

1) What are the key differences between discriminative and generative models?

a) Discriminative models learn the decision boundary between classes and patterns that differentiate them. They estimate the probability $P(y|x)$, which is the probability of a particular label y , given the input data x . These models focus on distinguishing between different categories.

Generative models learn the distribution of the data itself by modeling the joint probability $P(x,y)$, which involves sampling data points from this distribution. After being trained on thousands of images of digits, this sampling could produce a new image of a digit.

2) What are some challenges associated with training and evaluating generative AI models?

a) **Computational cost:** High [computational power](#) and hardware requirements for training more complex models.

- **Training complexity:** Training generative models could be challenging and full of nuances.
- **Evaluation metrics:** It's challenging to quantitatively assess the quality and diversity of the model outputs.
- **Data requirements:** Generative models often require massive amounts of data with high quality and diversity. The collection of such data could be time-consuming and expensive.
- **Bias and fairness:** Unchecked models can amplify the [biases](#) present in the training data, leading to unfair outputs

3) How can you control the style or attributes of generated content using generative AI models?

a) **Prompt engineering:** Specify the desired output style by providing detailed prompts highlighting the style or the tone of the content generation. This is an effective and simple method in both text-to-text and text-to-image models. It is a much more effective method if you do it in alignment with the specific requirements or the documentation of the particular model in question.

b) **Temperature and sampling control:** The *temperature* parameter controls how random the outputs would be. Lower temperatures mean a more

conservative and predictable token selection, and higher temperature allows more creative generation. Other parameters such as *top-k* and *top-p* can also control how creatively the model selects possible next tokens while generating.

- c) **Style transfer (Images):** Another technique that can be used during inference for the models that support it, is to apply the style of one image (reference image) to an input image.
 - d) **Fine-tuning:** We can use a pretrained model and fine-tune it on a specific dataset containing the style or tone that is desired. This means training the model further on additional data to learn additional specific styles or attributes.
 - e) **Reinforcement learning:** We can guide the model to prefer certain outputs and steer away from other outputs by providing feedback. This feedback will be used to modify the model through reinforcement learning. Over time, the model will be aligned to the preferences of the users and/or preference datasets. An example of this, in the context of LLMs, is Reinforcement learning from human feedback (RLHF).
- 4) How can prompts be strategically designed to elicit desired behaviors or outputs from the model? What are some best practices for effective prompt engineering?

Prompting is important in directing LLMs to respond to specific tasks. Effective prompts can even mitigate the need for fine-tuning models by using techniques such as few-shot learning, task decomposition, and prompt templates.

Some **best practices** for effective, prompt engineering include:

1. **Be clear and concise:** Provide specific instructions so the model knows exactly what task you want it to perform. Be straightforward and to-the-point.
2. **Use examples:** For in-context learning, showing a few input-output pairs helps the model understand the task the way you would like.
3. **Break down complex tasks:** If the task is complicated, breaking it into smaller steps can improve the quality of the response.
4. **Set constraints or formats:** If you need a specific output style, format, or length, clearly state those requirements within the prompt.

Interview Questions

Easy Questions

Question 1

****Q:**** What is the difference between TensorFlow and PyTorch?

****A:**** TensorFlow is a symbolic math library used for machine learning applications such as neural networks. It allows for both high-level and low-level API. PyTorch, on the other hand, is an open-source machine learning library based on the Torch library, used for applications such as computer vision and natural language processing. PyTorch is known for its dynamic computation graph, which makes it more intuitive and easier to debug.

Question 2

****Q:**** What is pruning in model compression?

****A:**** Pruning is a technique used to reduce the size of a neural network by removing weights that are less important. This helps in reducing the model size and improving inference speed without significantly affecting the model's performance.

Question 3

****Q:**** What is a language model in NLP?

****A:**** A language model is a statistical model that calculates the probability of a sequence of words. It is used to predict the next word in a sentence, generate text, and improve the performance of various NLP tasks.

Question 4

****Q:**** What is the purpose of tokenization in NLP?

****A:**** Tokenization is the process of breaking down text into smaller units, such as words or subwords, called tokens. It is a crucial step in NLP as it prepares the text for further processing by machine learning models.

Question 5

****Q:**** What is data preprocessing in machine learning?

****A:**** Data preprocessing is the process of transforming raw data into a format that can be used by machine learning models. It includes steps such as data cleaning, normalization, transformation, and feature extraction to improve the quality and performance of the model.

Question 6

****Q:**** What is Hugging Face?

****A:**** Hugging Face is a company that provides open-source tools and libraries for natural language processing (NLP). Their most popular library, Transformers, allows users to easily access and use pre-trained models for various NLP tasks such as text classification, translation, and question answering.

Medium Questions

Question 1

****Q:**** Explain the concept of transfer learning in NLP.

****A:**** Transfer learning in NLP involves taking a pre-trained model on a large dataset and fine-tuning it on a smaller, task-specific dataset. This approach leverages the knowledge gained from the large dataset to improve performance on the smaller dataset, often leading to better results than training from scratch.

Question 2

****Q:**** How does quantization help in model compression?

****A:**** Quantization reduces the number of bits required to represent each weight in the model, typically from 32-bit floating-point to 8-bit integers. This reduces the model size and can speed up inference by allowing the use of more efficient hardware operations.

Question 3

****Q:**** What is the difference between supervised and unsupervised learning?

****A:**** Supervised learning involves training a model on labeled data, where the input data is paired with the correct output. Unsupervised learning, on the other hand, involves training a model on unlabeled data, where the model tries to find patterns and relationships within the data without any explicit guidance.

Question 4

****Q:**** What are the advantages of using BERT for NLP tasks?

****A:**** BERT (Bidirectional Encoder Representations from Transformers) captures context from both directions (left and right) in a sentence, making it more effective for understanding the meaning of words in context. This bidirectional approach leads to better performance on various NLP tasks such as question answering, text classification, and named entity recognition.

Question 5

****Q:**** How does pruning improve model performance?

****A:**** Pruning improves model performance by reducing the number of parameters in the neural network, which can lead to faster inference times and lower memory usage. By removing less important weights, the model becomes more efficient without significantly affecting its accuracy.

Question 6

****Q:**** How do you handle missing data in a dataset?

****A:**** Handling missing data can be done using several techniques, such as:

- Removing rows or columns with missing values.
- Imputing missing values using statistical methods like mean, median, or mode.
- Using algorithms that support missing values, such as decision trees.
- Applying advanced imputation techniques like K-Nearest Neighbors (KNN) or regression imputation.

Hard Questions

Question 1

****Q:**** Describe the architecture of GPT and how it differs from BERT.

****A:**** GPT (Generative Pre-trained Transformer) is a unidirectional transformer model that generates text by predicting the next word in a sequence. It is trained using a left-to-right context. BERT (Bidirectional Encoder Representations from Transformers), on the other hand, is a bidirectional model that considers both the left and right context

of a word during training. This allows BERT to capture more context and generally perform better on tasks like question answering and text classification.

Question 2

****Q:**** How would you approach a text classification problem in the Telecom domain?

****A:**** First, I would gather and preprocess the relevant data, including cleaning and tokenizing the text. Next, I would choose a suitable model architecture, such as BERT or another transformer-based model, and fine-tune it on the Telecom-specific dataset. I would also consider using domain-specific embeddings if available. Finally, I would evaluate the model using appropriate metrics and iterate on the preprocessing and model tuning steps to improve performance.

Question 3

****Q:**** Explain the concept of model distillation and its benefits.

****A:**** Model distillation is a technique where a smaller, simpler model (student) is trained to replicate the behavior of a larger, more complex model (teacher). The student model learns from the teacher model's predictions, which helps it achieve similar performance with reduced computational resources. The benefits of model distillation include faster inference times, reduced memory usage, and the ability to deploy models on resource-constrained devices.

Question 4

****Q:**** How do you handle imbalanced datasets in machine learning?

****A:**** Handling imbalanced datasets can be done using several techniques, such as:

- Resampling the dataset by oversampling the minority class or undersampling the majority class.
- Using different evaluation metrics like precision, recall, and F1-score instead of accuracy.
- Applying algorithms that are robust to imbalanced data, such as ensemble methods (e.g., Random Forest, Gradient Boosting).
- Using synthetic data generation techniques like SMOTE (Synthetic Minority Over-sampling Technique) to create balanced datasets.

Question 5

****Q:**** What are the benefits of using Hugging Face's Transformers library?

****A:**** The benefits of using Hugging Face's Transformers library include:

- Access to a wide range of pre-trained models for various NLP tasks.
- Easy-to-use API for model training, fine-tuning, and inference.
- Support for multiple deep learning frameworks, including TensorFlow and PyTorch.
- Active community and extensive documentation for troubleshooting and learning.

Question 6

****Q:**** Explain the process of text classification using a pre-trained model from Hugging Face.

****A:**** The process of text classification using a pre-trained model from Hugging Face involves the following steps:

1. Load a pre-trained model and tokenizer from the Hugging Face library.
2. Preprocess the input text data using the tokenizer.
3. Fine-tune the pre-trained model on the specific text classification dataset.
4. Evaluate the model's performance using appropriate metrics.
5. Use the fine-tuned model to make predictions on new text data.