1. A real estate company wants to develop a system that predicts house prices based on square footage, number of bedrooms, and location.

Q: Identify the problem type and outline the step-by-step logic to solve it.

### **Problem Type:**

This is a **Supervised Machine Learning** problem, specifically a **Regression** problem, because the target variable (house price) is continuous.

#### Step-by-Step Logic to Solve the Problem:

## 1. Problem Understanding

- Goal: Predict house prices based on input features like square footage, number of bedrooms, and location.
- Output: A numerical value (house price).

#### 2. Data Collection

- Gather historical data with the following features:
  - Square Footage (numerical)
  - Number of Bedrooms (numerical)
  - Location (categorical)
  - House Price (target variable numerical)

## 3. Data Preprocessing

- Handle missing values (e.g., fill or remove rows/columns).
- Encode categorical data like Location using techniques like:
  - Label Encoding (for ordinal locations)

- One-Hot Encoding (for nominal locations)
- Feature Scaling for numerical features (optional, depending on the algorithm):
  - Use Min-Max Scaling or Standardization.
- **Split the dataset** into training and testing sets (e.g., 80% train, 20% test).

### 4. Model Selection

- Choose regression models such as:
  - Linear Regression (baseline)
  - o Decision Tree Regressor
  - Random Forest Regressor
  - Gradient Boosting Regressor

## 5. Model Training

• Fit the model on the training data using features (X\_train) and target (y\_train).

#### 6. Model Evaluation

- Predict on the test set and evaluate using metrics like:
  - Mean Absolute Error (MAE)
  - Mean Squared Error (MSE)
  - Root Mean Squared Error (RMSE)
  - R<sup>2</sup> Score (Coefficient of Determination)

### 7. Model Tuning (Optional)

- Use **Hyperparameter Tuning** (e.g., Grid Search, Random Search).
- Apply Cross-Validation to validate model performance more robustly.

### 8. Deployment

- Deploy the trained model in a web/mobile app.
- Take user input for square footage, number of bedrooms, and location.
- Predict and display the estimated house price.

## 9. Monitoring and Maintenance

- Continuously collect new data.
- Retrain or update the model periodically to maintain accuracy.
- 2. A bank wants to build a model to detect fraudulent transactions by analyzing customer spending behavior and transaction history.

Q: Identify the problem type and outline the step-by-step logic to solve it.

#### **Problem Type:**

This is a **Supervised Machine Learning** problem, specifically a **Classification** problem because the target variable (fraudulent or not) is **categorical** (e.g., 0 for non-fraud, 1 for fraud).

## **Step-by-Step Logic to Solve the Problem:**

### 1. Problem Understanding

- **Goal:** Detect whether a transaction is fraudulent based on customer behavior and transaction history.
- Output: Binary classification (Fraud or Not Fraud).

### 2. Data Collection

- Gather historical transaction data with features such as:
  - Transaction amount
  - Time of transaction
  - Customer ID or demographics
  - Merchant location
  - Type of transaction (online, in-store)
  - Device used
  - Transaction history
  - Label: Fraudulent (1) or Not Fraudulent (0)

## 3. Data Preprocessing

- Handle missing values appropriately.
- Convert categorical variables using encoding techniques:
  - One-Hot Encoding (for nominal data like merchant type, device type)
  - Label Encoding (if the categories have an order)
- Feature scaling (e.g., StandardScaler or MinMaxScaler) for numerical data.
- Handle class imbalance (since fraud cases are usually rare):
  - Techniques like **SMOTE**, **Random Undersampling**, or using **class weights** in models.
- Split the dataset into training and test sets (e.g., 80% train, 20% test).

### 4. Model Selection

Choose classification algorithms such as:

- Logistic Regression
- Decision Tree Classifier
- Random Forest Classifier
- XGBoost or LightGBM
- Support Vector Machine (SVM)
- Neural Networks (for advanced systems)

## 5. Model Training

• Train the selected model on the training data (X\_train, y\_train).

### 6. Model Evaluation

Evaluate the model using classification metrics:

- Accuracy (if classes are balanced)
- Precision, Recall, F1-score (important for imbalanced data)
- Confusion Matrix
- ROC-AUC Curve (for probability-based classifiers)

## 7. Model Tuning

- Use **Hyperparameter Tuning** (Grid Search, Randomized Search)
- Apply **Cross-Validation** to improve generalization

### 8. Deployment

- Integrate the model into the bank's transaction system.
- For every new transaction, the model should:
  - Analyze the input features
  - Predict fraud probability or class
  - Trigger alerts for suspicious transactions

## 9. Monitoring & Maintenance

- Monitor false positives/negatives.
- Continuously retrain the model with new data.
- Use feedback loop to improve model accuracy.
- 3. A supermarket wants to segment its customers based on their shopping patterns to provide personalized promotions.

Q: Identify the problem type and outline the step-by-step logic to solve it.

#### **Problem Type:**

This is an **Unsupervised Machine Learning** problem, specifically a **Clustering** problem because there is **no target variable** — the goal is to group customers based on similar behavior.

### **Step-by-Step Logic to Solve the Problem:**

## 1. Problem Understanding

- Goal: Segment customers into groups based on shopping behavior.
- **Output:** Cluster labels (e.g., Group 1: Budget shoppers, Group 2: Premium customers, etc.)

#### 2. Data Collection

Gather customer data with features such as:

- Purchase frequency
- Average purchase value
- Time spent in store or on site
- Product categories purchased
- Recency (how recently a customer purchased)
- Demographic info (age, income level, etc.)

## 3. Data Preprocessing

- Handle missing values (e.g., fill or drop).
- Convert categorical variables using encoding (e.g., One-Hot Encoding).
- Normalize or standardize features (especially for distance-based algorithms like K-Means).
- Feature engineering (e.g., calculate RFM Recency, Frequency, Monetary value).

#### 4. Model Selection

Choose appropriate clustering algorithms:

- K-Means Clustering (most common for customer segmentation)
- Hierarchical Clustering
- **DBSCAN** (if the clusters are of different shapes/sizes or contain noise)

### 5. Determine the Optimal Number of Clusters

Use methods like:

- **Elbow Method** (plot WCSS vs number of clusters)
- Silhouette Score
- **Dendrogram** (for hierarchical clustering)

## 6. Model Training

- Apply the clustering algorithm to the preprocessed data.
- Assign each customer to a cluster.

## 7. Analyze and Interpret Clusters

- Profile each cluster:
  - O What makes this group unique?
  - For example: "Cluster 1 = Low spend, high frequency", "Cluster 2 = High spend, low frequency"
- Use visualization techniques:
  - o Pairplot, t-SNE, PCA for dimensionality reduction

## 8. Use Clusters for Targeted Marketing

- Design personalized promotions for each segment.
  - E.g., loyalty programs for high-spenders, discounts for infrequent shoppers

## 9. Monitoring and Updating

Update clustering periodically as customer behavior changes.

- Re-evaluate features and cluster performance over time.
- 4. A company wants to estimate an employee's salary based on their years of experience, job title, and education level.

Q: Identify the problem type and outline the step-by-step logic to solve it.

#### **Problem Type:**

This is a **Supervised Machine Learning** problem, specifically a **Regression** problem, because the target variable (salary) is **continuous**.

### **Step-by-Step Logic to Solve the Problem:**

### 1. Problem Understanding

- Goal: Predict an employee's salary based on:
  - Years of experience (numerical)
  - Job title (categorical)
  - Education level (categorical)
- Output: Estimated salary (a numerical value)

#### 2. Data Collection

Collect historical employee data with the following fields:

- Years of Experience (numerical)
- Job Title (categorical)
- Education Level (categorical)
- Salary (target variable numerical)

### 3. Data Preprocessing

- Handle missing values
- Encode categorical variables:
  - $\circ$  Job Title  $\rightarrow$  One-Hot Encoding or Ordinal Encoding
  - Education Level → Ordinal Encoding (e.g., High School < Bachelor's < Master's < PhD)</li>
- **Feature Scaling** (if needed for the algorithm, e.g., for Gradient Descent-based models)
- **Split data** into training and test sets (e.g., 80% train, 20% test)

#### 4. Model Selection

Choose regression algorithms such as:

- Linear Regression
- Decision Tree Regressor
- Random Forest Regressor
- Support Vector Regressor (SVR)
- Gradient Boosting Regressor

## 5. Model Training

• Train the model using the training data (X\_train, y\_train)

### 6. Model Evaluation

Evaluate model performance using metrics like:

- Mean Absolute Error (MAE)
- Mean Squared Error (MSE)
- Root Mean Squared Error (RMSE)
- R<sup>2</sup> Score (coefficient of determination)

## 7. Model Tuning (Optional)

- Use **Hyperparameter Tuning** (Grid Search, Randomized Search)
- Apply **Cross-Validation** for better generalization

### 8. Deployment

- Deploy the trained model into the HR system or web app
- Input: Employee's experience, job title, and education level
- Output: Predicted salary

## 9. Monitoring and Updating

- Periodically update the model with new employee data
- Monitor prediction accuracy and recalibration as necessary

5.An email provider wants to automatically classify incoming emails as spam or not spam based on their content and sender details.

Q: Identify the problem type and outline the step-by-step logic to solve it.

#### **Problem Type:**

This is a **Supervised Machine Learning** problem, specifically a **Binary Classification** problem because the output variable is **categorical** with two classes: **Spam** or **Not Spam**.

### **Step-by-Step Logic to Solve the Problem:**

## 1. Problem Understanding

- **Goal:** Classify emails as spam (1) or not spam (0) using email content and sender information.
- Output: A label (spam or not spam)

### 2. Data Collection

Collect labeled email data with features like:

- Email subject and body text
- Sender email address or domain
- Presence of suspicious links or attachments
- Frequency of certain keywords (e.g., "win", "free", "click here")
- Length of email
- Historical spam labels (0 or 1)

## 3. Data Preprocessing

- Text preprocessing:
  - Lowercasing
  - Remove punctuation, stop words, and special characters
  - Tokenization
  - Stemming/Lemmatization (optional)
- Convert text into numerical format using:
  - Bag of Words (BoW)

- TF-IDF (Term Frequency-Inverse Document Frequency)
- Word embeddings (e.g., Word2Vec, BERT for advanced models)
- Encode categorical variables like sender domain
- **Feature scaling** (if needed for certain algorithms)
- Split data into training and test sets (e.g., 80% train, 20% test)

#### 4. Model Selection

Choose classification algorithms such as:

- Naive Bayes (commonly used for spam detection)
- Logistic Regression
- Support Vector Machine (SVM)
- Random Forest Classifier
- Gradient Boosting Classifier
- Deep learning models (e.g., LSTM, BERT for advanced cases)

## 5. Model Training

• Train the model on the training data (X\_train, y\_train)

#### 6. Model Evaluation

Use classification metrics to evaluate performance:

- Accuracy
- Precision and Recall

- **F1-Score** (especially important for imbalanced data)
- Confusion Matrix
- ROC-AUC Score

### 7. Model Tuning

- Use Grid Search or Randomized Search for hyperparameter tuning
- Apply **Cross-Validation** for more reliable evaluation

### 8. Deployment

- Integrate the trained model into the email system
- Automatically classify new incoming emails as spam or not spam
- Move spam to a separate folder

### 9. Monitoring and Updating

- Regularly collect feedback (e.g., user marks email as "not spam")
- Retrain the model with new data to adapt to new spam techniques
- Monitor false positives and false negatives

6.A business wants to analyze customer reviews of its products and determine whether the sentiment is positive or negative.

Q: Identify the problem type and outline the step-by-step logic to solve it.

#### **Problem Type:**

This is a **Supervised Machine Learning** problem, specifically a **Binary Classification** problem because the goal is to classify sentiment as **Positive** or **Negative**.

## **Step-by-Step Logic to Solve the Problem:**

## 1. Problem Understanding

- Goal: Determine the sentiment (positive or negative) from customer review text.
- Output: A label (1 for positive, 0 for negative)

#### 2. Data Collection

Collect customer reviews data with:

- Text reviews (e.g., "Great product!", "Terrible quality")
- Sentiment label (positive/negative or 1/0) this is needed for training

#### Sources may include:

- Product reviews from websites
- Customer feedback forms
- Social media posts (if labeled)

## 3. Data Preprocessing

- Clean the text:
  - Convert to lowercase
  - Remove punctuation, numbers, stop words
  - Remove HTML tags or emojis (if applicable)
  - Tokenize text into words
  - Stemming/Lemmatization (optional)

- Convert text to numerical form:
  - Bag of Words (BoW)
  - TF-IDF (Term Frequency-Inverse Document Frequency)
  - Word Embeddings (e.g., Word2Vec, GloVe, or BERT for deep learning models)
- **Split data** into training and testing sets (e.g., 80/20)

#### 4. Model Selection

Choose a classification model suitable for text:

- Naive Bayes (great for text classification)
- Logistic Regression
- Support Vector Machine (SVM)
- Random Forest Classifier
- Deep Learning models (e.g., LSTM, GRU, BERT for more advanced solutions)

### 5. Model Training

• Train the model using the processed text data and labels.

#### 6. Model Evaluation

Evaluate model performance using:

- Accuracy
- **Precision**, **Recall**, **F1-Score** (especially if data is imbalanced)
- Confusion Matrix

7. Model Tuning						
Use Grid Search or Randomized Search for hyperparameter tuning.						
•	Apply <b>Cross-Validation</b> for better generalization.					
. De	ployment					
•	Integrate the sentiment analysis model into the business platform.					
•	Automatically analyze new reviews and tag them as positive or negative.					
•	Generate reports or visualizations based on customer sentiment.					
. Mc	onitoring and Updating					
•	Retrain the model regularly with new reviews.					
•	Monitor changes in customer sentiment over time.					
•	Adapt the model to new slang, product changes, or trends.					
laim	An insurance company wants to predict whether a customer is likely to file a in the next year based on their driving history and demographics.					

• ROC-AUC Score (for probability-based outputs)

## **Problem Type:**

This is a **Supervised Machine Learning** problem, specifically a **Binary Classification** problem, because the goal is to **predict whether a customer will file a claim (Yes/No or 1/0)**.

## **Step-by-Step Logic to Solve the Problem:**

## 1. Problem Understanding

- Goal: Predict whether a customer will file an insurance claim in the next year.
- Output: Binary label 1 (Will file a claim), ∅ (Will not file a claim)

#### 2. Data Collection

Collect historical customer data with relevant features such as:

- Driving history:
  - o Number of past claims
  - Accident history
  - Traffic violations
  - o Mileage driven per year
- Demographics:
  - o Age
  - Gender
  - Marital status
  - Location
  - o Vehicle type
- Target variable:
  - Filed claim next year (Yes/No)

## 3. Data Preprocessing

- Handle missing values (e.g., imputation or deletion)
- Encode categorical features using:
  - One-Hot Encoding (for non-ordinal categories like gender or vehicle type)
  - Ordinal Encoding (for ordered categories)
- Feature scaling (for numerical data if needed, e.g., StandardScaler)
- Class imbalance handling if majority of customers didn't file claims:
  - SMOTE (Synthetic Minority Oversampling Technique)
  - o Class weights in models
- **Split data** into training and test sets (e.g., 80/20)

#### 4. Model Selection

Use classification algorithms such as:

- Logistic Regression
- Decision Tree Classifier
- Random Forest Classifier
- XGBoost or LightGBM
- Support Vector Machine (SVM)
- Neural Network (for more complex problems)

## 5. Model Training

• Train the model using the training data (X\_train, y\_train)

### 6. Model Evaluation

	Evaluate	using	classification	metrics
--	----------	-------	----------------	---------

- Accuracy
- Precision and Recall
- F1-score
- Confusion Matrix
- **ROC-AUC Curve** (for probability-based evaluation)

## 7. Model Tuning

- Apply **Grid Search** or **Randomized Search** to optimize hyperparameters.
- Use **Cross-Validation** to improve generalization.

## 8. Deployment

- Integrate the model into the insurance company's system.
- For every customer, use their profile to predict the likelihood of filing a claim.
- Adjust premiums or trigger further review based on risk levels.

### 9. Monitoring and Updating

- Periodically retrain the model with new customer data.
- Monitor false positives/negatives.
- Use customer feedback or new outcomes to refine the model.
- 8. A streaming platform wants to recommend movies to users by grouping them based on their viewing preferences and watch history.

#### Q: Identify the problem type and outline the step-by-step logic to solve it.

#### **Problem Type:**

This is an **Unsupervised Machine Learning** problem, specifically a **Clustering** or **Recommendation System** problem (depending on the approach).

- If the platform wants to **group users** → it's a **Clustering** problem.
- If the platform wants to recommend movies → it's a Recommendation System (often combining unsupervised and collaborative filtering techniques).

## **Step-by-Step Logic to Solve the Problem:**

Option A: User Clustering for Personalized Recommendations

(If the goal is to group users by behavior first)

## 1. Problem Understanding

- **Goal:** Group users based on viewing preferences and watch history to recommend personalized content.
- Output: User segments or clusters (e.g., Action Lovers, Drama Fans, etc.)

#### 2. Data Collection

Collect data such as:

- Movies watched
- Genres preferred
- Watch time per genre
- Viewing frequency
- Ratings given

• Time of watching (e.g., night/day)

# 3. Data Preprocessing

- Handle missing values
- Normalize numerical features (e.g., watch time)
- Encode categorical features (e.g., genre preferences)
- Create a **user-feature matrix** (rows = users, columns = behavior features)

### 4. Model Selection

Apply clustering algorithms:

- K-Means Clustering
- Hierarchical Clustering
- **DBSCAN** (for irregular clusters)

## 5. Determine Optimal Clusters

Use:

- Elbow Method to find the best value of K
- Silhouette Score for cluster quality

## 6. Cluster Analysis

- Label clusters based on viewing patterns (e.g., Cluster 1: Sci-Fi Fans)
- Use visualizations (PCA, t-SNE) for analysis

### 7. Content Recommendation

- Recommend popular or similar content within the user's cluster
- Use collaborative filtering within clusters for better suggestions

# ✓ Option B: Build a Recommendation System Directly

## 1. Use Collaborative Filtering

- User-Based: Recommend movies watched by similar users.
- **Item-Based:** Recommend movies similar to the ones a user has already watched.

## 2. Use Content-Based Filtering

 Recommend movies with similar genres, actors, directors, or tags as previously watched items.

## 3. Use Matrix Factorization Techniques

- SVD (Singular Value Decomposition)
- ALS (Alternating Least Squares)

## 4. Use Hybrid Models

• Combine collaborative + content-based methods (used by Netflix, Amazon, etc.)

#### 5. Model Evaluation

- Use metrics like:
  - o Precision@K, Recall@K

- **RMSE** (for rating prediction)
- **AUC** (for binary relevance)
- Use A/B testing for real-world impact

### 6. Deployment

- Show top-N movie recommendations to each user
- Update recommendations dynamically as user behavior changes

## 7. Monitoring and Updating

- Continuously gather user feedback (watch behavior, ratings)
- Retrain/update models regularly for evolving preferences
- 9. A hospital wants to predict the recovery time of patients after surgery based on their age, medical history, and lifestyle habits.
- Q: Identify the problem type and outline the step-by-step logic to solve it.

### **Problem Type:**

This is a **Supervised Machine Learning** problem, specifically a **Regression** problem, because the target variable — **recovery time (in days/weeks)** — is **continuous**.

## **Step-by-Step Logic to Solve the Problem:**

### 1. Problem Understanding

- Goal: Predict how long it will take for a patient to recover after surgery.
- Output: A continuous value (e.g., number of days)

### 2. Data Collection

Gather patient data, such as:

- Demographics:
  - Age
  - o Gender

### Medical history:

- Pre-existing conditions (e.g., diabetes, hypertension)
- o Previous surgeries
- Medication usage

### Lifestyle habits:

- Smoking status
- Alcohol consumption
- Physical activity level
- o Diet and sleep patterns

### • Surgery-related features:

- Type of surgery
- o Surgery duration
- o Post-operative complications

### • Target variable:

Recovery time (in days/weeks)

# 3. Data Preprocessing

- Handle missing values (e.g., mean/mode imputation or removing rows)
- Encode categorical variables (e.g., One-Hot Encoding for surgery type)
- Normalize or scale numerical features (e.g., age, BMI)
- Outlier detection and treatment (especially in recovery time)
- Split the data into training and testing sets (e.g., 80/20 split)

#### 4. Model Selection

Choose regression models:

- Linear Regression
- Decision Tree Regressor
- Random Forest Regressor
- Gradient Boosting Regressor (e.g., XGBoost, LightGBM)
- Support Vector Regressor
- Neural Networks (for complex data patterns)

## 5. Model Training

• Fit the chosen model(s) on the training dataset.

### 6. Model Evaluation

Evaluate using regression metrics:

- Mean Absolute Error (MAE)
- Mean Squared Error (MSE)
- Root Mean Squared Error (RMSE)

7. Model Tuning					
Use <b>Grid Search</b> or <b>Randomized Search</b> to find optimal hyperparameters.					
Apply <b>Cross-Validation</b> to avoid overfitting.					
ployment					
Integrate the trained model into the hospital's patient management system.					
Predict recovery time for new patients after surgery.					
Use predictions to plan follow-up care, hospital stays, and resource allocation.					
Continuously update the model with new patient recovery data.  Track prediction accuracy over time.  Improve model based on feedback from doctors and outcomes.					
A university wants to predict a student's final exam score based on study					

• R<sup>2</sup> Score (coefficient of determination)

### **Step-by-Step Logic to Solve the Problem:**

### 1. Problem Understanding

- **Goal:** Predict a student's final exam score.
- Output: A numeric value (e.g., marks out of 100)

#### 2. Data Collection

Gather student data such as:

- Study hours (daily or weekly)
- Attendance percentage
- Past academic performance (e.g., previous test scores, GPA)
- Optional features:
  - o Participation in class
  - Use of online learning resources
  - Assignment grades

## 3. Data Preprocessing

- Handle missing values (e.g., fill with mean or median)
- Convert categorical data (if any) using encoding techniques
- Scale or normalize numerical features (especially study hours or GPA)
- Detect and handle outliers (e.g., very high or low scores)
- Split the dataset into training and testing sets (e.g., 80/20 split)

### 4. Model Selection

Choose appropriate regression algorithms:

- Linear Regression
- Ridge or Lasso Regression (for regularization)
- Decision Tree Regressor
- Random Forest Regressor
- Gradient Boosting Regressor (like XGBoost)
- Support Vector Regressor (SVR)

## 5. Model Training

- Train the selected model(s) on the training data.
- Tune features for better performance if needed.

#### 6. Model Evaluation

Use regression metrics to evaluate model performance:

- Mean Absolute Error (MAE)
- Mean Squared Error (MSE)
- Root Mean Squared Error (RMSE)
- R<sup>2</sup> Score (for variance explanation)

## 7. Model Tuning

- Use techniques like Grid Search CV or Random Search CV to find the best hyperparameters.
- Apply Cross-Validation to ensure the model generalizes well.

# 8. Deployment

- Deploy the model in the university system to estimate expected exam scores.
- Use results to:
  - o Provide early academic support
  - Suggest study improvements
  - o Inform parents/teachers

## 9. Monitoring and Updating

- Continuously collect new student data and exam results.
- Retrain the model periodically to maintain accuracy.
- Improve the model with additional features (e.g., student motivation, learning styles).