# Software Salary Prediction Using Machine Learning

Team Id: LTVIP2025TMID41628

By:

Team Leader: Shaik Manisha

Team member: Goli Varun Kumar

Team member: Mandala Pujitha

Team member: Bathula Kalyan Babu

## Abstract

This project focuses on predicting salaries based on various features such as company name, company rating, job title, salaries reported ,employment status and location using machine learning algorithms. The model is trained on a dataset of salary records and is integrated with a web-based interface developed using Flask. The aim is to provide an intelligent system that can estimate salaries for different roles in the industry, assisting both recruiters and job seekers.

## Table of Contents

## 1. Introduction

In the modern job market, salary estimation plays an essential role for employees and employers alike. With the advancement of Artificial Intelligence (AI) and Machine Learning (ML), data-driven predictions can now be made with high accuracy. This project implements a salary prediction system that utilizes ML algorithms trained on real-world data. The web interface allows users to input their details and receive an estimated salary instantly.

## 2. Problem Statement

Determining fair and accurate salaries is often challenging due to the diversity of job roles, industries, and experience levels. Manual methods or generic online tools may not consider all influencing factors. This project addresses this issue by creating a machine learning-based model that predicts salaries using multiple relevant features, providing more personalized and data-driven results.

## 3. Objectives

- To collect and preprocess a dataset containing various job-related features.

- To train and evaluate ML models for accurate salary prediction.

- To develop a Flask web application for user interaction.

- To visualize the dataset and model performance for better understanding.

## 4. Existing System vs Proposed System

Existing System:

Existing salary estimation methods rely mainly on static data or manual entry systems, which are often inaccurate. These systems do not utilize machine learning and thus lack predictive capabilities.
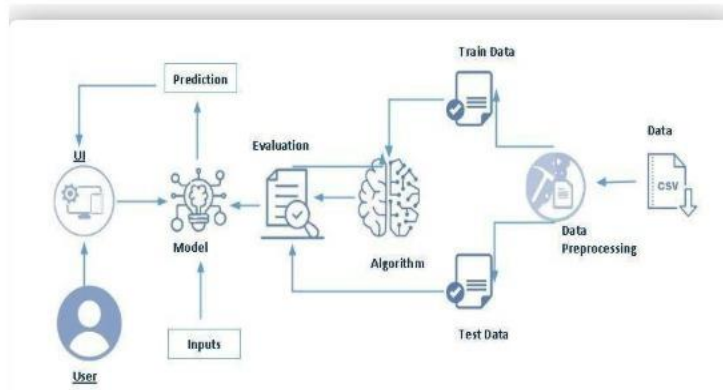
Proposed System:

The proposed system leverages machine learning techniques to predict salaries dynamically. It considers multiple features simultaneously and improves accuracy through training and evaluation.

## 5. System Architecture

The system consists of three primary components:
1. Data Collection and Preprocessing
2. Model Training and Evaluation
3. Flask Web Interface for Deployment

**Technical Architecture:**



## 6. Modules Description

The project is divided into the following modules:

1. Data Preprocessing Module

2. Model Training Module

3. Web Application Module

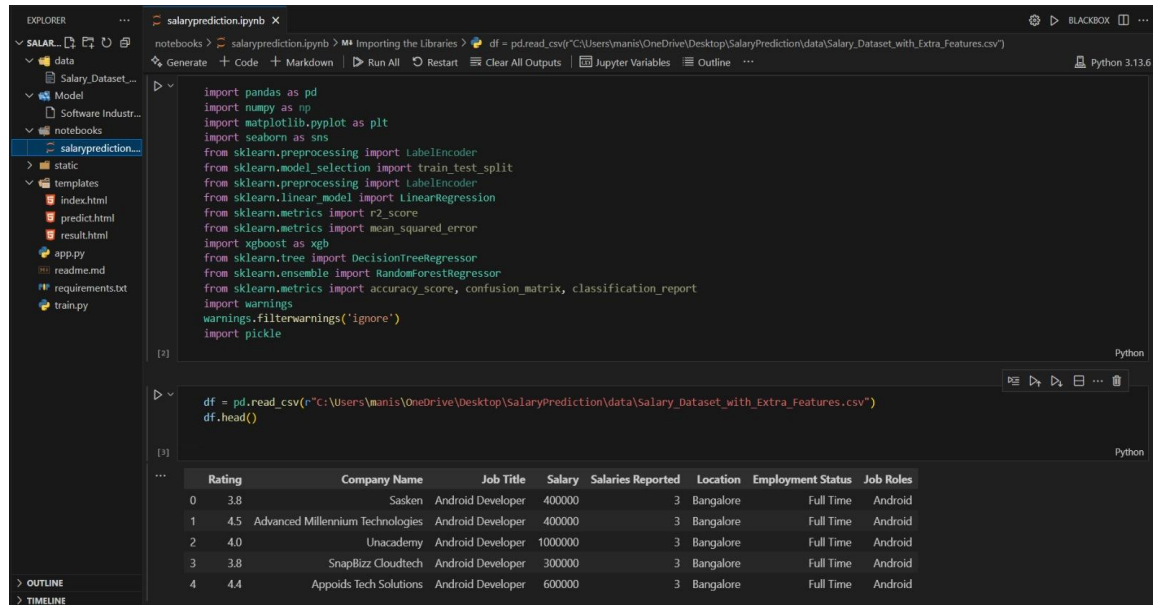4. Result and Visualization Module

## 7. Dataset Description

The dataset used in this project 'Salary_Dataset_with_Extra_Features.csv' contains features like company name, rating, job title, salaries reported, location, employment status and job roles These attributes significantly influence the salary outcome. The data was cleaned, encoded, and standardized before model training.

```
Rating,Company Name,Job Title,Salary,Salaries Reported,Location,Employment Status,Job Roles
3.8,Sasken,Android Developer,400000,3,Bangalore,Full Time,Android
4.5,Advanced Millennium Technologies,Android Developer,400000,3,Bangalore,Full Time,Android
4.0,Unacademy,Android Developer,1000000,3,Bangalore,Full Time,Android
3.8,SnapBizz Cloudtech,Android Developer,300000,3,Bangalore,Full Time,Android
4.4,Appoids Tech Solutions,Android Developer,600000,3,Bangalore,Full Time,Android
4.2,Freelancer,Android Developer,100000,3,Bangalore,Full Time,Android
3.7,SQUARE N CUBE,Android Developer,192000,3,Bangalore,Full Time,Android
3.1,Samsung R&D Institute India - Bangalore,Android Developer,400000,3,Bangalore,Full Time,Android
3.7,DXMinds Technologies,Android Developer,300000,3,Bangalore,Full Time,Android
3.6,Endeavour Software Technologies,Android Developer,600000,3,Bangalore,Full Time,Android
3.6,Craft Silicon,Android Developer,300000,3,Bangalore,Full Time,Android
3.9,Baronford & Associates,Android Developer,240000,2,Bangalore,Full Time,Android
3.7,Wibmo,Android Developer,900000,2,Bangalore,Full Time,Android
4.8,Retail Pulse,Android Developer - Intern,24000,2,Bangalore,Intern,Android
3.9,Bookmyshow,Android Developer,600000,2,Bangalore,Full Time,Android
3.9,Knowledge Flex,Android Developer,228000,2,Bangalore,Full Time,Android
3.6,Novopay Solutions,Android Developer,600000,2,Bangalore,Full Time,Android
3.7,WealthEngine,Android Developer,360000,2,Bangalore,Full Time,Android
4.0,J.P. Morgan,Android Developer,1000000,2,Bangalore,Full Time,Android
3.6,Acviss,Android Developer,500000,2,Bangalore,Full Time,Android
4.1,Fresher,Android Developer,408000,2,Bangalore,Full Time,Android
4.2,MedOnGo,Android Developer,300000,2,Bangalore,Full Time,Android
4.0,Nuclei,Android Developer,800000,2,Bangalore,Full Time,Android
4.4,eSecForte Technologies,Android Developer,228000,2,Bangalore,Full Time,Android
4.3,Moveinsync Technology Solutions,Android Developer,100000,2,Bangalore,Full Time,Android
3.6,Tech Mahindra,Android Developer,500000,2,Bangalore,Full Time,Android
4.0,ThiDiff Technologies,Android Developer,200000,2,Bangalore,Full Time,Android
4.9,Retranz Infolabs,Android Developer,500000,2,Bangalore,Full Time,Android
4.0,FicusLot,Android Developer,228000,2,Bangalore,Full Time,Android
3.7,KrazyBee,Android Developer,708000,2,Bangalore,Full Time,Android
5.0,powerplay app,Android Developer - Intern,396000,2,Bangalore,Intern,Android
4.2,Dcoder,Android Developer,700000,2,Bangalore,Full Time,Android
4.6,Masai School,Android Developer,500000,2,Bangalore,Full Time,Android
3.9,Integra Micro Software Services (P),Android Developer,500000,2,Bangalore,Full Time,Android
3.3,DocsApp,Android Developer,900000,2,Bangalore,Full Time,Android
5.0,Vispara Technosoft,Android Developer,200000,2,Bangalore,Full Time,Android
4.3,Vmoksha Technologies,Android Developer,228000,2,Bangalore,Full Time,Android
2.9,Artoo,Android Developer,400000,2,Bangalore,Full Time,Android
```
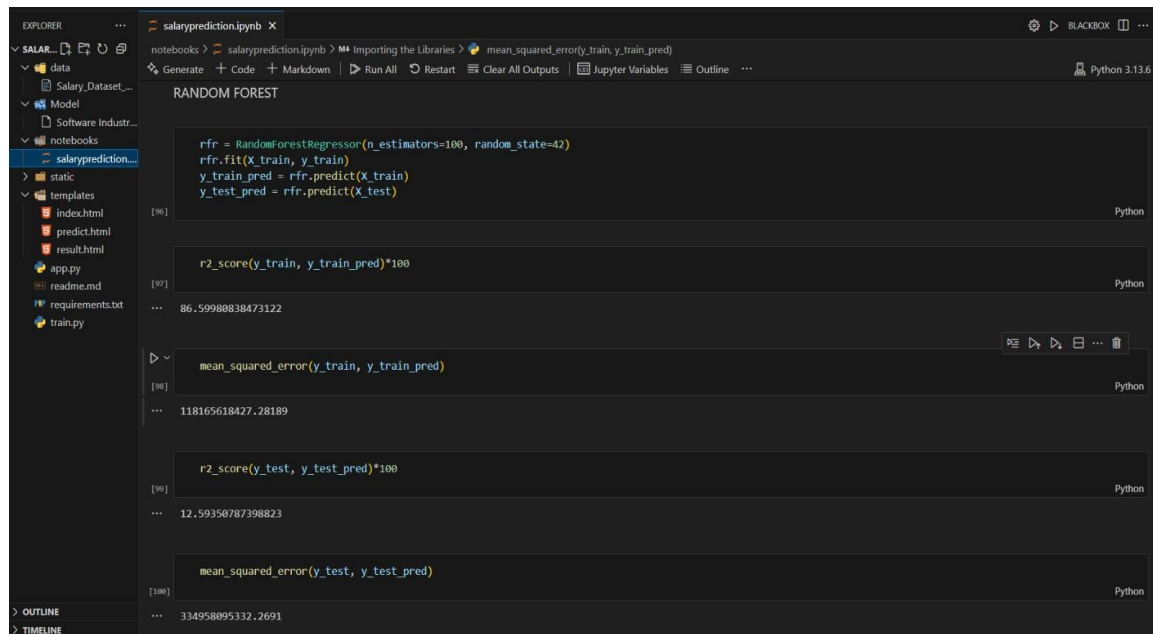
## 8. Model Training and Evaluation

The training process involves loading the dataset, preprocessing, splitting into training and testing sets, and applying algorithms such as Random Forest Regressor for prediction. The model's performance was evaluated based on metrics like Mean Squared Error (MSE) and $R^2$ Score.

## 9. Implementation

The implementation was carried out using Python. Flask framework was used for backend web development, and HTML/CSS for frontend design. The trained ML model was integrated into Flask using Pickle serialization. Users can enter their job-related details on the web page to get predicted salary values.

## 10. Results

The system successfully predicts salary values based on user input. The results show that the Random Forest model provides accurate predictions compared to other algorithms tested. The web interface allows seamless interaction with the model.

```
TESTING THE MODEL (RANDOM FOREST)

        rfr.predict([[0.7,6422,461,1,0,1,3]])
[111]                                                                      Python
...    array([5812240.])


        rfr.predict([[2.5,5116,709,1,9,1,3]])
[112]                                                                      Python
...    array([4269880.])


        rfr.predict([[1.2,4718,1071,1,0,1,5]])
[113]                                                                      Python
...    array([981640.])


  ▷     rfr.predict([[1.2,3412,8942,1,7,1,2]])
[114]                                                                      Python
...    array([982040.])
```

```
LINEAR REGRESSION

        reg = LinearRegression()
        reg.fit(X_train, y_train)
        y_train_pred = reg.predict(X_train)
        y_test_pred = reg.predict(X_test)
[106]                                                                      Python


        r2_score(y_train, y_train_pred)*100
[107]                                                                      Python
...    1.8362672309302885


        mean_squared_error(y_train, y_train_pred)
[108]                                                                      Python
...    865627785245.2856


  ▷     r2_score(y_test, y_test_pred)*100
[109]                                                                      Python
...    4.118277837531159
```

## 11. Conclusion and Future Scope

This project demonstrates how machine learning can be applied to real-world problems like salary prediction. The model achieves good accuracy and helps users estimate expected salaries effectively. Future enhancements may include expanding the dataset, adding more features and integrating deep learning methods for improved predictions.

### 9.1 Demonstration

This section provides visual evidence of the developed project, including source code snippets and the web application interface.

Figure 1: Screenshot of train.py (Model Training Code)

```python
1   import pandas as pd
2   from sklearn.model_selection import train_test_split
3   from sklearn.preprocessing import LabelEncoder, StandardScaler
4   from sklearn.ensemble import RandomForestRegressor
5   import pickle
6   import os
7
8   # Load dataset
9   df = pd.read_csv('data/Salary_Dataset_with_Extra_Features.csv')
10
11  # Fill missing values if any
12  df.fillna(method='ffill', inplace=True)
13
14  # Select features and target (now including Salary Reported)
15  features = ['Rating', 'Company Name', 'Job Title', 'Location',
16              'Employment Status', 'Job Roles', 'Salaries Reported']
17  target = 'Salary'
18
19  X = df[features]
20  y = df[target]
21
22  # Encode categorical variables
23  label_encoders = {}
24  for col in ['Company Name', 'Job Title', 'Location', 'Employment Status', 'Job Roles']:
25      le = LabelEncoder()
26      X[col] = le.fit_transform(X[col])
27      label_encoders[col] = le
28
29  # Scale numeric features (Rating and Salary Reported)
30  scaler = StandardScaler()
31  X[['Rating', 'Salaries Reported']] = scaler.fit_transform(X[['Rating', 'Salaries Reported']])
32
33  # Split data into training and testing sets
34  X_train, X_test, y_train, y_test = train_test_split(
35      X, y, test_size=0.2, random_state=42
36  )
37
38  # Train RandomForest model
39  model = RandomForestRegressor(n_estimators=100, random_state=42)
40  model.fit(X_train, y_train)
41
42  # Ensure Model directory exists
43  os.makedirs('Model', exist_ok=True)
44
45  # Save model, encoders, and scaler
46  with open('Model/Software Industry Salary Prediction.pkl', 'wb') as f:
47      pickle.dump({
48          'model': model,
49          'label_encoders': label_encoders,
50          'scaler': scaler
51      }, f)
52
53  print("Model trained and saved successfully!")
54
```
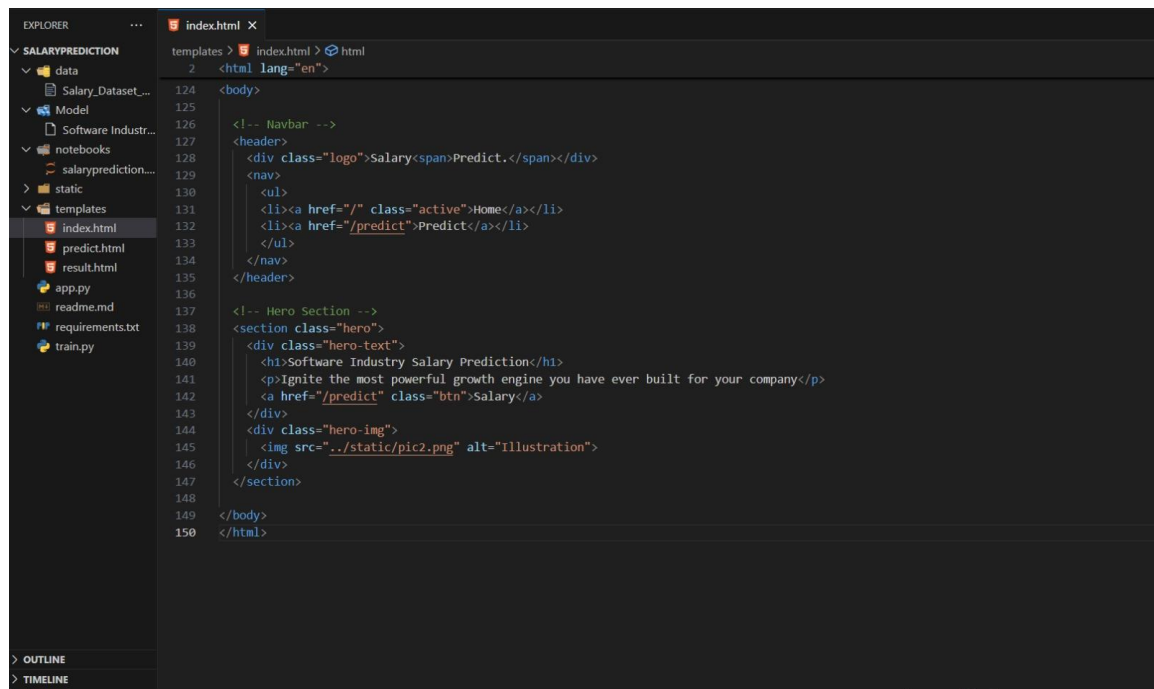
Figure 2: Screenshot of app.py (Flask backend code)



```python
from flask import Flask, render_template, request
import pickle
import numpy as np

app = Flask(__name__)

# Load the saved dictionary (model + label_encoders + scaler)
with open("Model/Software Industry Salary Prediction.pkl", "rb") as f:
    saved = pickle.load(f)

model = saved["model"]
label_encoders = saved["label_encoders"]
scaler = saved["scaler"]

def encode_label(le, value):
    if value in le.classes_:
        return le.transform([value])[0]
    else:
        # For unseen categories, return a default value (0)
        return 0

@app.route("/")
def index():
    return render_template("index.html")   # Landing page


@app.route("/predict", methods=["GET", "POST"])
def predict():
    if request.method == "POST":
        try:
            # Collect form data
            company_name = request.form["company_name"]
            job_title = request.form["job_title"]
            location = request.form["location"]
            job_roles = request.form["job_roles"]
            employment_status = request.form["employment_status"]
            salaries_reported = int(request.form["salaries_reported"])
            rating = float(request.form["rating"])
```



```python
    def predict():

            # Encode categorical variables using saved LabelEncoders
            company_name_encoded = encode_label(label_encoders['Company Name'], company_name)
            job_title_encoded = encode_label(label_encoders['Job Title'], job_title)
            location_encoded = encode_label(label_encoders['Location'], location)
            job_roles_encoded = encode_label(label_encoders['Job Roles'], job_roles)
            employment_status_encoded = encode_label(label_encoders['Employment Status'], employment_status)

            # Scale numeric features (Rating and Salaries Reported)
            scaled = scaler.transform([[rating, salaries_reported]])
            rating_scaled = scaled[0][0]
            salaries_scaled = scaled[0][1]

            # Prepare final input array
            input_features = np.array([[rating_scaled, company_name_encoded, job_title_encoded,
                                        location_encoded, employment_status_encoded, job_roles_encoded,
                                        salaries_scaled]])

            # Predict salary
            prediction = model.predict(input_features)[0]
            prediction = round(prediction, 2)

            return render_template("result.html", prediction=prediction)

        except Exception as e:
            return f"Error: {e}"

    # GET request → show the prediction form
    return render_template("predict.html")


if __name__ == "__main__":
    app.run(debug=True)
```

Figure 3: Screenshot of index.html.



Figure 4: Screenshot of predict.html.

Figure 5: Screenshot of result.html.



Figure 6: Web Application – Index Page.



Figure 7: Web Application – Predict Page.

Figure 8: Web Application – Result Page.



GitHub Repository: https://github.com/shaikmanisha38/SalaryPrediction.git

Demo Video:
https://drive.google.com/file/d/1XVYv6Aa8gRUFIhFP4Z3to7w_SqF9t1yq/view?usp=drive_link

## 12. References

1. Python Documentation - https://docs.python.org/3/

2. Scikit-learn Library - https://scikit-learn.org/

3. Flask Framework - https://flask.palletsprojects.com/

4. Dataset Source - Internal Project Dataset