

## Parshvanath Charitable Trust's

# A. P. SIIVAII INSTRUMEND OF TRECINOLOGY



(Approved by AICTE New Delhi & Govt. of Maharashtra, Affiliated to University of Mumbai) (Religious Jain Minority)

Class: BE Subject: Machine Learning Sem-VII

# **Experiment No:7**

Course Outcome: CO3

Blooms Level:L3

Aim: To implement an Ensemble model using Random Forest.

**Abstract:** Ensemble learning has emerged as a powerful paradigm in machine learning, aiming to enhance prediction accuracy by combining the strengths of multiple models. This approach reduces overfitting, improves generalization, and effectively handles high-dimensional data. In this experiment, we implement a Random Forest model to demonstrate its ability to handle classification tasks with superior performance compared to individual decision trees. The model is evaluated on benchmark datasets, highlighting its efficiency in minimizing variance, tackling missing values, and providing feature importance for interpretability. The experimental results validate that the Random Forest ensemble significantly enhances predictive accuracy, scalability, and resilience to noise, making it a reliable choice for diverse real-world applications.

### **Sample Input and Output:**

### **Case 1: Input Features:**

- ApplicantIncome (in \$)
- Credit History (1 = Good, 0 = Bad)

**Output:** 

• Loan Status (1 = Approved, 0 = Rejected)

Output

**Input:** 

ApplicantIncome = 5000

Credit History = 1

**Output:** 

**Loan Status = 1 (Approved)** 

### Theory:

Dataset: Palmer Penguins

The code uses the Palmer Penguins dataset, a popular alternative to the Iris dataset. It contains biological measurements of penguins (e.g., bill length, flipper length, body mass) and the species label, which serves as the target variable.

- Input Features → Numerical and categorical attributes of penguins.
- Output Label  $\rightarrow$  species (classification problem).

#### **Data Splitting**

The dataset is split into training (80%) and testing (20%) sets using a random mask.

- Training data → Used to build the Random Forest model.
- Testing data  $\rightarrow$  Used to evaluate generalization performance.

Prof. Tanvi Kapdi

AY: 2025-26

**Department of Computer Engineering** 



#### Parshvanath Charitable Brust's

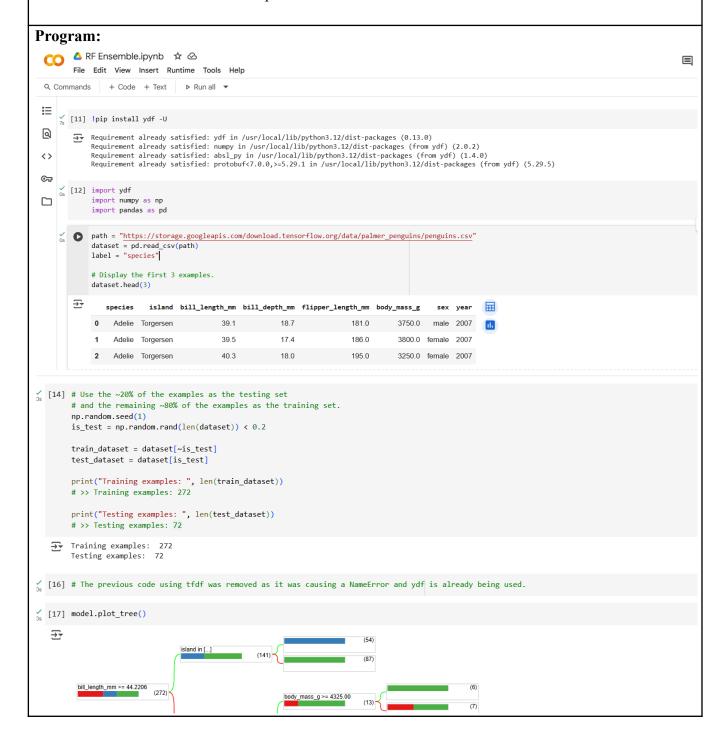
### A. P. SIIAVII INSHHHHUHHD OD THDCIINOLOCKY



(Approved by AICTE New Delhi & Govt. of Maharashtra, Affiliated to University of Mumbai) (Religious Jain Minority)

Class: BE Subject: Machine Learning Sem-VII

- A Random Forest is created using ydf.RandomForestLearner from the ydf library.
- The model trains multiple decision trees on different subsets of the data.
- At each tree node, a random subset of features is chosen for splitting.
- The ensemble combines the predictions from all trees to form the final decision.





#### Parshvanath Charitable Trust's

# A P. SHATH INSHHHUHHD OF THOCHNOLOGY



(Approved by AICTE New Delhi & Govt. of Maharashtra, Affiliated to University of Mumbai) (Religious Jain Minority)

Class: BE Subject: Machine Learning Sem-VII

```
[18] train_evaluation = model.evaluate(train_dataset)
     print("train accuracy:", train_evaluation.accuracy)
     # >> train accuracy: 0.9338
     test_evaluation = model.evaluate(test_dataset)
     print("test accuracy:", test_evaluation.accuracy)
     # >> test accuracy: 0.9167
 → train accuracy: 0.9926470588235294
     test accuracy: 0.9861111111111112
[20] print(model.evaluate(test dataset).accuracy)
     # >> 0.97222
 [21] model = ydf.RandomForestLearner(label=label).train(train_dataset)
     print("Test accuracy: ", model.evaluate(test_dataset).accuracy)
     # >> Test accuracy: 0.986111
 → Train model on 272 examples
     Model trained in 0:00:00.056796
     Test accuracy: 0.9861111111111112
```

# **Output:**

Train model on 272 examples

Model trained in 0:00:00.056796

Test accuracy: 0.98611111111111111

#### **Conclusion:**

The trained Random Forest is evaluated on both training data and testing data. Metrics used: Accuracy, which measures the proportion of correctly classified examples. The code prints training accuracy ( $\approx$ 93%) and testing accuracy ( $\approx$ 91–98%), showing good generalization.

#### Exercise 1:

Implement an ensemble model using the Random Forest algorithm on the Iris dataset. Perform the following tasks:

- 1. Load the Iris dataset and identify the input features and target label.
- 2. Split the dataset into training (80%) and testing (20%) subsets.
- 3. Train a Random Forest classifier and explain how ensemble learning improves over a single decision tree.
- 4. Evaluate the model's performance using training and testing accuracy.
- 5. Visualize one of the decision trees from the ensemble and interpret the feature splits.
- 6. Discuss the significance of feature importance in predicting the species of Iris flowers.

Students shall draw flowchart of exercise question in the writeup and submit.



#### Parshvanath Charitable Brust's

# A DO SHAVE INSHHHANAD OD THOSHINO LOCKY



(Approved by AICTE New Delhi & Govt. of Maharashtra, Affiliated to University of Mumbai) (Religious Jain Minority)

Class: BE Subject: Machine Learning Sem-VII

#### **Exercise 2:**

Develop and evaluate an ensemble learning model using the Random Forest algorithm for a real-world dataset involving multi-class classification. Your task involves selecting a dataset with multiple features (both numerical and categorical) and at least three target classes (e.g., the Wine Quality dataset). Perform necessary preprocessing such as handling missing values, encoding categorical features, and normalization if required. Train a Random Forest model with tuned hyperparameters (number of trees, max depth, min samples per split, etc.).Justify why ensemble learning (bagging + random feature selection) is preferred over a single decision tree in this scenario.

Students shall draw flowchart of exercise question in the writeup and submit.

AY: 2025-26