

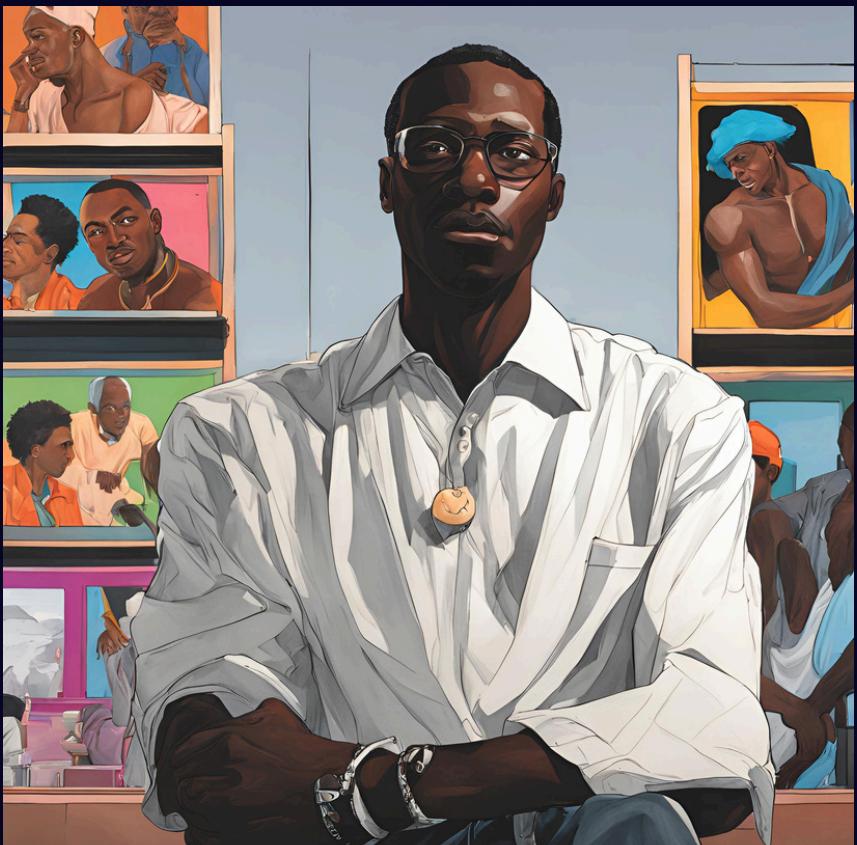


PROJECT_2

Regression analysis on
Garment Workers productivity
dataset

Tomorrow

MEET OUR TEAM



Jaiden Flint



Tuan Huynh



Kavita Gopal



Why is Productivity so important to the manufacturing industry?

Productivity reduces per unit manufacturing cost by effectively utilising all resources and reducing waste. Lower per-unit cost means an increase in profit level.

An economy is able to produce—and consume—increasingly more goods and services for the same amount of work.

Productivity is important to workers, consumers, business leaders, and analysts involved in policy making and government statisticians.

Project Overview

This project aims to predict the productivity of garment employees based on features such as working hours, incentives, overtime, and more. The goal is to develop machine learning models that can accurately predict the continuous productivity score for garment workers, helping factories understand key drivers of productivity and optimize their operations.

OUR PROJECT

We used the UCI's ML repository to get our dataset.

The screenshot shows the UCI Machine Learning Repository page for the 'Productivity Prediction of Garment Employees' dataset. The page includes the dataset title, a brief description, dataset characteristics, feature types, and detailed information sections like 'Dataset Information' and 'Additional Information'. On the right side, there are download links (CSV, ZIP, and Python), citation statistics (1 citation, 27941 views), DOI (10.24432/C51S6D), and a license section (Creative Commons Attribution 4.0 International). The URL in the browser bar is archive.ics.uci.edu/dataset/597/productivity+prediction+of+garment+employees.

Productivity Prediction of Garment Employees
Donated on 8/2/2020

This dataset includes important attributes of the garment manufacturing process and the productivity of the employees which had been collected manually and also been validated by the industry experts.

Dataset Characteristics	Subject Area	Associated Tasks
Multivariate, Time-Series	Business	Classification, Regression
Feature Type	# Instances	# Features
Integer, Real	1197	14

Dataset Information

Additional Information

The Garment Industry is one of the key examples of the industrial globalization of this modern era. It is a highly labour-intensive industry with lots of manual processes. Satisfying the huge global demand for garment products is mostly dependent on the production and delivery performance of the employees in the garment manufacturing companies. So, it is highly desirable among the decision makers in the garments industry to track, analyse and predict the productivity performance of the working teams in their factories. This dataset can be used for regression purpose by predicting the productivity range (0-1) or for classification purpose by transforming the productivity range (0-1) into different classes.

SHOW LESS ^

Has Missing Values?
Yes

DOWNLOAD (92.9 KB)
IMPORT IN PYTHON
CITE
1 citations
27941 views
DOI
10.24432/C51S6D
License
This dataset is licensed under a Creative Commons Attribution 4.0 International (CC BY 4.0) license.
This allows for the sharing and adaptation of the datasets for any purpose, provided that the appropriate credit is given.

Dataset

- dataset Garment Worker Productivity
 - consists of 1,197 rows and 15 features
- date
 - quarter
 - department
 - day
 - team
 - targeted_productivity
 - smv (Standard Minute Value)
 - wip (Work In Progress)
 - over_time
 - incentive
 - idle_time
 - idle_men
 - no_of_style_change
 - no_of_workers
 - actual_productivity
- target variable: actual_productivity
 - continuous primary variable
 - we aim to predict using various machine learning models.

OUR PROJECT

ES

Explortory Data Analysis (EDA)

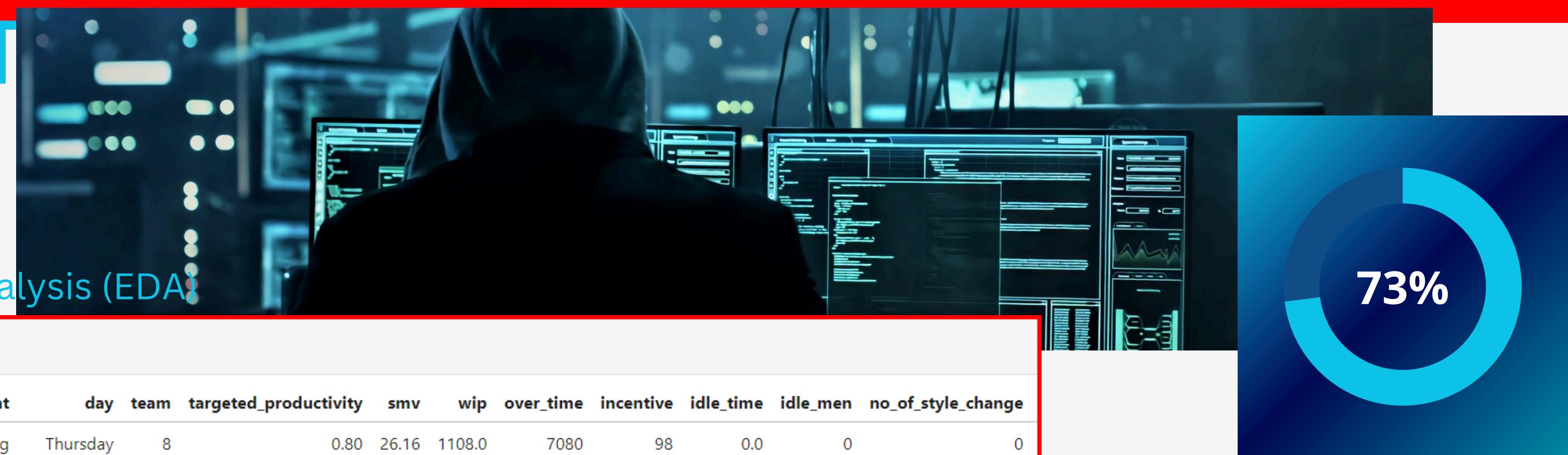
```
# Display the updated DataFrame  
updated_grpr_df
```

	date	quarter	department	day	team	targeted_productivity	smv	wip	over_time	incentive	idle_time	idle_men	no_of_style_change
0	1/1/2015	Quarter1	sewing	Thursday	8	0.80	26.16	11080.0	7080	98	0.0	0	0
1	1/1/2015	Quarter1	finishing	Thursday	1	0.75	3.94	NaN	960	0	0.0	0	0
2	1/1/2015	Quarter1	sewing	Thursday	11	0.80	11.41	9680.0	3660	50	0.0	0	0
3	1/1/2015	Quarter1	sewing	Thursday	12	0.80	11.41	9680.0	3660	50	0.0	0	0
4	1/1/2015	Quarter1	sewing	Thursday	6	0.80	25.90	11700.0	1920	50	0.0	0	0
...
1192	3/11/2015	Quarter2	finishing	Wednesday	10	0.75	2.90	NaN	960	0	0.0	0	0
1193	3/11/2015	Quarter2	finishing	Wednesday	8	0.70	3.90	NaN	960	0	0.0	0	0
1194	3/11/2015	Quarter2	finishing	Wednesday	7	0.65	3.90	NaN	960	0	0.0	0	0
1195	3/11/2015	Quarter2	finishing	Wednesday	9	0.75	2.90	NaN	1800	0	0.0	0	0
1196	3/11/2015	Quarter2	finishing	Wednesday	6	0.70	2.90	NaN	720	0	0.0	0	0

1197 rows × 15 columns

```
# Print the list of columns to inspect the actual column names  
print(cols)
```

```
['date', 'quarter', 'department', 'day', 'team', 'targeted_productivity', 'smv', 'wip', 'over_time', 'incentive', 'idle_time', 'idle_men', 'no  
of style change', 'no of workers', 'actual productivity', 'target achieved']
```



OUR PROJECT

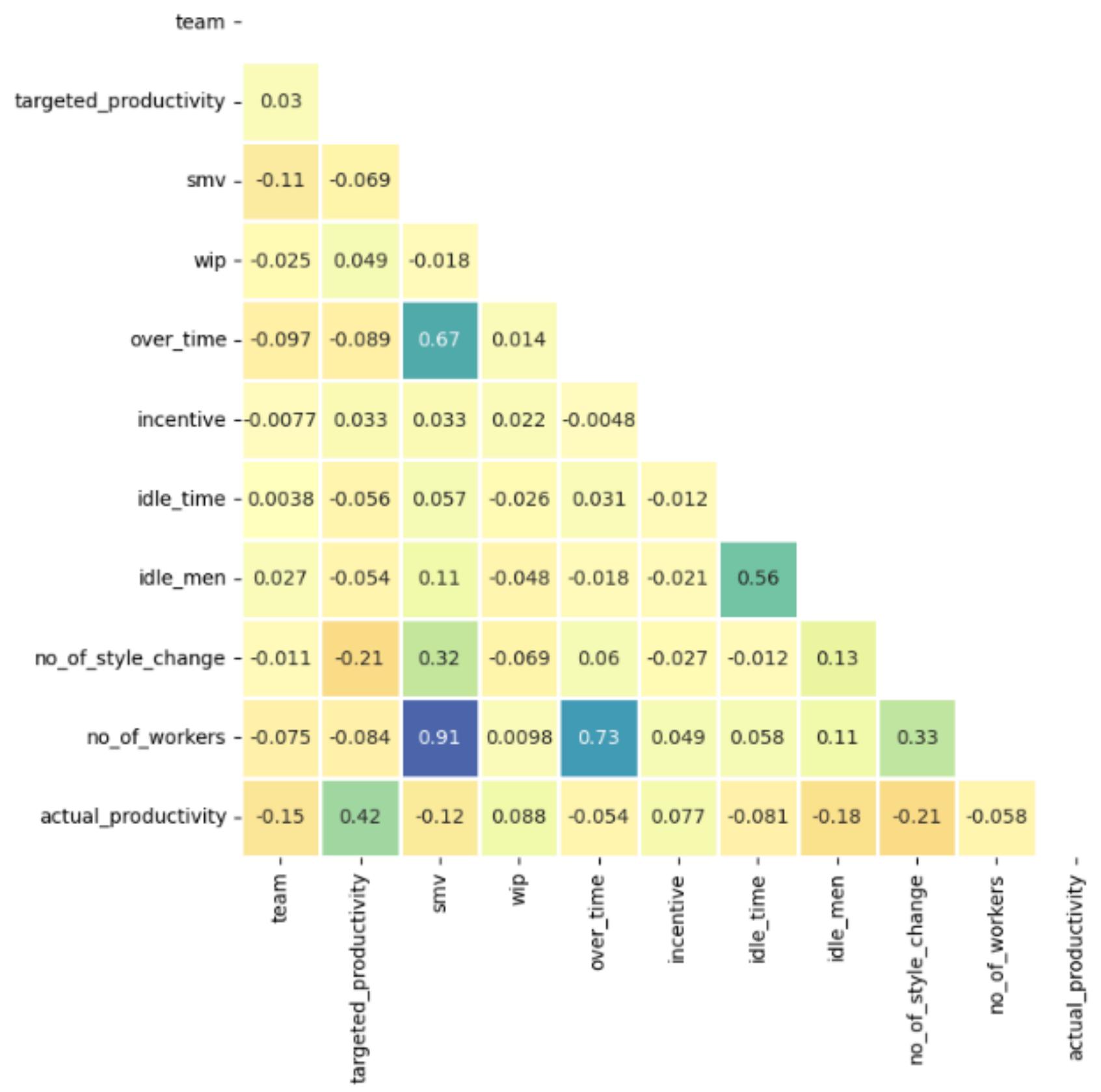
OUR VISION

```
# Create a list of teams in the desired order
team_order = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]

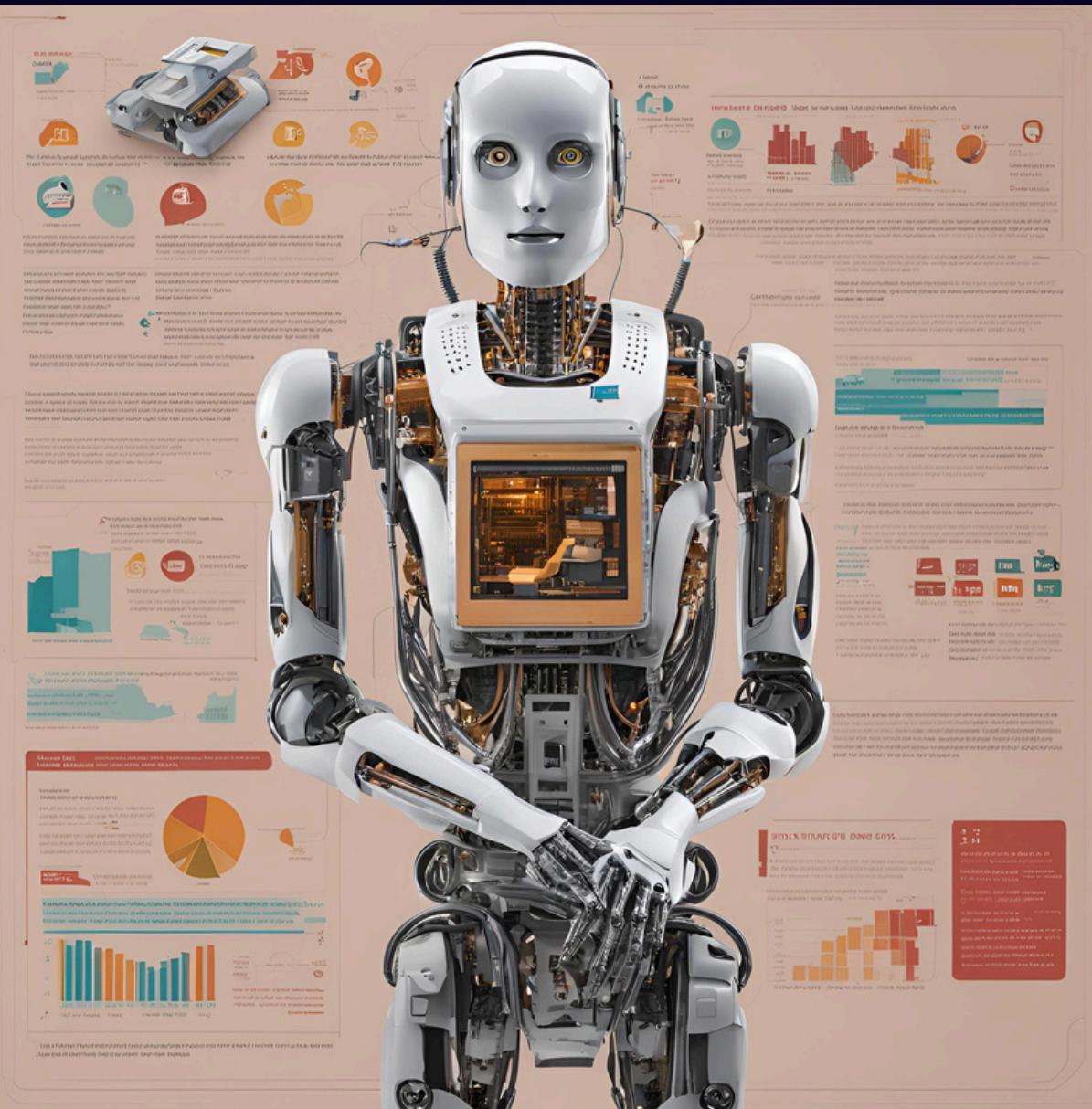
# Create the box plot with the specified order
fig = px.box(updated_grpr_df, x='quarter', y='actual_productivity', color='day',
              facet_col='team', facet_col_wrap=4, category_orders={"team": team_order})
fig.show()
```



OUR PROJECT



Correlation between features in the dataset: more incentives did not correlate to productivity as shown in the plot. No strong correlation when the value is closer to zero. Correlations vary from -1 to +1. Closer to +1 means strong positive correlation whereas closer to -1 means strong negative correlation. The variables with strong correlations are probably the features for model building.



OUR VISION



Data Cleaning and Processing

01

Handling Missing Values

- wip

02

Feature Selection

- actual_productivity, targeted_productivity removed to avoid data leakage
- date col removed: no relevant data for prediction

03

Categorical Feature Encoding

- One-Hot Encoding (scikit-learn): applied to categorical variables (quarter, department, and day).
- Transformation of categorical variables into numerical values suitable for ML.

04

Scaling Features

- Missing values in the features were filled using the mean, and all numeric features were scaled for better model performance.



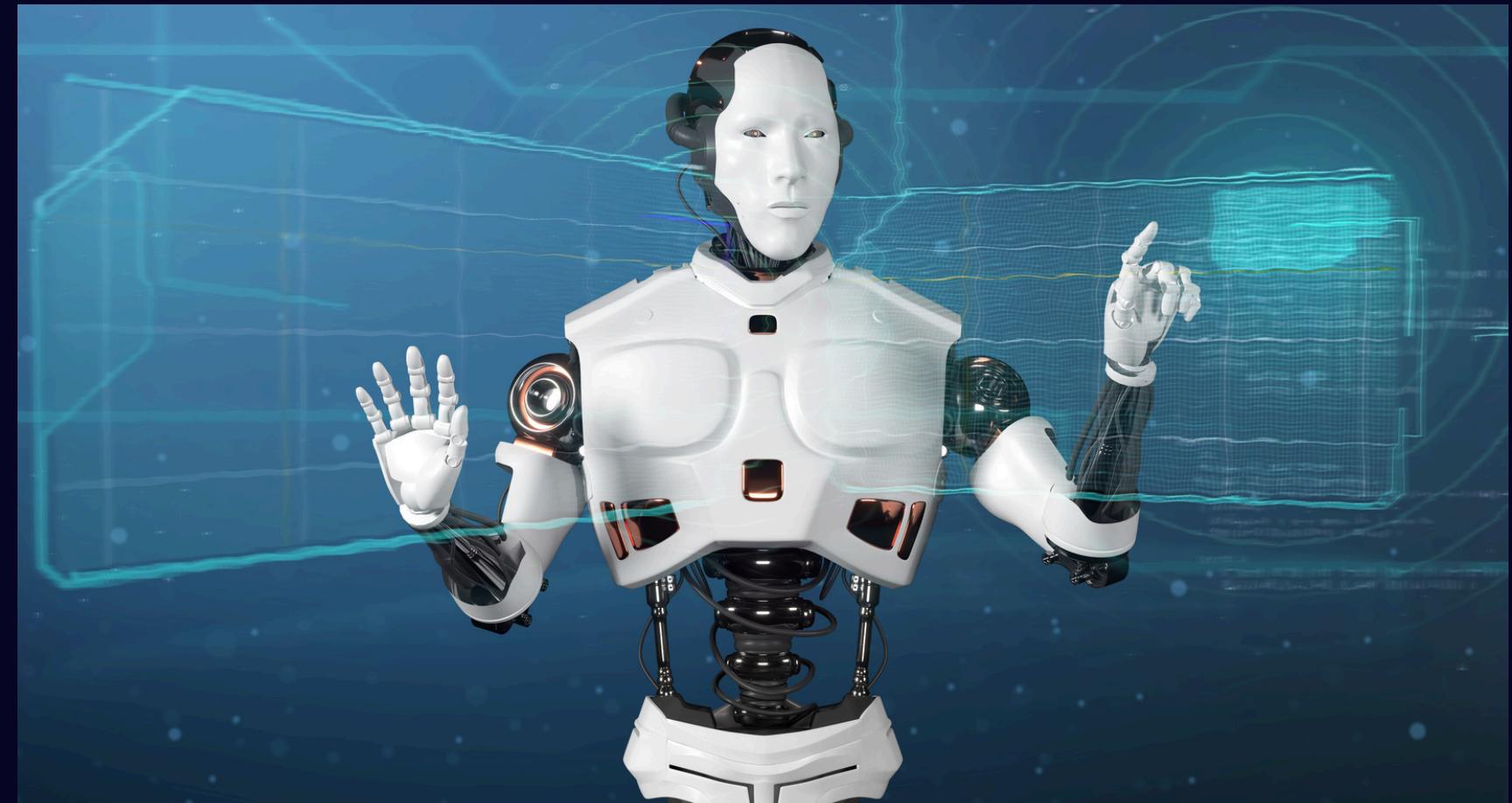
01

- Linear Regression
- Support Vector Regression
- KNeighbors Regression
- Gradient Boost Regression
- Random Forest Regression
- XgBoost Regression
- Decision Tree Regression
- Lasso Regression
- Elastic Net

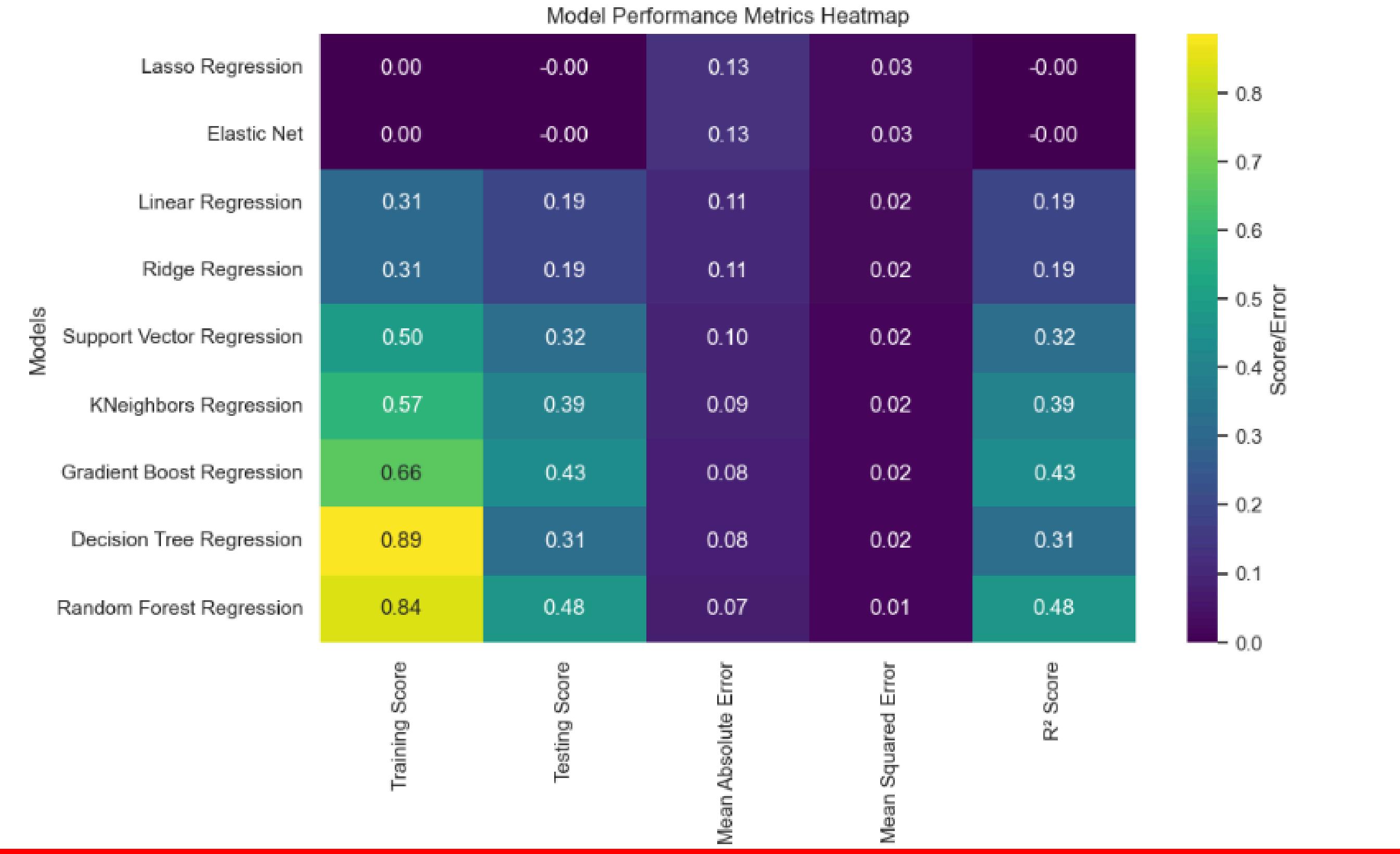
02

Choice of Metric?

- MAE is often preferred when you want a straightforward interpretation of average error.
- MSE is useful when you want to penalize larger errors more heavily.
- R^2 is useful for understanding the explanatory power of the model, but it should not be the sole metric for evaluation.



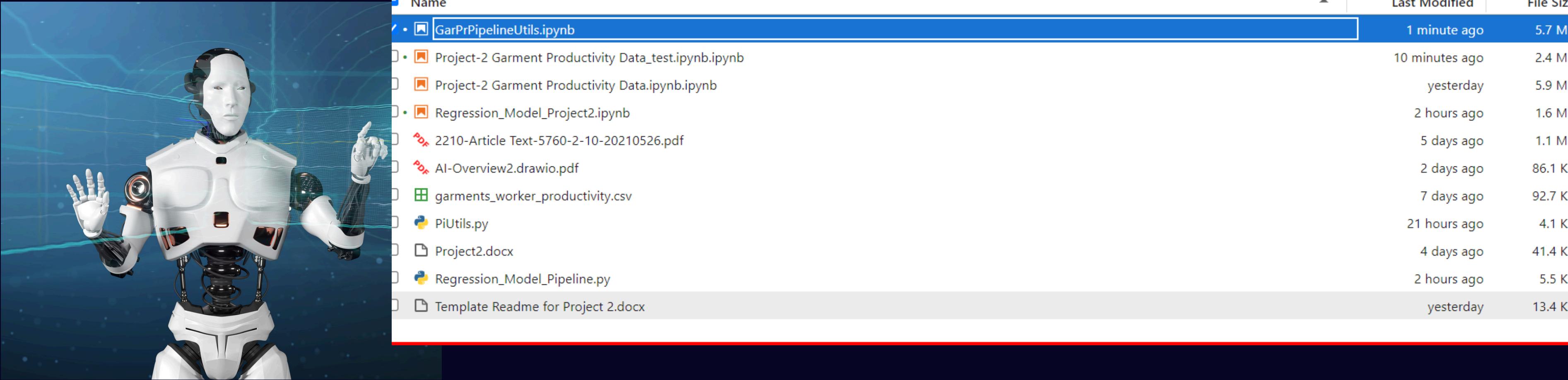
CRASH & BURN



The table is sorted descending based on the MAE which is often preferred when you want a straightforward interpretation of average error.

A lower MAE value indicates better model performance, as it means that the model's predictions are closer to the actual values on average. This is generally desirable in predictive modeling.

A higher MAE value suggests that the model's predictions are further away from the actual values, indicating poorer performance.



We had to start-over and rethink the strategy to clean the data for more meaningful model training!!!

A white humanoid robot stands on the left side of a table. The table has two main sections: 'Original' and 'Fitted'. The 'Original' section contains columns for 'Training Score', 'Testing Score', 'MAE', 'MSE', and 'R² Score'. The 'Fitted' section contains columns for 'Training Score', 'Testing Score', 'MAE', 'MSE', and 'R² Score'. Both sections have rows for various regression models: Lasso Regression, Elastic Net, Linear Regression, Ridge Regression, Support Vector Regression, KNeighbors Regression, Gradient Boost Regression, Decision Tree Regression, and Random Forest Regression. The entire table is enclosed in a red border.

	Original				Fitted					
	Training Score	Testing Score	MAE	MSE	R ² Score	Training Score	Testing Score	MAE	MSE	R ² Score
Lasso Regression	0.000000	-0.000685	0.128256	0.028178	-0.000685	0.254333	0.170077	0.109296	0.023370	0.170077
Elastic Net	0.000000	-0.000685	0.128256	0.028178	-0.000685	0.303058	0.193400	0.108250	0.022713	0.193400
Linear Regression	0.306007	0.188046	0.109552	0.022864	0.188046	0.306007	0.188046	0.109552	0.022864	0.188046
Ridge Regression	0.305997	0.188579	0.109464	0.022849	0.188579	0.305275	0.191695	0.108875	0.022761	0.191695
Support Vector Regression	0.496357	0.317175	0.095912	0.019228	0.317175	0.496357	0.317175	0.095912	0.019228	0.317175
KNeighbors Regression	0.574632	0.388269	0.086326	0.017226	0.388269	0.514701	0.323550	0.089639	0.019048	0.323550
Gradient Boost Regression	0.663522	0.434893	0.081950	0.015913	0.434893	0.737412	0.401290	0.082453	0.016859	0.401290
Decision Tree Regression	0.885893	0.314978	0.079652	0.019290	0.314978	0.748720	0.345696	0.082957	0.018425	0.345696
Random Forest Regression	0.840723	0.476603	0.074872	0.014738	0.476603	0.736627	0.454106	0.078817	0.015372	0.454106



- retnrun to the drawing board
- do something different in analysizing the data
- recreate the pipeline!!!
- try hyperparameter testing???



Metrics of the dataset.

First run through of the Regression Models.

Metrics of the dataset.

After Hyper Parameter Tuning.

```
X_train shape: (897, 10)
y_train shape: (897, 1)
X_test shape: (300, 10)
y_test shape: (300, 1)
```

REMEMBER NAN VALUES WERE FILLED WITH MEAN VALUES

	Training Score	Testing Score	Mean Absolute Error	Mean Squared Error	R ² Score
Lasso Regression	0.000000	-0.000685	0.128256	0.028178	-0.000685
Elastic Net	0.000000	-0.000685	0.128256	0.028178	-0.000685
Linear Regression	0.306007	0.188046	0.109552	0.022864	0.188046
Ridge Regression	0.305997	0.188579	0.109464	0.022849	0.188579
Support Vector Regression	0.496357	0.317175	0.095912	0.019228	0.317175
KNeighbors Regression	0.574632	0.388269	0.086326	0.017226	0.388269
Gradient Boost Regression	0.663522	0.434893	0.081950	0.015913	0.434893
Decision Tree Regression	0.885893	0.314978	0.079652	0.019290	0.314978
Random Forest Regression	0.840723	0.476603	0.074872	0.014738	0.476603

```
X_train shape: (897, 10)
y_train shape: (897, 1)
X_test shape: (300, 10)
y_test shape: (300, 1)
```

	Training Score	Testing Score	Mean Absolute Error	Mean Squared Error	R ² Score
Linear Regression	0.306007	0.188046	0.109552	0.022864	0.188046
Lasso Regression	0.254333	0.170077	0.189296	0.023370	0.170077
Ridge Regression	0.305275	0.191695	0.108875	0.022761	0.191695
Elastic Net	0.303058	0.193400	0.108250	0.022713	0.193400
Support Vector Regression	0.496357	0.317175	0.095912	0.019228	0.317175
KNeighbors Regression	0.514701	0.323550	0.089639	0.019048	0.323550
Decision Tree Regression	0.748720	0.345696	0.082957	0.018425	0.345696
Gradient Boost Regression	0.737412	0.401290	0.082453	0.016859	0.401290
Random Forest Regression	0.736627	0.454106	0.078817	0.015372	0.454106

PIPELINE.PY

```
# -----
# Hyper Parameter Tuning Using Grid Search
# -----  
  
def tune_models(pipelines, param_grids, X_train, y_train):  
    results = {}  
    for name, pipeline in pipelines.items():  
        print(f"Tuning {name}...")  
        grid_search = GridSearchCV(estimator=pipeline, param_grid=param_grids[name], cv=5, scoring='neg_mean_squared_error')  
        grid_search.fit(X_train, y_train)  
        results[name] = {  
            'best_params': grid_search.best_params_,  
            'best_score': grid_search.best_score_  
        }  
    return results
```

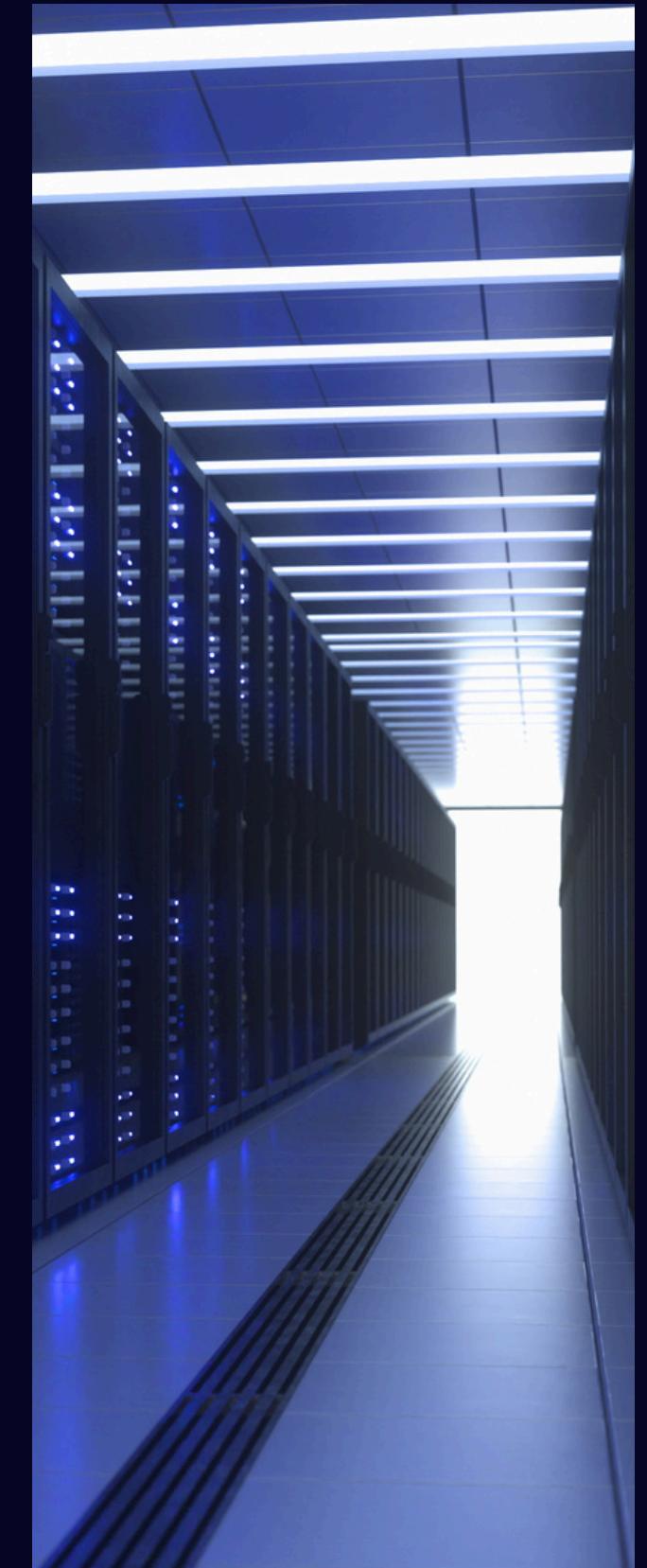
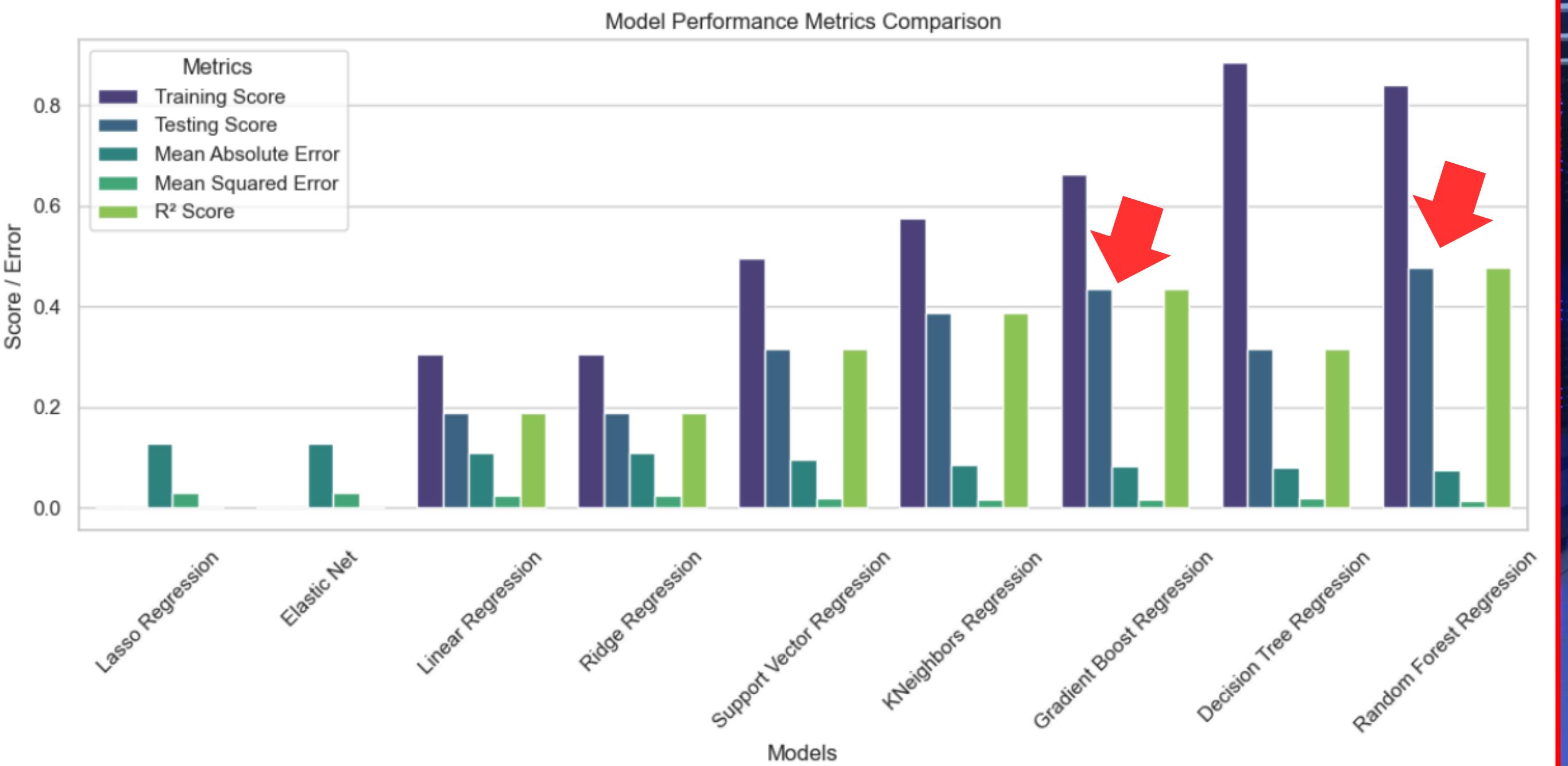
Used Grid Search CV for Hyper Parameter Tuning of the data.
Then re-ran the models with the best parameters.

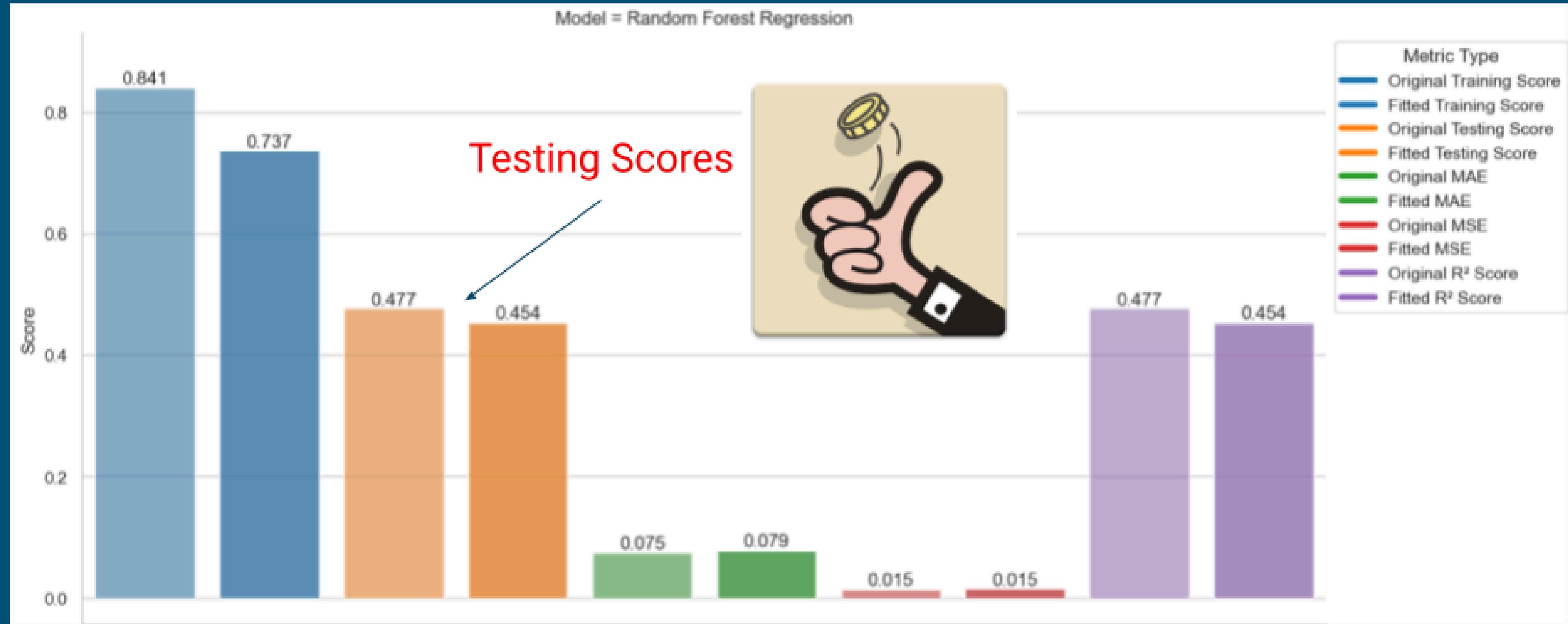
MAIN SCRIPT

```
fitted_models = {}  
  
for model_name, result in tuning_results.items():  
    best_params = result['best_params']  
    fitted_model = create_and_fit_model(model_name, best_params, X_train, y_train, pipelines)  
    fitted_models[model_name] = fitted_model # Store the fitted model  
  
fitted_metrics_df = evaluate_models(fitted_models, X_train, y_train, X_test, y_test)
```

ABOUT FUTURE TECHNOLOGY

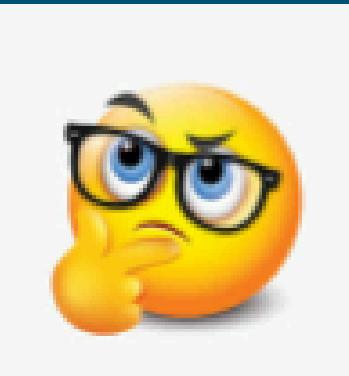
In this dataset the Tree models: Random Forest and Gradient Boost Regression Models had the best results for model training and testing for the first run through.





Scores don't look promising for trying to predict anything after training and testing the data.

EVEN AFTER HYPER
PARAMETER TUNING!!!



Time to rethink what else can we do or fix the dataset?

Feature Selection: Evaluate the features used in the model.

Data Quality: Check for missing values, outliers, or errors in the dataset.

Model Complexity: Consider whether the model is too simple (underfitting) or too complex (overfitting).

Bias-Variance Trade-off: Assess the bias and variance of your model. If your model has high bias, you may need a more complex model. If it has high variance, you may need to simplify the model or use techniques like regularization.

Cross-Validation: Implement cross-validation to ensure that your model's performance is consistent across different subsets of the data. This helps to validate that your model is not just performing well on a specific training set.

Evaluation Metrics: Use appropriate evaluation metrics that align with your goals. For example, if your data is imbalanced, accuracy might not be the best metric; consider using metrics like R-squared, adjusted R-squared, or mean squared error (MSE).

Hyperparameter Tuning: Experiment with different hyperparameters to find the optimal settings for your model.



Metrics of the dataset.

Dataset with NaN values filled with zeros.

Metrics of the dataset.

Dataset with NaN values being dropped.

```
X_train shape: (897, 10)
y_train shape: (897, 1)
X_test shape: (300, 10)
y_test shape: (300, 1)
```

	Training Score	Testing Score	Mean Absolute Error	Mean Squared Error	R ² Score
Lasso Regression	0.000000	-0.000685	0.128256	0.028178	-0.000685
Elastic Net	0.000000	-0.000685	0.128256	0.028178	-0.000685
Linear Regression	0.305653	0.187449	0.109557	0.022881	0.187449
Ridge Regression	0.305644	0.187981	0.109469	0.022866	0.187981
Support Vector Regression	0.491446	0.315339	0.096369	0.019279	0.315339
Decision Tree Regression	0.885893	0.168972	0.087337	0.023401	0.168972
KNeighbors Regression	0.574089	0.389248	0.086383	0.017198	0.389248
Gradient Boost Regression	0.667752	0.426369	0.082126	0.016153	0.426369
Random Forest Regression	0.838553	0.453764	0.076116	0.015382	0.453764

The dataset is modified with the NaN values being filled with 0

```
X_train shape: (518, 10)
y_train shape: (518, 1)
X_test shape: (173, 10)
y_test shape: (173, 1)
```

	Training Score	Testing Score	Mean Absolute Error	Mean Squared Error	R ² Score
Elastic Net	0.000000	-0.000209	0.113491	0.023086	-0.000209
Lasso Regression	0.224467	0.233420	0.098067	0.017694	0.233420
Support Vector Regression	0.826962	0.704936	0.060286	0.006811	0.704936
Decision Tree Regression	1.000000	0.661861	0.043695	0.007805	0.661861
Linear Regression	0.796404	0.828054	0.042139	0.003969	0.828054
Ridge Regression	0.796402	0.827945	0.042128	0.003971	0.827945
KNeighbors Regression	0.852334	0.838551	0.036810	0.003726	0.838551
Gradient Boost Regression	0.946769	0.851488	0.034835	0.003428	0.851488
Random Forest Regression	0.972617	0.830805	0.033738	0.003905	0.830805

The dataset is modified with the NaN values being dropped

NaN = 0

NaN = Drop

**THANK YOU
FOR YOUR PARTICIPATION!**

