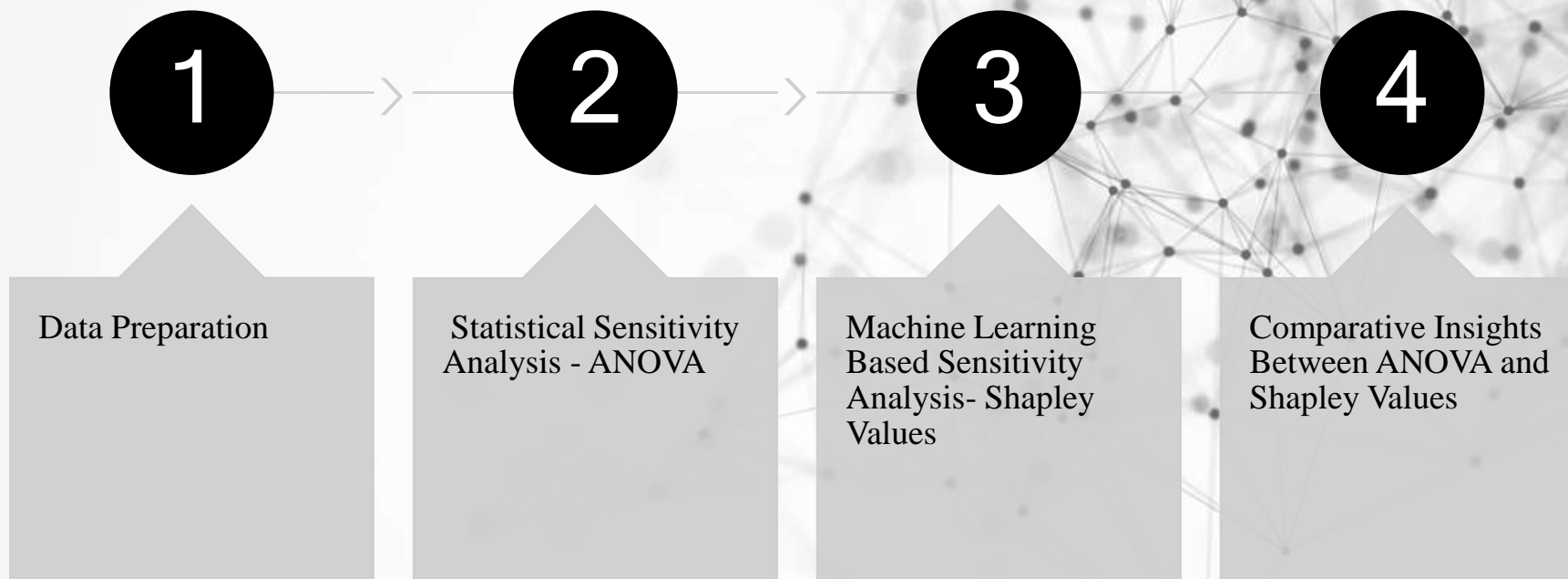




ENS 491

Machine Learning Exercise Part 2 Report

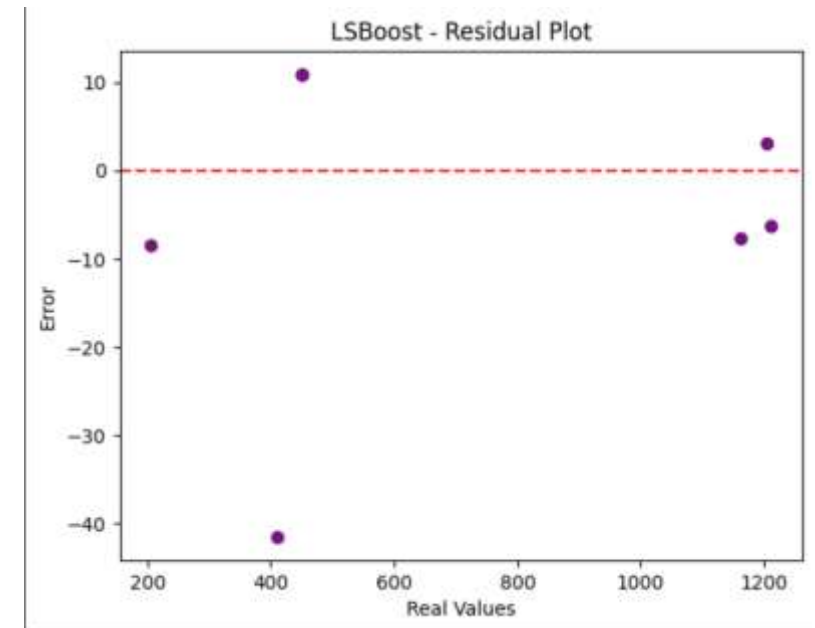
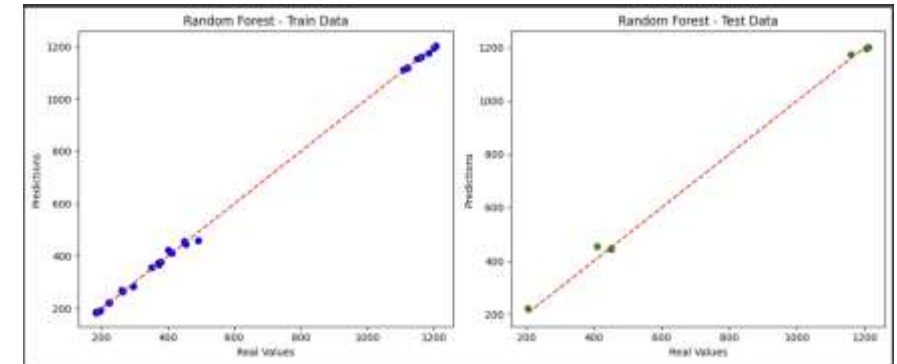
- Zeynep Derya Selçuk (29391)
- Darwin Serdar Dikici (30051)
- Selin Tıraş (2949)
- Efe Çelik (29268)
- Tunahan Arslan (29210)



Data Preparation

- Data preprocessing was done in Task 1.
- Normalization was done MinMaxScalar was used.
- LS Boost and Random Forest performance compared.
- You can find the necessary analyzes and steps in the ML exercise part 1 report.
- Necessary additions were made for Part 2, using what was made in Part 1.

⇒ Training set size: (38, 3)
Test set size: (7, 3)

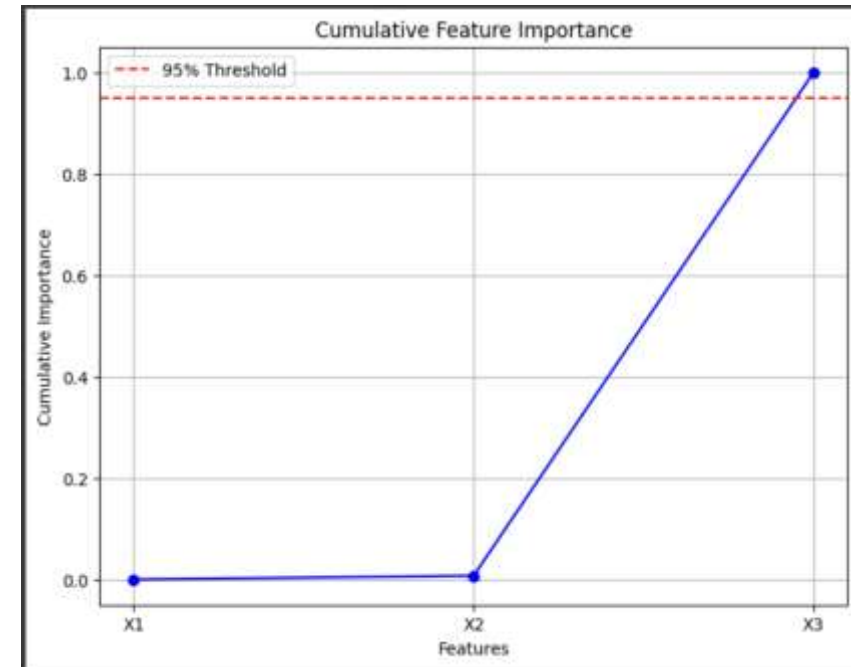


	Model	R^2 (Train)	R^2 (Test)	AME	RMSE
0	Random Forest	0.999636	0.997592	15.078571	20.121620
1	LSBoost	0.999569	0.998180	12.703682	17.494185

```

1
2 feature_importances = rf_model.feature_importances_
3
4
5 cumulative_importances = feature_importances.cumsum()
6
7
8 features = ['X1', 'X2', 'X3']
9
10
11 plt.figure(figsize=(8, 6))
12 plt.plot(features, cumulative_importances, marker='o', linestyle='-', color='b')
13 plt.title("Cumulative Feature Importance")
14 plt.xlabel("Features")
15 plt.ylabel("Cumulative Importance")
16 plt.grid()
17 plt.axhline(y=0.95, color='r', linestyle='--', label="95% Threshold")
18 plt.legend()
19 plt.show()
20

```



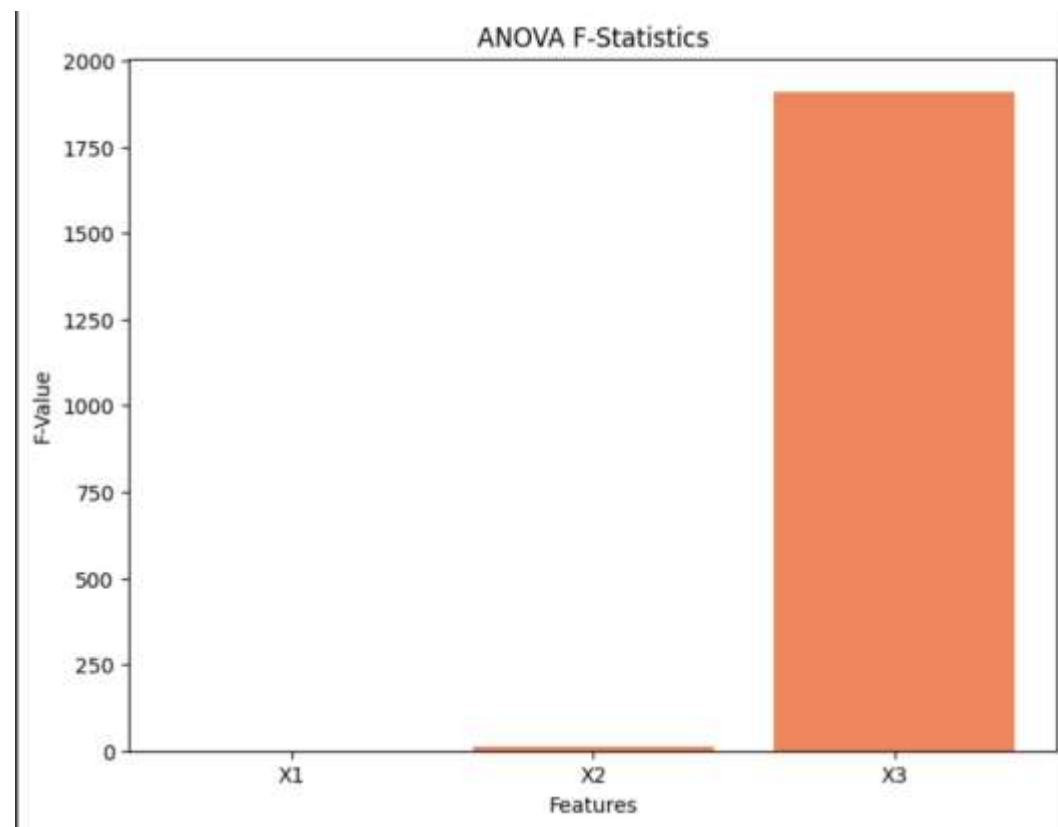
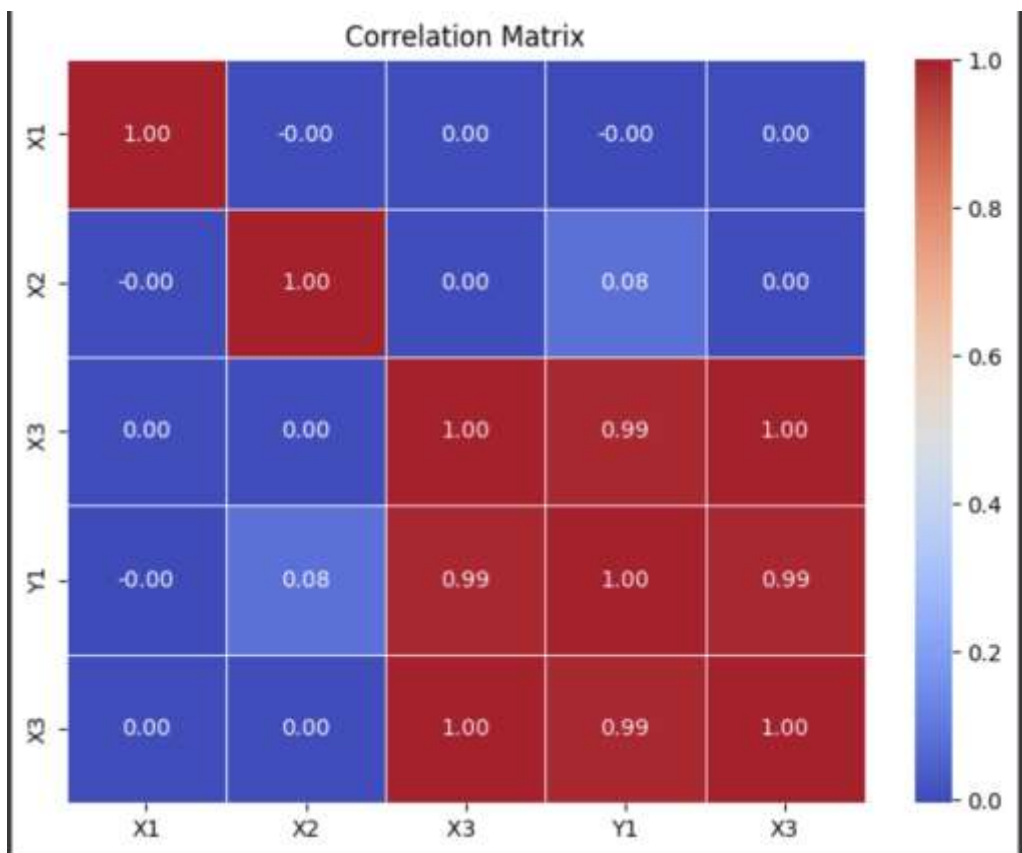
Statistical Sensitivity Analysis - ANOVA

- Conducted an Analysis of Variance (ANOVA) using the Ordinary Least Squares (OLS) method to evaluate the influence of each input feature (X1, X2, X3) on the target output (Y1).
- **Key Findings:** X3 was identified as the most significant feature with the highest F-statistic and lowest P-value.
- X2 showed moderate importance, while X1 had negligible influence.



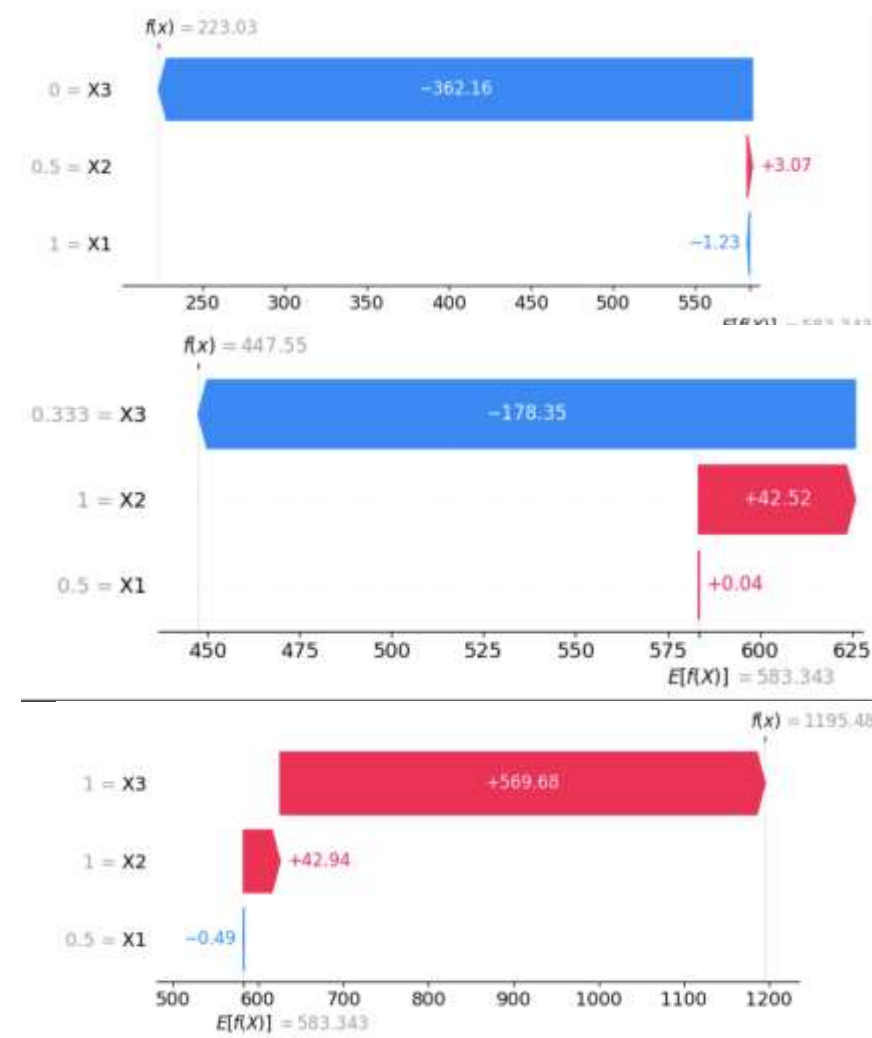
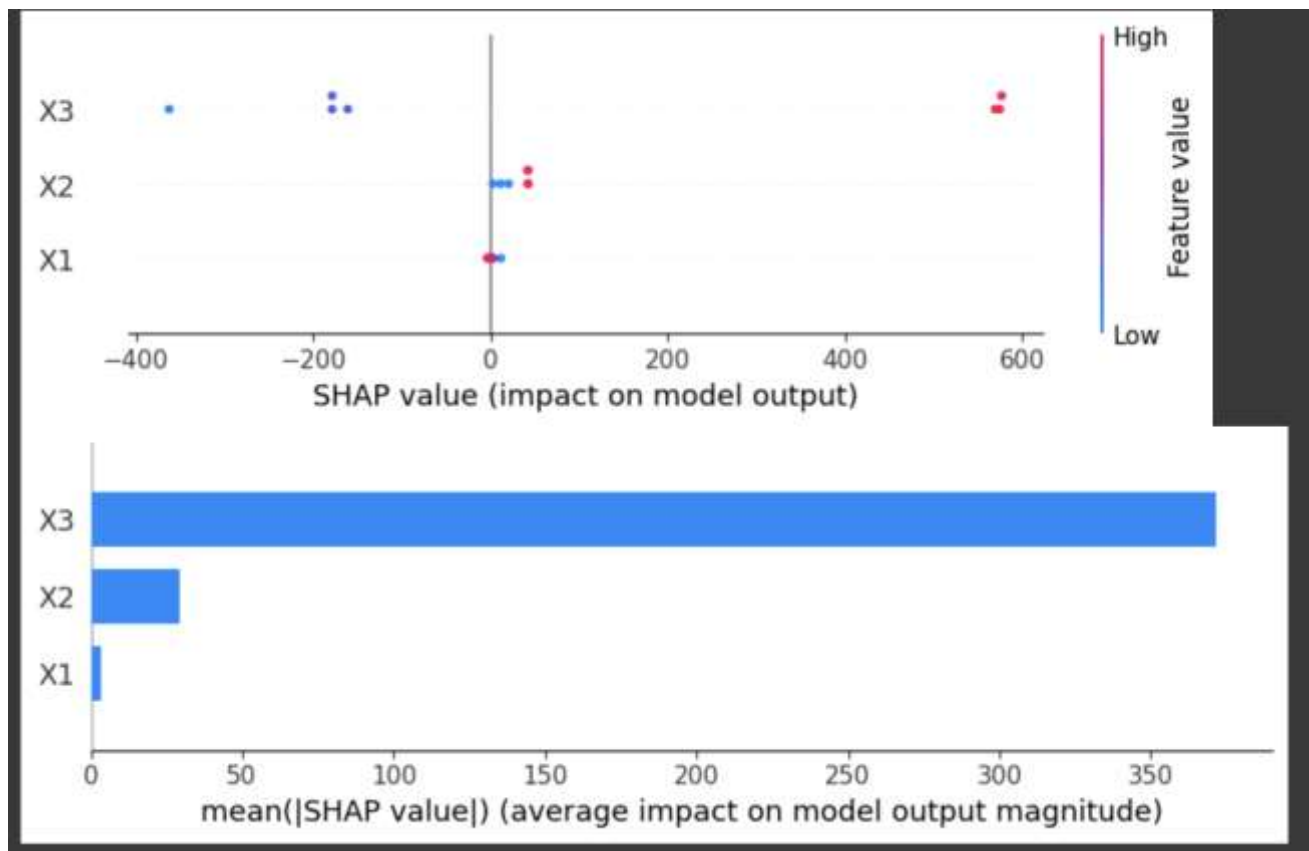
ANOVA results:

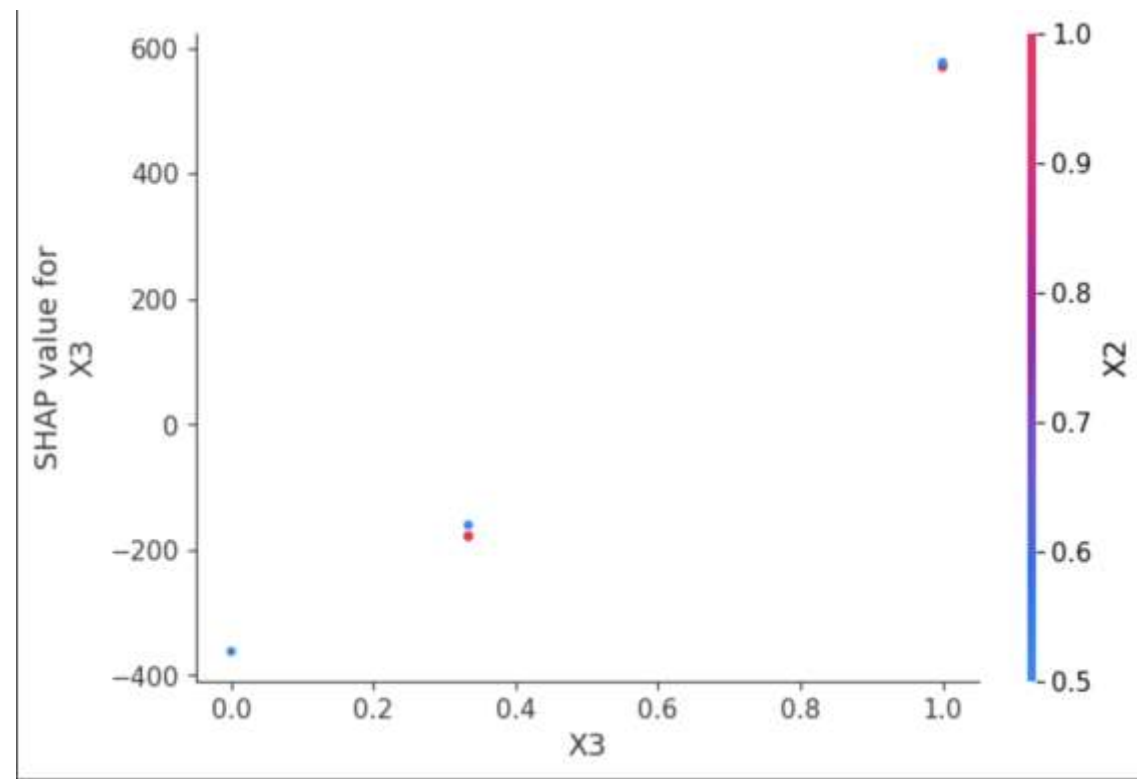
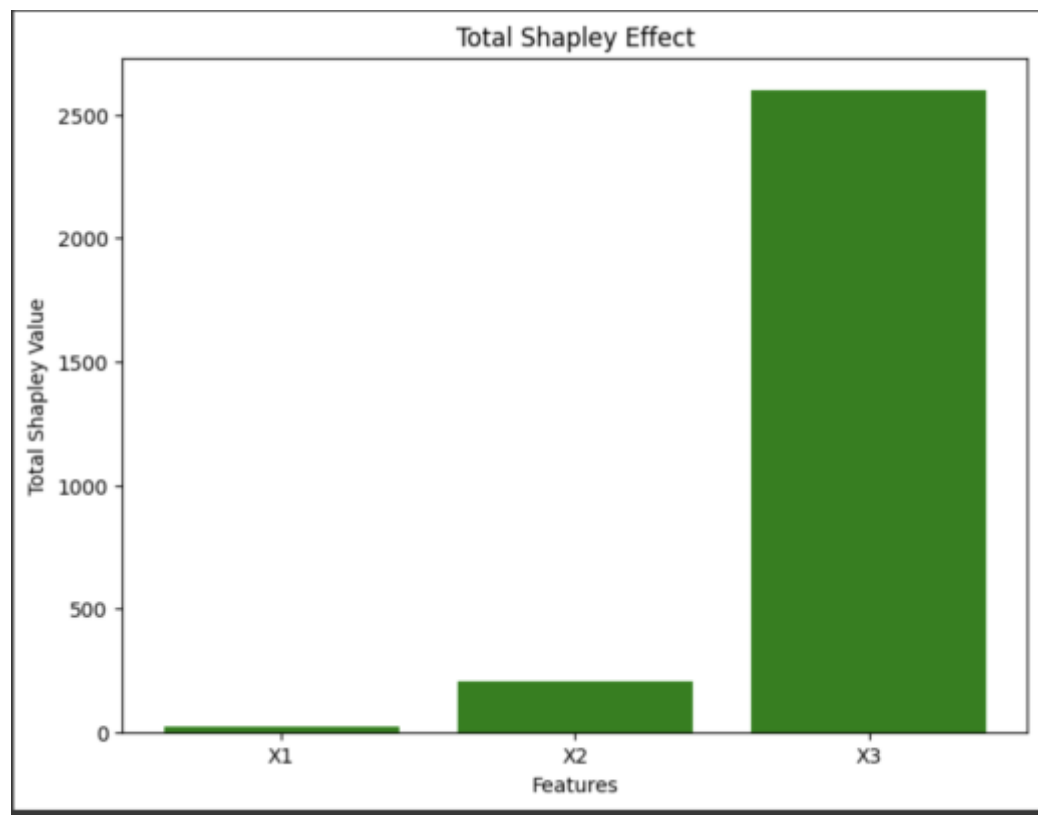
	df	sum_sq	mean_sq	F	PR(>F)
X1	1.0	1.369000e+02	1.369000e+02	0.036261	8.499172e-01
X2	1.0	5.216670e+04	5.216670e+04	13.817517	6.028574e-04
X3	1.0	7.211158e+06	7.211158e+06	1910.036406	5.115929e-36
Residual	41.0	1.547915e+05	3.775403e+03	NaN	NaN



Machine Learning Based Sensitivity Analysis- Shapley Values

- **Shapley Values:**Used shap.TreeExplainer to calculate Shapley values, providing a machine learning-based perspective on feature importance.
- **Steps:**
 - Calculated Shapley values for the test dataset.
 - Generated summary plots and bar graphs to visualize global feature importance.
 - Highlighted interactions between features (e.g., X3 and X2).
- **Interaction Analysis:**
 - Explored feature interactions using Shapley dependency plots, which revealed a strong relationship between X3 and X2.



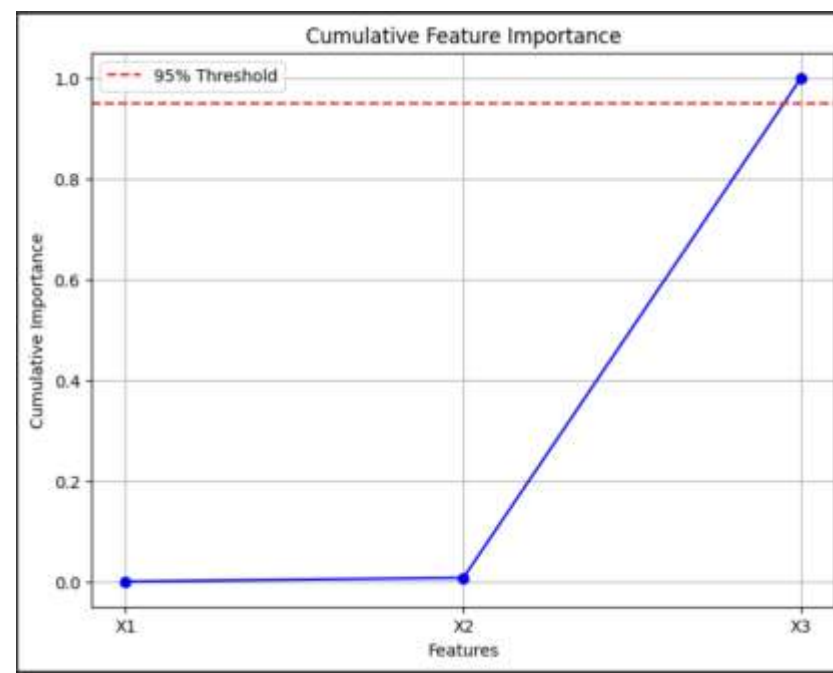
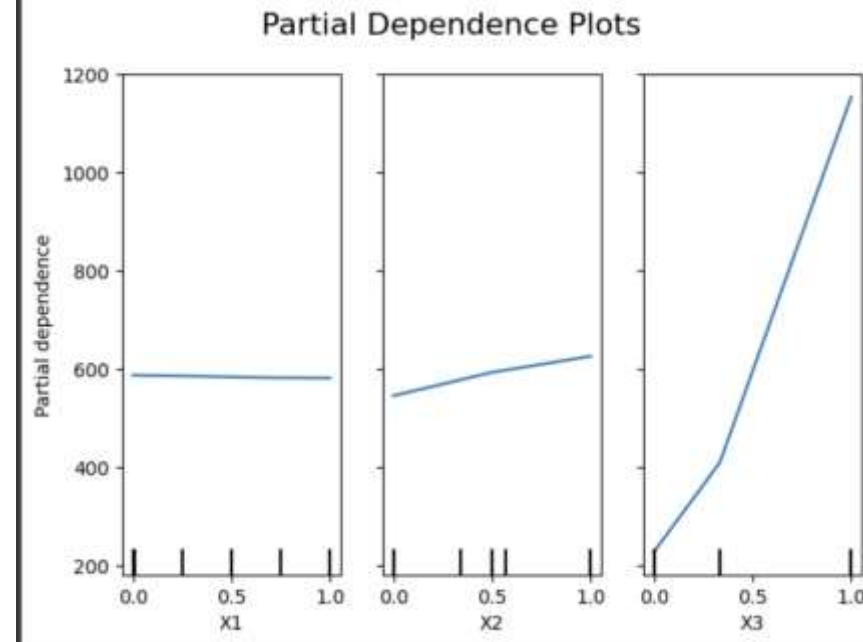
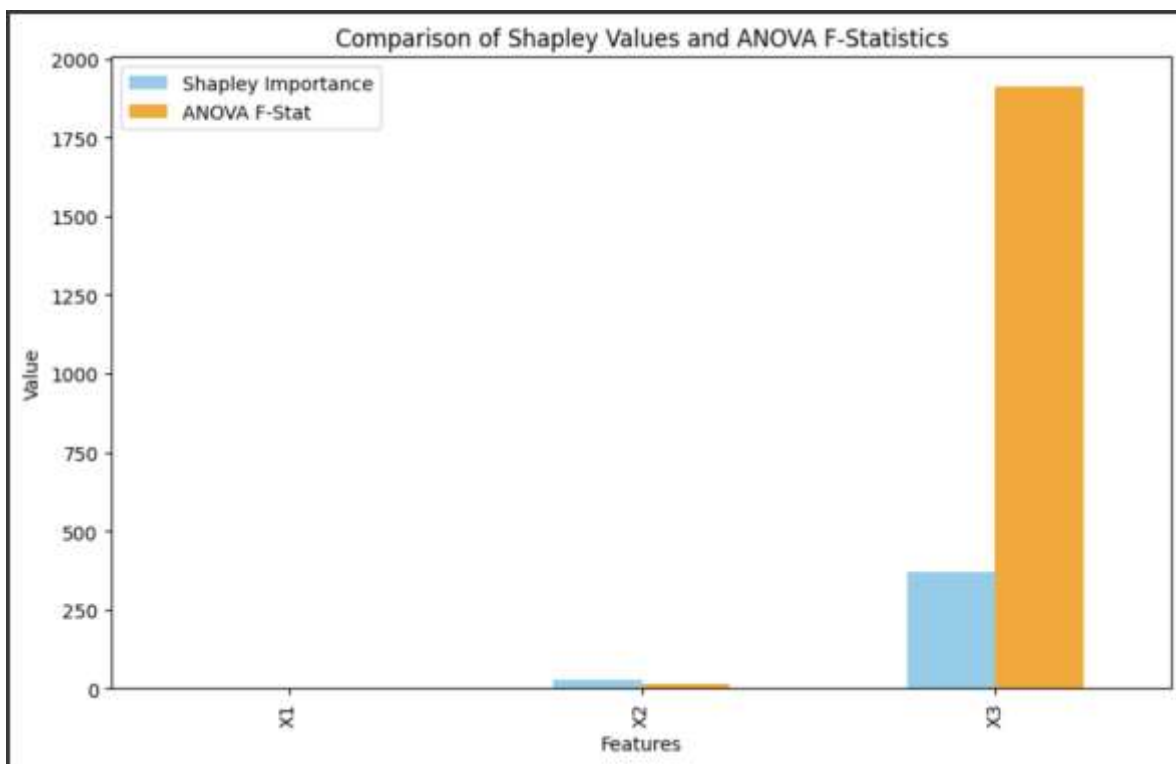


Comparative Insights Between ANOVA and Shapley Values

- **Comparison Between ANOVA and Shapley Values:** Both methods confirmed X3 as the most influential feature, followed by X2.
- Shapley analysis provided additional interpretability, capturing feature interactions that ANOVA could not.
- **Conclusion:** X3 is the dominant feature influencing the model's behavior, supported by both statistical and machine learning-based analyses.
- Shapley values provided more granular insights, particularly regarding feature interactions, making it a valuable tool for sensitivity analysis.

[→] Comparison of ANOVA and Shapley Results:

Feature	ANOVA P-Value	Shapley Mean Importance
0 X1	8.499172e-01	3.209077
1 X2	6.028574e-04	29.489981
2 X3	5.115929e-36	371.504010





END OF THE REPORT