# Title: Predicting Employee Performance Using Machine Learning

## Introduction

This project explores the use of supervised machine learning to predict employee performance scores based on workplace data. The goal is to provide HR departments with tools to identify top-performing employees, understand factors influencing performance, and assess potential for future growth.

## Problem Statement

The objective of this study is to create a predictive model that uses employee demographic details, workplace conditions, and engagement metrics to forecast performance scores. This model serves purposes such as:

- Recognizing high-performing individuals.
- Understanding key performance drivers.
- Supporting HR in talent development and strategic planning.

## Dataset Overview

- Dataset Size: 100,000 rows with 20 attributes.

**Features:**
- Demographic Information: Age, Gender, Education Level.
- Workplace Metrics: Department, Job Title, Years at Company, Monthly Salary, Team Size.
- Engagement Metrics: Training Hours, Employee Satisfaction Score, among others.
- Target Variable: Performance Score, measured on a scale from 1 to 5.

## Methodology

## Data Pre-processing

1. Handling Categorical Data:
   - Applied One-Hot Encoding to categorical variables such as Department and Gender.
   - Used Ordinal Encoding for Education Level to reflect hierarchy.
2. Scaling Numerical Features:
   - Normalized numerical variables using MinMaxScaler to ensure consistency.
3. Feature Selection:
   - Leveraged feature importance metrics from Random Forest to select influential features.
4. Data Cleaning:
   - Ensured the dataset was free of missing or invalid entries.
   - Verified distributions.

## Exploratory Data Analysis (EDA)

- Conducted histogram and box plot visualizations to understand the distribution of numerical variables.
- Analyzed correlations between features and the target variable to identify significant relationships.

## Model Selection and Training

1. Evaluated Models:
    - Random Forest Classifier
    - Gradient Boosting Classifier
    - XGBoost Classifier
2. Optimal Model:
    - Random Forest Classifier performed best with the following configuration:
    - 500 estimators and "sqrt" max features.
    - Achieved an accuracy of 99.95% on the validation set.
3. Cross-Validation:
    - Used GridSearchCV to fine-tune hyperparameters, employing a 2-fold cross-validation approach.

## Results

## Performance Metrics

- Confusion Matrix: Demonstrated excellent precision, recall, and F1-scores across all performance categories.
- Overall Scores:
- Macro Average F1-Score: 1.00
- Weighted Average F1-Score: 1.00

**The model's remarkable accuracy was largely due to the inclusion of the Monthly Salary feature. Without this feature, the accuracy remained significantly lower, highlighting its critical role in prediction**.

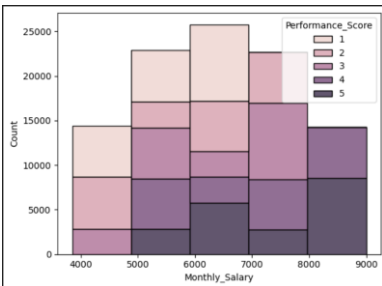## Feature Importance

# Top predictors of performance:

1. Monthly Salary: Contributed 46.36% to predictions.
2. Training Hours: Accounted for 4.73% of importance.
3. Projects Handled: 4.42% contribution.
4. Work Hours Per Week: 4.14% contribution.
5. Team Size and Experience: Additional significant factors.

| | Feature | Importance |
|---|---|---|
| 27 | remainder__Monthly_Salary | 0.463574 |
| 25 | remainder__Training_Hours | 0.047270 |
| 20 | remainder__Projects_Handled | 0.044188 |
| 17 | remainder__Age | 0.043194 |
| 19 | remainder__Work_Hours_Per_Week | 0.041438 |
| 21 | remainder__Overtime_Hours | 0.041412 |
| 24 | remainder__Team_Size | 0.037404 |
| 22 | remainder__Sick_Days | 0.035233 |
| 18 | remainder__Years_At_Company | 0.031664 |
| 23 | remainder__Remote_Work_Frequency | 0.023072 |

## Insights and Achievements

1. Key Insights:
   - Employees with higher salaries and more training hours tend to perform better.



   - Factors such as workload (e.g., hours worked and projects handled) significantly influence outcomes.
   - Employee satisfaction and promotion history are crucial for performance enhancement.
2. Notable Achievements:
   - Built a highly accurate predictive model for HR applications.
   - Provided actionable insights that support strategic decision-making.

## Challenges

- Data Bias: Historical biases in employee records may skew predictions and need to be addressed.
- Interpretability: Balancing model complexity with the need for transparent feature importance posed challenges.

## Future Work

1. Validate the model using real-time employee data to ensure practical applicability.

2. Expand the dataset with additional variables such as peer reviews and certifications.
3. Develop interactive dashboards to make predictions and insights more accessible to HR professionals.

## Conclusion

This project successfully implemented a machine learning solution to predict employee performance with remarkable accuracy. By offering data-driven insights, this model empowers HR teams to make informed decisions, enhancing productivity and workforce development.