

Data Collection and Preprocessing Phase

Date: 31 July 2025

SkillWallet ID: SWUID20240141492

Project Title: Employee Performance Prediction using Machine Learning

Maximum Marks: 2 Marks

Data Quality Report:

The Data Quality Report highlights issues in the dataset sourced from Kaggle (Garment Worker Productivity Dataset). It includes identified issues, their severity, and proposed resolution plans to ensure data consistency and quality for model development.

| Data Source | Data Quality Issue | Severity | Resolution Plan |
|----------------|---|----------|---|
| Kaggle Dataset | Missing values in <code>idle_time</code> , <code>idle_men</code> , and <code>incentive</code> columns | Moderate | Use mean/median imputation for numerical fields. |
| Kaggle Dataset | Presence of categorical variables such as <code>department</code> , <code>day</code> , and <code>quarter</code> | Moderate | Apply Label Encoding or One-Hot Encoding . |
| Kaggle Dataset | Outliers in <code>over_time</code> and <code>no_of_style_change</code> columns | High | Apply IQR (Interquartile Range) method to handle outliers. |
| Kaggle Dataset | Variations in <code>targeted_productivity</code> values due to scaling issues | Low | Normalize or standardize the feature. |
| Kaggle Dataset | Mixed data formats in <code>month</code> column (text and numbers) | Low | Convert all month values to a uniform numeric format . |