

Data Collection and Preprocessing Phase

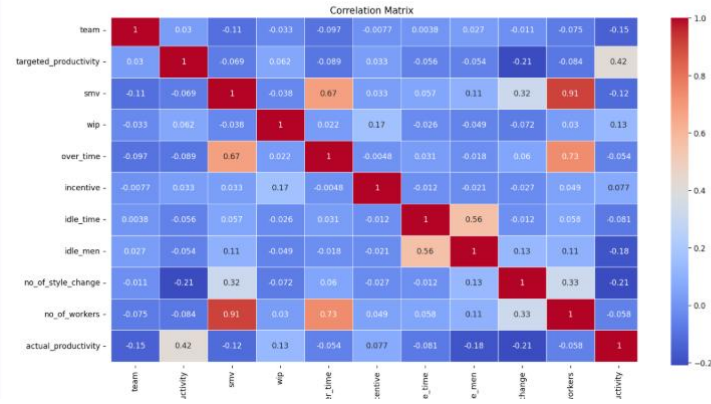
| | |
|---------------|--|
| Date | 26 June 2025 |
| Team ID | NONE |
| Project Title | Employee Performance Prediction using Machine Learning |

Data Exploration and Preprocessing Report

The dataset variables are analyzed to identify patterns and relationships affecting employee productivity. Using Python, preprocessing tasks such as handling missing values, label encoding categorical features, and exploratory data analysis (EDA) are performed. Data cleaning ensures high-quality input for machine learning models, while correlation analysis and visualization help uncover key drivers of performance. This phase forms a solid foundation for accurate productivity prediction using algorithms like Random Forest and XGBoost

| Section | Description |
|---------------------|--|
| Data Overview | <p><u>Dimension:</u> 1197 rows × 15 columns</p> <p><u>Descriptive statistics:</u></p> <pre>import pandas as pd df = pd.read_csv('data.csv') df.head() date quarter department ... no_of_style_change no_of_workers actual_productivity 0 1/1/2015 Quarter1 sweing ... 0 59.0 0.940725 1 1/1/2015 Quarter1 finishing ... 0 8.0 0.886500 2 1/1/2015 Quarter1 sweing ... 0 30.5 0.800570 3 1/1/2015 Quarter1 sweing ... 0 30.5 0.800570 4 1/1/2015 Quarter1 sweing ... 0 56.0 0.800382</pre> |
| Univariate Analysis | <p>Univariate analysis was performed to understand the distribution of individual features in the dataset. Key observations include:</p> <p>actual_productivity (Target Variable): The distribution is approximately normal with a slight right skew, centered around 0.77, indicating most employees perform close to average.</p> <p>over_time: Highly right-skewed — many employees work moderate overtime, while a few work significantly longer hours.</p> <p>idle_time: Most employees have low idle time, with peaks at 0 and 60 minutes, suggesting scheduled breaks or downtime.</p> <p>department: "sweing" department has the highest employee count, followed by "finishing".</p> <p>quarter: Q2 and Q4 show higher activity, likely due to seasonal demand.</p> |

These insights were visualized using histograms and bar plots in visualization.py to guide preprocessing and feature engineering decisions.



Outliers and Anomalies

Data Preprocessing Code Screenshots

Loading Data

```
(base) PS C:\Users\vansh\OneDrive\Desktop\employee_performance_ml> 8
Linear Regression R2 Score: 0.1681682566306545
Random Forest R2 Score: 0.44671974539154335
XGBoost R2 Score: 0.3538597397101051
Best model saved to model/best_model.pkl
(base) PS C:\Users\vansh\OneDrive\Desktop\employee_performance_ml>
```

Feature Engineering

Employee Productivity Prediction

| |
|-----------------------|
| Team No. |
| Targeted Productivity |
| SMV |
| Work In Progress |
| Over Time (in mins) |
| Incentive |
| Idle Time |
| Idle Men |
| No. of Style Change |
| No. of Workers |
| Q1 |
| Finishing |

Predict