# A Model Evaluation on LightGBM vs. Catboost:

## A Study on Nasdaq-listed Stock Closing Auction Price Prediction

JUNWEI SHEN & CHE ZHU

## Introduction

- Main Objective: This study will utilize two cutting-edge gradient boosting models, LightGBM and Catboost, to predict the future price movements of stocks through a synthetic index composed of NASDAQ-listed stocks.
- Major Approach: We take a methodical approach encompassing data preprocessing, feature engineering, and model training and tuning, and we strive to minimize the mean absolute error (MAE).
- Responsible AI development: model comparison analysis of the two gradient-boosting models.



*Fig 1: Nasdaq Stock Market*



*Fig 2: Correlation Heatmap of Features and Outcome:* There are several correlations between different variables; besides relatively highly correlated with bid_price, ask_price, bid_size, and reference_price, the target variable also has a medium-level correlation with time_id



*Fig 3: Different Price Measurements Change over Time:* the target line shows a significant divergence from the price lines; the target line fluctuating in a different scale implies that its behavior or scale is different from that of the prices

*Fig 4: LightGBM and Catboost Model Accuracy Results:* Compared to a variance of the target (59.797), global minimum (-103.030), and global maximum (113.180), both models have relatively small MAE; The Catboost model slightly outperforms the LightGBM model with an MAE of 4.892

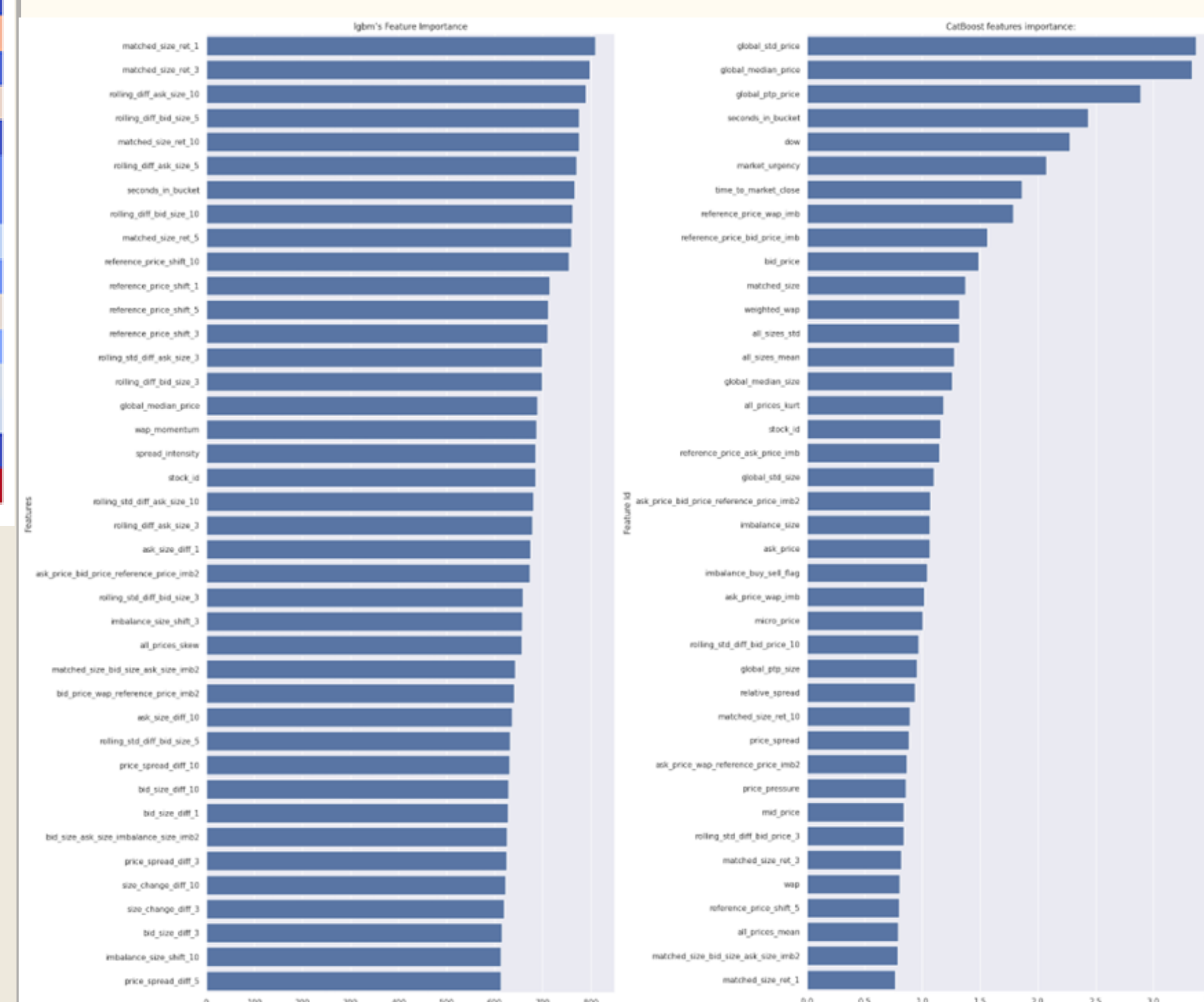| | |
|---|---|
| LightGBM Final MAE score | 4.986 |
| Catboost Final MAE score | 4.892 |
| Testing Set Target variance | 59.797 |
| Testing Set Target Range | (-103.030, 113.180) |



*Fig 5: LightGBM (Left) and Catboost(Right) Feature Importance Result:* Certain features are highly ranked in both results, implying that these features are strong predictors regardless of the model type. However, some features are given different significance levels in the ranking, this is the result of the algorithms' functionalitiess

## Model Comparsion

**CatBoost and LightGBM Overview:**

- Part of the gradient boosting family, known for efficiency, speed, and performance.
- Distinct in optimization and handling of categorical data.

**Performance Implications:**

- CatBoost: Better for datasets with correlated features (e.g., auction data), due to ordered boosting and symmetric trees techniques.
- LightGBM: Offers similar accuracy levels with less computing time, ideal for complex ML tasks with large datasets.
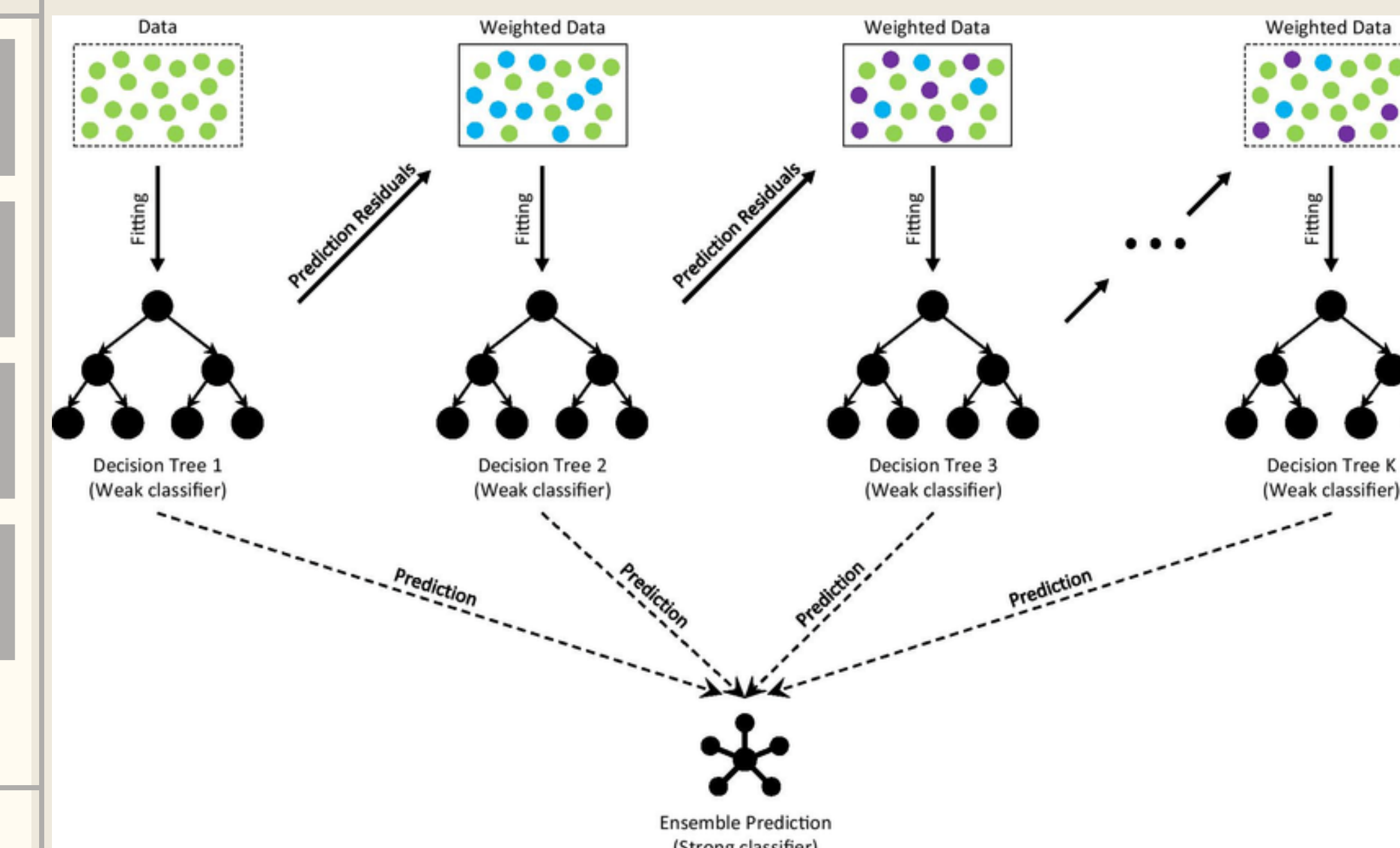


*Fig 6: The architecture of Gradient Boosting Decision Tree*

## Bias Consideration

Stock market prediction models, influenced by human actions and major events, face challenges in adapting to rapid changes and biases. Despite high data transparency, selective data use and historical assumptions can mislead and reduce generalization. It's crucial to account for these factors to ensure fair and accurate market analysis