

Data Collection and Preprocessing Phase

Date	19 July 2025
Project Title	Machine Learning Approach for Employee Performance Prediction
Maximum Marks	

Data Exploration and Preprocessing Report

Dataset variables will be statistically analyzed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong foundation for insights and predictions.

Section	Description																																																														
Data Overview	<u>Dimension:</u> 614 rows × 13 columns <u>Descriptive statistics:</u>																																																														
	<table><tr><th>targeted_productivity</th><th>smv</th><th>over_time</th><th>incentive</th><th>no_of_style_change</th><th>no_of_workers</th><th>actual_productivity</th></tr><tr><td>1197.000000</td><td>1197.000000</td><td>1197.000000</td><td>1197.000000</td><td>1197.000000</td><td>1197.000000</td><td>1197.000000</td></tr><tr><td>0.729632</td><td>15.062172</td><td>4567.460317</td><td>38.210526</td><td>0.150376</td><td>34.609858</td><td>0.735091</td></tr><tr><td>0.097891</td><td>10.943219</td><td>3348.823563</td><td>160.182643</td><td>0.427848</td><td>22.197687</td><td>0.174488</td></tr><tr><td>0.070000</td><td>2.900000</td><td>0.000000</td><td>0.000000</td><td>0.000000</td><td>2.000000</td><td>0.233705</td></tr><tr><td>0.700000</td><td>3.940000</td><td>1440.000000</td><td>0.000000</td><td>0.000000</td><td>9.000000</td><td>0.650307</td></tr><tr><td>0.750000</td><td>15.260000</td><td>3960.000000</td><td>0.000000</td><td>0.000000</td><td>34.000000</td><td>0.773333</td></tr><tr><td>0.800000</td><td>24.260000</td><td>6960.000000</td><td>50.000000</td><td>0.000000</td><td>57.000000</td><td>0.850253</td></tr><tr><td>0.800000</td><td>54.560000</td><td>25920.000000</td><td>3600.000000</td><td>2.000000</td><td>89.000000</td><td>1.120437</td></tr></table>	targeted_productivity	smv	over_time	incentive	no_of_style_change	no_of_workers	actual_productivity	1197.000000	1197.000000	1197.000000	1197.000000	1197.000000	1197.000000	1197.000000	0.729632	15.062172	4567.460317	38.210526	0.150376	34.609858	0.735091	0.097891	10.943219	3348.823563	160.182643	0.427848	22.197687	0.174488	0.070000	2.900000	0.000000	0.000000	0.000000	2.000000	0.233705	0.700000	3.940000	1440.000000	0.000000	0.000000	9.000000	0.650307	0.750000	15.260000	3960.000000	0.000000	0.000000	34.000000	0.773333	0.800000	24.260000	6960.000000	50.000000	0.000000	57.000000	0.850253	0.800000	54.560000	25920.000000	3600.000000	2.000000	89.000000
targeted_productivity	smv	over_time	incentive	no_of_style_change	no_of_workers	actual_productivity																																																									
1197.000000	1197.000000	1197.000000	1197.000000	1197.000000	1197.000000	1197.000000																																																									
0.729632	15.062172	4567.460317	38.210526	0.150376	34.609858	0.735091																																																									
0.097891	10.943219	3348.823563	160.182643	0.427848	22.197687	0.174488																																																									
0.070000	2.900000	0.000000	0.000000	0.000000	2.000000	0.233705																																																									
0.700000	3.940000	1440.000000	0.000000	0.000000	9.000000	0.650307																																																									
0.750000	15.260000	3960.000000	0.000000	0.000000	34.000000	0.773333																																																									
0.800000	24.260000	6960.000000	50.000000	0.000000	57.000000	0.850253																																																									
0.800000	54.560000	25920.000000	3600.000000	2.000000	89.000000	1.120437																																																									
Data Preprocessing Code Screenshots																																																															

Loading Data	<pre>df = pd.read_csv("garments_worker_productivity.csv") df</pre> <table><thead><tr><th></th><th>date</th><th>quarter</th><th>department</th><th>day</th><th>team</th><th>targeted_productivity</th><th>smv</th><th>wip</th><th>over_time</th></tr></thead><tbody><tr><td>0</td><td>1/1/2015</td><td>Quarter1</td><td>sweing</td><td>Thursday</td><td>8</td><td>0.80</td><td>26.16</td><td>1108.0</td><td>7080</td></tr><tr><td>1</td><td>1/1/2015</td><td>Quarter1</td><td>finishing</td><td>Thursday</td><td>1</td><td>0.75</td><td>3.94</td><td>NaN</td><td>960</td></tr><tr><td>2</td><td>1/1/2015</td><td>Quarter1</td><td>sweing</td><td>Thursday</td><td>11</td><td>0.80</td><td>11.41</td><td>968.0</td><td>3660</td></tr><tr><td>3</td><td>1/1/2015</td><td>Quarter1</td><td>sweing</td><td>Thursday</td><td>12</td><td>0.80</td><td>11.41</td><td>968.0</td><td>3660</td></tr><tr><td>4</td><td>1/1/2015</td><td>Quarter1</td><td>sweing</td><td>Thursday</td><td>6</td><td>0.80</td><td>25.90</td><td>1170.0</td><td>1920</td></tr></tbody></table>		date	quarter	department	day	team	targeted_productivity	smv	wip	over_time	0	1/1/2015	Quarter1	sweing	Thursday	8	0.80	26.16	1108.0	7080	1	1/1/2015	Quarter1	finishing	Thursday	1	0.75	3.94	NaN	960	2	1/1/2015	Quarter1	sweing	Thursday	11	0.80	11.41	968.0	3660	3	1/1/2015	Quarter1	sweing	Thursday	12	0.80	11.41	968.0	3660	4	1/1/2015	Quarter1	sweing	Thursday	6	0.80	25.90	1170.0	1920
	date	quarter	department	day	team	targeted_productivity	smv	wip	over_time																																																				
0	1/1/2015	Quarter1	sweing	Thursday	8	0.80	26.16	1108.0	7080																																																				
1	1/1/2015	Quarter1	finishing	Thursday	1	0.75	3.94	NaN	960																																																				
2	1/1/2015	Quarter1	sweing	Thursday	11	0.80	11.41	968.0	3660																																																				
3	1/1/2015	Quarter1	sweing	Thursday	12	0.80	11.41	968.0	3660																																																				
4	1/1/2015	Quarter1	sweing	Thursday	6	0.80	25.90	1170.0	1920																																																				
Handling Missing Value	<pre>df.isnull().sum()</pre> <table><tbody><tr><td></td><td>0</td></tr><tr><td>date</td><td>0</td></tr><tr><td>quarter</td><td>0</td></tr><tr><td>department</td><td>0</td></tr><tr><td>day</td><td>0</td></tr><tr><td>team</td><td>0</td></tr><tr><td>targeted_productivity</td><td>0</td></tr><tr><td>smv</td><td>0</td></tr><tr><td>wip</td><td>506</td></tr></tbody></table> <pre>df.drop('wip',axis=1, inplace=True)</pre>		0	date	0	quarter	0	department	0	day	0	team	0	targeted_productivity	0	smv	0	wip	506																																										
	0																																																												
date	0																																																												
quarter	0																																																												
department	0																																																												
day	0																																																												
team	0																																																												
targeted_productivity	0																																																												
smv	0																																																												
wip	506																																																												
Feature Engineering	<pre># handling date column df['date'] = pd.to_datetime(df['date']) df['month'] = df['date'].dt.month df.drop('date',axis=1,inplace=True) #handling department column category_mapping = { 'sweing':'sewing', 'finishing':'finishing' } df['department'] = df['department'].str.strip().map(category_mapping)</pre> <table><thead><tr><th></th><th>count</th></tr></thead><tbody><tr><td>department</td><td></td></tr><tr><td>sewing</td><td>691</td></tr><tr><td>finishing</td><td>506</td></tr></tbody></table>		count	department		sewing	691	finishing	506																																																				
	count																																																												
department																																																													
sewing	691																																																												
finishing	506																																																												
Label Encoding	<h3>3. Label Encoding for categorical columns</h3> <pre>le = MultiColumnLabelEncoder.MultiColumnLabelEncoder() df = le.fit_transform(df)</pre>																																																												
Save Processed Data	-																																																												

Correlaton Analysis

