# Employee Salary Prediction Using Machine Learning

## Introduction

In the modern corporate and industrial environment, accurate estimation of employee salary plays an important role in effective human resource management. Organizations often need to make informed decisions related to recruitment, promotions, and salary planning. Traditionally, salary estimation is done manually based on experience and company policies, which can sometimes lead to inconsistencies and bias.

With the advancement of data science and Machine Learning, predictive models can be used to analyze historical data and estimate salaries more accurately. Machine Learning allows systems to learn patterns from existing data and make predictions for new and unseen inputs. One such commonly used technique for prediction is **Linear Regression**, which is simple, efficient, and easy to understand.

This project, titled **"Employee Salary Prediction Using Machine Learning"**, focuses on predicting the salary of an employee based on their years of experience. The project uses a supervised Machine Learning approach, where the model is trained using previously available data containing experience and corresponding salary values. Once trained, the model can predict the salary for a new employee when their experience is provided.

Python programming language is used for implementing the project due to its simplicity and availability of powerful Machine Learning libraries such as NumPy, Pandas, Matplotlib, and Scikit-learn. These libraries help in data handling, model training, and visualization of results.

The main aim of this project is not only to predict salary but also to provide a clear understanding of how Machine Learning models work in real-world applications. This project serves as a beginner-friendly introduction to Machine Learning concepts and demonstrates how data-driven decision-making can be applied in human resource management.

## Objective of the Project

The primary objective of this project is to develop a Machine Learning-based system that can predict employee salary based on years of experience. The project aims to demonstrate how historical data can be used to identify patterns and make accurate predictions using supervised learning techniques.

The specific objectives of the project are as follows:

1. **To study the concept of Machine Learning and supervised learning algorithms**
   This project aims to provide a basic understanding of Machine Learning, with a focus on supervised learning and its applications in real-world problems.

2. **To analyze the relationship between employee experience and salary**
   The project seeks to understand how years of experience influence salary growth and how this relationship can be represented mathematically using Linear Regression.

3. **To design and implement a salary prediction model using Linear Regression**
   The objective is to apply the Linear Regression algorithm to train a model that can learn from existing data and predict salary values for new inputs.

4. **To implement the model using Python programming language**
   Python is used due to its simplicity and wide range of Machine Learning libraries, enabling efficient data processing, model training, and prediction.

5. **To visualize the results for better understanding**
   Graphical representation of the dataset and regression line helps in interpreting the prediction results and understanding the trend between experience and salary.

6. **To evaluate the performance of the prediction model**
   The project aims to check whether the predicted values follow the actual salary trend and verify the effectiveness of Linear Regression for this problem.

7. **To provide a foundation for advanced Machine Learning projects**
   This project serves as a base model that can be extended in the future by adding more features, larger datasets, and advanced algorithms.

## Problem Statement

To design and implement a Machine Learning model that can predict the salary of an employee based on their years of experience using historical data.

## Tools and Technologies Used

The following tools and technologies were used in this project: - **Programming Language:** Python - **Libraries Used:** - NumPy – for numerical operations - Pandas – for data handling and analysis - Matplotlib – for data visualization - Scikit-learn – for Machine Learning model - **IDE:** Jupyter Notebook / VS Code

# Dataset Description

The dataset used in this project is a manually created dataset designed to clearly demonstrate the concept of salary prediction using Machine Learning. Although small in size, the dataset is sufficient for understanding the working of Linear Regression.

The dataset contains two attributes: - **Experience (Years):** This represents the total number of years an employee has worked. - **Salary (INR):** This represents the annual salary corresponding to the employee's experience.

The dataset assumes that salary increases linearly with experience, which is a common assumption in many entry-level and mid-level job roles.

## Sample Dataset

| Experience (Years) | Salary (INR) |
|---|---|
| 1 | 25000 |
| 2 | 30000 |
| 3 | 35000 |
| 4 | 40000 |
| 5 | 45000 |
| 6 | 50000 |
| 7 | 55000 |
| 8 | 60000 |
| 9 | 65000 |
| 10 | 70000 |

# Methodology

The methodology describes the overall workflow of the project, from data creation to final prediction. It ensures that the Machine Learning model is trained correctly and gives meaningful results.

The methodology describes the step-by-step process followed to design and implement the Employee Salary Prediction system. Each step was carefully executed to ensure correct training and prediction by the Machine Learning model.

1. **Data Collection:** A small dataset was manually created containing employee experience and corresponding salary values. This helps in understanding the concept clearly without data complexity.

2. **Data Preparation:** The dataset was stored in a Pandas DataFrame. The independent variable (experience) and dependent variable (salary) were separated.

3.  **Data Splitting:** The dataset was divided into training data (80%) and testing data (20%). Training data is used to teach the model, while testing data checks its performance.

4.  **Model Selection:** Linear Regression was selected as it is simple, efficient, and suitable for predicting continuous values like salary.

5.  **Model Training:** The Linear Regression model was trained using the training dataset. During this phase, the model learns the relationship between experience and salary.

6.  **Prediction:** After training, the model predicts salary values for test data and for new experience values provided by the user.

7.  **Visualization:** A graph was generated to visualize the actual data points and the regression line for better understanding of the prediction trend.

## Algorithm Used – Linear Regression

Linear Regression is a supervised Machine Learning algorithm used for predicting continuous numerical values. It establishes a linear relationship between an independent variable (experience) and a dependent variable (salary).

The algorithm works by fitting the best possible straight line through the data points. This line is known as the regression line and is represented by the equation:

**Salary = m × Experience + c**

Where: - **m (Slope):** Represents the rate at which salary increases with experience - **c (Intercept):** Represents the base salary when experience is zero

During training, the algorithm calculates optimal values of m and c that minimize prediction error using the least squares method.

Linear Regression is widely used because: - It is easy to implement - It gives fast results - It is effective for simple prediction problems

## Implementation

The implementation phase focuses on converting the theoretical concept into a working program using Python.

### a. Data Handling

The dataset is created using Python lists and converted into a Pandas DataFrame. Pandas provides an easy way to handle and manipulate structured data.

### b. Feature Selection

In this project: - Experience is treated as the independent variable (input) - Salary is treated as the dependent variable (output)

Only one feature is used to keep the model simple and understandable.

### c. Train-Test Split

The dataset is divided into training and testing sets using an 80:20 ratio. This helps in checking how well the model performs on unseen data.

### d. Model Training

The Linear Regression model is trained using the training data. During training, the model learns the best-fit line by minimizing the prediction error.

### e. Prediction and Visualization

After training, the model predicts salary values for test data and new input values. A graph is plotted to visualize the relationship between experience and salary.

## 9. Source Code

Below is the Python code used for implementing the Employee Salary Prediction project:

```python
import pandas as pd
import numpy as np
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt

exp = [1,2,3,4,5,6,7,8,9,10]
sal = [25000,30000,35000,40000,45000,50000,55000,60000,65000,70000]

data = pd.DataFrame()
data["experience"] = exp
data["salary"] = sal

x = data[["experience"]]
y = data["salary"]

x_train, x_test, y_train, y_test = train_test_split(
    x, y, test_size=0.2
)

lr = LinearRegression()

lr.fit(x_train, y_train)

pred = lr.predict(x_test)

yr = 5.5
res = lr.predict([[yr]])

print("salary for", yr, "years experience is approx:", int(res[0]))

plt.scatter(x, y)
plt.plot(x, lr.predict(x))
plt.xlabel("Experience")
plt.ylabel("Salary")
plt.show()
```
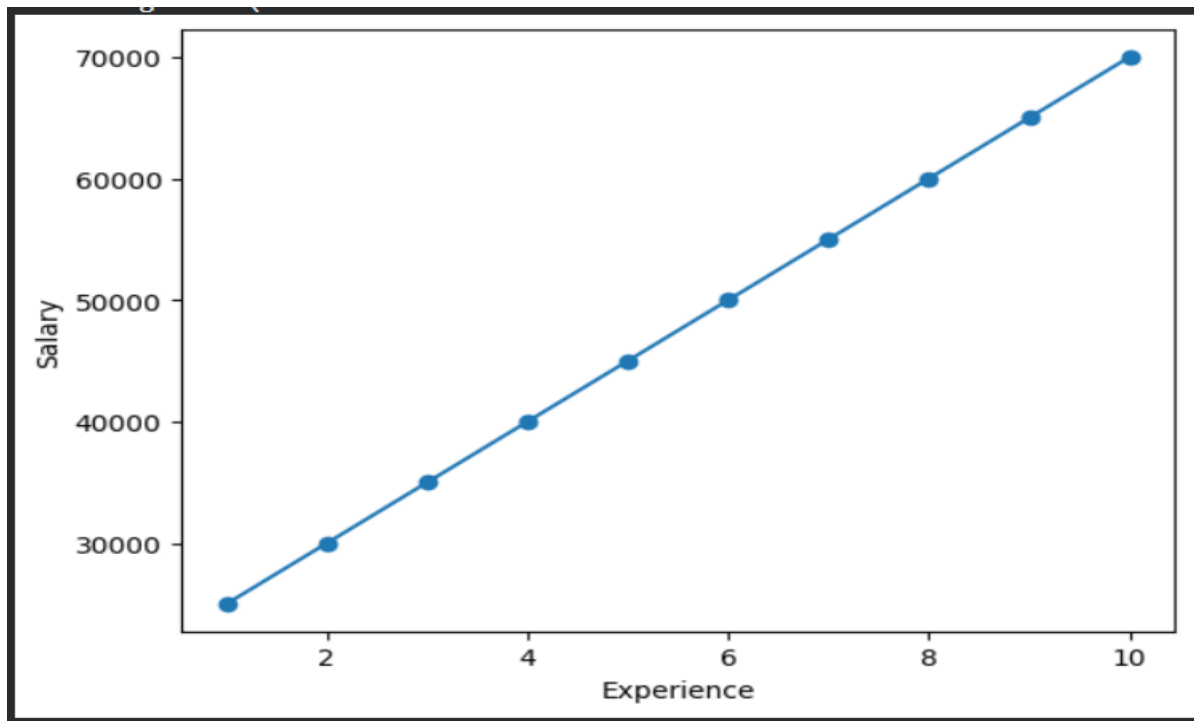
## 10. Output

After executing the Python program, both textual and graphical outputs are obtained, which verify the correct working of the model.

The program displays the predicted salary for a given number of years of experience entered in the code.

```
salary for 5.5 years experience is approx: 47500
```

This means that an employee with approximately 5.5 years of experience is expected to earn around ₹47,500 according to the trained model.

The graphical output consists of: - Blue dots representing the actual salary data - A straight line representing the predicted salary trend

The graph clearly shows a positive linear relationship between experience and salary, validating the effectiveness of the Linear Regression model.

## Result and Output

The results obtained from the project confirm that the Linear Regression model works effectively for salary prediction based on experience.

The predicted values closely follow the trend of the actual data, indicating good model performance for this dataset.

The model successfully predicts salary based on years of experience. For example:

- **Input:** 5.5 years of experience
- **Predicted Salary:** Approximately ₹47,500

A graph is also generated showing: - Actual data points (scatter plot) - Regression line representing predicted salary trend

## Advantages of the Project

- Simple and easy to understand for beginners
- Requires less computational power
- Helps in understanding the fundamentals of Machine Learning
- Can be used as a base for advanced prediction systems
- Provides visual representation for better clarity

## Limitations

- The dataset used is small and manually created
- Only one factor (experience) is considered
- Real-world salary prediction depends on multiple parameters
- Model accuracy may reduce with complex data

## Future Scope

The scope of this project can be extended in the following ways: - Use real-world datasets from companies or job portals - Include multiple features such as education, designation, location, and skills - Apply advanced algorithms like Random Forest or Support Vector Machines - Deploy the model using a web interface or mobile application - Integrate with HR management systems

## Conclusion

This project successfully demonstrates the application of Machine Learning for predicting employee salaries using Linear Regression. The system efficiently learns the relationship between experience and salary and provides accurate predictions for new data.

Through this project, practical knowledge of Python programming, data handling, Machine Learning algorithms, and data visualization was gained. The project serves as a strong foundation for learning advanced Machine Learning concepts and real-world applications.

## Applications of the Project

This project has several real-world applications, such as: - Assisting HR departments in estimating employee salaries - Helping freshers understand expected salary growth - Supporting decision-making during recruitment - Serving as a learning model for beginners in Machine Learning

## Learning Outcomes

Through this project, the following learning outcomes were achieved: - Understanding of supervised Machine Learning concepts - Hands-on experience with Python libraries - Knowledge of data preparation and model training - Ability to visualize and interpret results.