# Machine Learning Interview Questions

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

**Answer:**

- **Supervised Learning** uses labeled data. The algorithm learns to map inputs to outputs.
  *Example: Regression, Classification*
- **Unsupervised Learning** uses unlabeled data to identify hidden patterns.
  *Example: Clustering, Dimensionality Reduction*

---

## 2. Explain overfitting and underfitting. How can you prevent them?

**Answer:**

**Overfitting** happens when a model learns the training data too well, including noise and outliers. It performs excellently on training data but poorly on unseen data because it fails to generalize.

**Prevention:**

- Use cross-validation to test generalization
- Apply regularization (L1/L2) to reduce model complexity
- Prune decision trees
- Add more training data
- Use simpler models if needed

**Underfitting** occurs when a model is too simple to capture the underlying patterns in the data. It performs poorly on both training and test sets.

**Prevention:**

- Increase model complexity
- Improve feature engineering
- Train the model longer if necessary
- Reduce regularization

---

## 3. What is the bias-variance tradeoff?

**Answer:**

- **Bias**: Error due to simplistic assumptions in the model (underfitting).

- **Variance**: Error due to model complexity (overfitting).
   **Tradeoff:** Ideal model balances bias and variance to minimize total error.

---

## 4. How do personalized recommendation systems work in machine learning?

Personalized recommendation systems suggest relevant content or products to users by analyzing their past behavior, preferences, and similarities with other users or items.

There are three main types:

- Collaborative Filtering: This approach recommends items based on the preferences of similar users. For example, Netflix suggests movies by identifying users with viewing habits similar to yours.

- Content-Based Filtering: This method recommends items similar to what the user has already liked or interacted with. For example, Spotify suggests songs based on the genres or artists you frequently listen to.

- Hybrid Approach: This combines collaborative and content-based filtering to provide more accurate and robust recommendations. Platforms like Amazon often use hybrid models for personalized shopping experiences.

---

## 5. What is the difference between classification and regression?

**Answer:**

| Aspect | Classification | Regression |
|---|---|---|
| **Output Type** | Categorical (discrete classes) | Continuous (real-valued numbers) |
| **Goal** | Predict class labels | Predict numerical values |
| **Examples** | Spam detection, disease diagnosis, sentiment analysis | House price prediction, stock forecasting, car emissions prediction |
| **Algorithms** | Logistic Regression, SVM, Decision Tree Classifier | Linear Regression, Random Forest Regressor, XGBoost Regressor |

## 6. Explain the steps in a machine learning pipeline.

**Answer:**

1. Data Collection
2. Data Cleaning
3. Exploratory Data Analysis (EDA)
4. Feature Engineering
5. Model Selection
6. Training
7. Evaluation
8. Deployment & Monitoring

## 7. What is cross-validation? Why is it used?

**Answer:**
Cross-validation is a technique to assess model performance by splitting data into training and validation sets multiple times.
 **Use:** Reduces overfitting, ensures generalization.

---

## 8. What are precision, recall, F1 score, and accuracy? When do you use each?

**Answer:**
1. Accuracy

- Use: When classes are balanced and all errors matter equally.
- Explanation: Measures how often the model is correct overall.
- Formula: (TP + TN) / (TP + TN + FP + FN)

2. Precision

- Use: When false positives are more harmful (e.g., spam filter flagging real emails).
- Explanation: Of all predicted positives, how many are truly positive?
- Formula: TP / (TP + FP)

3. Recall (Sensitivity)

- Use: When false negatives are more harmful (e.g., missing a cancer diagnosis).
- Explanation: Of all actual positives, how many did the model correctly find?
- Formula: TP / (TP + FN)

4. F1 Score

- Use: When you need a balance between precision and recall (especially with imbalanced data).
- Explanation: Combines precision and recall into one metric using harmonic mean.
- Formula: 2 × (Precision × Recall) / (Precision + Recall)

---

## 9. What are the assumptions of a linear regression model?

**Answer:**

| Assumption | Description |
|---|---|
| **Linearity** | The relationship between input variables and the output is linear. |
| **Independence** | Residuals (errors) are independent from each other. |
| **Homoscedasticity** | Constant variance of residuals across all levels of input features. |
| **Normality** | Residuals should be normally distributed. |
| **No Multicollinearity** | Independent variables should not be highly correlated with each other. |

## 10. How does a decision tree work? What are entropy and information gain?

**Answer:**

- A **Decision Tree** splits data into branches based on feature values to make decisions. At each step, the algorithm selects the best feature that divides the data to achieve maximum **purity** in child nodes. Decision Trees are intuitive, handle both categorical and numerical data, and are the building blocks for powerful ensembles like Random Forest and Gradient Boosted Trees.

- **Entropy** is a measure of impurity or randomness in the dataset. Lower entropy = purer node. **Example:** If all data points in a node belong to one class, entropy = 0.

- **Information Gain** is the reduction in entropy after a dataset is split on a feature.
  - Formula: `Information Gain = Entropy(Parent) — Weighted Avg. Entropy(Children)`
  - Higher information gain = better feature for splitting.

---

## 11. What is regularization? Explain L1 vs L2.

**Answer:**

- Regularization adds a penalty to the loss function to avoid overfitting.
- **L1 (Lasso):** Shrinks some coefficients to 0 (feature selection).
- **L2 (Ridge):** Shrinks all coefficients but doesn't make them zero.

---

## 12. How does KNN work and how do you choose K?

**Answer:**

- K-Nearest Neighbors (KNN) classifies a data point based on the **majority class among its K-nearest neighbors** (for classification) or the **average of neighbors** (for regression)
- Choose K using cross-validation
- Small K = noise-sensitive; Large K = underfitting

---

## 13. What are the differences between bagging and boosting?

**Answer:**

- **Bagging:** Trains multiple models independently and averages predictions (e.g., Random Forest).
- **Boosting:** Trains models sequentially, each correcting the errors of the previous (e.g., XGBoost).

---

## 14. Explain how Random Forest works.

**Answer:**

- Ensemble of decision trees using bootstrapped samples and random feature selection.
- Reduce overfitting and improves accuracy.

---

## 15. What is PCA? How does it reduce dimensionality?

**Answer:**

- PCA (Principal Component Analysis) transforms data into fewer uncorrelated variables (principal components) that capture the most variance.
- Reduces dimensions while preserving information.

---

## 16. How do you handle imbalanced datasets?

**Answer:**

- Resampling: Over/Under Sampling
- Use metrics like ROC-AUC, F1
- Synthetic Data: SMOTE
- Algorithm-level solutions: Class weighting

---

## 17. What is the difference between ROC curve and Precision-Recall curve?

**Answer:**

- **ROC Curve:** Plots TPR vs FPR, good when classes are balanced.
- **PR Curve:** Plots precision vs recall, better for imbalanced data.

---

## 18. Explain the kernel trick in SVM.

**Answer:**
The kernel trick allows SVM to operate in a high-dimensional space without explicitly computing the transformation.
Common kernels: Linear, Polynomial, RBF

---

## 19. What are hyperparameters? How do you tune them?

**Answer:**
Hyperparameters are external settings of a model (e.g., learning rate, depth). They control how the model learns and performs.
Tuning methods:

- Grid Search
- Random Search
- Bayesian Optimization
- AutoML tools

---

## 20. Explain ensemble learning and different types of ensemble methods.

**Answer:**
Ensemble learning combines multiple models to improve performance.
Types:

- Bagging (Random Forest)
- Boosting (XGBoost, AdaBoost)
- Stacking (combining different algorithms)

---

## 21. What are some challenges in deploying ML models in production?

**Answer:**

- Data Drift
- Model Monitoring
- Version Control
- Scalability
- Latency
- Retraining pipelines

---

## 22. How do you handle data leakage in ML pipelines?

**Answer:**

- Ensure training data doesn't include information from the test set.
- Avoid target leakage (e.g., using future data).
- Proper feature engineering within cross-validation.

---

## 23. What is the role of feature engineering in machine learning? Give examples.

**Answer:**
Feature engineering transforms raw data into useful features.
Examples:

- Encoding categorical variables
- Creating interaction features
- Binning, scaling, log transformation

---

## 24. How would you explain your ML model to a non-technical stakeholder?

**Answer:**

- Focus on **business outcomes** and **impact**
- Use simple analogies (e.g., decision tree = flowchart)
- Avoid jargon
- Show visuals (charts, confusion matrix)

---

## 25. What are the different types of clustering in machine learning?

## Answer:

There are several types of clustering methods, each with its own approach:

| Clustering Type | Description | Example Algorithm |
|---|---|---|
| **Partitioning Clustering** | Divides data into distinct, non-overlapping clusters. | K-Means, K-Medoids |
| **Hierarchical Clustering** | Builds a tree of clusters (dendrogram) using a bottom-up or top-down approach. | Agglomerative, Divisive |
| **Density-Based Clustering** | Forms clusters based on dense regions of data separated by low-density areas. | DBSCAN, OPTICS |
| **Grid-Based Clustering** | Divides data space into a grid and forms clusters based on dense cells. | STING, CLIQUE |
| **Model-Based Clustering** | Assumes data is generated from a mixture of underlying probability distributions. | Gaussian Mixture Models (GMM) |

---