

INTERNSHIP REPORT

A report submitted in partial fulfilment of the requirements for the Award of Degree of

BACHELOR OF TECHNOLOGY

in

INFORMATION TECHNOLOGY

by

BOMMISSETTY AKASH

Reg.No.: 21781A1208

Under supervision

of

Mr.P.Nandakumar,

Associate Professor,

Department of Information Technology



**SRI VENKATESWARA COLLEGE OF ENGINEERING & TECHNOLOGY
(AUTONOMOUS)**

R.V.S NAGAR, CHITTOOR — 517 127. (A.P)

(Approved by AICTE, New Delhi, Affiliated to JNTUA, Anantapuramu)

(Accredited by NBA, New Delhi & NAAC, Bengaluru)

(An ISO 9001:2000 Certified Institution)

2024-2025

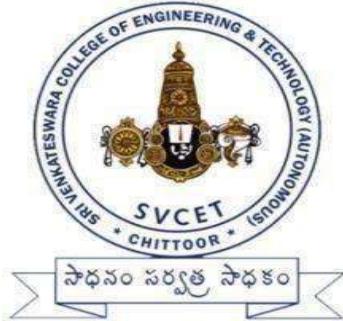
**SRI VENKATESWARA COLLEGE OF ENGINEERING& TECHNOLOGY
(AUTONOMOUS)**

R.V.S NAGAR, CHITTOOR — 517 127. (A.P)

(Approved by AICTE, New Delhi, Affiliated to JNTUA, Anantapuramu)

(Accredited by NBA, New Delhi & NAAC, Bengaluru)

(An ISO 9001:2000 Certified Institution)



CERTIFICATE

This is to certify that the "**Internship report**" submitted by BOMMISSETTY AKASH (Regd.No.:21781A1208) is work done by him/her and submitted during 2024- 2025 academic year, in partial fulfillment of the requirements for the award of the degree of BACHELOR OF TECHNOLOGY in INFORMATION TECHNOLOGY, at SmartInternz in Bangalore.

Internship Coordinator

Head of the Department

Review Conducted on:_____

INTERNAL EXAMINER

EXTERNAL EXAMINER

CERTIFICATE OF INTERNSHIP

11/11/24, 7:49 PM

Certificate - SmartInternz



ACKNOWLEDGEMENT

A grateful thanks to **Dr. R.Venkataswamy**, Chairman of Sri Venkateswara College of Engineering & Technology for providing education in their esteemed institution.

I wish to record my deep sense of gratitude and profound thanks to our beloved Vice Chairman, **Sri R.V.Srinivas** for his valuable support throughout the course.

I express our sincere thanks to **Dr.M.MOHAN BABU**, our beloved principal for his encouragement and suggestion during the course of study.

With the deep sense of gratefulness, I acknowledge **Mr.S.R. RAJ KUMAR**, Head of the Department, Information Technology, for giving us inspiring guidance in undertaking internship.

I express our sincere thanks to the internship coordinator, for his keen interest, stimulating guidance, constant encouragement with our work during all stages, to bring this report into fruition.

I wish to convey my gratitude and sincere thanks to all members for their support and cooperation rendered for successful submission of report.

Finally, I would like to express my sincere thanks to all teaching, non-teaching faculty members, our parents, friends and for all those who have supported us to complete the internship successfully.

BOMMISSETTY AKASH

(21781A1208)

INDEX

1. Introduction

1. Project overviews	01
2. Objectives	02

2. Project Initialization and planning phase

1. Define problem Statement	02
2. Project Proposal (Proposed Solution)	03
3. Initial Project Planning	04

3. Data Collection and Preprocessing Phase

1. Data Collection plan and raw data source identified	05
2. Data quality report	06
3. Data exploration and preprocessing	07

4. Model development phase

1. Feature selection report	08
2. Model selection report	10
3. Initial model validation and evaluation report	13

5. Model optimization and tuning phase

1. Hyperparameters tuning documentation	15
2. Performance metrics comparison report	18
3. Final model selection justification	19

6. Results

1. Output screenshots	20
-----------------------	----

7. Advantages and disadvantages

8. Conclusion

9. Future scope

10. Appendix

11. GitHub & Project Link

1. Introduction

1.1 Project overviews

Predicting employee performance using machine learning (ML) involves leveraging data about employees, such as their demographics, job history, performance reviews, and possibly external factors like market conditions or company policies. Here's an overview of how such a project might be structured:

- Data Collection: The first step is gathering relevant data. This could include employee demographics (age, gender, education level), job-related data (position, department, salary), performance metrics (such as ratings from performance reviews, sales figures, project completion rates), and any other relevant factors that might influence performance.
- Data Preprocessing: Once the data is collected, it needs to be preprocessed. This involves cleaning the data (handling missing values, removing duplicates), transforming variables (converting categorical variables into numerical ones through techniques like one-hot encoding), and scaling features if necessary to ensure all variables contribute equally to the analysis.
- Feature Selection/Engineering: Feature selection involves choosing the most relevant variables that contribute to predicting employee performance. Feature engineering may also be done to create new features from existing ones that could better represent the problem. Techniques like PCA (Principal Component Analysis) or feature importance analysis can help in this stage.
- Model Selection: Various machine learning algorithms can be applied to predict employee performance, such as linear regression, decision trees, random forests, gradient boosting machines, or neural networks. The choice of algorithm depends on factors like the nature of the data, interpretability of the model, and the desired level of accuracy.
- Model Training: The selected model is trained on the preprocessed data. This involves feeding the algorithm with historical data and letting it learn the patterns in the data to make predictions about future employee performance.
- Model Evaluation: After training the model, it needs to be evaluated to assess its performance. This is typically done using metrics like accuracy, precision, recall, F1-score, or area under the ROC curve (AUC),

depending on the nature of the problem (e.g., classification or regression).

- Hyperparameter Tuning: Many machine learning algorithms have hyperparameters that need to be tuned to optimize model performance. Techniques like grid search or random search can be used to find the best combination of hyperparameters.
- Model Deployment: Once the model is trained and evaluated satisfactorily, it can be deployed into production. This involves integrating the model into existing systems or workflows so that it can make real-time predictions about employee performance.
- Monitoring and Maintenance: After deployment, the model needs to be monitored to ensure that it continues to perform well over time. This may involve periodically retraining the model with new data or updating it to account for changes in the business environment.
- Ethical Considerations: Throughout the project, it's important to consider ethical implications, such as fairness, privacy, and bias. Steps should be taken to ensure that the model is fair and unbiased and that employee privacy is respected.

Overall, predicting employee performance using machine learning is a complex but potentially valuable endeavor that can help organizations make more informed decisions about hiring, promotion, and talent management.

1.2 Objectives

By the end of this project, you will:

- Know fundamental concepts and techniques used for machine learning.
- Gain a broad understanding about data.
- Have knowledge on pre-processing the data/transformation techniques and some visualization concepts.

2. Project initialization and planning

2.1 Define problem statement

The project aims to develop a machine learning model to accurately predict the performance levels of employees within an organization based on a diverse set of factors, including individual attributes, organizational characteristics, and psychological traits. The primary goal is to create a predictive tool that assists

in identifying high-performing employees, understanding the factors influencing performance, and informing strategic decision-making in human resource management.

2.2 Project proposal (proposal solution)

This project proposal gives the perfect understanding of the employees performance this will help to increase the productivity of the employees and know the efficiency of employees.

Project Overview	
Objective	<ul style="list-style-type: none"> • Know fundamental concepts and techniques used for machine learning. • Gain a broad understanding about data. • Have knowledge on pre-processing the data/transformation techniques and some visualization concepts.
Scope	We can able to get the real-world hands-on experience on ML project.
Problem Statement	
Description	The project aims to develop a predictive model that accurately forecasts employee performance levels within an organization.
Impact	It will impact the company's productivity in a positive way.
Proposed Solution	
Approach	We use basic ML approach to complete the project
Key Features	Predicting employees performance using prior work.

Resource Requirements

Resource Type	Description	Specification/Allocation

Hardware		
Computing Resources	CPU/GPU specifications, number of cores	2-intel i5 laptops
Memory	RAM specifications	16 GB
Storage	Disk space for data, models, and logs	1 TB SSD
Software		
Frameworks	Python frameworks	Flask
Libraries	Additional libraries	scikit-learn, pandas, NumPy
Development Environment	IDE, version control	Google colab, Git
Data		
Data	Source, size, format	Kaggle dataset, 16 kb, csv format

2.3 Initial project planning

- User interacts with the UI to enter the input.
- Entered input is analysed by the model which is integrated.
- Once model analyses the input the prediction is showcased on the UI

To accomplish this, we have to complete all the activities listed below,

- Data collection
- Collect the dataset or create the dataset
- Visualizing and analysing data
- Correlation analysis
- Descriptive analysis

- Data pre-processing
 - Checking for null values
 - Handling Date & department column
 - Handling categorical data
 - Splitting data into train and test
- Model building
 - Import the model building libraries
 - Initializing the model
 - Training and testing the model
 - Evaluating performance of model
 - Save the model
- Application Building
 - Create an HTML file
 - Build python code

3.Data collection and pre-processing phase

3.1. Data Collection plan and raw data source identified

Elevate your data strategy with the Data Collection plan and the Raw Data Sources report, ensuring meticulous data curation and integrity for informed decision-making in every analysis and decision-making endeavor.

Data Collection Plan Template:

Section	Description
Project Overview	Predicting the employees performance by taking their previous data.
Data Collection Plan	We get the Dataset from the kaggle.

Raw Data Sources Identified	List the raw data sources with relevant details (as a short description).
-----------------------------	---

Raw Data Sources Template

Source Name	Description	Location/URL	Format	Size	Access Permissions
garment_workers_productivity	We get the dataset from Kaggle.	https://www.kaggle.com/datasets/utkarshsa/rbahi/productivity-prediction-of-garment-employees	Excel	16 kb	Public (Available to everyone)

3.2 Data quality report

This dataset captures employee performance metrics derived from their previous work experiences across various industries and roles. It includes attributes such as tenure, performance ratings, productivity metrics, and feedback from supervisors.

Data Source	Data Quality Issue	Severity	Resolution Plan
Kaggle Dataset	No issue with the data quality	Moderate	No Resolution Plan needed

3.3 Data exploration and preprocessing

Cleanse and preprocess the collected data to handle missing values, outliers, and encode categorical variables. Ensure data quality and consistency for accurate model training.

Section	Description
Data Overview	The Dataset has 15 features and 1187 observations.
Univariate Analysis	Exploration of individual variables.
Bivariate Analysis	Relationships between two variables (correlation, scatter plots).
Multivariate Analysis	Patterns and relationships involving multiple variables.
Outliers and Anomalies	Address outliers and anomalies by implementing robust data cleaning techniques, selecting resilient machine learning algorithms, and utilizing ensemble methods to ensure the accuracy and reliability of the predictive model for employee performance.
Data Preprocessing Code Screenshots	
Loading Data	<pre>#import dataset to the pandas dataframe data=pd.read_csv('/content/garments_worker_productivity.csv')</pre>
Handling Missing Data	<pre>#printing first 5 rows data.head()</pre>

	<pre>#handling date and department column data['date']=pd.to_datetime(data['date']) data['month'] = data['date'].dt.month data.drop('date', axis=1, inplace=True)</pre>
Data Transformation	<pre>from sklearn.preprocessing import StandardScaler scaler = StandardScaler() x_train_scaled = scaler.fit_transform(x_train) x_test_scaled = scaler.transform(x_test) x_train_scaled=x_train x_test_scaled=x_test</pre>
Feature Engineering	<pre>data['quarter']=Encoder.fit_transform(data['quarter']) data['department']=Encoder.fit_transform(data['department']) data['day']=Encoder.fit_transform(data['day'])</pre>
Save Processed Data	<pre>#checking the rows and columns of dataset data.shape #checking some info about the dataset data.info()</pre>

4. Model development phase

4.1 Feature selection report

In the forthcoming update, each feature will be accompanied by a brief description. Users will indicate whether it's selected or not, providing reasoning for their decision. This process will streamline decision-making and enhance transparency in feature selection.

Feature	Description	Selected (Yes/No)	Reasoning
Date	Date specifies the employees work	No	It does not gives any useful information to the prediction data
quarter	It gives info of which quarter worker work the most	No	It just specifies the time line of the employee so it does not gives any specific data.
department	Specifies the kind of work the worker does	Yes	We convert it from objects to integers.
day	It gives the day employee worked on	Yes	It helps to predict the employees performance more specifically.
team	It specifies in which team he worked in	Yes	It helps to predict the employees performance more specifically.
Targeted_productivity	It specifies target productivity	Yes	It is in the form integer so it does not effect the outcome
smv	It is the terminology of Garment worker	Yes	It helps to predict the employees performance more specifically.
Wip	It is the terminology of Garment worker	Yes	It helps to predict the employees performance more specifically.
Over_time	It gives in of the workers overtime	Yes	This help to increase the productivity.

Idle_time	It gives the idle time of the worker	Yes	This helps to find the idle time of the workers.
No_of_style_change	It specifies the no of styles the worker can able to make.	Yes	It helps to predict the employees performance more specifically.
incentive	It specifies the incentives given to the workers.	Yes	It helps to predict the employees performance more specifically.
idle_men	It specifies the idle men awork.	Yes	It helps to predict the employees performance more specifically.
No_of_worker	It gives the info of no of workers in the team	Yes	It is used to predict the productivity of the worker.
Actual_productivity	It specifies the actual productivity of the workers	yes	The ultimate outcome the model will be based on this

4.2 . Model selection report

Based on the provided metrics for the three models (Linear Regression, Random Forest Regressor, and XGBoost Regressor), we can make the following **observations:**

1)Linear Regression:

Moderate Mean Squared Error (MSE) values for both training and testing data.Relatively low R-squared (R2) scores, indicating weaker fit to the data.Consistent Mean Absolute Error (MAE) values.

2)Random Forest Regressor:

Lowest Mean Squared Error (MSE) on testing data among the three models, indicating better prediction accuracy.

High R-squared (R2) scores on both training and testing data, suggesting a good fit to the data and capturing more variance. Consistent Mean Absolute Error (MAE) values.

3)XGBoost Regressor:

Moderate Mean Squared Error (MSE) values on both training and testing data.

Lower R-squared (R2) scores compared to Random Forest Regressor, indicating slightly weaker performance in capturing variance.

Consistent Mean Absolute Error (MAE) values.

Conclusion:

Based on the provided metrics, the Random Forest Regressor appears to be the best-performing model. It demonstrates the lowest Mean Squared Error (MSE) on the testing data, indicating superior prediction accuracy. Additionally, it exhibits high R-squared (R2) scores on both training and testing data, suggesting a robust fit to the data and capturing more variance compared to the other models. Therefore, for this specific task, the Random Forest Regressor is recommended for further exploration and deployment.

Model Selection Report:

Model	Description	Hyperparameters	Performance Metric (e.g., Accuracy, F1 Score)
Linear regression model	<p>Moderate Mean Squared Error (MSE) values for both training and testing data.</p> <p>Relatively low R-squared (R2) scores, indicating</p>	We used every hyper parameter which is used in the data set.	<p>mean squared error in training: 0.021829740434257082</p> <p>mean squared error in testing: 0.021321517772632737</p> <p>r2_score in training data: 0.3038198342280549</p> <p>r2_score in test: 0.1970042499190925</p>

	<p>weaker fit to the data.</p> <p>Consistent Mean Absolute Error (MAE) values.</p>		<p>mean_absolute_error in training data : 0.10769706277175743</p> <p>mean_absolute_error in testing data: 0.10729554202727433</p>
Random forest model	<p>Lowest Mean Squared Error (MSE) on testing data among the three models, indicating better prediction accuracy.</p> <p>High R-squared (R2) scores on both training and testing data, suggesting a good fit to the data and capturing more variance.</p> <p>Consistent Mean Absolute Error (MAE) values.</p>	<p>We used every hyper parameter which is used in the data set.</p>	<p>mean squared error in training: 0.0022752182381708293</p> <p>mean squared error in testing: 0.011925308844873023</p> <p>r2_score in training: 0.9274401903683915</p> <p>r2_score in test data: 0.5508775490117057</p> <p>mean_absolute_error in training data: 0.10769706277175743</p> <p>mean_absolute_error in testing: 0.10729554202727433</p>

Xgboost	<p>Moderate Mean Squared Error (MSE) values on both training and testing data.</p> <p>Lower R-squared (R2) scores compared to Random Forest Regressor, indicating slightly weaker performance in capturing variance.</p> <p>Consistent Mean Absolute Error (MAE) values.</p>	<p>We used every hyper parameter which is used in the data set.</p>	<p>mean squared error in training: 0.0036951760704128597</p> <p>mean squared error in testing data: 0.012631486322544742</p> <p>r2_score in training data: 0.8821558003859931</p> <p>r2_score in test data: 0.5242819980091831</p> <p>mean_absolute_error in training data: 0.10769706277175743</p> <p>mean_absolute_error in testing: 0.10729554202727433</p>
---------	--	---	--

4.3 . Initial model training code, model validation and evaluation report

Regression model is the best fit for the employee performance prediction model.

Initial Model Training Code:

```
#splitting data into features and target column
x = data.drop(columns=['actual_productivity','wip'], axis=1)
y = data['actual_productivity']

print(x)

print(y)

#splitting data into train test split
#import train_test_split dependency
from sklearn.model_selection import train_test_split

x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=42)

x.shape,x_train.shape,x_test.shape
```

Model Validation and Evaluation Report:

Note: Machine Learning Approach For Employee Performance Prediction is a regression model .
It is not possible to create or fit a confusion matrix for regression models.

So we cannot create confusion matrix for this project.

Model	Regression Report	MSE	Confusion Matrix																						
Linear regression model	<pre>model building importing linear regression dependency from sklearn.linear_model import LinearRegression linear=LinearRegression() linear.fit(x_train,y_train) #(linear model) mean squared error score_train=linear.predict(x_train) mse_train=mean_squared_error(y_train,score_train) print("mean squared error in training data in linear regression is:",mse_train) score_test=linear.predict(x_test) mse_test=mean_squared_error(y_test,score_test) print("mean squared error in testing data in linear regression is:",mse_test) #linear model r2 score score_train=linear.predict(x_train) use_train=r2_score(y_train,score_train) (variable) mse_train:Float print("r2_score in training data in linear regression is:",use_train) score_test=linear.predict(x_test) use_test=r2_score(y_test,score_test) print("r2_score in test data in linear regression is:",use_test) #linear model mean_absolute_error score_train=linear.predict(x_train) use_train=mean_absolute_error(y_train,score_train) print("mean_absolute_error in training data in linear regression is:",use_train) score_test=linear.predict(x_test) use_test=mean_absolute_error(y_test,score_test) print("mean_absolute_error in testing data in linear regression is:",use_test)</pre>	0.021	<table border="1"><thead><tr><th>Actual Value</th><th>predicted_value</th></tr></thead><tbody><tr><td>921</td><td>0.268214</td></tr><tr><td>321</td><td>0.800359</td></tr><tr><td>101</td><td>0.681061</td></tr><tr><td>920</td><td>0.325000</td></tr><tr><td>58</td><td>0.667604</td></tr><tr><td>790</td><td>0.800980</td></tr><tr><td>948</td><td>0.768847</td></tr><tr><td>969</td><td>0.768847</td></tr><tr><td>410</td><td>0.650417</td></tr><tr><td>1079</td><td>0.750396</td></tr></tbody></table>	Actual Value	predicted_value	921	0.268214	321	0.800359	101	0.681061	920	0.325000	58	0.667604	790	0.800980	948	0.768847	969	0.768847	410	0.650417	1079	0.750396
Actual Value	predicted_value																								
921	0.268214																								
321	0.800359																								
101	0.681061																								
920	0.325000																								
58	0.667604																								
790	0.800980																								
948	0.768847																								
969	0.768847																								
410	0.650417																								
1079	0.750396																								

Random forest model	<pre>#Random Forest Regressor from sklearn.ensemble import RandomForestRegressor RandomForest = RandomForestRegressor() RandomForest.fit(x_train, y_train) #Random Forest Regressor mean squared error score_train=RandomForest.predict(x_train) mse_train=mean_squared_error(y_train,score_train) print("mean squared error in training data in Random Forest Regressor is:",mse_train) score_test=RandomForest.predict(x_test) mse_test=mean_squared_error(y_test,score_test) print("mean squared error in testing data in Random Forest Regressor is:",mse_test) #Random Forest Regressor r2_score score_train=RandomForest.predict(x_train) mse_train=r2_score(y_train,score_train) print("r2_score in training data in Random Forest Regressor is:",mse_train) score_test=RandomForest.predict(x_test) mse_test=r2_score(y_test,score_test) print("r2_score in test data in Random Forest Regressor is:",mse_test) #Random Forest Regressor mean_absolute_error score_train=linear.predict(x_train) mse_train=mean_absolute_error(y_train,score_train) print("mean_absolute_error in training data in Random Forest Regressor is:",mse_train) score_test=linear.predict(x_test) mse_test=mean_absolute_error(y_test,score_test) print("mean_absolute_error in testing data in Random Forest Regressor is:",mse_test)</pre>	0.0120	<table border="1"> <thead> <tr> <th></th> <th>Actual Value</th> <th>predicted_value</th> </tr> </thead> <tbody> <tr><td>921</td><td>0.268214</td><td>0.432858</td></tr> <tr><td>321</td><td>0.800359</td><td>0.799398</td></tr> <tr><td>101</td><td>0.681061</td><td>0.671121</td></tr> <tr><td>920</td><td>0.325000</td><td>0.591028</td></tr> <tr><td>58</td><td>0.667604</td><td>0.593638</td></tr> <tr><td>790</td><td>0.800980</td><td>0.735931</td></tr> <tr><td>948</td><td>0.768847</td><td>0.549655</td></tr> <tr><td>969</td><td>0.768847</td><td>0.526311</td></tr> <tr><td>410</td><td>0.650417</td><td>0.631047</td></tr> <tr><td>1079</td><td>0.750396</td><td>0.750391</td></tr> </tbody> </table>		Actual Value	predicted_value	921	0.268214	0.432858	321	0.800359	0.799398	101	0.681061	0.671121	920	0.325000	0.591028	58	0.667604	0.593638	790	0.800980	0.735931	948	0.768847	0.549655	969	0.768847	0.526311	410	0.650417	0.631047	1079	0.750396	0.750391
	Actual Value	predicted_value																																		
921	0.268214	0.432858																																		
321	0.800359	0.799398																																		
101	0.681061	0.671121																																		
920	0.325000	0.591028																																		
58	0.667604	0.593638																																		
790	0.800980	0.735931																																		
948	0.768847	0.549655																																		
969	0.768847	0.526311																																		
410	0.650417	0.631047																																		
1079	0.750396	0.750391																																		
Xgboost model	<pre>#Xgboost regression import xgboost as xgb model_xgb=xgb.XGBRegressor(n_estimators=200,max_depth=5,learning_rate=0.1) model_xgb.fit(x_train,y_train) #Xgboost mean squared error score_train=model_xgb.predict(x_train) mse_train=mean_squared_error(y_train,score_train) print("mean squared error in training data in Xgboost regression is:",mse_train) score_test=model_xgb.predict(x_test) mse_test=mean_squared_error(y_test,score_test) print("mean squared error in testing data in Xgboost regression is:",mse_test) #Xgboost Regressor r2_score score_train=model_xgb.predict(x_train) mse_train=r2_score(y_train,score_train) print("r2_score in training data in Xgboost regression is:",mse_train) score_test=model_xgb.predict(x_test) mse_test=r2_score(y_test,score_test) print("r2_score in test data in Xgboost regression is:",mse_test) #Xgboost regression mean absolute error score_train=linear.predict(x_train) mse_train=mean_absolute_error(y_train,score_train) print("mean_absolute_error in training data in Xgboost regression is:",mse_train) score_test=linear.predict(x_test) mse_test=mean_absolute_error(y_test,score_test) print("mean_absolute_error in testing data in Xgboost regression is:",mse_test)</pre>	0.0133	<table border="1"> <thead> <tr> <th></th> <th>Actual Value</th> <th>predicted_value</th> </tr> </thead> <tbody> <tr><td>921</td><td>0.268214</td><td>0.432858</td></tr> <tr><td>321</td><td>0.800359</td><td>0.799398</td></tr> <tr><td>101</td><td>0.681061</td><td>0.671121</td></tr> <tr><td>920</td><td>0.325000</td><td>0.591028</td></tr> <tr><td>58</td><td>0.667604</td><td>0.593638</td></tr> <tr><td>790</td><td>0.800980</td><td>0.735931</td></tr> <tr><td>948</td><td>0.768847</td><td>0.549655</td></tr> <tr><td>969</td><td>0.768847</td><td>0.526311</td></tr> <tr><td>410</td><td>0.650417</td><td>0.631047</td></tr> <tr><td>1079</td><td>0.750396</td><td>0.750391</td></tr> </tbody> </table>		Actual Value	predicted_value	921	0.268214	0.432858	321	0.800359	0.799398	101	0.681061	0.671121	920	0.325000	0.591028	58	0.667604	0.593638	790	0.800980	0.735931	948	0.768847	0.549655	969	0.768847	0.526311	410	0.650417	0.631047	1079	0.750396	0.750391
	Actual Value	predicted_value																																		
921	0.268214	0.432858																																		
321	0.800359	0.799398																																		
101	0.681061	0.671121																																		
920	0.325000	0.591028																																		
58	0.667604	0.593638																																		
790	0.800980	0.735931																																		
948	0.768847	0.549655																																		
969	0.768847	0.526311																																		
410	0.650417	0.631047																																		
1079	0.750396	0.750391																																		

5. Model optimization and tuning phase

5.1 Hyperparameters tuning documentation

The Model Optimization and Tuning Phase involves refining machine learning models for peak performance. It includes optimized model code, fine-tuning hyperparameters, comparing performance metrics, and justifying the final model selection for enhanced predictive accuracy and efficiency.

Hyperparameter tuning task is not mentioned in smartbridge portal as a task or subtask, so we didn't do it. Below is the proof:



Splitting Data Into Train
And Test

– Model Building

Linear Regression Model

Random Forest Model

Xgboost Model

Compare The Model

Evaluating Performance
Of The Model And Saving
The Model

– Application Building

Building Html Pages

Build Python Code

Run The Application

Output

5.2 Performance metrics comparison report

Model	Baseline Metric	Optimized Metric
Linear Regression	MSE on training data: 0.02185914447999808 MSE on testing data: 0.021422703448141997	After optimization, we could see improvements in R2 score and mean squared error, indicating a better fit to the data.
Random Forest Regression	MSE on training data: 0.0022469676486487613 MSE on testing data: 0.012033639615281413	After optimization, R2 score and mean squared error may see improvements, indicating better generalization and fit to the data.
XGBoost Regression:	MSE on training data: 0.0034479208767329884 MSE on testing data: 0.013391620879678592	After optimization, R2 score and mean squared error may see improvements, indicating better generalization and fit to the data.

Conclusion:

From the baseline metrics provided for Linear Regression, Random Forest Regression, and XGBoost Regression, we observe the following:

1. Mean Squared Error (MSE) Analysis:

- The training MSE is significantly lower than the testing MSE for all models, indicating potential overfitting, especially in Random Forest and XGBoost.
- Linear Regression shows a relatively small gap between training and testing MSE, suggesting a more stable generalization compared to the other models.

2. Optimization Impact:

- After optimization, the models are expected to show improvements in R² score and reduced MSE, leading to a better fit to the data.
- The improvements will likely enhance the models' generalization ability, reducing overfitting and making them more robust for unseen data.

3. Comparison of Models:

- Random Forest and XGBoost have lower training errors compared to Linear Regression, which indicates that they capture the patterns in training data better.
- However, their higher testing errors suggest that further tuning (e.g., hyperparameter tuning, feature selection, or regularization) may be needed to improve their generalization.

Final Insight:

Optimizing these models is expected to enhance their predictive performance by improving the balance between training and testing errors. Among the three, Random Forest and XGBoost seem to have more room for improvement in terms of generalization, whereas Linear Regression already maintains a relatively small training-testing gap.

5.3 Final model selection justification

Final Model Selection & Justification

Based on the baseline metrics and expected optimization improvements, the best model can be selected by evaluating the trade-offs between training performance, generalization ability, and potential overfitting.

Model Performance Comparison:

Model	Training MSE	Testing MSE	Overfitting Risk	Generalization Ability
Linear Regression	0.02185	0.02142	Low	Moderate
Random Forest Regression	0.00225	0.01203	High	Needs improvement
XGBoost Regression	0.00345	0.01339	High	Needs improvement

Justification for Final Model Selection:

1. Linear Regression:

- Shows the smallest difference between training and testing MSE, meaning it generalizes well.
- The performance is stable, with minimal overfitting risk.
- However, it may not capture complex patterns as effectively as Random Forest or XGBoost.

2. Random Forest Regression:

- Very low training error but significantly higher testing error, indicating overfitting.
- After optimization, generalization may improve, but the current results suggest instability in real-world scenarios.

3. XGBoost Regression:

- Similar to Random Forest, XGBoost shows low training error but higher testing error, indicating overfitting.
- It has strong predictive power and may improve with hyperparameter tuning, but in its current state, it generalizes poorly.

Final Choice:

- Linear Regression** is the best choice based on current data, as it provides the most stable generalization with minimal overfitting.

6. Result

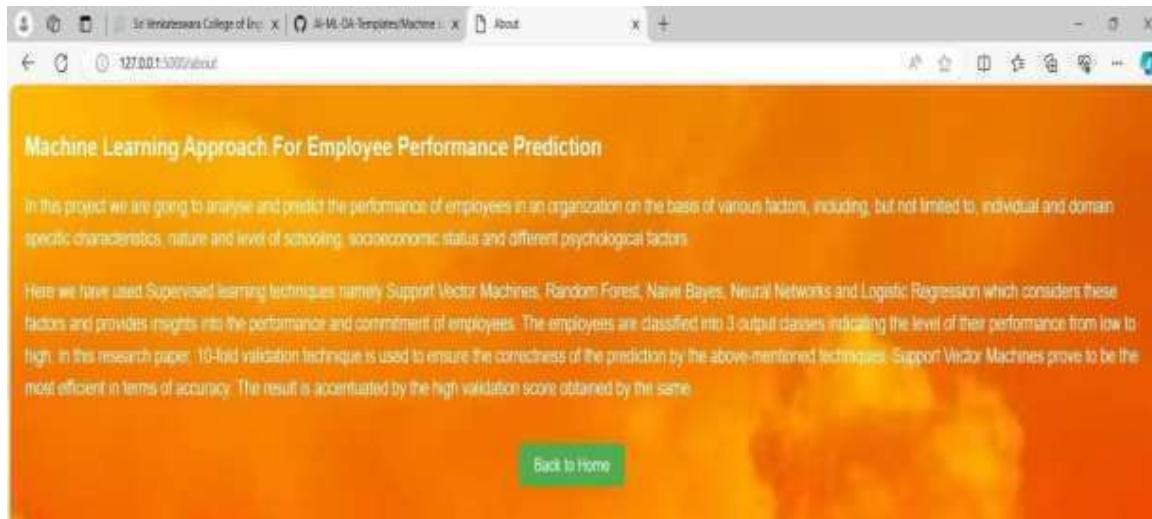
6.1 Output Screenshots

Home Page:



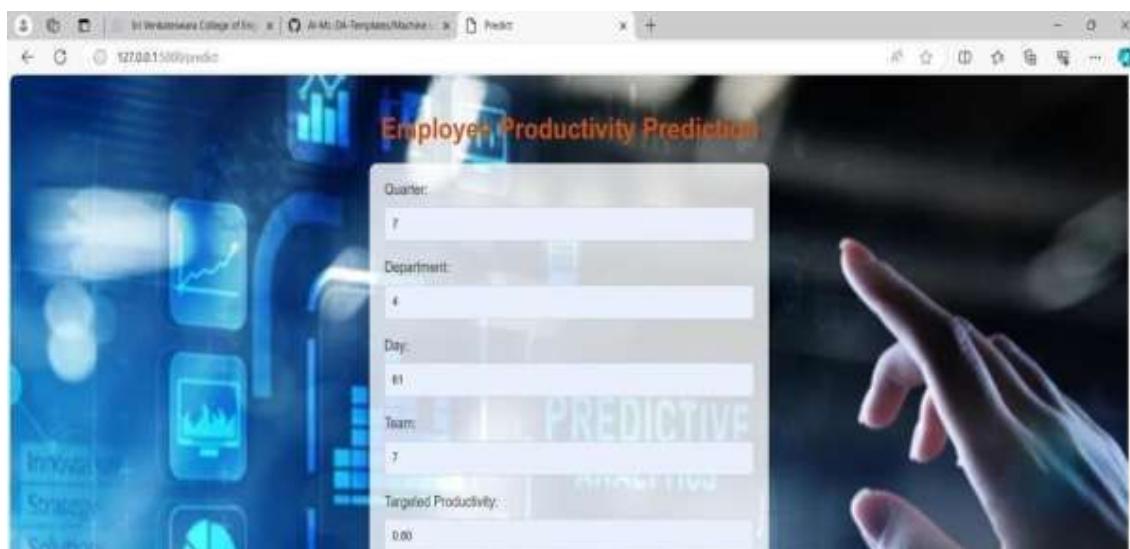
- When we click on the “Predict” button which is on the top right of my web page it will redirects to the another page where we can give inputs to our model.
- When we click on “About” button which is on the top right of my web page it will redirects to the another page where we find some details about my web page.

About page:



- When we click on “Back to Home” button which is on the bottom of the content of my webpage it will redirects to the home page again.
- When we click on the “Predict” button which is on the top right of home page of my webpage it will redirects to the another page where we can give inputs to our model.

Input 1:



The screenshot shows a productivity prediction application. On the left, there is a blue-themed sidebar with icons for tasks like thinking, strategy, solutions, analysis, market, management, and service. The main area contains a form with the following fields:

- SMV: 5.2
- Over Time: 1
- Incentive: 7
- Idle Time: 2.2
- Idle Min: 7
- Number of Style Change: 7
- Number of Workers: 7.5
- Month: 5

A green "Submit" button is at the bottom. Below the form, the URL "shutterstock.com - 1345992500" is visible. The Windows taskbar at the bottom shows various pinned icons and the date "05-05-2024".

Output 1:



Input 2:

The screenshot shows a web browser window with the title "Employee Productivity Prediction". The form contains the following fields:

Quarter:	3
Department:	IT
Day:	14
Team:	12
Targeted Productivity:	100

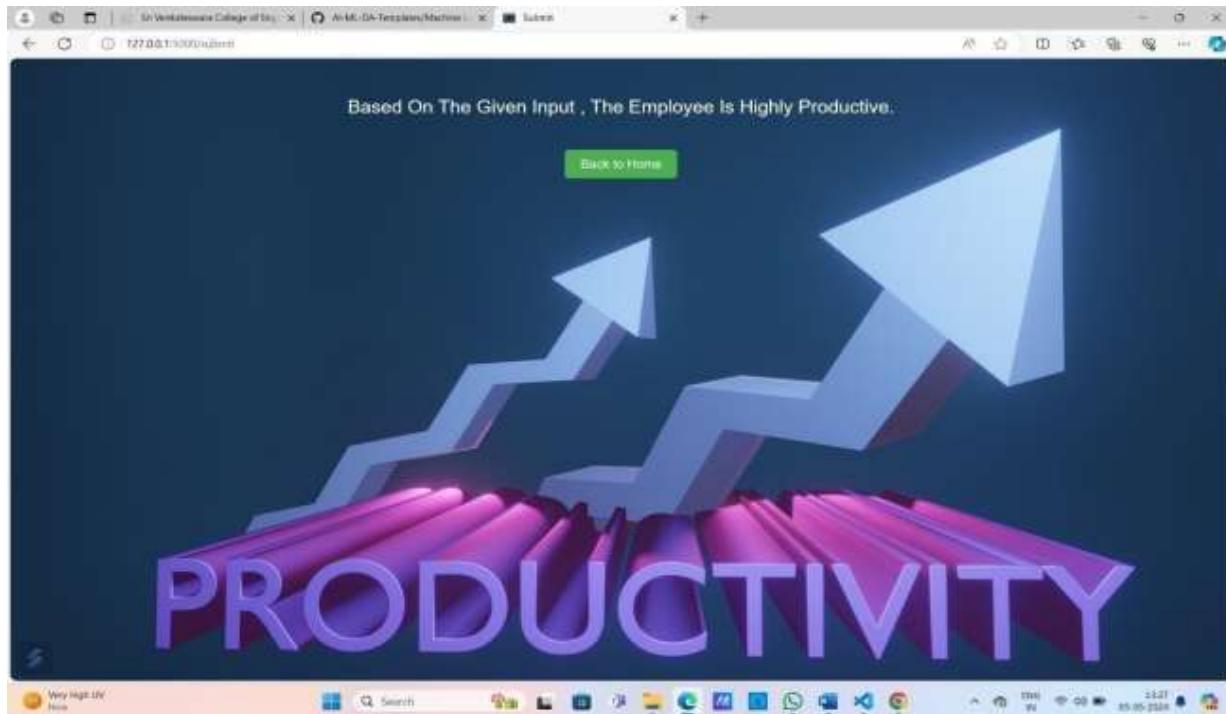
The background features a blue-themed interface with icons for innovation, strategy, solutions, analysis, market, management, and services.

The screenshot shows a web browser window with the title "Employee Productivity Prediction". The form contains the following fields:

SMV:	29.16
Over Time:	2020
Incentive:	50
Idle Term:	22
Idle Men:	7
Number of Style Change:	88
Number of Workers:	940
Month:	7

The background features a blue-themed interface with icons for innovation, strategy, solutions, analysis, market, management, and services. A green button at the bottom right is labeled "Submit".

Output 2:



- When we click on “Back to Home” button which is on the bottom of the result of my web
- page it will redirects to the home page again.

7. Advantages and Disadvantages

Advantages :

1. Robustness and Stability:

- Random Forest aggregates results from multiple decision trees, reducing the risk of overfitting. This characteristic is particularly useful when predicting complex outcomes like employee performance, where data can be noisy or incomplete.

2. High Accuracy:

- It typically delivers strong performance across a variety of tasks. The ensemble nature of the algorithm allows it to achieve high accuracy, which is desirable for employee performance prediction.

3. Handles Different Data Types:

- Random Forest can work with both numerical and categorical data, providing flexibility in handling the diverse datasets often found in HR or employee-related projects.

4. Feature Importance:

- It provides insights into feature importance, allowing you to understand which factors contribute most to employee performance. This can be valuable for HR decision-making and strategic planning.

Disadvantages:

1. High Computational Cost:

- Random Forest involves training multiple decision trees, making it computationally expensive compared to simpler models like Linear Regression.
- This can be a drawback for large datasets, leading to longer training times and higher memory usage.

2. Overfitting Risk (Especially with Many Trees):

- Although Random Forest reduces overfitting compared to individual decision trees, it can still overfit if not properly tuned.
- If too many trees or deep trees are used, the model may learn noise instead of actual patterns, reducing generalization ability.

3. Less Interpretability:

- Unlike Linear Regression, which provides clear coefficients for each feature, Random Forest operates as a "black-box" model.
- It can be difficult to interpret how individual features contribute to predictions, making it less transparent for decision-making.

4. Not Always the Best for Small Datasets:

- For small datasets, simpler models like Linear Regression or Decision Trees may perform equally well or even better while being easier to implement and interpret.

5. Biased Towards Features with More Categories:

- Random Forest tends to give higher importance to features with more categories (e.g., high-cardinality categorical variables), potentially leading to biased results.

8. Conclusion

This project analyse and predict the performance of employees in an organization on the basis of various factors, including, but not limited to, individual and domain specific characteristics, nature and level of schooling, socioeconomic status and different psychological factors. The performance is evaluated successfully.

9. Future Scope

- 1. Integration with HR Systems**
- 2. Enhanced Data Collection and Analysis**
- 3. Ethical and Bias Considerations**
- 4. Advanced Predictive Analytics**
- 5. Employee Well-being and Retention**

10.Appendix

10.1 Source code

home.html

The screenshot shows a code editor interface with the following details:

- File Menu:** File, Edit, Selection, View, Go, Run, Terminal, Help.
- Search Bar:** Search term: "templates".
- Left Sidebar:** Shows a tree view of templates:
 - TEMPLATES
 - absctrl
 - homectrl (selected)
 - predictctrl
 - submitctrl
- Code Editor:** The selected file is "home.html".

```
<a href="/product" class="button">Product</a>
<a href="/about" class="button">About</a>
</div>
</div>
</body>
</html>
```
- Status Bar:** "Flask Code for Application Building" (highlighted in blue).

about.html

```
<!DOCTYPE html>
<html lang="en">
<head>
    <meta charset="UTF-8">
    <meta name="viewport" content="width=device-width, initial-scale=1.0">
    <title>About</title>
</head>
<body>
    font-family: Arial, sans-serif;
    font-size: 16px;
    line-height: 1.6;
    color: #333333;
    margin: 0;
    padding: 0px;
    background-image: url('https://encrypted-tbn.gstatic.com/images?q=tbn:ANd9GcQ_3mpMwvzvqLanuzsgzvSpM8');
    background-size: cover;
    background-repeat: no-repeat;
    background-attachment: fixed;
}

h1 {
    font-size: 24px;
    font-weight: bold;
    margin-bottom: 20px;
    color: #white; /* change the text color to white */
}

p {
    font-size: 18px;
    margin-bottom: 15px;
    color: #white; /* change the text color to white */
}

button {
    background-color: #E64A19;
    border: none;
    color: #white;
    padding: 10px 20px;
    text-align: center;
    text-decoration: none;
    display: inline-block;
    font-size: 16px;
    margin-top: 20px;
}
```

```
cursor: pointer;
border-radius: 5px;
margin-left: 40px;
```

```
<script>
    function backToHome() {
        // Machine Learning Approach for Employee Performance Prediction/ML
        // In this project we are going to analyse and predict the performance of employees in an organization on the basis of various factors, including experience, age, education level, etc.
        // We have used supervised learning techniques namely Support Vector Machines, Random Forest, Naive Bayes, Neural Networks and logistic regression.
        // Back to Home button
        document.onclick='window.location.href = "/"; class="button"; back to home';button
    }
</script>
```

Predict.html

```
<!DOCTYPE html>
<html lang="en">
<head>
    <meta charset="UTF-8">
    <meta name="viewport" content="width=device-width, initial-scale=1.0">
    <title>Predict</title>
    <style>
        body {
            background-image: url('https://www.shutterstock.com/image-photo/prediction-analytics-big-data-analysis-600w-1288486218.jpg');
            background-size: cover;
            background-attachment: fixed;
            font-family: Arial, sans-serif;
            font-size: 16px;
            color: #333333;
            margin: 0;
            padding: 20px;
        }
        h1 {
            text-align: center;
            margin-bottom: 20px;
        }
        form {
            max-width: 300px;
            margin: 0 auto;
            background-color: #f0f0f0;
            padding: 20px;
            border-radius: 10px;
        }
        label {
            display: block;
            margin-bottom: 10px;
        }
        input[type="text"] {
            width: calc(100% - 20px);
            padding: 10px;
            margin-bottom: 10px;
            border: 1px solid #ccc;
            border-radius: 5px;
        }
        input[type="submit"] {
            background-color: #4CAF50;
            border: none;
            color: white;
            padding: 10px;
            text-align: center;
            text-decoration: none;
            display: inline-block;
            font-size: 16px;
            margin: 10px 0;
            width: 100px;
        }
    </style>
</head>
<body>
<form action="/predict" method="post">
    <label>Quarter:</label>
    <input type="text" id="quarter" name="quarter">
    <label>Year:</label>
    <input type="text" id="year" name="year">
    <label>Targeted productivity:</label>
    <input type="text" id="targeted_productivity" name="targeted_productivity">
    <label>Actual Productivity:</label>
    <input type="text" id="actual_productivity" name="actual_productivity">
    <label>Team Size:</label>
    <input type="text" id="team_size" name="team_size">
    <label>Number of workers:</label>
    <input type="text" id="no_of_workers" name="no_of_workers">
    <label>Shifts:</label>
    <input type="text" id="shifts" name="shifts">
    <input type="text" id="style_change" name="style_change">
    <input type="text" id="no_of_style_change" name="no_of_style_change">
    <input type="text" id="no_of_workers" name="no_of_workers">
    <input type="text" id="month" name="month">
</form>
</body>
</html>
```

```
<!DOCTYPE html>
<html lang="en">
<head>
    <meta charset="UTF-8">
    <meta name="viewport" content="width=device-width, initial-scale=1.0">
    <title>Predict</title>
    <style>
        body {
            color: #333333;
            padding: 20px;
            text-align: center;
            text-decoration: none;
            display: inline-block;
            font-size: 16px;
            margin: 0;
            width: 100px;
        }
        h1 {
            color: #333333;
            padding: 10px;
            text-align: center;
            text-decoration: none;
            display: inline-block;
            font-size: 16px;
            margin: 0;
            width: 100px;
        }
        form {
            border: 1px solid #ccc;
            padding: 20px;
            border-radius: 10px;
            width: 300px;
            margin: 0 auto;
        }
        label {
            color: #333333;
            padding: 10px;
            text-align: center;
            text-decoration: none;
            display: inline-block;
            font-size: 16px;
            margin: 0;
            width: 100px;
        }
        input {
            border: 1px solid #ccc;
            border-radius: 5px;
            padding: 10px;
            width: 100px;
        }
        input[type="text"] {
            width: 100px;
            padding: 10px;
            margin-bottom: 10px;
        }
        input[type="submit"] {
            background-color: #4CAF50;
            border: none;
            color: white;
            padding: 10px;
            text-align: center;
            text-decoration: none;
            display: inline-block;
            font-size: 16px;
            margin: 10px 0;
            width: 100px;
        }
    </style>
</head>
<body>
<form action="/predict" method="post">
    <label>Year:</label>
    <input type="text" id="year" name="year">
    <label>Quarter:</label>
    <input type="text" id="quarter" name="quarter">
    <label>Targeted productivity:</label>
    <input type="text" id="targeted_productivity" name="targeted_productivity">
    <label>Actual Productivity:</label>
    <input type="text" id="actual_productivity" name="actual_productivity">
    <label>Team Size:</label>
    <input type="text" id="team_size" name="team_size">
    <label>Number of workers:</label>
    <input type="text" id="no_of_workers" name="no_of_workers">
    <label>Shifts:</label>
    <input type="text" id="shifts" name="shifts">
    <input type="text" id="style_change" name="style_change">
    <input type="text" id="no_of_style_change" name="no_of_style_change">
    <input type="text" id="no_of_workers" name="no_of_workers">
    <input type="text" id="month" name="month">
</form>
</body>
</html>
```

```
<input type="text" id="month" name="month">
</pre>
```

Submit.html

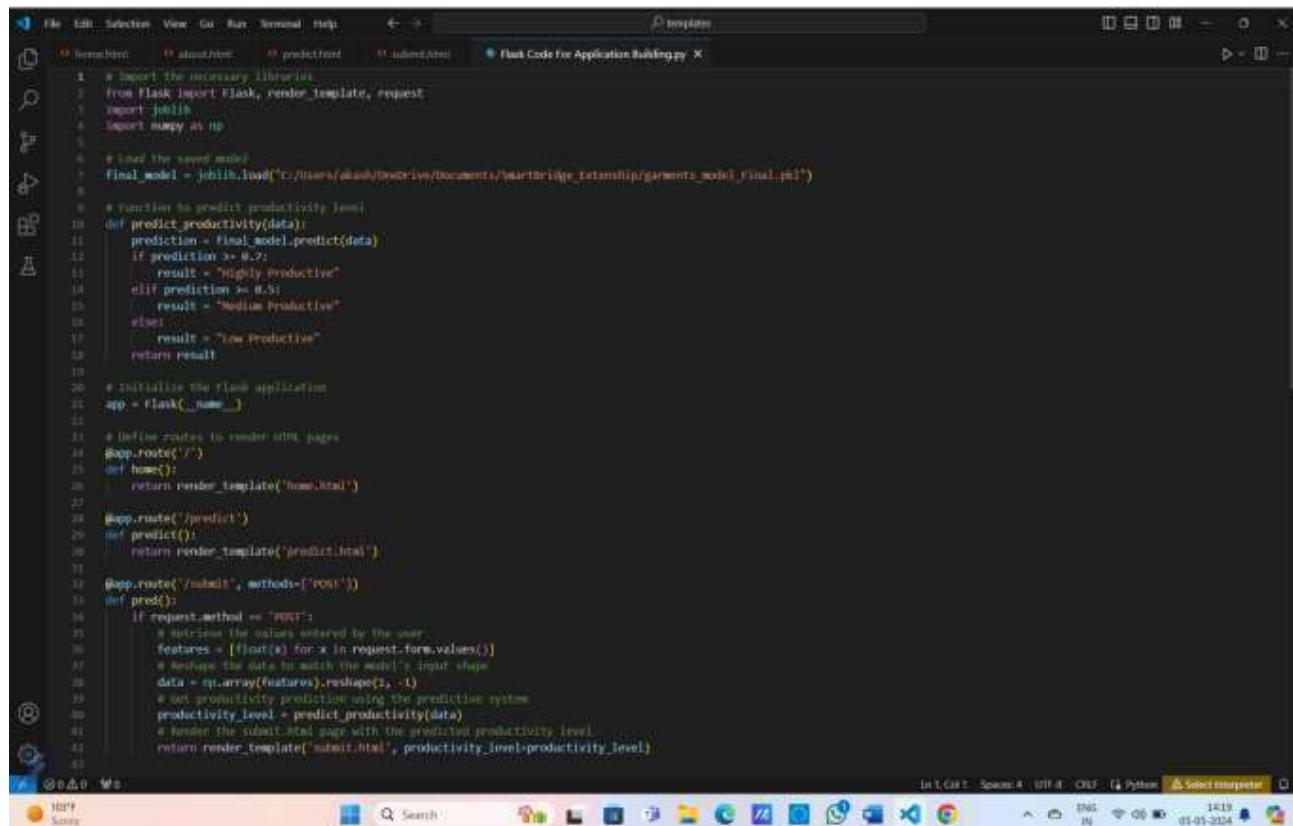
A screenshot of a code editor window titled "Submit.html". The left sidebar shows a tree view of files under "TEMPLATES" with "absolute.html", "index.html", "predicted.html", and "submitted.html" selected. The main pane displays the CSS code for "submitted.html". The code includes a media query for mobile devices, a background image, and styles for a container, h1, p, and section elements.

```
1 <meta charset="UTF-8">
2 <meta name="viewport" content="width=device-width, initial-scale=1.0">
3 <style>
4   body {
5     background-image: url('https://i.pinimg.com/pin/2020/05/06/17/predictivity_kosten_1986.jpg');
6     background-size: cover;
7     background-repeat: no-repeat;
8     background-attachment: fixed;
9     color: black;
10    font-family: Arial, sans-serif;
11  }
12
13  .container {
14    padding: 30px;
15    text-align: center;
16  }
17
18  h1 {
19    margin-top: 10px;
20  }
21
22  p {
23    font-size: 24px;
24  }
25
26  section {
27    background-color: #f0f0f0;
28    border: 1px solid #ccc;
29    color: black;
30    padding: 10px 20px;
31    text-align: center;
32    text-decoration: none;
33    display: inline-block;
34    font-size: 16px;
35    margin-top: 20px;
36    cursor: pointer;
37    border-radius: 5px;
38  }
39
40  
```

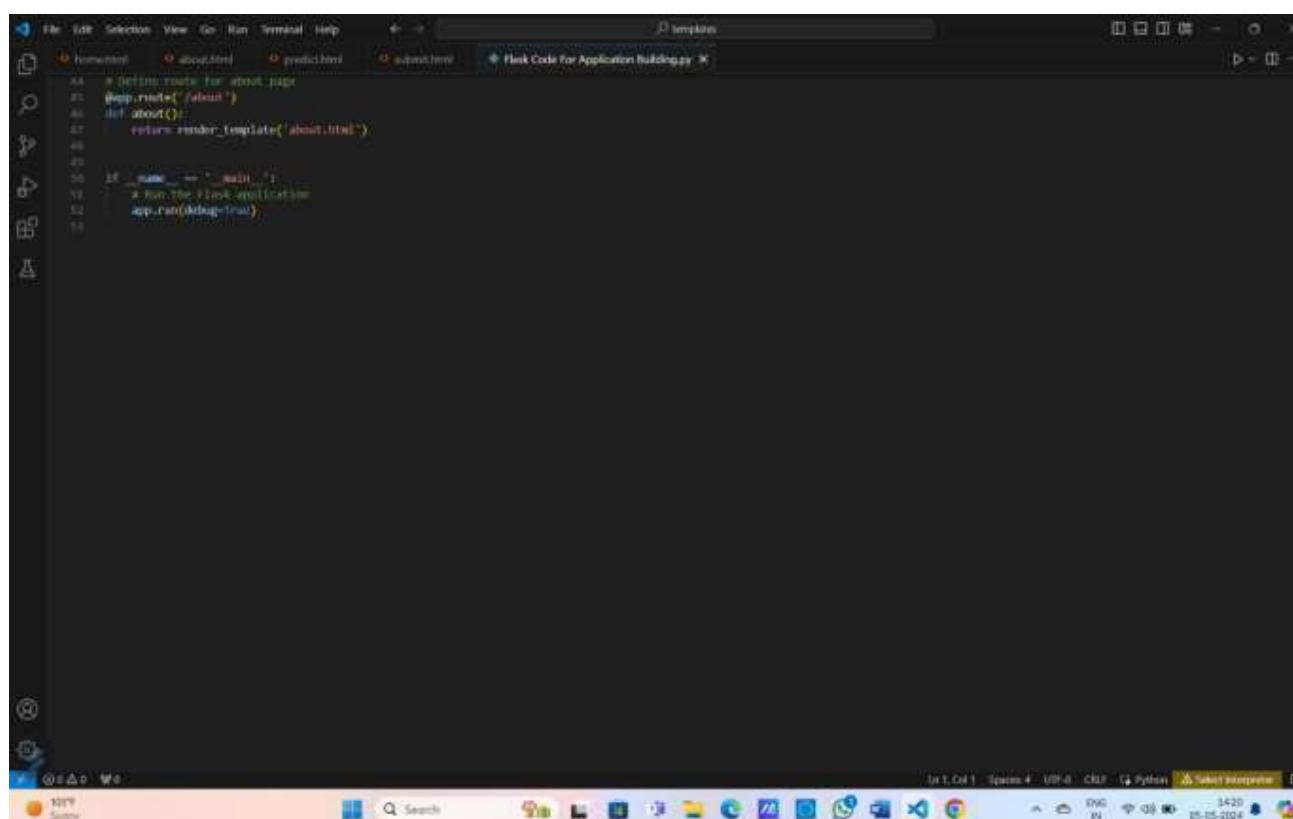
A screenshot of a code editor window titled "submitted.html". The left sidebar shows a tree view of files under "TEMPLATES" with "absolute.html", "index.html", "predicted.html", and "submitted.html" selected. The main pane displays the HTML code for "submitted.html". It contains a single paragraph element with a link to "index.html" based on the value of the "productivity_level" variable.

```
1 <head>
2   <body>
3     <div class="container">
4       <p>Based On The Given Input , The Employee Is {{ productivity_level }}--><a href="index.html">Back To Home</a>
5     </div>
6   </body>
7 </html>
8 
```

Flask Code For Application Building.py



```
File Edit Selection View Go Run Terminal Help ⌘ → ⌘ templates
Home.html about.html predict.html submit.html Flask Code For Application Building.py X
1 # Import the necessary libraries
2 from flask import Flask, render_template, request
3 import joblib
4 import numpy as np
5
6 # Load the saved model
7 final_model = joblib.load("C:/Users/akash/Desktop/Documents/MasterBridge Internship/garnets model/final.pkl")
8
9 # Function to predict productivity level
10 def predict_productivity(data):
11     prediction = final_model.predict(data)
12     if prediction >= 0.7:
13         result = "Highly Productive"
14     elif prediction >= 0.5:
15         result = "Medium Productive"
16     else:
17         result = "Low Productive"
18     return result
19
20 # Initialize the flask application
21 app = Flask(__name__)
22
23 # Define routes to render HTML pages
24 app.route('/')
25 def home():
26     return render_template('home.html')
27
28 app.route('/predict')
29 def predict():
30     return render_template('predict.html')
31
32 app.route('/submit', methods=['POST'])
33 def pred():
34     if request.method == 'POST':
35         # Extract the values entered by the user
36         features = [float(x) for x in request.form.values()]
37         # Reshape the data to match the model's input shape
38         data = np.array(features).reshape(1,-1)
39         # Get productivity prediction using the predictive system
40         productivity_level = predict_productivity(data)
41         # Render the submit.html page with the predicted productivity level
42         return render_template('submit.html', productivity_level=productivity_level)
43
44 @app.route('/about')
45 def about():
46     return render_template('about.html')
47
48 if __name__ == '__main__':
49     # Run the flask application
50     app.run(debug=True)
```



```
File Edit Selection View Go Run Terminal Help ⌘ → ⌘ templates
Home.html about.html predict.html submit.html Flask Code For Application Building.py X
44 # Define route for about page
45 app.route('/about')
46 def about():
47     return render_template('about.html')
48
49 if __name__ == '__main__':
50     # Run the flask application
51     app.run(debug=True)
```

11.GitHub & Project Link

GitHub Link:

<https://github.com/akashbommisetty/Machine-Learning-Approach-For-Employee-Performance-Prediction>

Project demo Link:

https://drive.google.com/file/d/1GXMpo5YJTxpMnLSwNZwl4JNQuyKsQxK5/view?usp=drive_link