# Predicting Supersymmetric Particles Report

Introduction:

The goal of this data analysis project is to predict whether a process produces supersymmetric particles or not. The dataset contains both training and test data, and the variable of interest is the "Signal" column, which represents whether the process produces supersymmetric particles (1) or not (0). The objective is to build a predictive model that can accurately classify the test data based on the provided features.

Methods:

The methods used in this analysis include Lasso and Ridge regression. The main objective of these methods is classification, where the target variable is binary (0 or 1). Lasso and Ridge are regularization techniques that add penalty terms to the regression model, helping to prevent overfitting and select relevant features. Lasso regression performs both feature selection and parameter shrinkage by adding an L1 penalty term, which can lead to some coefficients being exactly zero. Ridge regression, on the other hand, uses an L2 penalty term, which helps to reduce the impact of less important features without necessarily eliminating them entirely.

Data:

The dataset consists of features and the target variable "Signal." The provided features were numeric, and there were no categorical variables to deal with, simplifying the preprocessing steps. The initial data cleaning process involved removing the first column (assumed to be an index or identifier) from both the training and test datasets. Additionally, the missing values in the test dataset were filled with the corresponding column means. After combining the train and test datasets for feature engineering, missing values in the combined dataset were replaced with column means.

Analyze:

The analysis started with the preprocessing steps, which included handling missing values and preparing the data for model training. Lasso and Ridge regression models were trained using the training dataset, with the best lambda values selected through cross-validation. The final Lasso and Ridge models were then fitted using the chosen lambda values. The root mean squared error (RMSE) was computed to evaluate the performance of the models on the training data.

The final Lasso model had an RMSE of approximately 0.4248, while the final Ridge model had an RMSE of around 0.4254. These relatively low RMSE values indicate that the models were able to capture the underlying patterns in the data well. It is important to note that the objective here was classification, not regression, and thus the RMSE values do not represent prediction errors in the traditional sense.

Conclusion

In conclusion, this analysis aimed to predict whether a process produces supersymmetric particles or not using Lasso and Ridge regression methods. Both models performed well, achieving low RMSE values on the training data. The Kaggle leaderboard score of approximately 0.7570.