

# Clinical Trial Success Prediction

## Business Problem

Pharmaceutical companies face significant financial risks due to the high failure rate of clinical trials, with more than 85% of new drug candidates failing to progress through the development pipeline. A predictive system capable of estimating the likelihood of trial success based on design features could transform decision-making processes, improve R&D efficiency, and reduce sunk costs.

## Background and History

Historically, clinical trial planning has depended on experience and precedent. The emergence of the AACT database provides structured access to metadata on hundreds of thousands of trials, enabling data-driven forecasting. Several studies have noted that sponsor type, trial phase, and design characteristics are correlated with trial outcomes.

## Data Explanation

The dataset comes from the Aggregate Analysis of ClinicalTrials.gov (AACT) and includes structured tables for study metadata, sponsors, interventions, design elements, outcomes, and eligibility criteria. Key fields used include phase, study type, enrollment, masking, allocation, intervention type, agency class, outcome type, gender eligibility, and age ranges.

## Methods

The analysis involved EDA, missing value imputation, one-hot encoding for categorical features, and feature engineering. Three models were tested: Logistic Regression (with scaled inputs), Random Forest, and XGBoost. Model performance was measured using AUC-ROC, precision, recall, and F1-score. SHAP was used to interpret feature importance.

## Analysis

Random Forest achieved the best results with an AUC of 0.80 and an F1-score of 0.71 for successful trials. XGBoost underperformed in recall, and Logistic Regression failed to converge meaningfully, performing only slightly better than random. SHAP analysis revealed enrollment size, sponsor type, presence of primary outcomes, and intervention type as key features driving model output.

## **Conclusion**

Predictive modeling based on trial design metadata can offer useful insights for trial planning and investment decisions. Random Forest provided the best performance in this project, indicating that decision-tree-based models perform best in this domain.

## **Assumptions**

This analysis assumes that the presence of submitted results is an acceptable proxy for trial success. It is also assumed that missing values in the dataset are missing at random and that imputing them using statistical methods does not introduce significant bias. Additionally, it is presumed that trials failing to report results are functionally equivalent to unsuccessful trials, although some may have succeeded but not reported.

## **Limitations**

A limitation is the use of result submission as a proxy for success, which may not reflect actual clinical efficacy. There is also a significant amount of missing data in critical fields such as phase and participant age, which may affect the performance of the model. The model was evaluated on a held-out test set derived from an 80/20 train-test split. However, no validation was performed on an external dataset beyond the AACT database snapshot used in this analysis.

## **Challenges**

The analysis had challenges related to data quality and consistency. High numbers of missing values in fields like age, masking, and allocation complicated preprocessing. Combining and aligning multiple AACT tables required indexing and merging logic. The use of one-hot encoding increased feature dimensionality, which can impact model performance and interpretability. Maintaining a balance between model complexity and explainability was a concern throughout the analysis.

## **Future Uses / Additional Applications**

This modeling approach could be extended to support trial risk scoring systems for investors or funding bodies. It could also aid in trial design optimization tools that suggest more successful configurations. Additionally, the model might assist in prioritizing funding within therapeutic areas or predicting the likelihood of

regulatory approval. Integration with real-time trial databases could enhance ongoing training.

## **Recommendations**

The model should be used as a decision support tool rather than as a replacement for expert judgment. In future iterations, unstructured protocol documents and trial objectives could be incorporated to capture additional features. Therapeutic area-specific models could increase precision in certain domains. Moreover, outputs should be probabilistic to allow decision-makers to calibrate thresholds based on specific risk tolerances.

## **Implementation Plan**

A production-ready pipeline could be deployed using a dashboard (e.g., Streamlit) to allow users to input trial metadata and receive a success probability along with SHAP-based interpretability. The system should be retrained periodically as the AACT database is updated.

## **Ethical Assessment**

Ethical considerations include the risk of deprioritizing trials in underrepresented diseases due to data-driven biases. There is potential for reinforcing existing inequalities in research funding if historical trends are not corrected for. To mitigate these risks, transparency should be a major goal by using interpretability tools such as SHAP. Model recommendations should be used alongside human review and domain knowledge.

## **Possible Questions**

1. How did you define clinical trial success, and why was result submission used as a proxy?
2. What factors led to Random Forest outperforming XGBoost and Logistic Regression in this context?
3. How were trials with missing or ambiguous design characteristics handled during preprocessing?
4. Can this model be applied at the proposal stage, or does it require finalized trial registration data?
5. To what extent does the model reflect or reinforce existing biases in clinical research funding?

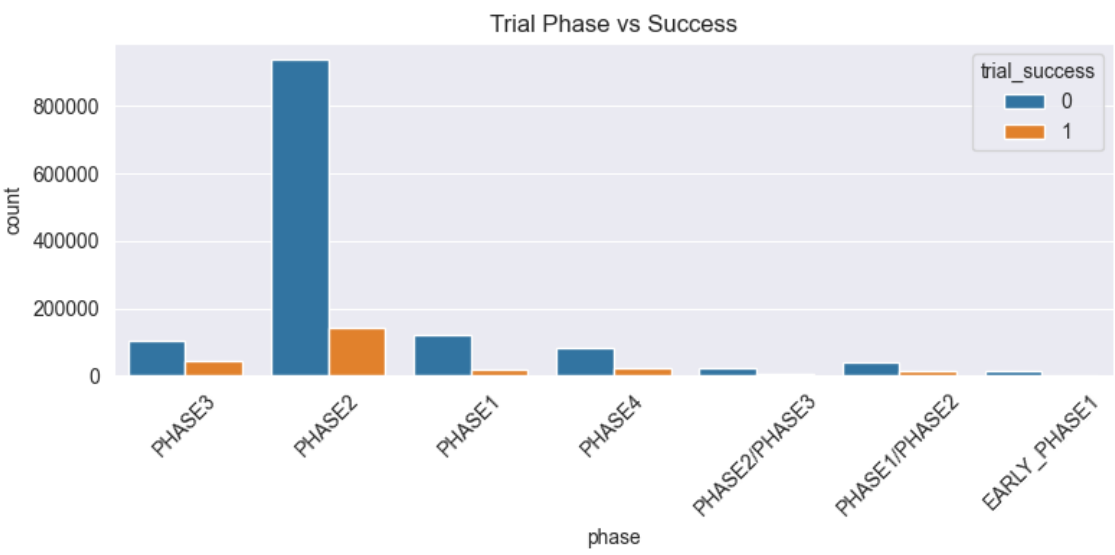
6. How would the model's predictions generalize to international trials or those outside of ClinicalTrials.gov?
7. What validation strategies could be used in future work to improve generalizability beyond this dataset?
8. Could this model be adapted for use in specific therapeutic areas such as oncology or neurology?
9. What safeguards are in place to ensure that eligibility and demographic fields do not introduce bias?
10. How could this model contribute to identifying low-quality or noncompliant trials early on?

## Appendix: Illustrations

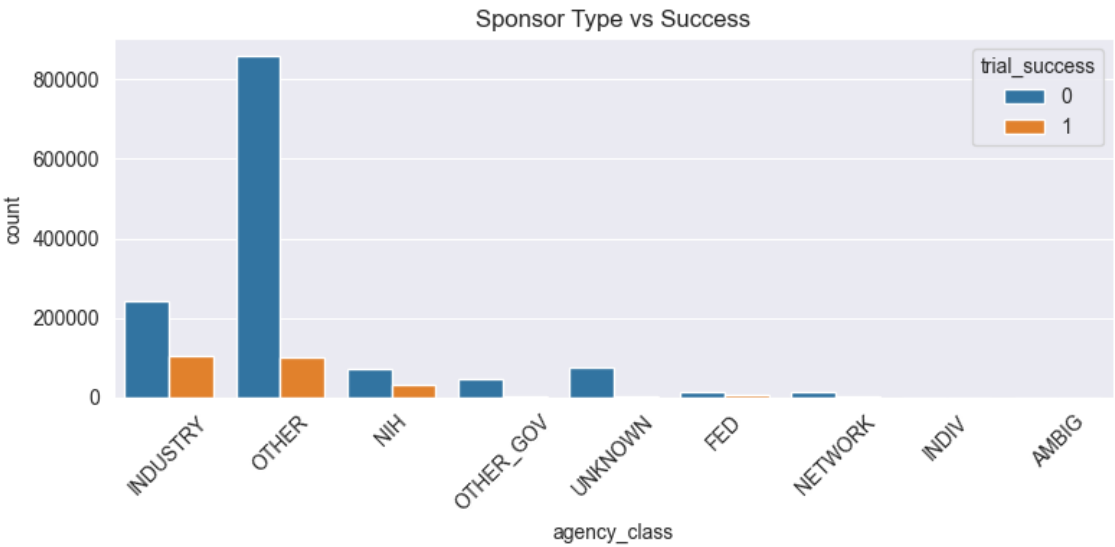
Eda Trial Success Distribution



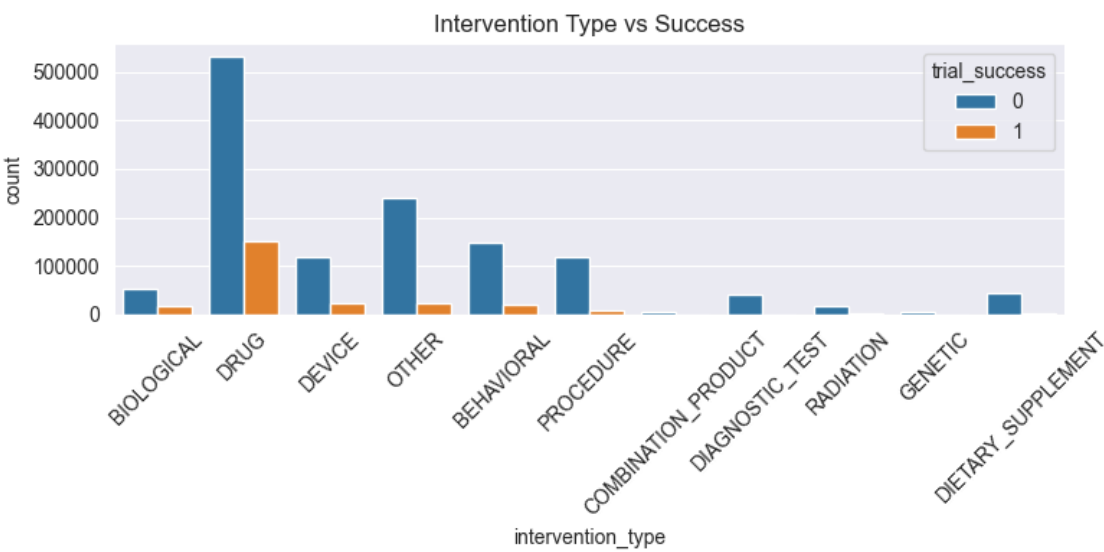
Eda Phase Vs Success



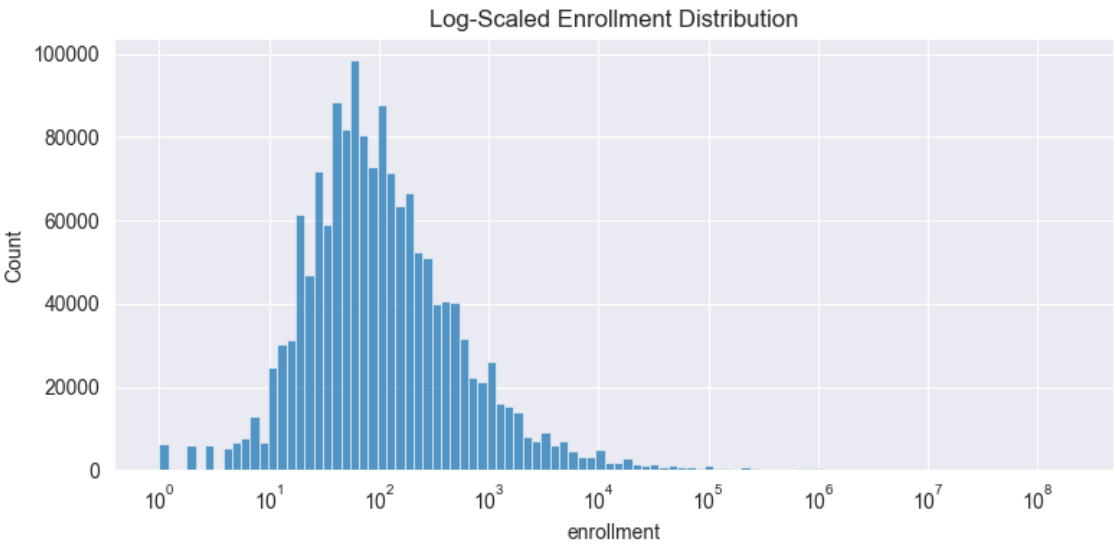
Eda Sponsor Type Vs Success



Eda Intervention Type Vs Success

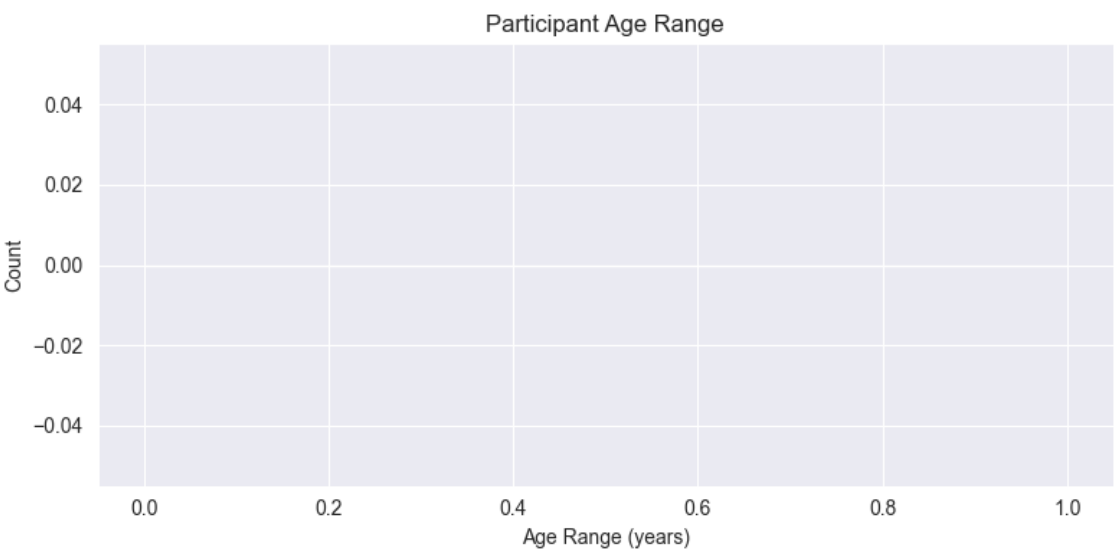


Eda Log Enrollment Distribution

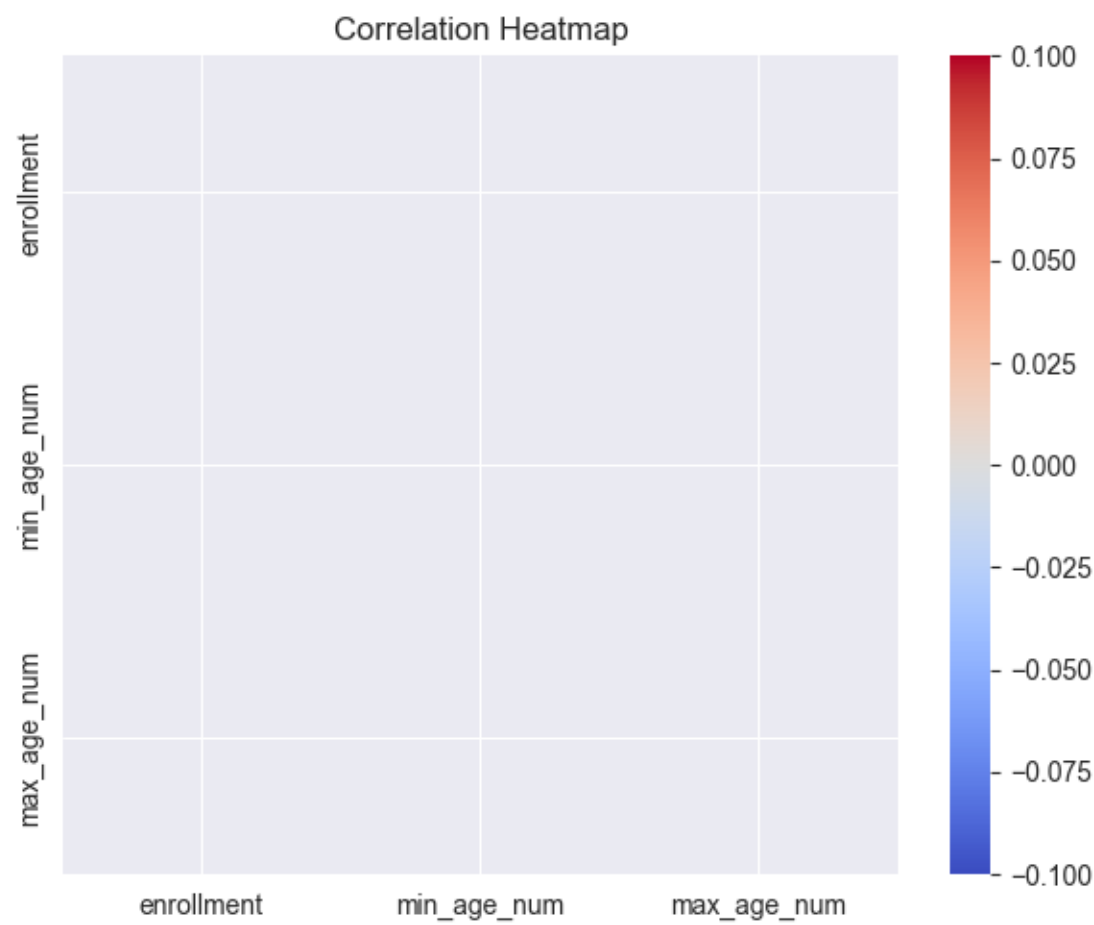




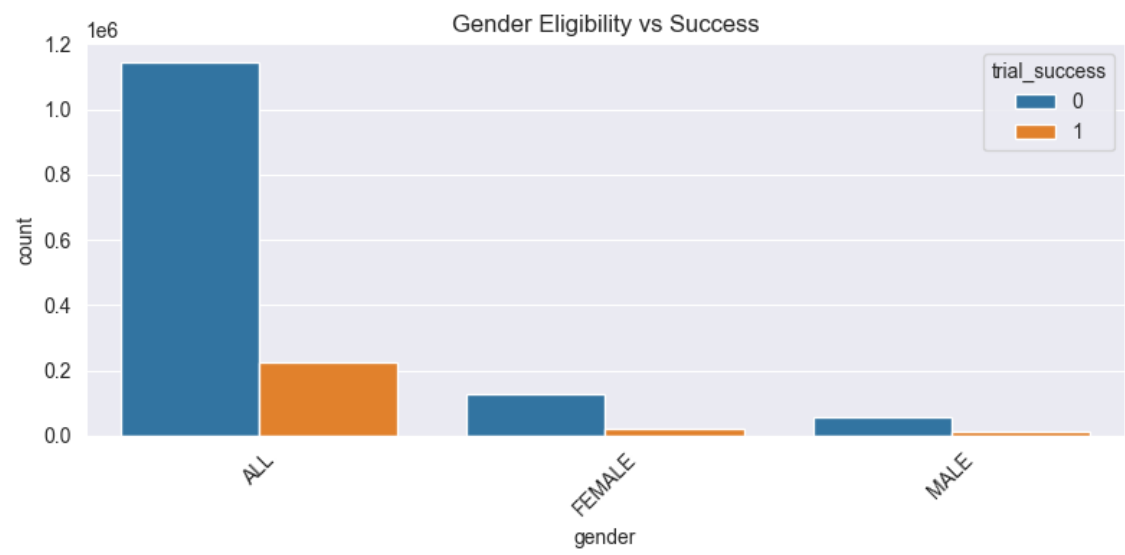
Eda Age Range Distribution



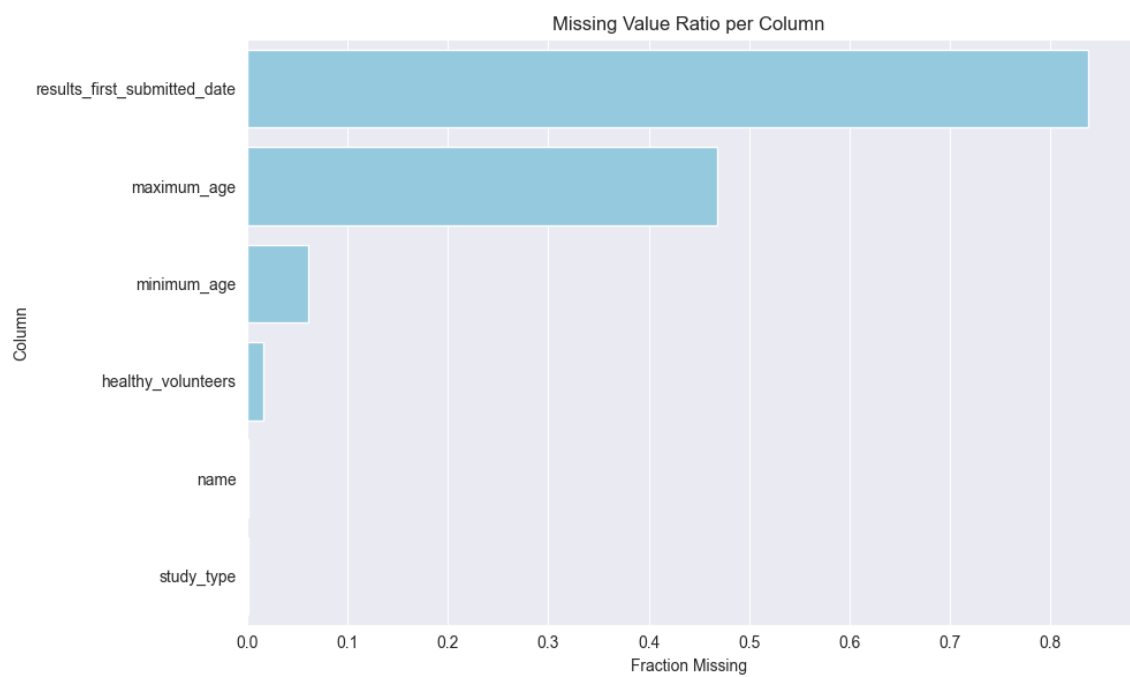
Eda Correlation Heatmap



Eda Gender Vs Success



# Eda Missing Values Barplot



Shap Summary Bar

