# Data Collection and Preprocessing Phase

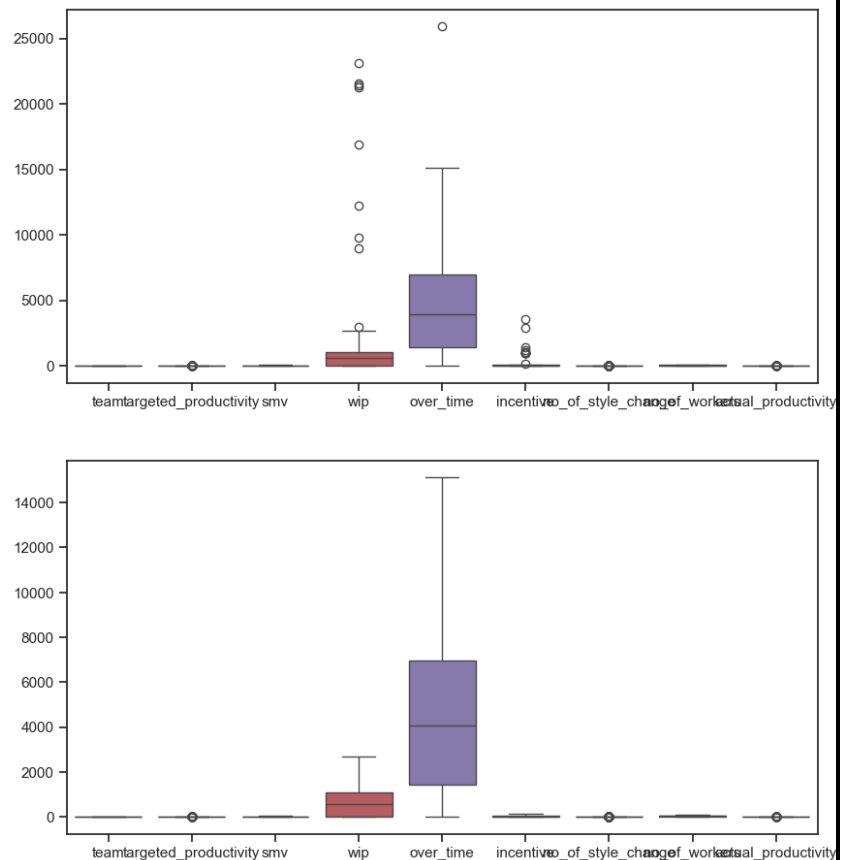| Date | 15 March 2024 |
|---|---|
| Team ID | xxxxxx |
| Project Title | xxxxxx |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing Template**

Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

| Section | Description |
|---|---|
| Data Overview |  |
| Univariate Analysis |  |
| Bivariate Analysis | |

Correlation Heatmap


Scatter Matrix

| Outliers and Anomalies |  |
| --- | --- |

## Data Preprocessing Code Screenshots

| Loading Data |  |
| --- | --- |

```
df = pd.read_csv("productivity.csv")
df.head()
✓ 0.0s                                                                                    Pytho
```

| quarter | department | day | team | targeted_productivity | smv | wip | over_time | incentive | idle_time | idle_men | no_of_style_change | no_of_workers | actual_productivity |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Quarter1 | sweing | Thursday | 8 | 0.80 | 26.16 | 1108.0 | 7080 | 98 | 0.0 | 0 | 0 | 59.0 | 0.94072! |
| Quarter1 | finishing | Thursday | 1 | 0.75 | 3.94 | NaN | 960 | 0 | 0.0 | 0 | 0 | 8.0 | 0.886500 |
| Quarter1 | sweing | Thursday | 11 | 0.80 | 11.41 | 968.0 | 3660 | 50 | 0.0 | 0 | 0 | 30.5 | 0.80057! |
| Quarter1 | sweing | Thursday | 12 | 0.80 | 11.41 | 968.0 | 3660 | 50 | 0.0 | 0 | 0 | 30.5 | 0.80057! |
| Quarter1 | sweing | Thursday | 6 | 0.80 | 25.90 | 1170.0 | 1920 | 50 | 0.0 | 0 | 0 | 56.0 | 0.80038. |

| Handling Missing Data | |
| --- | --- |

```
df.isnull().sum()
✓ 0.0s

quarter                    0
day                        0
team                       0
targeted_productivity      0
smv                        0
wip                      506
over_time                  0
incentive                  0
no_of_style_change         0
no_of_workers              0
actual_productivity        0
dtype: int64
```

```
    df['wip'].fillna(0,inplace=True)
    df.isnull().sum()
✓ 0.0s

C:\Users\hp\AppData\Local\Temp\ipykernel_17788\818197784.py:1
The behavior will change in pandas 3.0. This inplace method wi

For example, when doing 'df[col].method(value, inplace=True)'

    df['wip'].fillna(0,inplace=True)
```

```
quarter                   0
day                       0
team                      0
targeted_productivity     0
smv                       0
wip                       0
over_time                 0
incentive                 0
no_of_style_change        0
no_of_workers             0
actual_productivity       0
dtype: int64
```

## Data Transformation

```
df_encoded = pd.get_dummies(df, columns=["quarter", "day"])
features= ['targeted_productivity', 'smv', 'wip', 'over_time', 'incentive', 'no_of_style_change', 'no_of_workers', 'actual_productivity']
scaler = MinMaxScaler()
df_encoded[features] = scaler.fit_transform(df_encoded[features])
print(df_encoded.head())
✓ 0.0s

   team  targeted_productivity       smv       wip  over_time  incentive  \
0     8               1.000000  0.450252  0.410675   0.468254   0.823529
1     1               0.931507  0.020132  0.000000   0.063492   0.000000
2    11               1.000000  0.164731  0.358784   0.242063   0.420168
3    12               1.000000  0.164731  0.358784   0.242063   0.420168
4     6               1.000000  0.445219  0.433655   0.126984   0.420168

   no_of_style_change  no_of_workers  actual_productivity  quarter_Quarter1  \
0                 0.0       0.655172             0.797332              True
1                 0.0       0.068966             0.736180              True
2                 0.0       0.327586             0.639274              True
3                 0.0       0.327586             0.639274              True
4                 0.0       0.620690             0.639062              True

   quarter_Quarter2  quarter_Quarter3  quarter_Quarter4  quarter_Quarter5  \
0             False             False             False             False
1             False             False             False             False
2             False             False             False             False
3             False             False             False             False
4             False             False             False             False

   day_Monday  day_Saturday  day_Sunday  day_Thursday  day_Tuesday  \
0       False         False       False          True        False
1       False         False       False          True        False
2       False         False       False          True        False
...
1             False
2             False
3             False
```

| Save Processed Data | ```
df_encoded.to_csv('processed_productivity.csv', index=False)
✓  0.0s
``` |