom model by using AWS Key Management Service (AWS KMS).

**Answer: B**

**Explanation:**

The correct answer is **B. Mask the confidential data in the inference responses by using dynamic data masking.**

Here's a detailed justification:

The primary goal is to prevent the exposure of confidential data in the model's output, not to simply secure the model itself or its training data after the fact (though those are important security considerations). The model has already been trained with potentially sensitive data. Options A, C, and D don't directly address this specific problem of preventing confidential information from appearing in the model's inference outputs.

**Why B is correct:** Dynamic data masking focuses on modifying the output data at runtime, i.e., when the model is generating responses. This allows you to identify and redact or replace confidential elements in the generated text before it's presented to the user. This ensures that even if the model has learned some confidential data, it will not expose it to the outside world. This does not require re-training the model (expensive) and works to protect data at output.

**Why A is incorrect:** While retraining with a cleaned dataset is a good practice for long-term model improvement and data privacy, it's a time-consuming and resource-intensive approach. It also doesn't provide immediate protection against existing risks associated with the current model. There is no guarantee that a future model will not learn some level of confidential data based on surrounding context.

**Why C is incorrect:** Encryption of data in inference responses with Amazon SageMaker does not mask or redact sensitive information; instead, it transforms the data into an unreadable format without the appropriate decryption key. It prevents unauthorized access to the entire response but wouldn't allow specific elements within the response (that aren't confidential) to be viewed. The original confidential information would still be present, just encrypted. This is often not the desired outcome.

**Why D is incorrect:** Encrypting the entire custom model itself with KMS protects the model artifact from unauthorized access. However, it does not prevent the model from generating confidential data in its responses. Encryption alone doesn't modify the output of the model, it protects it. The confidential data would still be generated.

In summary, dynamic data masking is the most direct and effective technique to prevent the existing trained model from generating responses that contain confidential data. Other options involve more drastic measures or simply do not solve the problem.

**Supporting resources:**

While AWS doesn't offer a specific "dynamic data masking" service for Bedrock, the concept is widely used in data governance and security. The functionality of achieving dynamic data masking for Bedrock would involve building custom logic to analyze generated text responses.

For more about secure data handling and governance, refer to AWS's security best practices: https://aws.amazon.com/security/
Explore AWS KMS for encryption basics: https://aws.amazon.com/kms/

## Question: 79                                                      CertyIQ

A company has built a solution by using generative AI. The solution uses large language models (LLMs) to translate training manuals from English into other languages. The company wants to evaluate the accuracy of the solution by examining the text generated for the manuals.
Which model evaluation strategy meets these requirements?

    A.Bilingual Evaluation Understudy (BLEU)

    B.Root mean squared error (RMSE)

    C.Recall-Oriented Understudy for Gisting Evaluation (ROUGE)

    D.F1 score

**Answer: A**

**Explanation:**

The correct answer is A, Bilingual Evaluation Understudy (BLEU). Here's a detailed justification:

BLEU is specifically designed for evaluating the quality of machine-translated text. It works by comparing the generated translation to one or more reference translations, calculating a score based on n-gram precision. In essence, BLEU measures how many words and phrases in the machine-generated text are also present in the human-created reference translations. A higher BLEU score indicates better similarity between the generated and reference texts, thus better translation accuracy. In the context of translating training manuals, the company needs a metric to assess how well the LLM is capturing the meaning and nuances of the original English manuals in the target languages. BLEU effectively serves this purpose.

Option B, Root Mean Squared Error (RMSE), is used for evaluating regression models, where the goal is to predict a numerical value. It measures the average magnitude of the error between predicted and actual values. This is irrelevant for text translation, where the output is text, not numerical data.

Option C, Recall-Oriented Understudy for Gisting Evaluation (ROUGE), is primarily used for evaluating summarization tasks. While it does compare generated text to reference text, it focuses on recall (how much of the reference text is present in the generated text) rather than precision like BLEU. While ROUGE could provide some insights, BLEU is more directly suited for translation quality assessment.

Option D, F1 score, is a metric used for evaluating classification models. It is the harmonic mean of precision and recall, providing a balanced measure of the model's performance. Like RMSE, F1 score isn't appropriate for evaluating text translation because the output is not a classification.

In summary, BLEU is the most appropriate metric because it directly addresses the task of evaluating the quality of machine-generated translations by comparing them to reference translations, making it ideal for assessing the accuracy of the company's LLM-powered training manual translation solution.

Further Research:

**BLEU:** https://aclanthology.org/P02-1040/ (Original BLEU paper)
**Evaluating Text Output (Google):** https://developers.google.com/machine-learning/crash-course/classification/check-your-work (Although focused on classification, it provides context on different evaluation metrics.)

**Question: 80**

A large retailer receives thousands of customer support inquiries about products every day. The customer support inquiries need to be processed and responded to quickly. The company wants to implement Agents for Amazon Bedrock.
What are the key benefits of using Amazon Bedrock agents that could help this retailer?

    A.Generation of custom foundation models (FMs) to predict customer needs

    B.Automation of repetitive tasks and orchestration of complex workflows

    C.Automatically calling multiple foundation models (FMs) and consolidating the results

    D.Selecting the foundation model (FM) based on predefined criteria and metrics

**Answer: B**

**Explanation:**

The correct answer is **B. Automation of repetitive tasks and orchestration of complex workflows.**

Here's why:

Agents for Amazon Bedrock are designed to automate tasks and streamline complex workflows. In the context of customer support inquiries, this translates to agents being able to automatically categorize inquiries, extract relevant information, query databases for product details or customer history, and even draft responses. This automation significantly reduces the workload on human agents, allowing them to focus on more complex or nuanced issues.

Option A, Generation of custom foundation models (FMs) to predict customer needs, is not the primary function of Agents for Amazon Bedrock. While Bedrock allows for customization of FMs, Agents are about orchestrating actions using these models, not building them from scratch.

Option C, Automatically calling multiple foundation models (FMs) and consolidating the results, is a capability of Bedrock, but not specifically the key benefit provided by Agents. Agents use the models available through Bedrock.

Option D, Selecting the foundation model (FM) based on predefined criteria and metrics, is also functionality within Bedrock, but is not as central to the purpose of agents. The primary value of agents is in automating complex operations, which includes choosing, invoking and managing FMs.

In summary, Agents for Amazon Bedrock directly address the retailer's need to quickly process and respond to a high volume of customer support inquiries by automating repetitive tasks involved in understanding the request, finding solutions, and delivering responses. This automated orchestration optimizes the workflow, reducing response times and improving customer satisfaction.

For further research:

**AWS documentation on Agents for Amazon Bedrock:** https://aws.amazon.com/bedrock/agents/
**Blog post on Agents for Amazon Bedrock:** https://aws.amazon.com/blogs/aws/agents-for-amazon-bedrock-build-agents-that-complete-tasks-for-you/

# Thank you

Thank you for being so interested in the premium exam material.
I'm glad to hear that you found it informative and helpful.

## But Wait

I wanted to let you know that there is more content available in the full version.
The full paper contains additional sections and information that you may find helpful,
and I encourage you to download it to get a more comprehensive and detailed view of
all the subject matter.

**Download Full Version Now**

**Future is Secured**

100% Pass Guarantee

**24/7 Customer Support**

Mail us - certyiqofficial@gmail.com

**Free Updates**

Lifetime Free Updates!

Total: **334 Questions**

Link: https://certyiq.com/papers/amazon/aws-certified-ai-practitioner-aif-c01