<u>Predicting Startup Funding Success- A Machine Learning Approach</u>

Below are the question answers that my audience might ask during the presentation. I did not include them in the presentation due to time constraints.

1. What makes the "Startup Investments Crunchbase" dataset suitable for this analysis?
➢ The dataset is comprehensive, covering approximately 54,000 startups with detailed information on funding rounds, investors, and industry classifications. This richness makes it a robust foundation for predicting funding likelihood.

2. Why were Logistic Regression and Decision Tree models chosen for this study?
➢ These models were selected due to their proven effectiveness and interpretability in classification tasks, making them suitable for predicting the binary outcome of whether a startup will receive funding.

3. How did you handle missing data and outliers in your dataset?
➢ Missing data was addressed through imputation or exclusion based on the context, while outliers, especially in total funding, were managed using the Interquartile Range (IQR) method and visual analysis to ensure data integrity.

4. What are the critical predictors of a startup's funding status that you determined from this project?
➢ Key predictors include the number of funding rounds, the amount of venture funding received, and the market sector the startup operates in. These factors were identified through exploratory data analysis and correlation studies.

5. Why did Logistic Regression outperform the Decision Tree model?
➢ Logistic Regression provided higher accuracy and precision, with fewer misclassifications of non-funded startups. This indicates its better suitability for the dataset and the prediction task.

6. Can you explain the significance of the models' F1 scores?
➢ The F1 score, being the harmonic mean of precision and recall, measures a model's accuracy in terms of its balance between predicting true positives and minimizing false positives and negatives. A high F1 score, as seen in both models, indicates effective prediction capabilities.

7. What are the limitations of your study?
➢ Limitations include the dataset's scope, which may not capture all variables influencing funding decisions, and the complexity of the startup funding process, which is difficult to fully model.

8. How did you ensure the models' predictive power and avoid overfitting?
➢ Cross-validation was employed to test the models on multiple subsets of the data, ensuring generalizability and robustness of the predictive power while avoiding overfitting.

9. What ethical considerations were taken into account in your analysis?
➢ The study emphasizes the importance of maintaining the privacy and security of startup data, mitigating biases in the models, and transparently communicating limitations and uncertainties to ensure fairness and innovation.

10. Why is model interpretability important in your study, and how does it affect the choice between Logistic Regression and Decision Trees?
➢ Interpretability is crucial for stakeholders to understand the decision-making process of the model, allowing for informed strategic planning. Logistic Regression was preferred for its balance of accuracy and interpretability, but the choice can vary based on specific needs, such as the importance of reducing false positives or the need for a more straightforward explanation of predictions.