

# Fundamentals of Machine Learning (4341603) - Summer 2025 Solution

Milav Dabgar

May 17, 2025

## Question 1(a) [3 marks]

Define machine Learning. Give any two applications of machine learning.

### Solution

Machine Learning is a subset of artificial intelligence that enables computers to learn and make decisions from data without being explicitly programmed for every task.

#### Applications:

- **Email spam detection:** Automatically identifies and filters spam emails
- **Recommendation systems:** Suggests products on e-commerce sites like Amazon

Table 1. ML vs Traditional Programming

Traditional Programming	Machine Learning
Input data + Program $\rightarrow$ Output	Input data + Output $\rightarrow$ Program
Rules are explicitly coded	Rules are learned from data

### Mnemonic

“ML = Make Learning from data”

## Question 1(b) [4 marks]

Define: Under fitting and overfitting.

### Solution

**Underfitting** occurs when a model is too simple to capture underlying patterns in data, resulting in poor performance on both training and test data.

**Overfitting** occurs when a model learns training data too well, including noise, causing poor performance on new unseen data.

Table 2. Comparison

Aspect	Underfitting	Overfitting
Training accuracy	Low	High
Test accuracy	Low	Low
Model complexity	Too simple	Too complex
Solution	Increase complexity	Reduce complexity

**Mnemonic**

“Under = Under-performs, Over = Over-learns”

**Question 1(c) [7 marks]**

Describe different types of machine learning with suitable example.

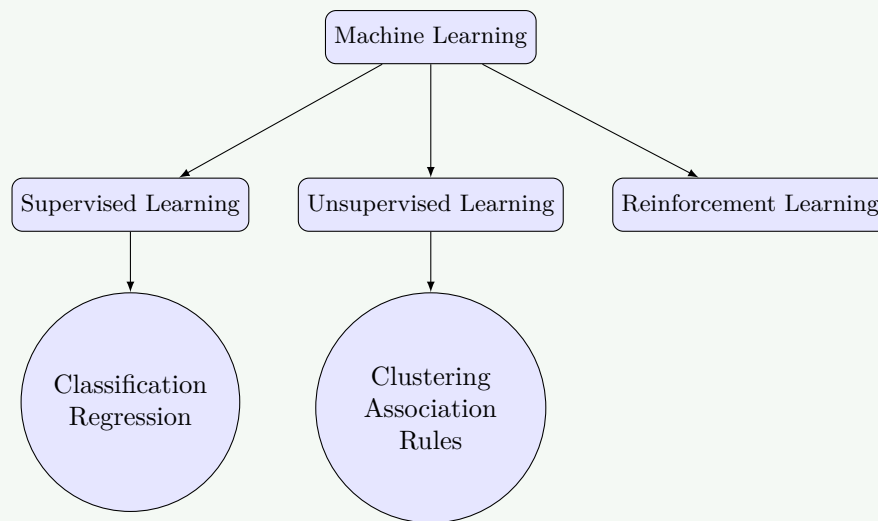
**Solution****Table 3.** Types of Machine Learning

Type	Description	Example
<b>Supervised</b>	Uses labeled training data	Email classification
<b>Unsupervised</b>	No labeled data, finds patterns	Customer segmentation
<b>Reinforcement</b>	Learns through rewards/penalties	Game playing AI

**Supervised Learning** uses input-output pairs to train models. The algorithm learns from examples to predict outcomes for new data.

**Unsupervised Learning** discovers hidden patterns in data without target labels. It groups similar data points together.

**Reinforcement Learning** trains agents to make decisions by rewarding good actions and penalizing bad ones.

**Figure 1.** Types of Machine Learning**Mnemonic**

“Super Un-supervised Reinforces learning”

**Question 1(c) OR [7 marks]**

Describe different tools and technology used in the field machine learning.

**Solution****Table 4.** ML Tools and Technologies

Category	Tools	Purpose
<b>Programming</b>	Python, R	Core development
<b>Libraries</b>	Scikit-learn, TensorFlow	Model building
<b>Data Processing</b>	Pandas, NumPy	Data manipulation
<b>Visualization</b>	Matplotlib, Seaborn	Data plotting

**Python** is the most popular language due to its simplicity and extensive libraries.  
**Scikit-learn** provides simple tools for data mining and analysis, perfect for beginners.  
**TensorFlow** and **PyTorch** are advanced frameworks for deep learning applications.  
**Jupyter Notebook** offers interactive development environment for experimentation.



**Figure 2.** ML Tools Workflow

### Mnemonic

“Python Pandas Scikit Tensor Jupyter”

## Question 2(a) [3 marks]

Give the difference between Qualitative data and Quantitative data.

### Solution

**Table 5.** Qualitative vs Quantitative Data

Qualitative Data	Quantitative Data
<b>Non-numerical</b> categories	<b>Numerical</b> values
Colors, names, grades	Height, weight, price
Cannot be measured	Can be measured

**Qualitative data** describes qualities or characteristics that cannot be measured numerically.

**Quantitative data** represents measurable quantities expressed as numbers.

### Mnemonic

“Quality = Categories, Quantity = Numbers”

## Question 2(b) [4 marks]

Find the mean and median for the following data: 3,4,5,5,7,8,9,11,12,14.

### Solution

**Given data:** 3, 4, 5, 5, 7, 8, 9, 11, 12, 14

**Mean calculation:**

- Sum =  $3 + 4 + 5 + 5 + 7 + 8 + 9 + 11 + 12 + 14 = 78$
- Count = 10 numbers
- Mean** =  $78/10 = 7.8$

**Median calculation:**

- Data is already sorted

- For 10 numbers: Median = (5th + 6th value)/2
- **Median** =  $(7 + 8)/2 = 7.5$

**Table 6.** Results

Measure	Value
Mean	7.8
Median	7.5

**Mnemonic**

“Mean = Average, Median = Middle”

**Question 2(c) [7 marks]**

Describe machine learning activities in detail.

**Solution****Table 7.** Machine Learning Activities

Activity	Description	Example
<b>Data Collection</b>	Gathering relevant data	Survey responses
<b>Data Preprocessing</b>	Cleaning and preparing data	Removing duplicates
<b>Feature Selection</b>	Choosing important variables	Age, income for loans
<b>Model Training</b>	Teaching algorithm patterns	Feeding training data
<b>Model Evaluation</b>	Testing model performance	Accuracy measurement

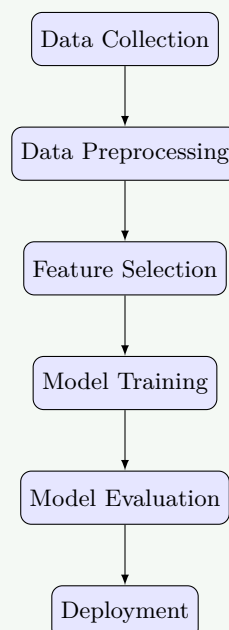
**Data Collection** involves gathering information from various sources like databases, sensors, or surveys.

**Data Preprocessing** includes cleaning, transforming, and organizing raw data for analysis.

**Feature Selection** identifies the most relevant variables that contribute to predictions.

**Model Training** uses algorithms to learn patterns from prepared training data.

**Model Evaluation** tests how well the trained model performs on new, unseen data.



**Figure 3.** Machine Learning Activities Flow**Mnemonic**

“Collect Process Feature Train Evaluate Deploy”

**Question 2(a) OR [3 marks]**

Give the difference between predictive model and descriptive model.

**Solution****Table 8.** Predictive vs Descriptive Models

Predictive Model	Descriptive Model
<b>Forecasts</b> future outcomes	<b>Explains</b> current patterns
Uses supervised learning	Uses unsupervised learning
Stock price prediction	Customer segmentation

**Predictive models** use historical data to make predictions about future events or unknown outcomes.

**Descriptive models** analyze existing data to understand current patterns and relationships.

**Mnemonic**

“Predict = Future, Describe = Present”

**Question 2(b) OR [4 marks]**

Classify the following using appropriate data type: hair color, gender, blood group type, time of day.

**Solution****Table 9.** Data Type Classification

Data	Type	Reason
<b>Hair color</b>	Nominal	Categories with no order
<b>Gender</b>	Nominal	Categories with no order
<b>Blood group</b>	Nominal	Categories with no order
<b>Time of day</b>	Continuous	Measurable quantity

**Nominal data** represents categories without any natural ordering.

**Continuous data** can take any value within a range and is measurable.

**Mnemonic**

“Names = Nominal, Numbers = Numerical”

**Question 2(c) OR [7 marks]**

Explain various methods used in data pre-processing.

## Solution

Table 10. Data Preprocessing Methods

Method	Purpose	Example
<b>Data Cleaning</b>	Remove errors and inconsistencies	Fix typos, remove duplicates
<b>Data Integration</b>	Combine multiple sources	Merge customer databases
<b>Data Transformation</b>	Convert to suitable format	Normalize values 0-1
<b>Data Reduction</b>	Reduce dataset size	Select important features

**Data Cleaning** removes or corrects erroneous, incomplete, or irrelevant data.

**Data Integration** combines data from multiple sources into a unified dataset.

**Data Transformation** converts data into appropriate formats for analysis.

**Data Reduction** decreases dataset size while maintaining information quality.

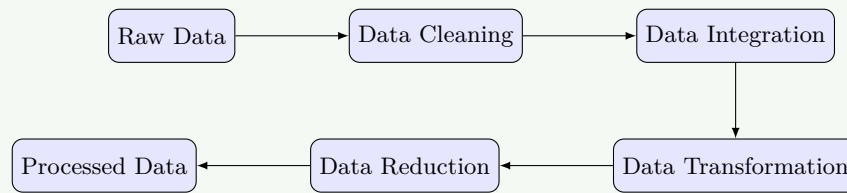


Figure 4. Data Preprocessing Pipeline

## Mnemonic

“Clean Integrate Transform Reduce”

## Question 3(a) [3 marks]

Give difference between classification and regression.

## Solution

Table 11. Classification vs Regression

Classification	Regression
<b>Discrete</b> output	<b>Continuous</b> output
Predicts categories	Predicts numerical values
Email: spam/not spam	House price prediction

**Classification** predicts discrete categories or classes from input data.

**Regression** predicts continuous numerical values from input data.

## Mnemonic

“Class = Categories, Regress = Real numbers”

## Question 3(b) [4 marks]

Write confusion matrix using appropriate example. Calculate accuracy and error rate for it.

**Solution****Example: Email Classification****Table 12.** Confusion Matrix

	Predicted Spam	Predicted Not Spam
Actual Spam	85 (TP)	15 (FN)
Actual Not Spam	10 (FP)	90 (TN)

**Calculations:**

- **Accuracy** =  $(TP + TN) / (TP + TN + FP + FN) = (85 + 90) / 200 = 87.5\%$
- **Error Rate** =  $(FP + FN) / (TP + TN + FP + FN) = (10 + 15) / 200 = 12.5\%$

**Key Terms:**

- **TP:** True Positive - Correctly predicted spam
- **TN:** True Negative - Correctly predicted not spam

**Mnemonic**

“True Positive True Negative = Correct predictions”

**Question 3(c) [7 marks]**

Explain KNN algorithm in detail.

**Solution**

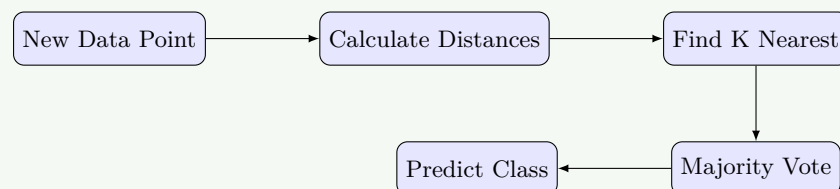
**K-Nearest Neighbors (KNN)** is a simple classification algorithm that classifies data points based on the majority class of their K nearest neighbors.

**Table 13.** KNN Algorithm Steps

Step	Description	Example
<b>Choose K</b>	Select number of neighbors	K=3
<b>Calculate Distance</b>	Find distance to all points	Euclidean distance
<b>Find Neighbors</b>	Identify K closest points	3 nearest points
<b>Vote</b>	Majority class wins	2 cats, 1 dog → cat

**Working Process:**

1. **Calculate distances** between test point and all training points
2. **Sort distances** and select K nearest neighbors
3. **Count votes** from each class among neighbors
4. **Assign class** with majority votes

**Figure 5.** KNN Process Flow**Advantages:**

- **Simple to implement** and understand
- **No training required** - lazy learning algorithm

**Mnemonic**

“K Nearest Neighbors Vote for classification”

**Question 3(a) OR [3 marks]**

Give any three applications of multiple linear regression.

**Solution**

**Applications of Multiple Linear Regression:**

Table 14. Applications		
Application	Variables	Purpose
House Price Prediction	Size, location, age	Estimate property value
Sales Forecasting	Advertising, season, price	Predict revenue
Medical Diagnosis	Symptoms, age, history	Risk assessment

**Multiple Linear Regression** uses multiple input variables to predict a continuous output variable.

**Mnemonic**

“Multiple inputs, One output”

**Question 3(b) OR [4 marks]**

Explain bagging, boosting and stacking in detail.

**Solution**

**Table 15. Ensemble Methods**

Method	Approach	Example
<b>Bagging</b>	Parallel training, average results	Random Forest
<b>Boosting</b>	Sequential training, learn from errors	AdaBoost
<b>Stacking</b>	Meta-learner combines models	Neural network combiner

**Bagging** trains multiple models on different data subsets and averages predictions.

**Boosting** trains models sequentially, each learning from previous model's mistakes.

**Stacking** uses a meta-model to learn how to combine predictions from base models.

**Mnemonic**

“Bag parallel, Boost sequential, Stack meta”

**Question 3(c) OR [7 marks]**

Explain single linear regression with its application.



### Solution

**Single Linear Regression** finds the best straight line relationship between one input variable (X) and one output variable (Y).

**Formula:**  $Y = a + bX$

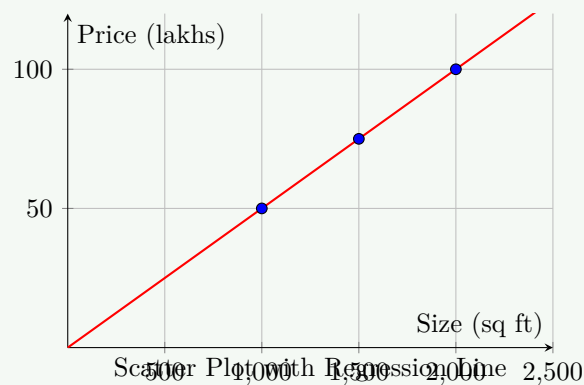
- **a:** Y-intercept
- **b:** Slope of line

**Table 16.** Application Example - House Price vs Size

House Size (sq ft)	Price (lakhs)
1000	50
1500	75
2000	100

#### Working Process:

1. **Collect data** with input-output pairs
2. **Plot points** on scatter graph
3. **Find best line** that minimizes error
4. **Make predictions** using line equation



**Figure 6.** Linear Regression Visualization

#### Applications:

- **Sales vs Advertising:** More ads → More sales
- **Temperature vs Ice cream sales:** Hot weather → More sales

### Mnemonic

“One X predicts One Y with a line”

## Question 4(a) [3 marks]

Define the following: (1) support (2) confidence.

### Solution

**Support** measures how frequently an itemset appears in the dataset.

**Confidence** measures how often items in consequent appear when antecedent is present.

**Table 17.** Definitions

Measure	Formula	Example
<b>Support</b>	$\text{Count}(\text{itemset}) / \text{Total transactions}$	Bread appears in 60% transactions
<b>Confidence</b>	$\text{Support}(\text{AUB}) / \text{Support}(\text{A})$	80% who buy bread also buy butter

**Support** = Frequency of occurrence

**Confidence** = Reliability of rule

#### Mnemonic

“Support = How often, Confidence = How reliable”

## Question 4(b) [4 marks]

Explain applications of unsupervised learning.

#### Solution

**Table 18.** Unsupervised Learning Applications

Application	Purpose	Example
<b>Customer Segmentation</b>	Group similar customers	Marketing campaigns
<b>Data Compression</b>	Reduce data size	Image compression
<b>Anomaly Detection</b>	Find unusual patterns	Fraud detection
<b>Recommendation Systems</b>	Suggest similar items	Music recommendations

**Customer Segmentation** groups customers with similar buying behavior for targeted marketing.

**Data Compression** reduces storage space by finding patterns and removing redundancy.

**Anomaly Detection** identifies unusual patterns that may indicate fraud or errors.

#### Mnemonic

“Segment Compress Detect Recommend”

## Question 4(c) [7 marks]

Write and explain apriori algorithm with suitable example.

#### Solution

**Apriori Algorithm** finds frequent itemsets and generates association rules for market basket analysis.

**Table 19.** Algorithm Steps

Step	Description	Example
<b>Find frequent 1-itemsets</b>	Count individual items	{Bread}:4, {Milk}:3
<b>Generate 2-itemsets</b>	Combine frequent items	{Bread,Milk}:2
<b>Apply minimum support</b>	Filter infrequent sets	Keep if support $\geq 50\%$
<b>Generate rules</b>	Create if-then rules	Bread $\rightarrow$ Milk

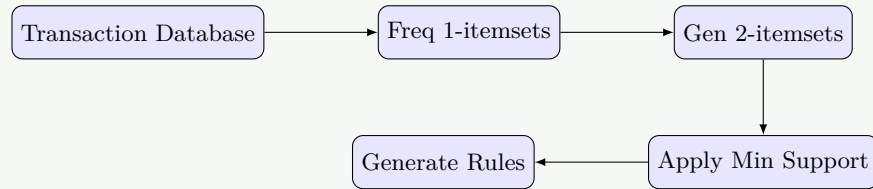
#### Example Dataset:

- Transaction 1: {Bread, Milk, Eggs}
- Transaction 2: {Bread, Milk}

- Transaction 3: {Bread, Eggs}
- Transaction 4: {Milk, Eggs}

**Working Process:**

1. **Scan database** to count item frequencies
2. **Generate candidate itemsets** of increasing size
3. **Prune infrequent itemsets** below minimum support
4. **Generate association rules** from frequent itemsets

**Figure 7.** Apriori Algorithm Steps**Mnemonic**

“A-priori knowledge helps find frequent patterns”

**Question 4(a) OR [3 marks]**

List out the difference between clustering and classification.

**Solution****Table 20.** Clustering vs Classification

Clustering	Classification
<b>Unsupervised</b> learning	<b>Supervised</b> learning
No labeled data	Uses labeled training data
Groups similar data	Assigns predefined labels

**Clustering** discovers hidden groups in unlabeled data.

**Classification** assigns new data to known categories using trained models.

**Mnemonic**

“Cluster = Groups unknown, Classify = Labels known”

**Question 4(b) OR [4 marks]**

Explain the clustering process in detail.

**Solution****Table 21.** Clustering Process Steps

Step	Description	Purpose
<b>Data Preparation</b>	Clean and normalize data	Ensure quality input
<b>Distance Metric</b>	Choose similarity measure	Euclidean, Manhattan
<b>Algorithm Selection</b>	Pick clustering method	K-means, Hierarchical
<b>Cluster Validation</b>	Evaluate cluster quality	Silhouette score

**Clustering Process** groups similar data points together based on their characteristics.  
**Key decisions include choosing the number of clusters and appropriate distance metrics.**  
**Validation ensures clusters are meaningful and well-separated.**

### Mnemonic

“Prepare Distance Algorithm Validate”

## Question 4(c) OR [7 marks]

Write and explain K-means clustering algorithm with suitable example.

### Solution

**K-means** partitions data into K clusters by minimizing within-cluster sum of squares.

**Table 22.** Algorithm Steps

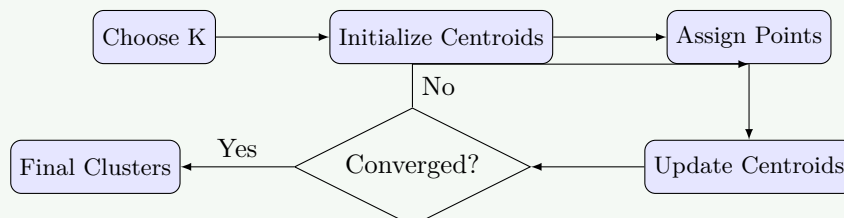
Step	Description	Example
<b>Initialize centroids</b>	Random K center points	C1(2,3), C2(8,7)
<b>Assign points</b>	Each point to nearest centroid	Point(1,2) → C1
<b>Update centroids</b>	Mean of assigned points	New C1(1.5, 2.5)
<b>Repeat</b>	Until centroids stop moving	Convergence

### Example: Customer Income vs Age

- Customer 1: (Income=30k, Age=25)
- Customer 2: (Income=35k, Age=30)
- Customer 3: (Income=70k, Age=45)
- Customer 4: (Income=75k, Age=50)

### Working Process:

1. **Choose K=2** clusters for young/old customers
2. **Initialize centroids** randomly
3. **Calculate distances** from each customer to centroids
4. **Assign customers** to nearest centroid
5. **Update centroid positions** to center of assigned customers
6. **Repeat until stable**



**Figure 8.** K-Means Logic

**Mnemonic**

“K centroids Mean their assigned points”

**Question 5(a) [3 marks]**

List the applications of matplotlib.

**Solution**

**Table 23.** Matplotlib Applications

Application	Purpose	Example
<b>Data Visualization</b>	Create charts and graphs	Bar charts, histograms
<b>Scientific Plotting</b>	Research presentations	Mathematical functions
<b>Dashboard Creation</b>	Interactive displays	Business metrics

**Matplotlib** is Python’s primary plotting library for creating static, animated, and interactive visualizations.

**Mnemonic**

“Mat-plot-lib = Math Plotting Library”

**Question 5(b) [4 marks]**

Write down code to plot a vertical line and horizontal line using matplotlib.

**Solution**

```

1  import matplotlib.pyplot as plt
2
3  # Create figure
4  plt.figure(figsize=(8, 6))
5
6  # Plot vertical line at x=3
7  plt.axvline(x=3, color='red', linestyle='--', label='Vertical Line')
8
9  # Plot horizontal line at y=2
10 plt.axhline(y=2, color='blue', linestyle='-', label='Horizontal Line')
11
12 # Add labels and title
13 plt.xlabel('X-axis')
14 plt.ylabel('Y-axis')
15 plt.title('Vertical and Horizontal Lines')
16 plt.legend()
17 plt.grid(True)
18 plt.show()

```

**Key Functions:**

- `axvline()`: Creates vertical line
- `axhline()`: Creates horizontal line

**Mnemonic**

“axvline = Vertical, axhline = Horizontal”

## Question 5(c) [7 marks]

Explain features and applications of Scikit-Learn.

### Solution

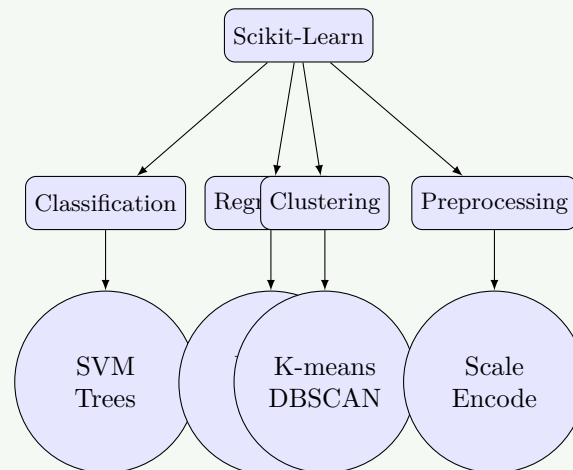
**Table 24.** Scikit-Learn Features

Feature	Description	Example
<b>Simple API</b>	Easy to use interface	fit(), predict()
<b>Multiple Algorithms</b>	Various ML methods	SVM, Random Forest
<b>Data Preprocessing</b>	Built-in data tools	StandardScaler
<b>Model Evaluation</b>	Performance metrics	accuracy_score

**Scikit-Learn** is Python's most popular machine learning library providing simple tools for data analysis.

#### Applications:

- **Classification:** Email spam detection
- **Regression:** House price prediction
- **Clustering:** Customer segmentation
- **Dimensionality Reduction:** Data visualization



**Figure 9.** Scikit-Learn Components

### Mnemonic

“Scikit = Science Kit for machine learning”

## Question 5(a) OR [3 marks]

Give the purpose of NumPy in machine learning.

### Solution

**Table 25.** NumPy Purpose in ML

Purpose	Description	Benefit
<b>Numerical Computing</b>	Fast array operations	Efficient calculations
<b>Foundation Library</b>	Base for other libraries	Pandas, Scikit-learn use it
<b>Mathematical Functions</b>	Built-in math operations	Statistics, linear algebra

**NumPy** provides the foundation for numerical computing in Python machine learning applications.  
**Essential for handling large datasets and performing mathematical operations efficiently.**

#### Mnemonic

“Num-Py = Numerical Python”

## Question 5(b) OR [4 marks]

Write down steps to import csv file in pandas.

#### Solution

```

1 import pandas as pd
2
3 # Step 1: Import pandas library
4 # Step 2: Use read_csv() function
5 data = pd.read_csv('filename.csv')
6
7 # Step 3: Display first few rows
8 print(data.head())
9
10 # Optional: Specify parameters
11 data = pd.read_csv('file.csv',
12                   delimiter=',',
13                   header=0,
14                   index_col=0)
```

#### Steps:

1. **Import pandas** library
2. **Use read\_csv()** function with filename
3. **Verify data** with head() method

#### Mnemonic

“Import Read Verify”

## Question 5(c) OR [7 marks]

Explain features and applications of Pandas.

#### Solution

**Table 26.** Pandas Features

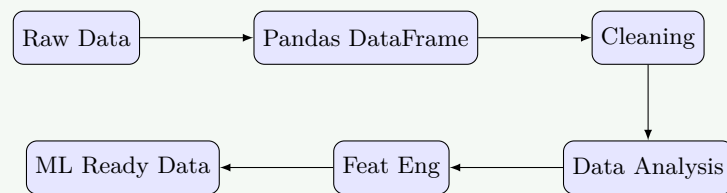
Feature	Description	Example
<b>Data Structures</b>	DataFrame and Series	Tabular data handling
<b>Data I/O</b>	Read/write multiple formats	CSV, Excel, JSON
<b>Data Cleaning</b>	Handle missing values	dropna(), fillna()
<b>Data Analysis</b>	Statistical operations	groupby(), describe()

**Pandas** is the primary data manipulation library in Python for machine learning projects.

#### Key Capabilities:

- **Data Loading** from various file formats

- **Data Cleaning** and preprocessing operations
- **Data Transformation** and reshaping
- **Statistical Analysis** and aggregation



**Figure 10.** Pandas Workflow

### Mnemonic

“Pandas = Panel Data for analysis”