

Linear Regression Analysis Report

Garments Worker Productivity Prediction

Student: Vaishnavi Pasumarthi, Akshay Nagarajan
NetID: VXP210093, AHN210000
Due Date: September 21, 2025

Dataset and Methodology

This analysis uses the Garments Worker Productivity dataset from the UCI Machine Learning Repository containing 1197 instances (some null values) with 14 feature variables. The target variable `actual_productivity` represents worker performance in garment manufacturing. The Standard Deviation for this feature is 0.174488.

Data preprocessing included handling missing values, extracting date-based features (month and weekday), creating ratio and per-worker features, standardizing numeric features, generating polynomial interaction terms for selected variables, and one-hot encoding categorical variables, all combined into a unified pipeline for train/test data preparation.

Model Development

Model 1: SGD Regression

Hyperparameter tuning was performed using 3-fold cross-validation to optimize the `SGDRegressor` model:

Parameter	Optimal Value	Parameter	Optimal Value
Loss	Huber	Learning Rate	Optimal
Penalty	L2	Eta ₀	0.001
Alpha	0.001	Max Iterations	8000
L1 Ratio	0.2	Epsilon	0.03

The Huber loss function was selected for robustness to outliers, with L2 regularization to prevent overfitting.

Model 2: OLS Regression

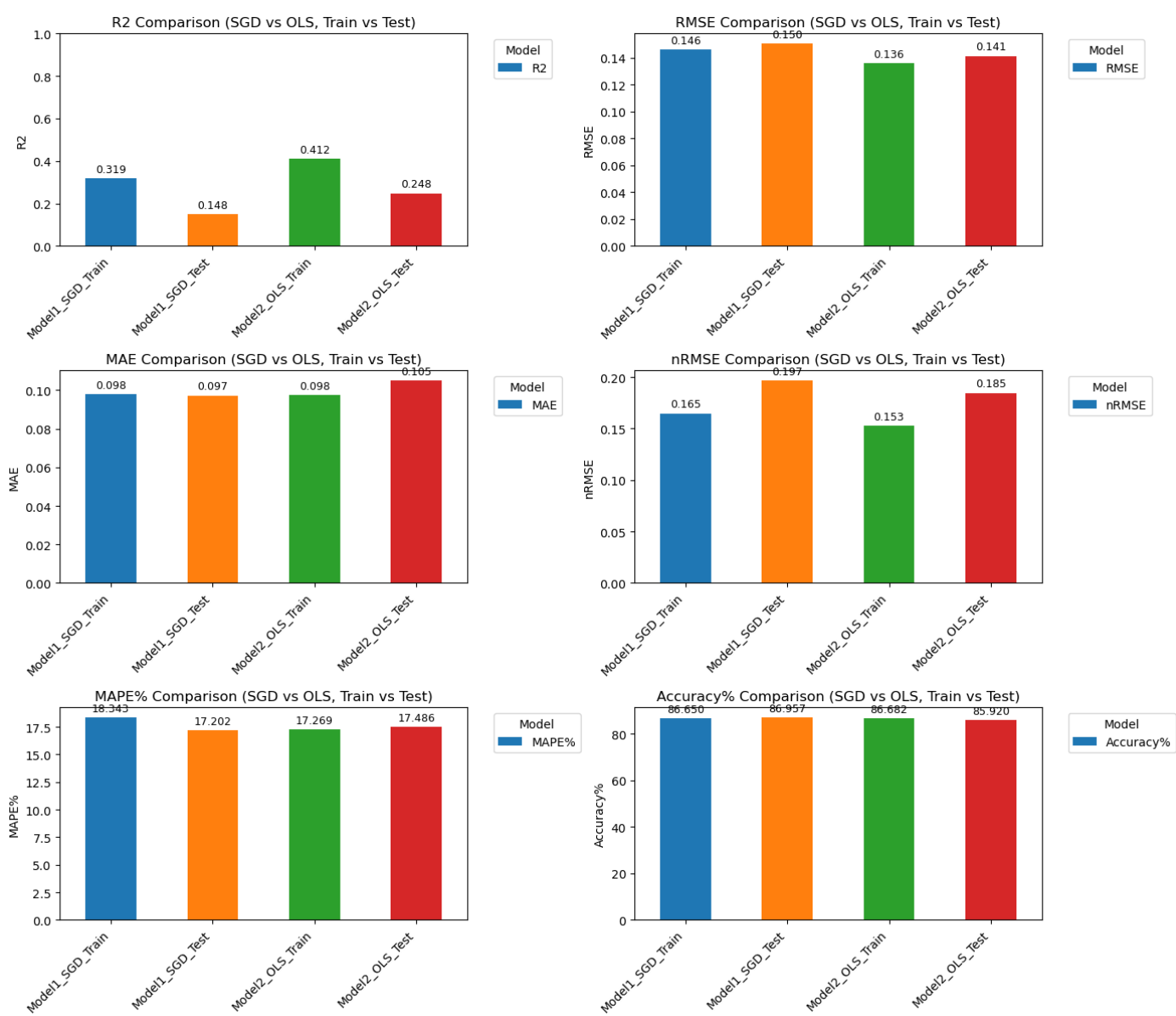
Ordinary Least Squares regression was implemented using `statsmodels` library, fitting a linear model with all 46 predictors plus intercept term. Here is the top half of the results table:

Metric	Value	Metric	Value
Dependent Variable	<code>actual_productivity</code>	Model	OLS
Method	Least Squares	Date	Sun, 21 Sep 2025
Time	20:18:12	No. Observations	957
Df Residuals	914	Df Model	42
R ²	0.412	Adj. R ²	0.385
F-statistic	15.28	Prob (F-statistic)	6.07e-79
Log-Likelihood	553.30	AIC	-1021.0
BIC	-811.5	Covariance Type	nonrobust

Results Analysis

Performance Comparison

Metric	SGD Train	SGD Test	OLS Train	OLS Test
R ²	0.319	0.148	0.412	0.248
RMSE	0.146	0.150	0.136	0.141
MAE	0.098	0.097	0.098	0.105
nRMSE	0.165	0.197	0.153	0.185
MAPE%	18.34	17.20	17.27	17.49
Accuracy%	86.65	86.96	86.68	85.92



Model Performance Analysis

R² Analysis: The OLS model shows superior performance with test R² of 0.248 compared to SGD's 0.148, indicating OLS explains 67% more variance in the test data. Both models exhibit overfitting, with training R² substantially higher than test R² (OLS: 0.412 vs 0.248; SGD: 0.319 vs 0.148).

Error Metrics: OLS demonstrates lower RMSE values (0.141 test) compared to SGD (0.150 test), suggesting better overall prediction accuracy. MAE values are similar between models (~0.10), indicating comparable absolute error magnitudes. The nRMSE values show OLS has better normalized error performance.

Percentage Metrics: MAPE values are consistent across models (17-18%), while accuracy percentages remain high (85-87%) for both approaches, suggesting reasonable predictive capability despite moderate R² values.

Statistical Analysis (OLS Model)

Overall Model Significance:

- **F-statistic:** 15.28 with p-value = 6.07e-79 (highly significant)
- **R²:** 0.412 (explains 41.2% of training variance)
- **Adjusted R²:** 0.385 (accounts for number of predictors)

Model Comparison and Interpretation

SGD vs OLS Performance: The OLS model consistently outperforms SGD across most metrics, particularly in R² values. This suggests that the linear relationships in the data are better captured through direct least squares estimation rather than gradient-based optimization with the chosen hyperparameters.

Overfitting Assessment: Both models show evidence of overfitting, with substantial drops from training to test performance. The OLS model maintains better generalization despite having more parameters.

Conclusions

The comparative analysis demonstrates that OLS regression outperforms SGD regression for this dataset, achieving superior predictive accuracy and lower error rates. However, both models face limitations including overfitting tendencies and moderate overall explanatory power (best test R² = 0.248). This is primarily low Standard Deviation, making the model estimates close to the overall mean. Due to this, both models are very limited by the variability of the data. Both models achieve reasonable practical accuracy (85-87%) despite the unexplained variability in worker productivity that may require additional features or non-linear modeling approaches.