

Technical interview questions with sample answers tailored for Junior AI Specialist roles:

1. What is the difference between supervised and unsupervised learning?

Answer:

- **Supervised Learning:** The model is trained on labeled data, meaning the input comes with the corresponding output. Example algorithms: Linear Regression, Decision Trees.
 - **Unsupervised Learning:** The model works with unlabeled data to find patterns or structure in the data. Example algorithms: K-Means Clustering, Principal Component Analysis.
-

2. Explain overfitting and underfitting in machine learning. How can you prevent them?

Answer:

- **Overfitting:** The model performs well on training data but poorly on test data due to excessive complexity.
 - **Underfitting:** The model performs poorly on both training and test data due to insufficient complexity.
 - **Prevention:** Use techniques like cross-validation, regularization (L1, L2), pruning decision trees, and ensuring the dataset size is sufficient.
-

3. What is a confusion matrix?

Answer: A confusion matrix is a table used to evaluate the performance of a classification model. It includes:

- **True Positives (TP):** Correctly predicted positives.
- **True Negatives (TN):** Correctly predicted negatives.
- **False Positives (FP):** Incorrectly predicted as positive.
- **False Negatives (FN):** Incorrectly predicted as negative.

From this, you can calculate metrics like accuracy, precision, recall, and F1-score.

4. What is gradient descent, and why is it important?

Answer: Gradient descent is an optimization algorithm used to minimize the cost function in machine learning models. It iteratively adjusts the model parameters (weights) in the direction of the negative gradient to find the optimal values that reduce prediction errors.

5. Can you explain the difference between AI, Machine Learning, and Deep Learning?

Answer:

- **AI (Artificial Intelligence):** A broad field focused on building systems that simulate human intelligence.
 - **Machine Learning:** A subset of AI that uses algorithms to learn from data.
 - **Deep Learning:** A subset of machine learning that uses neural networks with multiple layers to model complex patterns in data.
-

6. What is a neural network?

Answer: A neural network is a machine learning model inspired by the human brain, consisting of layers of nodes:

- **Input Layer:** Accepts the input data.
 - **Hidden Layers:** Processes the data using weights and activation functions.
 - **Output Layer:** Produces the prediction or result.
-

7. What are some commonly used activation functions in neural networks?

Answer:

- **Sigmoid:** Outputs values between 0 and 1, often used in binary classification.
 - **ReLU (Rectified Linear Unit):** Sets negative values to 0, commonly used in deep learning models.
 - **Softmax:** Converts logits into probabilities for multi-class classification.
-

8. Explain the role of regularization in machine learning.

Answer: Regularization techniques prevent overfitting by penalizing large weights in the model. Types include:

- **L1 Regularization (Lasso):** Adds the absolute value of weights to the loss function.
 - **L2 Regularization (Ridge):** Adds the squared value of weights to the loss function.
-

9. How would you handle missing data in a dataset?

Answer:

- Remove rows or columns with missing values (if the amount of missing data is small).
- Impute missing values using statistical methods like mean, median, or mode.
- Use advanced techniques like K-Nearest Neighbors (KNN) or model-based imputation.

10. What is cross-validation, and why is it used?

Answer: Cross-validation splits the dataset into training and validation sets to evaluate the model's performance. Common techniques include:

- **K-Fold Cross-Validation:** Splits the data into k subsets, trains on k-1 subsets, and tests on the remaining one.
 - **Leave-One-Out Cross-Validation:** Uses one observation as the test set and the rest for training.
-

11. Explain the concept of feature scaling.

Answer: Feature scaling normalizes or standardizes data to ensure all features contribute equally to the model. Common methods:

- **Normalization:** Rescales values to a range of 0 to 1.
 - **Standardization:** Rescales data to have a mean of 0 and a standard deviation of 1.
-

12. What are hyperparameters, and how do you optimize them?

Answer: Hyperparameters are model configurations set before training (e.g., learning rate, number of hidden layers). Optimization methods include:

- **Grid Search:** Tests all combinations of hyperparameter values.
 - **Random Search:** Tests a random selection of hyperparameters.
 - **Bayesian Optimization:** Uses probabilistic models to find the best set of hyperparameters.
-

13. How would you evaluate the performance of a regression model?

Answer: Common evaluation metrics:

- **Mean Absolute Error (MAE):** Average of absolute differences between predicted and actual values.
 - **Mean Squared Error (MSE):** Average of squared differences between predicted and actual values.
 - **R² Score:** Proportion of variance explained by the model.
-

14. What is the purpose of dimensionality reduction?

Answer: Dimensionality reduction simplifies datasets by reducing the number of features while retaining essential information. Techniques include:

- **PCA (Principal Component Analysis):** Projects data onto lower dimensions.
 - **t-SNE:** Visualizes high-dimensional data in 2D or 3D space.
-

15. Explain Natural Language Processing (NLP).

Answer: NLP is a field of AI focused on enabling machines to understand and generate human language. Common tasks:

- **Tokenization:** Splitting text into smaller units (e.g., words, sentences).
- **Sentiment Analysis:** Identifying the sentiment of text.
-