

# Machine Learning Interview Questions - Complete Study Guide

## I. GENERAL AI & ML CONCEPTS

### What is the difference between AI, Machine Learning, and Deep Learning?

**Artificial Intelligence (AI)** is the broadest concept - it's any technique that enables machines to mimic human intelligence, including rule-based systems, expert systems, and machine learning.

**Machine Learning (ML)** is a subset of AI where algorithms learn patterns from data without being explicitly programmed for every scenario. Instead of hard-coding rules, the system learns from examples.

**Deep Learning (DL)** is a subset of ML that uses neural networks with multiple layers (typically 3+ hidden layers) to automatically learn hierarchical representations of data.

Think of it as nested circles: AI contains ML, and ML contains Deep Learning.

### What are the main types of Machine Learning?

1. **Supervised Learning:** Learning with labeled examples (input-output pairs)
2. **Unsupervised Learning:** Finding patterns in data without labels
3. **Semi-supervised Learning:** Uses both labeled and unlabeled data
4. **Reinforcement Learning:** Learning through interaction with an environment using rewards/penalties

### Explain the difference between supervised and unsupervised learning

#### Supervised Learning:

- Uses labeled training data (input-output pairs)
- Goal: Learn a mapping function from input to output
- Examples: Classification (spam detection), Regression (price prediction)
- Algorithms: Linear Regression, SVM, Random Forest, Neural Networks

#### Unsupervised Learning:

- Uses only input data without corresponding outputs
- Goal: Discover hidden patterns or structures in data
- Examples: Clustering (customer segmentation), Dimensionality reduction (PCA)
- Algorithms: K-means, Hierarchical clustering, DBSCAN, PCA

### What is overfitting and underfitting?

#### Overfitting:

- Model learns training data too well, including noise

- High training accuracy, poor test accuracy
- Model is too complex for the amount of data
- Solutions: Regularization, more data, cross-validation, early stopping

### **Underfitting:**

- Model is too simple to capture underlying patterns
- Poor performance on both training and test data
- Solutions: More complex model, better features, reduce regularization

## **How do you evaluate a machine learning model?**

### **For Classification:**

- Accuracy, Precision, Recall, F1-score
- ROC curve and AUC
- Confusion matrix
- Cross-validation scores

### **For Regression:**

- Mean Squared Error (MSE)
- Root Mean Squared Error (RMSE)
- Mean Absolute Error (MAE)
- R-squared (coefficient of determination)

### **General Practices:**

- Train/validation/test split
- Cross-validation
- Learning curves
- Feature importance analysis

## **What is the bias-variance trade-off?**

**Bias:** Error from overly simplistic assumptions. High bias leads to underfitting.

**Variance:** Error from sensitivity to small fluctuations in training data. High variance leads to overfitting.

### **Trade-off:**

- Simple models: High bias, low variance
- Complex models: Low bias, high variance

- Goal: Find the sweet spot that minimizes total error ( $\text{bias}^2 + \text{variance} + \text{irreducible error}$ )

## Explain cross-validation and why it's important

Cross-validation splits data into multiple folds to get a more robust estimate of model performance:

### K-Fold Cross-Validation:

1. Split data into k equal parts
2. Train on k-1 folds, test on remaining fold
3. Repeat k times, each fold serves as test set once
4. Average the k performance scores

### Importance:

- Reduces overfitting to a particular train/test split
- Better utilizes limited data
- Provides confidence intervals for performance metrics
- Helps in hyperparameter tuning

## What are precision, recall, F1-score, and accuracy?

**Accuracy:**  $(TP + TN) / (TP + TN + FP + FN)$  - Overall correctness

**Precision:**  $TP / (TP + FP)$  - Of positive predictions, how many were correct?

**Recall (Sensitivity):**  $TP / (TP + FN)$  - Of actual positives, how many were found?

**F1-Score:**  $2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$  - Harmonic mean of precision and recall

*TP=True Positives, TN=True Negatives, FP=False Positives, FN=False Negatives*

## What is the difference between classification and regression?

### Classification:

- Predicts discrete categories or classes
- Output is categorical (spam/not spam, cat/dog/bird)
- Evaluation: Accuracy, precision, recall, F1-score
- Examples: Image recognition, sentiment analysis

### Regression:

- Predicts continuous numerical values
- Output is a real number (price, temperature, age)
- Evaluation: MSE, RMSE, MAE,  $R^2$

- Examples: Stock price prediction, sales forecasting

## What are some real-world applications of AI/ML?

- **Healthcare:** Medical diagnosis, drug discovery, personalized treatment
- **Finance:** Fraud detection, algorithmic trading, credit scoring
- **Transportation:** Autonomous vehicles, route optimization
- **Technology:** Recommendation systems, search engines, voice assistants
- **Retail:** Demand forecasting, price optimization, customer segmentation
- **Manufacturing:** Predictive maintenance, quality control, supply chain optimization

## II. DATA PREPROCESSING & FEATURE ENGINEERING

### How do you handle missing data in a dataset?

#### Strategies:

1. **Remove:** Delete rows/columns with missing values (if < 5% missing)
2. **Imputation:**
  - Mean/Median/Mode for numerical/categorical data
  - Forward/Backward fill for time series
  - KNN imputation
  - Multiple imputation
3. **Prediction:** Use other features to predict missing values
4. **Indicator variables:** Create binary flags for missingness

**Choice depends on:** Amount of missing data, type of missingness (MCAR, MAR, MNAR), and domain knowledge.

### What is normalization and standardization?

#### Normalization (Min-Max Scaling):

- Scales data to [0,1] range
- Formula:  $(x - \min) / (\max - \min)$
- Preserves relationships between data points
- Sensitive to outliers

#### Standardization (Z-score):

- Centers data around mean=0, std=1
- Formula:  $(x - \mu) / \sigma$

- Less sensitive to outliers
- Assumes normal distribution

**When to use:** Standardization for algorithms assuming normal distribution (SVM, Neural Networks).  
Normalization when you need bounded values.

## What is one-hot encoding? When would you use it?

One-hot encoding converts categorical variables into binary vectors where only one element is 1 (hot) and others are 0.

**Example:** Color: [Red, Blue, Green] becomes:

- Red: [1, 0, 0]
- Blue: [0, 1, 0]
- Green: [0, 0, 1]

**When to use:**

- Nominal categorical variables (no order)
- Algorithms that can't handle categorical data directly
- When you want to avoid imposing artificial ordering

**Avoid when:** High cardinality categories (creates too many features), ordinal data (use label encoding instead).

## How do you handle categorical variables?

1. **Label Encoding:** Assign integers to categories (for ordinal data)
2. **One-Hot Encoding:** Create binary columns for each category
3. **Target Encoding:** Replace categories with target variable statistics
4. **Binary Encoding:** Convert to binary representation
5. **Frequency Encoding:** Replace with frequency of occurrence
6. **Embedding:** Learn dense representations (for high cardinality)

## What is feature selection and why is it important?

Feature selection chooses the most relevant features for model training.

**Why important:**

- Reduces overfitting
- Improves model interpretability
- Decreases training time

- Reduces storage requirements
- Removes noise and irrelevant information

### **Methods:**

- **Filter:** Statistical tests, correlation analysis
- **Wrapper:** Forward/backward selection, recursive feature elimination
- **Embedded:** LASSO regularization, tree-based feature importance

## **What is dimensionality reduction? Explain PCA.**

**Dimensionality Reduction:** Technique to reduce the number of features while preserving important information.

### **Principal Component Analysis (PCA):**

- Finds directions (principal components) of maximum variance
- Projects data onto lower-dimensional space
- Components are orthogonal and ordered by variance explained
- Unsupervised technique

### **Steps:**

1. Standardize data
2. Compute covariance matrix
3. Find eigenvalues and eigenvectors
4. Select top k components
5. Transform data

**Use cases:** Visualization, noise reduction, preprocessing for other algorithms.

## **What is the curse of dimensionality?**

As the number of dimensions increases:

- Data becomes sparse in high-dimensional space
- Distance metrics become less meaningful
- Required sample size grows exponentially
- Visualization becomes impossible
- Computational complexity increases

**Solutions:** Dimensionality reduction, feature selection, regularization, domain expertise for feature engineering.

## What is feature scaling and when should it be applied?

Feature scaling ensures all features contribute equally to distance-based algorithms.

### When to apply:

- Algorithms using distance metrics (KNN, SVM, Neural Networks)
- Gradient descent optimization
- When features have different units/scales

### When not needed:

- Tree-based algorithms (Random Forest, Decision Trees)
- Naive Bayes
- Algorithms that don't use distance metrics

## How do you deal with imbalanced datasets?

### Techniques:

#### 1. Resampling:

- Oversampling minority class (SMOTE)
- Undersampling majority class
- Combination approaches

#### 2. Algorithm-level:

- Class weights
- Cost-sensitive learning
- Ensemble methods

#### 3. Evaluation:

- Use precision, recall, F1-score instead of accuracy
- ROC-AUC, Precision-Recall curves

#### 4. Data-level:

- Collect more minority class data
- Generate synthetic examples

## Explain the role of EDA (Exploratory Data Analysis) in ML

EDA is the critical first step in any ML project:

### Purposes:

- Understand data distribution and patterns

- Identify outliers and anomalies
- Discover relationships between variables
- Detect missing values and data quality issues
- Guide feature engineering decisions
- Inform algorithm selection

#### **Techniques:**

- Summary statistics
- Data visualization (histograms, box plots, scatter plots)
- Correlation analysis
- Distribution analysis

### **III. MACHINE LEARNING ALGORITHMS**

#### **How does the Decision Tree algorithm work?**

Decision Trees make predictions by learning simple decision rules inferred from data features.

#### **Algorithm:**

1. Start with entire dataset at root
2. Find best feature and split point that maximizes information gain
3. Split data into subsets based on feature value
4. Recursively repeat for each subset
5. Stop when stopping criteria met (max depth, min samples, pure nodes)

#### **Splitting Criteria:**

- **Classification:** Gini impurity, Information gain (entropy)
- **Regression:** Mean Squared Error

**Advantages:** Interpretable, handles both numerical and categorical data, no feature scaling needed

**Disadvantages:** Prone to overfitting, unstable (small data changes affect tree structure)

#### **What is the difference between bagging and boosting?**

#### **Bagging (Bootstrap Aggregating):**

- Trains multiple models in parallel on different bootstrap samples
- Combines predictions by averaging (regression) or voting (classification)
- Reduces variance, helps with overfitting
- Example: Random Forest



## Boosting:

- Trains models sequentially, each correcting previous model's errors
- Combines weak learners to create strong learner
- Reduces bias, can lead to overfitting if not careful
- Examples: AdaBoost, Gradient Boosting, XGBoost

## Explain how the K-Nearest Neighbors algorithm works

KNN is a lazy learning algorithm that makes predictions based on the k closest training examples.

### Algorithm:

1. Choose value of k
2. Calculate distance between query point and all training points
3. Find k nearest neighbors
4. For classification: majority vote among k neighbors
5. For regression: average of k neighbors' values

**Distance Metrics:** Euclidean, Manhattan, Minkowski, Hamming

**Advantages:** Simple, no assumptions about data, works well with small datasets **Disadvantages:** Computationally expensive, sensitive to irrelevant features, requires feature scaling

## What is the intuition behind Support Vector Machines?

SVM finds the optimal hyperplane that separates classes with maximum margin.

### Key Concepts:

- **Hyperplane:** Decision boundary separating classes
- **Support Vectors:** Data points closest to hyperplane
- **Margin:** Distance between hyperplane and nearest points from each class
- **Kernel Trick:** Maps data to higher dimensions where it becomes linearly separable

### Types of Kernels:

- Linear: For linearly separable data
- RBF (Gaussian): For non-linear data
- Polynomial: For specific polynomial relationships

**Advantages:** Effective in high dimensions, memory efficient, versatile with kernels **Disadvantages:** Poor performance on large datasets, sensitive to feature scaling

## How does Naive Bayes classifier work?

Naive Bayes applies Bayes' theorem with the "naive" assumption that features are conditionally independent.

**Bayes' Theorem:**  $P(A|B) = P(B|A) \times P(A) / P(B)$

**For classification:**  $P(\text{class}|\text{features}) = P(\text{features}|\text{class}) \times P(\text{class}) / P(\text{features})$

### Types:

- **Gaussian:** For continuous features (assumes normal distribution)
- **Multinomial:** For discrete features (text classification)
- **Bernoulli:** For binary features

**Advantages:** Fast, works well with small datasets, handles multiple classes naturally **Disadvantages:** Independence assumption often violated, poor probability estimates

## What is the difference between Random Forest and XGBoost?

### Random Forest:

- Ensemble of decision trees using bagging
- Trees trained independently in parallel
- Uses bootstrap sampling and random feature selection
- Averages predictions across trees
- Less prone to overfitting

### XGBoost (Extreme Gradient Boosting):

- Ensemble using gradient boosting
- Trees trained sequentially, each correcting previous errors
- Uses regularization to prevent overfitting
- More complex hyperparameter tuning
- Often achieves better performance but requires more careful tuning

## How do gradient descent and stochastic gradient descent differ?

### Gradient Descent (Batch GD):

- Uses entire dataset to compute gradient
- Updates parameters once per epoch
- More stable convergence
- Computationally expensive for large datasets

## Stochastic Gradient Descent (SGD):

- Uses single sample to compute gradient
- Updates parameters after each sample
- Faster but noisier convergence
- Can escape local minima due to noise

**Mini-batch GD:** Compromise using small batches (typically 32-512 samples)

## Explain logistic regression and where it's used

Logistic regression uses the logistic function to model the probability of binary outcomes.

**Logistic Function:**  $\sigma(z) = 1 / (1 + e^{(-z)})$  **Linear Combination:**  $z = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$

### Key Points:

- Output is probability between 0 and 1
- Uses maximum likelihood estimation
- Decision boundary at probability = 0.5
- Can be extended to multi-class (softmax regression)

### Use Cases:

- Binary classification problems
- Medical diagnosis
- Marketing response prediction
- Email spam detection

## What are ensemble models?

Ensemble models combine multiple individual models to create a stronger predictor.

### Types:

1. **Bagging:** Parallel training with bootstrap sampling (Random Forest)
2. **Boosting:** Sequential training correcting previous errors (XGBoost)
3. **Stacking:** Train meta-model on predictions of base models
4. **Voting:** Simple majority vote or averaging

### Benefits:

- Reduced overfitting
- Better generalization

- Improved accuracy
- More robust predictions

## What is the difference between L1 and L2 regularization?

### L1 Regularization (Lasso):

- Adds sum of absolute values of parameters
- Penalty:  $\lambda \sum |\beta_i|$
- Promotes sparsity (sets some coefficients to exactly zero)
- Performs feature selection
- Less stable with correlated features

### L2 Regularization (Ridge):

- Adds sum of squared parameters
- Penalty:  $\lambda \sum \beta_i^2$
- Shrinks coefficients toward zero but doesn't eliminate them
- Handles multicollinearity better
- More stable solution

**Elastic Net:** Combines L1 and L2 regularization

## IV. DEEP LEARNING & NEURAL NETWORKS

### What is a perceptron?

A perceptron is the simplest neural network unit, consisting of:

- Input features ( $x_1, x_2, \dots, x_n$ )
- Weights ( $w_1, w_2, \dots, w_n$ )
- Bias term ( $b$ )
- Activation function

**Output:**  $f(\sum w_i x_i + b)$

### Limitations:

- Can only learn linearly separable patterns
- Cannot solve XOR problem
- Single layer perceptron is very limited

**Multi-layer perceptrons (MLPs) overcome these limitations by stacking multiple layers.**

## How do activation functions like ReLU, Sigmoid, and Tanh work?

**Sigmoid:**  $\sigma(x) = 1/(1 + e^{(-x)})$

- Output: (0, 1)
- Problems: Vanishing gradients, not zero-centered
- Use: Binary classification output layer

**Tanh:**  $\tanh(x) = (e^x - e^{(-x)})/(e^x + e^{(-x)})$

- Output: (-1, 1)
- Zero-centered, still has vanishing gradient problem
- Use: Hidden layers in RNNs

**ReLU:**  $f(x) = \max(0, x)$

- Output:  $[0, \infty)$
- Advantages: Computationally efficient, reduces vanishing gradients
- Problems: Dead neurons (always output 0)
- Use: Most common in hidden layers

**Variants:** Leaky ReLU, ELU, Swish address ReLU's limitations

## What are epochs, batch size, and learning rate?

**Epoch:** One complete pass through the entire training dataset

**Batch Size:** Number of samples processed before updating model parameters

- Larger batches: More stable gradients, better hardware utilization
- Smaller batches: More frequent updates, better generalization

**Learning Rate:** Step size for parameter updates during optimization

- Too high: May overshoot optimal solution
- Too low: Slow convergence, may get stuck in local minima
- Common values: 0.001, 0.01, 0.1

## What is the vanishing gradient problem?

During backpropagation in deep networks, gradients become exponentially smaller as they propagate backward through layers.

**Causes:**

- Activation functions with small derivatives (sigmoid, tanh)

- Deep network architecture
- Poor weight initialization

### **Consequences:**

- Early layers learn very slowly
- Network fails to capture long-range dependencies

### **Solutions:**

- Better activation functions (ReLU)
- Proper weight initialization (Xavier, He initialization)
- Skip connections (ResNet)
- Batch normalization
- LSTM/GRU for sequential data

## **What is the difference between CNN and RNN?**

### **Convolutional Neural Networks (CNN):**

- Designed for spatial data (images, 2D/3D data)
- Uses convolution operations with filters/kernels
- Translation invariant
- Parameter sharing across spatial locations
- Architecture: Convolution → Pooling → Fully Connected
- Applications: Image recognition, computer vision

### **Recurrent Neural Networks (RNN):**

- Designed for sequential data (text, time series)
- Has memory/hidden state from previous time steps
- Can handle variable-length sequences
- Architecture includes recurrent connections
- Applications: Language modeling, machine translation, time series

## **How does an LSTM work and where is it used?**

**Long Short-Term Memory (LSTM)** addresses RNN's vanishing gradient problem using gating mechanisms:

### **Gates:**

1. **Forget Gate:** Decides what information to discard from cell state

2. **Input Gate:** Decides what new information to store
3. **Output Gate:** Controls what parts of cell state to output

**Cell State:** Carries information across time steps with minimal modification

**Process:**

1. Forget irrelevant information
2. Decide what new information to store
3. Update cell state
4. Output filtered version of cell state

**Applications:**

- Machine translation
- Speech recognition
- Time series prediction
- Sentiment analysis

## **What are convolutional layers in CNN?**

Convolutional layers apply filters (kernels) across input data to detect features.

**Key Concepts:**

- **Filter/Kernel:** Small matrix that slides across input
- **Stride:** Step size of filter movement
- **Padding:** Adding zeros around input borders
- **Feature Maps:** Output of convolution operation

**Operations:**

1. Element-wise multiplication between filter and input region
2. Sum all products
3. Apply activation function
4. Slide filter to next position

**Benefits:**

- Parameter sharing reduces overfitting
- Translation invariance
- Hierarchical feature learning
- Spatial relationship preservation

## What is transfer learning?

Transfer learning uses pre-trained models as starting points for new, related tasks.

### Approaches:

1. **Feature Extraction:** Use pre-trained model as fixed feature extractor
2. **Fine-tuning:** Update pre-trained weights for new task
3. **Using Pre-trained Models:** Adapt architecture for specific needs

### Benefits:

- Faster training
- Better performance with limited data
- Leverages learned representations
- Reduces computational requirements

**Common Pre-trained Models:** ResNet, VGG, BERT, GPT

## What is dropout and why is it used?

Dropout randomly sets a fraction of neurons to zero during training.

### How it works:

- During training: Randomly "drop" neurons with probability  $p$
- During inference: Use all neurons but scale outputs by  $(1-p)$

### Benefits:

- Prevents overfitting
- Reduces co-adaptation between neurons
- Acts as ensemble method
- Improves generalization

**Typical dropout rates:** 0.2-0.5 for hidden layers, 0.5-0.8 for input layers

## How do you prevent overfitting in deep learning models?

### Regularization Techniques:

1. **Dropout:** Randomly disable neurons during training
2. **L1/L2 Regularization:** Add penalty terms to loss function
3. **Batch Normalization:** Normalize layer inputs
4. **Early Stopping:** Stop training when validation loss increases



**Data Techniques:** 5. **Data Augmentation:** Artificially increase dataset size 6. **More Training Data:**

Collect additional samples

**Architecture Techniques:** 7. **Simpler Models:** Reduce parameters/complexity 8. **Cross-validation:** Better model selection

**Other Techniques:** 9. **Transfer Learning:** Use pre-trained models 10. **Ensemble Methods:** Combine multiple models

## V. NATURAL LANGUAGE PROCESSING

### What is tokenization in NLP?

Tokenization breaks text into smaller units (tokens) for processing.

#### Types:

- **Word Tokenization:** Split by spaces/punctuation
- **Sentence Tokenization:** Split into sentences
- **Subword Tokenization:** Split words into subparts (BPE, WordPiece)
- **Character Tokenization:** Individual characters as tokens

#### Challenges:

- Contractions (don't → do, n't)
- Punctuation handling
- Multiple languages
- Out-of-vocabulary words

**Tools:** NLTK, spaCy, Hugging Face tokenizers

### How do word embeddings like Word2Vec or GloVe work?

Word embeddings represent words as dense vectors in continuous space where similar words have similar representations.

#### Word2Vec:

- **Skip-gram:** Predict context words from target word
- **CBOW:** Predict target word from context
- Uses neural network with single hidden layer
- Captures semantic and syntactic relationships

#### GloVe (Global Vectors):

- Combines global matrix factorization and local context window methods
- Uses word co-occurrence statistics
- Optimizes ratio of co-occurrence probabilities

**Benefits:**

- Captures semantic similarity
- Reduces dimensionality
- Works with machine learning algorithms

**What is the difference between stemming and lemmatization?****Stemming:**

- Removes suffixes to get root form
- Rule-based approach (Porter, Snowball stemmers)
- Fast but may create non-words
- Example: running → run, studies → studi

**Lemmatization:**

- Reduces words to dictionary form (lemma)
- Uses vocabulary and morphological analysis
- Slower but more accurate
- Example: running → run, studies → study, better → good

**When to use:**

- Stemming: Speed is important, approximate matching okay
- Lemmatization: Accuracy is crucial, meaning preservation needed

**What is TF-IDF and why is it used?**

**Term Frequency-Inverse Document Frequency** measures word importance in document relative to collection of documents.

**Formula:**

- $TF(t,d) = (\text{Number of times term } t \text{ appears in document } d) / (\text{Total terms in } d)$
- $IDF(t,D) = \log(\text{Total documents} / \text{Documents containing term } t)$
- $TF-IDF(t,d,D) = TF(t,d) \times IDF(t,D)$

**Intuition:**

- High TF-IDF: Word appears frequently in document but rarely in corpus
- Low TF-IDF: Common words (the, and, or) or rare words

### Applications:

- Information retrieval
- Text mining
- Feature extraction for ML
- Document similarity

## What are Transformers in NLP?

Transformers use attention mechanisms to process sequences in parallel, revolutionizing NLP.

### Key Components:

- **Self-Attention:** Each position attends to all positions in sequence
- **Multi-Head Attention:** Multiple attention mechanisms in parallel
- **Position Encoding:** Add positional information to embeddings
- **Feed-Forward Networks:** Process attended representations

### Advantages:

- Parallelizable (faster training)
- Captures long-range dependencies
- No recurrence needed
- State-of-the-art performance

**Famous Models:** BERT, GPT, T5, RoBERTa

## VI. MODEL DEPLOYMENT & MLOps

### How do you save and load a trained model in Python?

#### Scikit-learn:

```
python

import joblib
# Save
joblib.dump(model, 'model.pkl')
# Load
model = joblib.load('model.pkl')
```

#### TensorFlow/Keras:

```
python
```

```
# Save
```

```
model.save('model.h5')
```

```
# Load
```

```
model = tf.keras.models.load_model('model.h5')
```

## PyTorch:

```
python
```

```
# Save
```

```
torch.save(model.state_dict(), 'model.pth')
```

```
# Load
```

```
model.load_state_dict(torch.load('model.pth'))
```

## Best Practices:

- Save preprocessing pipelines with models
- Version control for models
- Include metadata (training date, performance metrics)
- Use model registries for production

## What is model drift and how do you monitor it?

**Model Drift:** Degradation in model performance over time due to changes in data or environment.

### Types:

1. **Data Drift:** Input data distribution changes
2. **Concept Drift:** Relationship between inputs and outputs changes
3. **Label Drift:** Target distribution changes

### Monitoring Methods:

- Statistical tests (KS test, Chi-square)
- Performance metrics tracking
- Data distribution monitoring
- Feature importance changes
- Prediction distribution analysis

### Solutions:

- Regular model retraining

- Online learning
- Ensemble approaches
- A/B testing for model updates

## How do you deploy a machine learning model as an API?

### Steps:

1. **Wrap model in API framework** (Flask, FastAPI, Django)
2. **Create endpoints** for predictions
3. **Handle input validation** and preprocessing
4. **Return structured responses** (JSON)
5. **Add error handling** and logging
6. **Containerize** with Docker
7. **Deploy** to cloud platform

### Example with Flask:

```
python

from flask import Flask, request, jsonify
import joblib

app = Flask(__name__)
model = joblib.load('model.pkl')

@app.route('/predict', methods=['POST'])
def predict():
    data = request.json
    prediction = model.predict([data['features']])
    return jsonify({'prediction': prediction[0]})
```

## What are common tools for deploying ML models?

### Web Frameworks:

- **Flask:** Lightweight, simple for prototypes
- **FastAPI:** Modern, fast, automatic API documentation
- **Django:** Full-featured, good for complex applications

### Containerization:

- **Docker:** Package application with dependencies
- **Kubernetes:** Container orchestration for scaling

## Cloud Platforms:

- **AWS:** SageMaker, Lambda, EC2, ECS
- **Google Cloud:** AI Platform, Cloud Run, App Engine
- **Azure:** Machine Learning Studio, Container Instances

## Specialized Tools:

- **Streamlit:** Quick web apps for data science
- **Gradio:** Easy ML model interfaces
- **MLflow:** Model lifecycle management
- **Seldon:** Kubernetes-native ML deployments

## Explain the CI/CD pipeline in MLOps

**Continuous Integration/Continuous Deployment** for ML systems includes both code and model versioning.

### CI Pipeline:

1. **Code commit** triggers pipeline
2. **Automated testing** (unit tests, integration tests)
3. **Data validation** (schema checks, quality tests)
4. **Model training** on new data
5. **Model validation** (performance, bias checks)
6. **Model versioning** and artifact storage

### CD Pipeline:

1. **Model approval** (manual or automated gates)
2. **Staging deployment** for testing
3. **A/B testing** or canary deployment
4. **Production deployment**
5. **Monitoring** and alerting
6. **Rollback** capabilities

### Tools:

- **Version Control:** Git, DVC (Data Version Control)
- **CI/CD:** Jenkins, GitHub Actions, GitLab CI
- **Model Registry:** MLflow, Weights & Biases

- **Monitoring:** Evidently, Great Expectations
- **Orchestration:** Airflow, Kubeflow, Prefect

#### **Key Differences from Software CI/CD:**

- Data dependencies and versioning
- Model performance monitoring
- Gradual rollouts (A/B testing)
- Data drift detection
- Model retraining triggers