# Opening

I'm going to walk you through the detailed journey of my recent project, an exploration that spanned from meticulous data preprocessing to the strategic application of various machine learning models. This endeavor was not just about employing techniques but understanding their impact and the reasoning behind their selection.

# Data Preprocessing: A Careful Start

The foundation of any data science project lies in its data. I initiated this process with outlier detection to ensure the robustness of our models against anomalies that could skew predictions. Following this, feature scaling was applied, a crucial step to neutralize the effect of differing scales across variables, allowing each feature to contribute equally to the analysis. A correlation analysis provided insights into potential relationships and redundancies among features, guiding my feature selection process. Notably, the decision to remove the last column due to significant missing values was made to preserve the integrity and quality of our dataset, prioritizing meaningful data over volume.

# Initial Modeling with Logistic Regression

Choosing logistic regression as the initial model was strategic for its simplicity and interpretability. It served as a baseline, offering a clear view of the data's linear separability. However, the initial accuracy of 54.80% highlighted the model's limitations, nudging me towards more sophisticated methods to capture the dataset's complexity.

# Enhancement with Polynomial Features

The introduction of polynomial features was a pivotal moment. This technique expands the feature space by considering not only the original features but also their interactions and higher-order terms. It's akin to adding dimensions to our analysis, allowing us to capture

non-linear relationships. The improvement to 59.20% in accuracy was a testament to the complex nature of our data and the necessity of looking beyond linear assumptions. Polynomial features essentially provided the model with a richer, more nuanced language to describe the data's underlying patterns.

## Learning Curves and the Battle Against Overfitting

Upon examining the learning curves, a significant gap between training and validation performance hinted at overfitting. This is where regularization techniques came into play. L1 regularization, also known as Lasso, not only penalizes the complexity of the model but has the added benefit of feature selection by driving less important feature coefficients to zero. This dual action of simplification and selection brought our validation accuracy up to 66.80%, illustrating the power of L1 in enhancing generalization by focusing the model on the most salient features.

## Random Forest: A Game Changer

Despite the improvements, I ventured into random forest models, attracted by their ability to handle complex interactions and their intrinsic resistance to overfitting through ensemble learning. The initial performance, while not stellar, set the stage for a deeper analysis. By conducting a feature importance analysis and progressively focusing on the top 10 and then the top 6 features, I was able to dramatically refine the model's focus and efficiency, resulting in a remarkable accuracy of 92%. This leap in performance underscored the value of feature selection in eliminating noise and concentrating the model's learning on the most predictive aspects of the data.

## Neural Networks: The Final Frontier

Intrigued by the potential of neural networks to model even more complex relationships through their layered structure, I embarked on this final modeling journey. My initial foray, using a straightforward architecture, yielded modest results. However, by integrating L2 regularization, I aimed to combat overfitting by penalizing large weights, encouraging the

model to distribute its learning across a broader set of features rather than relying too heavily on any single one. Furthermore, focusing on the top six features, as identified through our previous analysis, and employing dropout for additional regularization, allowed us to achieve an impressive validation accuracy of 89%. This combination of techniques highlighted the nuanced balance between model complexity, feature selection, and regularization in achieving optimal performance.

## Summary of Best Results

To summarize, our journey through the realms of logistic regression, random forest, and neural networks revealed the intricate dance of machine learning techniques:

Logistic Regression with Polynomial Features and L1 Regularization culminated in a fine-tuned model achieving a cross-validation score of 78.80%.

Random Forest, after rigorous feature importance analysis and selection, delivered an exceptional accuracy of 92%.

Neural Network, optimized with L2 regularization and a focused feature set, reached a validation accuracy of 89%.

Each step of this journey, from the initial data preprocessing to the strategic application of various models, was guided by a continuous learning process. The enhancements made at each stage were not arbitrary but rooted in a deep understanding of the data and the tools at our disposal.

Closing

In closing, this project was a rich learning experience, underscoring the importance of a methodical approach, strategic thinking, and the continuous interplay between theory and practice in the field of data science. Thank you for your attention, and I look forward to any questions or discussions.