

Assignment Code: DA-AG-015

Boosting Techniques | **Assignment**

Instructions: Carefully read each question. Use Google Docs, Microsoft Word, or a similar tool to create a document where you type out each question along with its answer. Save the document as a PDF, and then upload it to the LMS. Please do not zip or archive the files before uploading them. Each question carries 20 marks.

Total Marks: 200

Question 1: What is Boosting in Machine Learning? Explain how it improves weak learners.

Answer:

Question 2: What is the difference between AdaBoost and Gradient Boosting in terms of how models are trained?

Answer:

Question 3: How does regularization help in XGBoost?

Answer:

Question 4: Why is CatBoost considered efficient for handling categorical data?

Answer:

Question 5: What are some real-world applications where boosting techniques are preferred over bagging methods?

Answer:

Datasets:

- Use `sklearn.datasets.load_breast_cancer()` for classification tasks.
- Use `sklearn.datasets.fetch_california_housing()` for regression tasks.

Question 6: Write a Python program to:

- Train an AdaBoost Classifier on the Breast Cancer dataset
- Print the model accuracy

(Include your Python code and output in the code box below.)

Answer:

Question 7: Write a Python program to:

- Train a Gradient Boosting Regressor on the California Housing dataset
- Evaluate performance using R-squared score

(Include your Python code and output in the code box below.)

Answer:

Question 8: Write a Python program to:

- Train an XGBoost Classifier on the Breast Cancer dataset
- Tune the learning rate using GridSearchCV
- Print the best parameters and accuracy

(Include your Python code and output in the code box below.)

Answer:

Question 9: Write a Python program to:

- Train a CatBoost Classifier
- Plot the confusion matrix using **seaborn**

(Include your Python code and output in the code box below.)

Answer:

Question 10: You're working for a FinTech company trying to predict loan default using customer demographics and transaction behavior.

The dataset is imbalanced, contains missing values, and has both numeric and categorical features.

Describe your step-by-step data science pipeline using boosting techniques:

- Data preprocessing & handling missing/categorical values
- Choice between AdaBoost, XGBoost, or CatBoost
- Hyperparameter tuning strategy
- Evaluation metrics you'd choose and why
- How the business would benefit from your model

(Include your Python code and output in the code box below.)

Answer: