

**DATA ANALYTICS**

**BRANCH:- OPEN ELECTIVE**

**SEMESTER – FIFTH**

***These important questions have been prepared using your previous exam papers (PYQs), verified concepts, and additional reference from trusted online academic sources. For deeper understanding, please refer to your class notes as well.***

 **For more study materials, notes, important questions, and updates, visit –**

[DiplomaWallah.in](https://DiplomaWallah.in)

 **To join our official WhatsApp group for free updates, contact: [CLICK HERE TO JOIN](#)**

---

**1 HIGH & LONG IMPORTANT QUESTIONS**

1. **Define Artificial Intelligence (AI).** Explain the different **types of AI** (Narrow, General, Super AI) and discuss the major **applications of AI** in the real world (Healthcare, Finance, E-commerce).
2. **What is Machine Learning (ML)?** Explain the **Machine Learning Workflow/Pipeline** (Data Collection → Preprocessing → Model Training → Evaluation → Deployment) with a neat block diagram.
3. **Difference between Supervised, Unsupervised, and Reinforcement Learning.** Give examples for each.
4. **What is Data Preprocessing?** Why is it important? Explain the steps to handle **Missing Values** and **Outliers** in a dataset.
5. **Explain Linear Regression.** What are its assumptions? Differentiate between Simple Linear Regression and Multiple Linear Regression.
6. **What is a Decision Tree?** Explain key terms like **Root Node, Leaf Node, Entropy, and Information Gain.** How do we prevent overfitting in Decision Trees (Pruning)?
7. **Define Deep Learning.** Explain the architecture of an **Artificial Neural Network (ANN)** including Input Layer, Hidden Layer, and Output Layer. How is it different from biological neurons?
8. **What is Cloud Computing in AI?** Explain the difference between **IaaS, PaaS, and SaaS** service models.
9. **Explain the Confusion Matrix** and how to calculate the following metrics:
  - Accuracy
  - Precision
  - Recall
  - F1-Score
10. **What is Natural Language Processing (NLP)?** Explain the steps of text processing: **Tokenization, Stemming, Lemmatization, and Stop Word Removal.**
11. **Explain the K-Means Clustering algorithm.** How do we decide the number of clusters (K)?

---

12. **What is Docker?** Explain the architecture of Docker (Client, Daemon, Images, Containers) and why containers are used in AI deployment.

---

## **2** IMPORTANT & SHORT QUESTIONS

1. **Differentiate between:**
  - o AI vs. Machine Learning vs. Deep Learning.
2. **Python Libraries:** Briefly explain the use of the following libraries in AI:
  - o **NumPy** (for numerical data)
  - o **Pandas** (for data manipulation)
  - o **Matplotlib/Seaborn** (for visualization)
  - o **Scikit-learn** (for ML models)
3. **What is Big Data?** Explain the **5 Vs of Big Data** (Volume, Velocity, Variety, Veracity, Value).
4. **Activation Functions:** Define and draw the graph for:
  - o Sigmoid
  - o ReLU (Rectified Linear Unit)
  - o Tanh
5. **What is Version Control System (VCS)?** Why is **Git** and **GitHub** useful for AI projects?
6. **Data Visualization:** Explain the difference between a **Histogram** and a **Scatter Plot**. Why do we visualize data?
7. **What is a Hypothesis?** Explain Null Hypothesis and P-value in short.
8. **Define Sentiment Analysis.** Give one real-world business example of it.
9. **What is MLOps?** Why is it needed in the industry?
10. **Explain SVM (Support Vector Machine).** What is a Hyperplane and Support Vector?

---

## **3** "AA BHI SAKTA HAI" QUESTIONS (20–30% Probability)

*(Tricky or new syllabus topics. Read them once if you want to score Top grades.)*

1. **Ethics in AI:** Discuss the ethical challenges in AI like **Bias, Privacy, and Job displacement**.
2. **Ensemble Learning:** Write a short note on **Bagging (Random Forest)** and **Boosting**. Why are they better than single models?
3. **Bayes' Theorem:** State the theorem and its application in Naive Bayes Classification.
4. **Dimensionality Reduction:** What is **PCA (Principal Component Analysis)** and why is it used to reduce data dimensions?
5. **TF-IDF:** What does Term Frequency-Inverse Document Frequency mean in NLP?
6. **Maths for AI:** Briefly explain the concept of **Eigenvalues and Eigenvectors** (Linear Algebra) and their use in AI.
7. **Data Integration:** What are the challenges in integrating data from different sources?

8. **Backpropagation:** Briefly explain how a Neural Network "learns" by updating weights to reduce loss.

---

 **Analyst Tip for JUT Exam:**

Since this subject (AI&ML) has a lot of "practical" topics like Python and Docker in the syllabus, **do not write code** in the exam unless specifically asked (which is rare in theory exams). Instead, **describe the concept**.

- **Example:** Instead of writing code to clean data, write: *"Data cleaning involves removing null values using methods like mean imputation..."*
- **Diagrams are key:** Draw the Neural Network layers, the ML pipeline blocks, and the Confusion Matrix table. This guarantees good marks in JUT.

**QUICK REVISE**

 **1. AI Fundamentals & Setup (Weeks 1, 2)**

Topic	Quick Note
<b>What is AI?</b>	Making computers smart like humans to solve problems, learn, and decide.
<b>Types of AI</b>	<b>Weak/Narrow AI:</b> Does <i>one</i> specific task (e.g., Google Search, Alexa). <b>Strong/General AI:</b> Can do <i>any</i> intellectual task (future goal).
<b>AI SDLC</b>	The process is <b>data-centric</b> and involves continuous <b>monitoring</b> and <b>retraining</b> of the model, unlike traditional software development.
<b>Ethics in AI</b>	<b>Fairness</b> (avoiding bias), <b>Privacy</b> (protecting data), and <b>Accountability</b> (fixing errors).
<b>Version Control (Git)</b>	A system to track and manage changes in your code/data over time, allowing easy rollback and teamwork.

 **2. ML, Cloud & Data Basics (Week 3, 5)**

## Diploma wallah

Topic	Quick Note
<b>Machine Learning</b>	Teaching computers to learn patterns from data <b>without explicit programming</b> .
<b>ML Types</b>	<b>Supervised:</b> Learning from labeled examples (Input + Correct Answer). <b>Unsupervised:</b> Finding hidden groups/patterns in unlabeled data.
<b>ML Pipeline</b>	Data Collection > Data Preprocessing > Model Training > Model Testing > Deployment.
<b>Cloud Models</b>	<b>IaaS:</b> Renting the <b>hardware</b> (servers, storage). <b>PaaS:</b> Renting the <b>platform</b> (OS, tools to build). <b>SaaS:</b> Renting the <b>ready-made software</b> (Gmail, CRM).
<b>Big Data (5 Vs)</b>	<b>Volume</b> (Large size), <b>Velocity</b> (High speed), <b>Variety</b> (Different formats), <b>Veracity</b> (Accuracy), <b>Value</b> (Usefulness).
<b>NumPy &amp; Pandas</b>	<b>NumPy:</b> Used for fast calculation on arrays and matrices. <b>Pandas:</b> Used for manipulating table-like data (DataFrames).

## 3. Data Exploration & Math (Week 4, 6)

Topic	Quick Note
<b>Data Visualization</b>	Presenting data visually (graphs/charts) to easily find <b>patterns, trends, and outliers</b> .
<b>Central Limit Theorem</b>	Even if the original data is weird, the average of many samples tends to form a <b>Normal Distribution</b> (Bell Curve).
<b>Bayes' Theorem</b>	Formula to calculate the <b>probability of an event</b> based on information we already know (prior knowledge).
<b>Correlation</b>	Measures how strongly two variables are <b>related</b> (from -1 to +1).
<b>Vectors &amp; Matrices</b>	Data and parameters in AI are stored as <b>Vectors</b> (1D array) and <b>Matrices</b> (2D array) for efficient processing.
<b>Eigenvalues/vectors</b>	Special numbers/vectors that show the <b>direction of maximum variance</b> (used in PCA for dimension reduction).

## 4. Data Preprocessing (Week 7)

Topic	Quick Note
<b>Missing Values</b>	Data entries that are empty. Handled by <b>imputation</b> (filling with mean/median) or <b>deletion</b> .
<b>Outliers</b>	Data points that are far away from the rest of the data. Handled by <b>removal</b> or using <b>log transformation</b> .
<b>Normalization</b>	Scaling data values to a small, fixed range (usually <b>0 to 1</b> ).
<b>Standardization</b>	Scaling data so that the mean is <b>0</b> and the standard deviation is <b>1</b> (Z-Score).
<b>Data Reduction</b>	Making the data smaller (reducing rows or features) while keeping important information.

## 5. Regression & Evaluation (Week 8)

Topic	Quick Note
<b>Data Splitting</b>	Dividing data into <b>Training</b> (for learning), <b>Validation</b> (for tuning), and <b>Testing</b> (for final grade).
<b>Overfitting</b>	Model learns the <b>training data too perfectly</b> (memorizes noise), performs badly on new data.
<b>Underfitting</b>	Model is <b>too simple</b> , cannot capture the basic patterns of the data.
<b>Linear Regression</b>	Used to <b>predict a continuous value</b> (like price, temperature) based on input features using a straight line equation.
<b>R-Squared</b>	Measures how well the regression line <b>fits the observed data</b> (higher is better, close to 1).
<b>RMSE</b>	Measures the <b>average size of the errors</b> (difference between predicted and actual value).
<b>Cross-Validation</b>	Method (like K-Fold) to repeatedly train and test a model on different subsets of data to get a <b>reliable performance estimate</b> .

## 6. Classification & Ensemble (Week 9)

## Diploma wallah

Topic	Quick Note
<b>Classification</b>	Used to <b>predict a category/class</b> (like Yes/No, Spam/Not Spam).
<b>Decision Tree</b>	A flow-chart model that makes decisions by asking a series of <b>Yes/No questions</b> on features.
<b>Entropy</b>	Measures the <b>impurity</b> or randomness in a group of data.
<b>Information Gain</b>	The drop in entropy after a split; used to decide the <b>best feature</b> to split the tree on.
<b>Logistic Regression</b>	Used for <b>Binary Classification</b> (Yes/No), uses the Sigmoid function to output probability (0 to 1).
<b>Confusion Matrix</b>	A table summarizing the model's predictions: <b>TP, TN, FP, FN</b> .

## 7. Deep Learning, Unsupervised & MLOps (Week 10)

Topic	Quick Note
<b>Deep Learning</b>	Uses <b>Neural Networks</b> with many hidden layers to solve complex problems, especially with unstructured data (images, text).
<b>Perceptron</b>	The basic unit of a Neural Network. It takes weighted inputs, sums them up, and passes them through an <b>Activation Function</b> .
<b>Activation Function</b>	Introduces <b>non-linearity</b> (e.g., ReLU, Sigmoid) so the network can learn complex, non-straight line relationships.
<b>Backpropagation</b>	The main learning algorithm. It calculates the <b>error</b> and sends it backward to <b>adjust the weights</b> of the network layers.
<b>K-Means Clustering</b>	Unsupervised algorithm that groups data points into <b>K number of clusters</b> based on how close they are to the cluster center.
<b>Dimensionality Reduction (PCA)</b>	Reduces the number of features/columns by finding the most important ones, helping to <b>speed up training</b> .

## Diploma wallah

Topic	Quick Note
MLOps	The set of practices for <b>deploying, monitoring, and managing</b> ML models in a real-world production environment.

### 8. NLP & Deployment (Week 11, 12)

Topic	Quick Note
NLP	Giving computers the ability to <b>understand and process human language</b> (text and speech).
Tokenization	Breaking down text into small meaningful units (words or phrases).
Stemming	Chopping off the end of words to get a rough root word (e.g., <i>running</i> > <i>run</i> ).
Lemmatization	Getting the proper dictionary base word (lemma) using vocabulary rules (e.g., <i>better</i> > <i>good</i> ).
TF-IDF	A score that tells you how <b>important a word</b> is in a specific document compared to all other documents.
Sentiment Analysis	Determining the <b>mood or opinion</b> (positive, negative, neutral) of a piece of text.
Docker	A tool to package an application and all its dependencies into a single, reliable unit called a <b>container</b> .
Container	Ensures that the trained ML model runs <b>exactly the same way</b> on a developer's computer, a tester's computer, or the final server.

DIPLOMA WALLAH ( SWANGAM ❤ )