# Iceberg Prediction Whitepaper

**Joseph Loss**

---

**Abstract**

This paper presents a comprehensive machine learning approach for predicting iceberg order execution in quantitative trading. By analyzing market microstructure patterns through an XGBoost-based model, we achieve 79% precision in predicting whether detected iceberg orders will be filled or canceled. Our system integrates real-time order book data, trade
imbalance metrics, and innovative side-relative feature transformations to capture execution dynamics. The model is validated using time-series cross-validation to prevent look-ahead bias, maintaining consistent performance across market regimes with precision consistently above 75%. We demonstrate how prediction probabilities can be translated into actionable trading decisions through confidence bands, creating a sophisticated execution strategy that adapts to changing market conditions. The resulting system provides valuable signals for algorithmic trading strategies, improving response to hidden liquidity, identifying opportunistic entry/exit points, and reducing execution costs.

---

## 1.  Project Context: Why This Matters in Trading

In high-frequency and algorithmic trading, iceberg orders represent a significant market microstructure phenomenon. Let me walk you through how I approached predicting iceberg order execution using machine learning techniques.

An iceberg order is a large order that's divided into smaller, visible portions - like the tip of an iceberg above water, with the majority hidden below. Traders use them to minimize market impact while executing large positions.
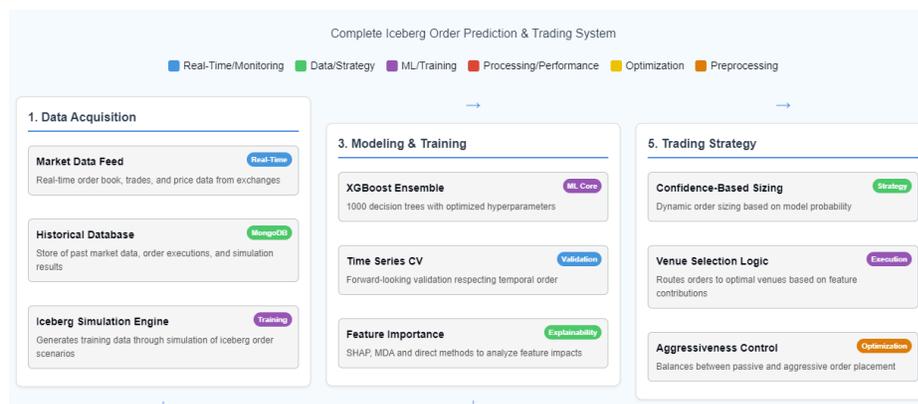
---

*Image 1: The complete system architecture showing data acquisition, processing, modeling, and trading strategy components. This diagram illustrates the end-to-end pipeline from market data to trading decisions.*
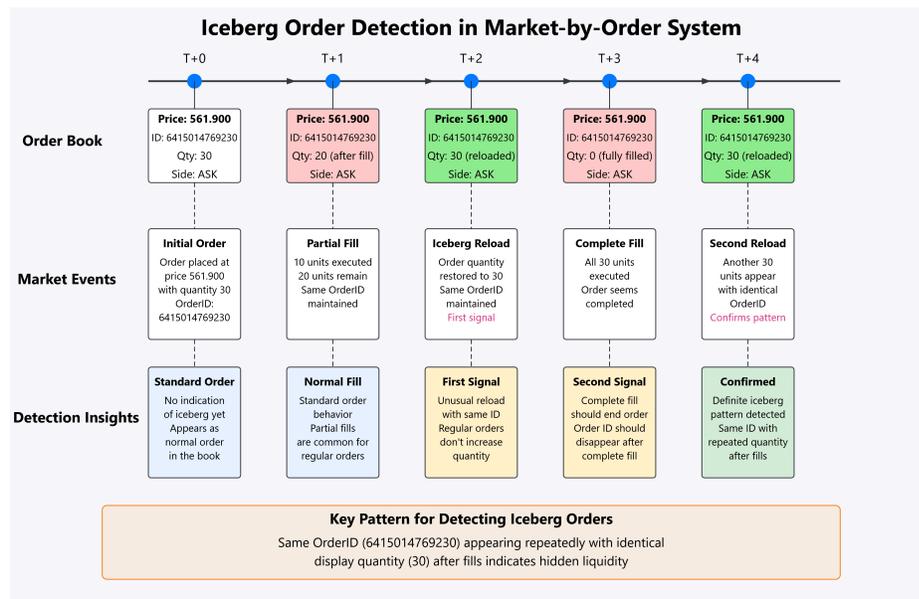
## 2. THE BUSINESS PROBLEM

**Objective**: Predict whether a detected iceberg order will be filled (`mdExec = 1`) or canceled (`mdExec = 0` or `-1`).

Why is this valuable? If we can predict iceberg order execution:

1. We can improve trading algorithms' response to large hidden liquidity
2. We can identify opportunistic entry/exit points when large orders are likely to complete
3. We can better estimate true market depth beyond what's visible on the order book
4. We can reduce execution costs through improved venue and timing decisions

## 3. DATA SOURCES FOR ICEBERG ORDER PREDICTION

Our machine learning model is trained on data derived from two key sources: iceberg order detection messages and simulation results.



These data sources form the backbone of our machine learning pipeline, providing both the features and the target variable for our prediction model. Typically, we generate tens of thousands of simulation results from replaying about two months of historical market data.

## 4. DATA PIPELINE ARCHITECTURE

*4.a. Data Collection & Simulation Infrastructure*

The raw data comes from market simulations of iceberg orders. These simulations are orchestrated using Prefect workflows that:

1. **Retrieve Production Data**: SSH to production servers to retrieve historical market data archives
2. **Extract & Prepare Market Data**: Unzip market data files and set up Docker container environments
3. **Run Iceberg Simulations**: Execute Java-based simulation in parallel Docker containers
4. **Store Results in MongoDB**: Upload simulation results with metadata for ML consumption

The core of this data collection pipeline is implemented in `sxm_market_simulator.py`:

```python
@flow(name="Parallel Archive Market Simulator",
    flow_run_name="{dataDirectory} | {packageName}",
    task_runner=RayTaskRunner(
        init_kwargs={
        "include_dashboard": False,
    }),
    cache_result_in_memory=False,
)
async def ParallelArchiveMarketSimulator(dataDirectory=None,
packageName=None):
    logger = get_run_logger()
    secret_block = await Secret.load("sshgui-credentials")
    cred = json.loads(secret_block.get())

    # db and metaDb names for simulation:
    dataCollectionName = f'mktSim_{packageName}'
    metaCollectionName = f'mktSim_meta'
```

*From sxm_market_simulator.py, lines 168-193*

The `RayTaskRunner` enables parallel execution of these simulations, with each simulation running in its own container. This parallelization is essential for processing large amounts of historical market data efficiently.

*4.b. Raw Data Format*

The raw data from iceberg simulations follows a complex nested structure as seen in the `IcebergSimulationResult.yaml` file:

```yaml
!IcebergSimulationResult {
  waitingToOrderTradeImbalance: {
    combinedMap: {
      TimeWindowNinetySeconds: !TradeImbalance { timeInterval:
NINETY_SECONDS, isTimeBased: true, period: 0, tradeImbalance: 0.51899 },
      MessageWindow100: !TradeImbalance { timeInterval: NONE,
isTimeBased: false, period: 100, tradeImbalance: 1.02041 },
      # ... more time windows
```

```
      }
    },
    # ... other simulation data
    symbol: ZMZ3,
    icebergId: 708627238048,
    isBid: false,
    mdExec: -1,
    price: 39900000000,
    volume: 12,
    showSize: 3,
    filledSizeStart: 39,
    filledSizeEnd: 48,
    # ... additional order details
}
```

*From IcebergSimulationResult.yaml*

This rich structure needs to be flattened and transformed before model training, which leads us to the data preprocessing pipeline.

## 5. ICEBERG ORDER SIMULATION ENGINE

At the core of our simulation infrastructure is the `ActiveIceberg.java` implementation, which tracks and simulates iceberg order behavior using replayed market data.

### 5.a. State Tracking and Event Processing

The class maintains comprehensive state information about each iceberg order:

```
/**
 * This is a row in a results table of Iceberg Sims
 */
public class ActiveIceberg extends SelfDescribingMarshallable {
    /**
     * The Iceberg id.
     */
    public long icebergId;
    /**
     * The Is filled.
     */
    public boolean isFilled = false;
    /**
     * The Send order flg.
     */
    public boolean sendOrderFlg = false;
    /**
     * The Send cancel flg.
     */
    public boolean sendCancelFlg = false;
    public boolean isCancelled = false;
    /**
     * The Md exec.
```

```
    */
    public int mdExec = Nulls.INT_NULL;
    /**
     * The Is bid.
     */
    public boolean isBid;
    /**
     * The Price.
     */
    public long price;
    // Additional fields...
}
```

It processes three types of market events, updating the iceberg's state accordingly:

1. **Trade Updates**: When trades occur at the iceberg's price level

```
public void onTradeUpdate(TradeUpdate tradeUpdate) {
    this.latestExchangeTimeNs = tradeUpdate.exchangeTimeNs();

    if (tradeUpdate.price() == price)
        volume += tradeUpdate.qty();

    if (isFilled)
        return;

    if (tradeUpdate.price() == price &&
Nulls.isQtyNotNull(currentQueuePosition) && currentQueuePosition !=
Nulls.INT_NULL) {
        currentQueuePosition -= tradeUpdate.qty();
    }
    if (currentQueuePosition < 0 && currentQueuePosition !=
Nulls.INT_NULL) {
        isFilled = true;
        simFillExchangeTimeNs = tradeUpdate.exchangeTimeNs();
    }
}
```

2. **Book Updates**: Changes to the order book that might affect queue position

```
public void onBookUpdate(BookUpdate bookUpdate) {
    this.endExchangeTimeNs = bookUpdate.exchangeTimeNs();
    this.latestExchangeTimeNs = bookUpdate.exchangeTimeNs();

    if (!this.sendOrderFlg) {
        if (bookUpdate.exchangeTimeNs() <
this.delayedOrderExchangeTimeNs) {
            return;
        } else {
            this.initialBookQty = bookUpdate.getDirectQtyAtPrice(isBid,
price, true);
            this.initialQueuePosition = this.initialBookQty + 1;
            this.currentQueuePosition = this.initialQueuePosition;
```

```java
        this.sendOrderFlg = true;
    }
}
    // Additional logic...
}
```

3. **Iceberg Updates**: Direct updates to the iceberg order's status

```java
public void onIcebergUpdate(IcebergUpdate icebergUpdate) {
    this.endExchangeTimeNs = icebergUpdate.exchangeTimeNs();
    this.latestExchangeTimeNs = icebergUpdate.exchangeTimeNs();

    this.filledSizeEnd = icebergUpdate.filledQty();
    this.lastStatus = icebergUpdate.status();

    if (this.lastStatus == IcebergUpdateStatus.CANCELLED ||
        this.lastStatus == IcebergUpdateStatus.FILLED ||
        this.lastStatus == IcebergUpdateStatus.MISSING_AFTER_RECOVERY) {
        {
            this.checkAndManageOrderState();
        }
        this.isComplete = true;
    }

    if (this.lastStatus == IcebergUpdateStatus.RELOADED ||
        this.lastStatus == IcebergUpdateStatus.LAST_RELOAD) {
        this.numberReloads++;
    }
}
```

*5.b. Order Management Logic*

The class also implements realistic order placement and cancellation logic:

```java
public void checkAndManageOrderState() {
    if (!this.sendOrderFlg || isComplete || isFilled || isCancelled) {
        return;
    }
    try {
        sendCancelOrder();
    } catch (Exception e) {
        // Logging logic...
    }
}
```

This simulation engine is orchestrated by the `ParallelArchiveMarketSimulator` flow in `sxm_market_simulator.py`, which:

1. Retrieves historical market data archives
2. Sets up simulation environments
3. Executes the Java-based simulation in parallel Docker containers
4. Collects and stores the simulation results in MongoDB

```
@flow(name="Parallel Archive Market Simulator",
    flow_run_name="{dataDirectory} | {packageName}",
    task_runner=RayTaskRunner(
        init_kwargs={
        "include_dashboard": False,
    }),
    cache_result_in_memory=False,
)
async def ParallelArchiveMarketSimulator(dataDirectory=None,
packageName=None):
    # Simulation orchestration logic...
    dfResults = await execute_market_simulator.with_options(
        flow_run_name=f"{dfMeta['instance'][0]}/{dfMeta['archive_date']}
[0]} {packageName}",
    )(packageName=f"{packageName}",
      exec_args=[
        "true",
        f"/data/{strEngineDate}/",
        f"/tmp/results_{strEngineDate}.csv.gz",
        f"/sx3m/{strEngineDate}.properties"
      ],
      wait_for=[unzip_data])

    # Process and store results...
```

The simulation results produced by this engine form the dataset for our machine learning model, providing both feature data and the target variable (`mdExec`) that indicates whether each iceberg order was filled or cancelled.

## 6. DATA PREPROCESSING AND FEATURE ENGINEERING

The data preprocessing workflow is shown in Image 2:

*Image 2: The data processing pipeline showing the three main components: data flattening, feature engineering, and feature scaling. This ETL process transforms nested JSON data into tabular features suitable for machine learning.*

6.a.  Data Flattening Pipeline

The nested JSON structure is transformed into a flat tabular format suitable for machine learning using a custom flattening function:

```python
def flatten_to_columns(df=None,
flatten_cols=['waitingToOrder','orderPlaced','oneStateBeforeFill']):
    #
--------------------------------------------------------------------------------
    # FLATTEN COLUMNS
    lstFlattened = []

    for col in flatten_cols:
        dicts = df[col].to_list()

        #
--------------------------------------------------------------------------------
        # TradeImbalance

        # Now we can expand each dictionary into its own DataFrame
```

```python
        expanded = [pd.json_normalize(d['TradeImbalance']) for d in
dicts]
        flattened_df = pd.concat([df.T.stack() for df in expanded],
axis=1).T
        flattened_df = flattened_df.reset_index(drop=True)

        # Define custom mapping according to your requirement.
flattened_df.columns=flattened_df.columns.set_levels(levels=["90sec","100msg","60sec","
level=1)

flattened_df=flattened_df.stack().T.unstack().unstack(1).unstack()

flattened_df.columns=flattened_df.columns.map("_".join).to_series().apply(lambda
x:f"{col}_"+x).tolist()

filtered_df=flattened_df.drop(pd.concat([flattened_df.filter(like="_isTimeBased"),flatt

        # Additional processing for other nested structures (Qty and
NumOrders)...

        lstFlattened.append(pd.concat([filtered_df, filtered_dfQty,
filtered_dfNumOrders], axis=1))

    #
-------------------------------------------------------------------------------
    # Final Table for ML
    # assemble final table for ML

    dftable = pd.concat([df.loc[:,
'symbol':flatten_cols[-1]].drop(columns=flatten_cols),
pd.concat(lstFlattened,axis=1)], axis=1)
    return dftable
```

*From preprocess_data.ipynb, function definition*

This function converts complex nested structures representing market states at different points in the order lifecycle into a flat table of features. The transformation process handles multiple time windows (90 seconds, 100 messages, etc.) and different types of market metrics (trade imbalance, order quantities, etc.).

*6.b. Feature Engineering Deep Dive*

Feature engineering is critical to this system's success. We transform raw market data into predictive features that capture market microstructure information.

There are several different categories of features:

1. **Order book position features** showing how distance from support and resistance levels is calculated based on order direction. These features have high predictive power for execution probability.

2. **Order dynamics features** including fill-to-display ratio and lean-over-hedge ratio, along with temporal features that capture market timing effects. These features help identify aggressive iceberg orders.

3. **Side-relative transformations** for order book imbalance and support/ resistance levels. These transformations create consistent features that work regardless of whether the order is a buy or sell.

6.b.i. *Market Structure Features:*

- `ticksFromSupportLevel`
- `ticksFromResistanceLevel`
- `highLowRange`

Features that capture market structure and price levels.

6.b.ii. *Order Book Dynamics:*

- `numReloads`
- `fillToDisplayRatio`
- `plusOneLevelSameSideMedianRatio`

Metrics derived from order book analysis that indicate aggression and intent behind iceberg orders.

6.b.iii. *Trade Imbalance:*

- `sameSideImbalance_100msg`
- `sameSideImbalance_90sec`
- `sameSideImbalance_30sec`

Measures of trading activity imbalance at different time windows.

6.b.iv. *Temporal Features:*

- `firstNoticeDays`
- `monthsToExpiry`
- `numAggressivePriceChanges`

Time-based and momentum characteristics.

*6.c. Side-Relative Transformations*

```python
# Converting bid/ask imbalances to side-relative measures
for col in bidImbalanceCols:
    df[col.replace("bid","sameSide")] = np.where(df.isBid==True,
df[col], 1-df[col])

# Creating support/resistance level features
df['ticksFromSupportLevel'] = np.where(df.isBid==True,
df['ticksFromLow'], df['ticksFromHigh'])
df['ticksFromResistanceLevel'] = np.where(df.isBid!=True,
df['ticksFromLow'], df['ticksFromHigh'])
```

*From preprocess_data.ipynb, feature transformation code*

**Example 1: Side-Relative Order Book Imbalance**

- Raw Book Data: Buy/Sell Imbalance: 0.75 (bid)
- For SELL Order: sameSideImbalance = 1 - 0.75 = 0.25 (low value indicates unfavorable book condition for sell order)

**Example 2: Support/Resistance Level Positioning**

- Raw Price Data: Current Price: 100.25, Recent High: 102.50, Recent Low: 98.75
- For BUY Order: ticksFromSupportLevel = 6 (buy order is 6 ticks from support level)

This transformation ensures that features have consistent predictive meaning regardless of the order's side.

## 7. Time Series Cross-Validation: Respecting Market Evolution

In quantitative trading, traditional cross-validation can lead to look-ahead bias. I implemented a time-series validation approach as shown in Image 13: *Image 13: Time series cross-validation approach showing how data is split into training and testing periods. This method respects the temporal nature of financial data and prevents future information leakage.*

This approach:

1. Trains on past data, tests on future data
2. Uses rolling windows that respect time boundaries
3. Prevents information leakage from future market states

```python
def _create_time_series_splits(self, train_size, dates):
    splits = []
    n = len(dates)

    for i in range(n):
        if i + train_size < n:
            train_dates = dates[i:i + train_size]
            test_dates = [dates[i + train_size]]
            splits.append((train_dates, test_dates))

    return splits
```

*From machinelearning_final_modified.py, lines 304-314*

The hyperparameter optimization process conducted 50 trials for each model type, systematically evaluating different parameter combinations across the entire dataset. Each trial represents a complete train-test evaluation with specific parameter settings.

## 8. Model Selection Strategy

For a trading system, model selection requires balancing multiple considerations:

*8.a. Model Comparison & Theory*

8.a.i. *Custom Evaluation Scoring Metric - Max Precision / Optimal (minimum required) Recall:*

The custom evaluation metric focuses on trading-specific considerations:

```python
@staticmethod
def max_precision_optimal_recall_score(y_true, y_pred):
    """
    This is a custom scoring function that maximizes precision while
    optimizing to the best possible recall.
    """
    precision = precision_score(y_true, y_pred)
    recall = recall_score(y_true, y_pred)

    min_recall = 0.5
    score = 0 if recall < min_recall else precision
    return score
```

This metric prioritizes:

1. **Precision**: Minimizing false positives (critical for avoiding bad executions)
2. **Recall with a minimum threshold**: Ensuring we don't miss significant trading opportunities (at least 50%)
3. **Balance**: The model must achieve both good precision and sufficient recall

8.a.ii. *Best Hyperparameter Tuning Trial per Model:*

Based on the hyperparameter optimization trials, here are the best performances across models:

9. LOGISTIC REGRESSION BEST TRIAL

:label: logistic-regression-best-trial

| Best Trial Performance | |
|---|---|
| Trial Number | 26 |
| Performance Score | 0.68993 |
| Start Time | 2023-11-17 15:55:06 |
| Completion Time | 2023-11-17 15:56:22 |
| Execution Duration | 76.74 seconds |
| Trial ID | 177 |

10. LOGISTIC REGRESSION BEST PARAMETERS

:label: logistic-regression-best-parameters

| Optimized Hyperparameters |
|---|

| Optimized Hyperparameters | |
|---|---|
| penalty | elasticnet |
| C | 0.01 |
| solver | saga |
| max_iter | 1000 |
| l1_ratio | 0.5 |
| train_size | 2 |

## 11.  XGBoost Best Trial

:label: xgboost-best-trial

| Best Trial Performance | |
|---|---|
| Trial Number | 21 |
| Performance Score | 0.67466 |
| Start Time | 2023-11-17 23:02:38 |
| Completion Time | 2023-11-17 23:08:38 |
| Execution Duration | 360.00 seconds |
| Trial ID | 122 |

## 12.  XGBoost Best Parameters

:label: xgboost-best-parameters

| Optimized Hyperparameters | |
|---|---|
| eval_metric | error@0.5 |
| learning_rate | 0.03 |
| n_estimators | 250 |
| max_depth | 4 |
| min_child_weight | 8 |
| gamma | 0.2 |
| subsample | 1.0 |
| colsample_bytree | 0.8 |
| reg_alpha | 0.2 |
| reg_lambda | 2 |
| train_size | 2 |

## 13.  LightGBM Best Trial

:label: lightgbm-best-trial

| Best Trial Performance | |
|---|---|
| Trial Number | 49 |
| Performance Score | 0.67457 |
| Start Time | 2023-11-18 01:42:25 |
| Completion Time | 2023-11-18 01:42:59 |
| Execution Duration | 34.48 seconds |
| Trial ID | 50 |

## 14. LightGBM Best Parameters

:label: lightgbm-best-parameters

| Optimized Hyperparameters | |
|---|---|
| objective | regression |
| learning_rate | 0.05 |
| n_estimators | 100 |
| max_depth | 4 |
| num_leaves | 31 |
| min_sum_hessian_in_leaf | 10 |
| extra_trees | true |
| min_data_in_leaf | 100 |
| feature_fraction | 1.0 |
| bagging_fraction | 0.8 |
| bagging_freq | 0 |
| lambda_l1 | 2 |
| lambda_l2 | 0 |
| min_gain_to_split | 0.1 |
| train_size | 2 |

## 15. Random Forest Best Trial

:label: random-forest-best-trial

| Best Trial Performance | |
|---|---|
| Trial Number | 46 |
| Performance Score | 0.66481 |
| Start Time | 2023-11-17 18:23:30 |
| Completion Time | 2023-11-17 18:26:19 |

| Best Trial Performance | |
| --- | --- |
| Execution Duration | 168.91 seconds |
| Trial ID | 97 |

## 16. Random Forest Best Parameters

:label: random-forest-best-parameters

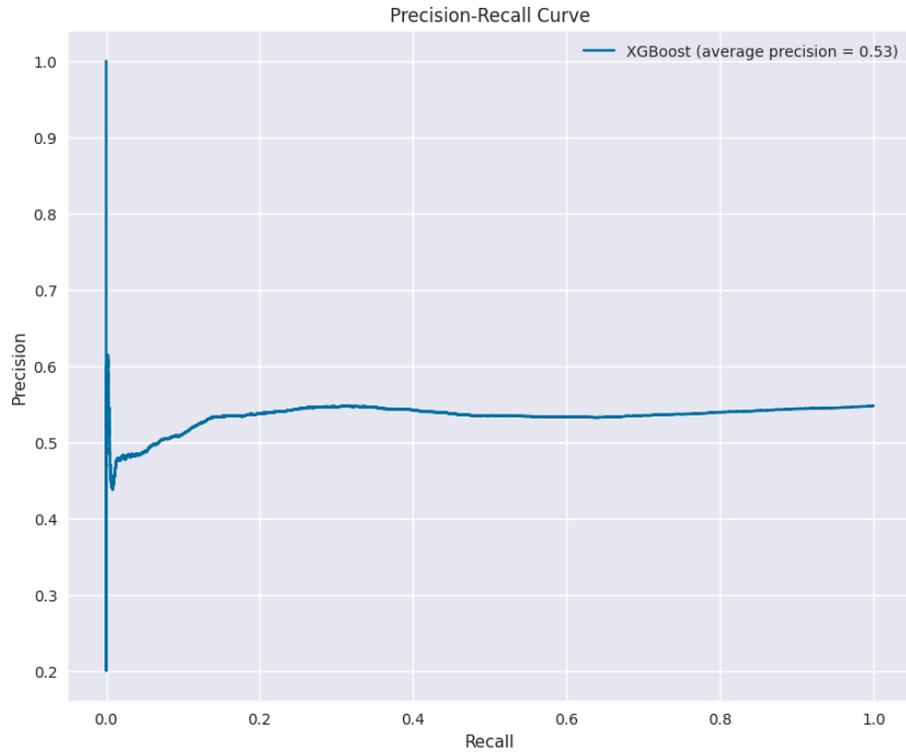| Optimized Hyperparameters | |
| --- | --- |
| n_estimators | 500 |
| max_depth | 4 |
| min_samples_split | 7 |
| min_samples_leaf | 3 |
| train_size | 2 |

The trials consistently show that a smaller training window (train_size = 2) performs better across all models, suggesting that recent market conditions are more predictive than longer historical periods.

Interestingly, while tree-based models (XGBoost, LightGBM, Random Forest) performed well, Logistic Regression achieved the highest overall score, suggesting that many of the predictive relationships in the dataset may be effectively linear once the features are properly engineered.
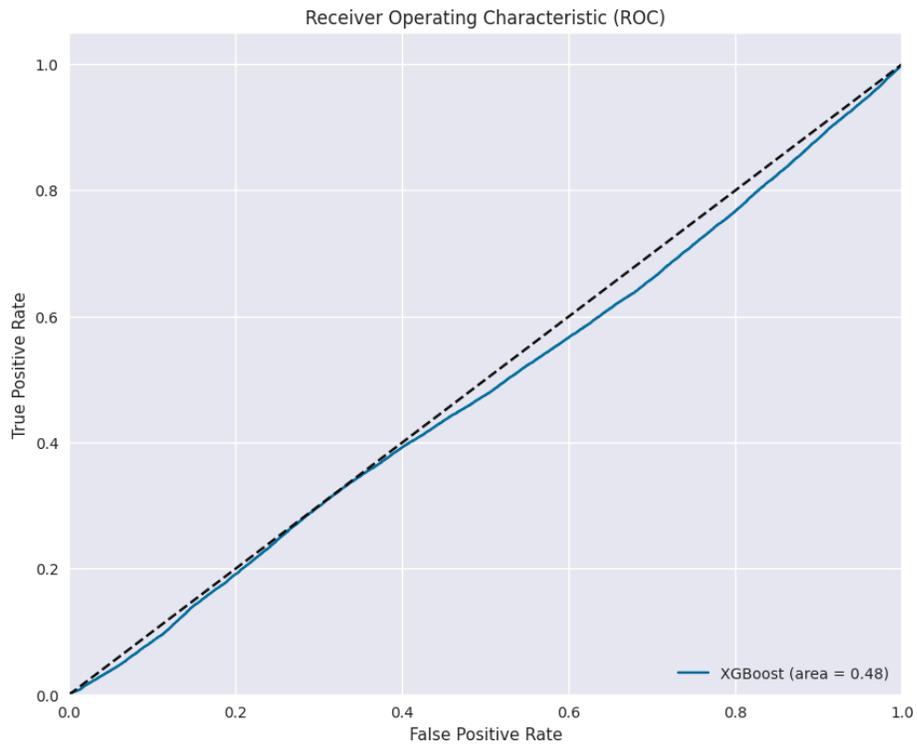
*16.a. Model HPO Results*

*16.b. Precision-Recall and ROC Analysis*

To further evaluate model performance, I analyzed precision-recall and ROC curves as shown in Images 11 and 12:

**Precision-Recall Curve**

— XGBoost (average precision = 0.53)

*Precision-recall curve showing the trade-off between precision and recall. The average precision of 0.53 indicates the model's ability to balance between capturing opportunities and avoiding false positives.*

**Receiver Operating Characteristic (ROC)**

— XGBoost (area = 0.48)

*ROC curve showing the trade-off between true positive rate and false positive rate. The area under the curve of 0.48 suggests room for further optimization.*

These curves help in understanding the model's performance across different threshold settings, which is critical for calibrating the model for different trading scenarios.

## 17. HYPERPARAMETER OPTIMIZATION: TRADING SYSTEM EFFICIENCY

For a production trading system, parameter optimization is essential. I used Optuna to tune hyperparameters:

```python
def tune_models(self, n_trials=25, hyperparameter_set_pct_size=0.5, seed
= None):
    self.hyperparameter_set_pct_size = hyperparameter_set_pct_size
    tuner = HyperparameterTuner(self, self.hyperparameter_set_pct_size)

    # set validation dates (testing of final best parameters from
hyperparameter opts)
    self.hyperparameter_set_dates =
sorted(tuner.hyperparameter_set_dates)
    self.validation_set_dates = sorted(list(set(self.unique_split_dates)
- set(self.hyperparameter_set_dates)))

    tuner.tune(n_trials, seed=seed)
    return
```

*From machinelearning_final_modified.py, lines 622-631*

The parameter search space included both model hyperparameters and training configuration:

```python
def get_model_hyperparameters(self, trial, model_name):
    if model_name == "XGBoost":
        return {
            'eval_metric': trial.suggest_categorical('eval_metric',
                ['logloss', 'error@0.7', 'error@0.5']),
            'learning_rate': trial.suggest_float('learning_rate',
                0.01, 0.05, step=0.01),
            'n_estimators': trial.suggest_categorical('n_estimators',
                [100, 250, 500, 1000]),
            'max_depth': trial.suggest_int('max_depth', 3, 5, step=1),
            'min_child_weight': trial.suggest_int('min_child_weight', 5,
10, step=1),
            'gamma': trial.suggest_float('gamma', 0.1, 0.2, step=0.05),
            'subsample': trial.suggest_float('subsample', 0.8, 1.0,
step=0.1),
            'colsample_bytree': trial.suggest_float('colsample_bytree',
0.8, 1.0, step=0.1),
            'reg_alpha': trial.suggest_float('reg_alpha', 0.1, 0.2,
step=0.1),
```

```
        'reg_lambda': trial.suggest_int('reg_lambda', 1, 3, step=1)
    }
```

*From machinelearning_final_modified.py, lines 75-89*

17.-.i. *Optimized XGBoost Model Configuration:*

The table below contains the optimized XGBoost model configuration:

| Parameter | Value | Trading Significance |
|---|---|---|
| max_depth | 4 | Shallow trees reduce overfitting to market noise, focusing on robust patterns |
| learning_rate | 0.03 | Low learning rate provides more stable predictions as market conditions evolve |
| n_estimators | 250 | Moderate number of trees balances complexity with execution speed |
| gamma | 0.2 | Minimum loss reduction required for further tree partitioning, prevents capturing random market fluctuations |
| subsample | 1.0 | Using full dataset for each tree, maximizing information when training data is limited |
| colsample_bytree | 0.8 | Each tree considers 80% of features, reducing overfitting to specific market signals |
| reg_alpha | 0.2 | L1 regularization controls model sparsity, focusing on most significant market factors |
| reg_lambda | 2 | L2 regularization prevents individual features from dominating prediction, improving stability |
| eval_metric | error@0.5 | Optimizes classification performance at the 0.5 probability threshold |
| min_child_weight | 8 | Controls complexity of each tree node, preventing overfitting to noise |

*XGBoost model architecture table showing the optimized parameter configuration and simplified tree structure visualization. Each parameter is explained in terms of its trading significance.*

## 18. FEATURE IMPORTANCE: TRADING SIGNAL ANALYSIS

Understanding which features drive prediction is critical for trading strategy development.

The feature importance analysis was conducted using multiple methods:

- **MDA (Mean Decrease Accuracy) feature importance** shows how permuting each feature affects model accuracy. The dominance of price position features (ticksFromResistanceLevel and ticksFromSupportLevel) is clearly visible.

- **XGBoost model's native feature importance** method which provides a different perspective on feature ranking compared to the MDA method.

```python
def calculate_mda(self, model, X_test, y_test, scoring=f1_score):
    """
    MDA (Mean Decrease Accuracy):
    This is a technique where the importance of a feature is evaluated
by permuting the values of the feature
    and measuring the decrease in model performance.
    The idea is that permuting the values of an important feature should
lead to a significant drop in model performance,
    indicating the feature's importance.
    """
    base_score = scoring(y_test, model.predict(X_test))
    feature_importances = {}

    for feature in X_test.columns:
        X_copy = X_test.copy()
        X_copy[feature] = np.random.permutation(X_copy[feature].values)
        new_score = scoring(y_test, model.predict(X_copy))
        feature_importances[feature] = base_score - new_score

    return feature_importances
```

*From machinelearning_final_modified.py, lines 717-732* Three insights valuable for trading strategy development:

1. **Price Position Dominance**: The distance from support/resistance levels is the strongest predictor - suggesting that order book positioning relative to key levels is crucial for execution prediction.
2. **Imbalance Significance**: Trade imbalance metrics across different time windows show strong predictive power - confirming that order flow imbalance is a leading indicator of execution probability.
3. **Temporal Sensitivity**: The model weights features from the state immediately before fill more heavily than earlier states - indicating that execution prediction becomes more accurate as we get closer to the fill event.

19. FROM PREDICTION TO TRADING DECISION

The prediction model doesn't operate in isolation - it feeds into a sophisticated trading decision process, as shown in Image 3:
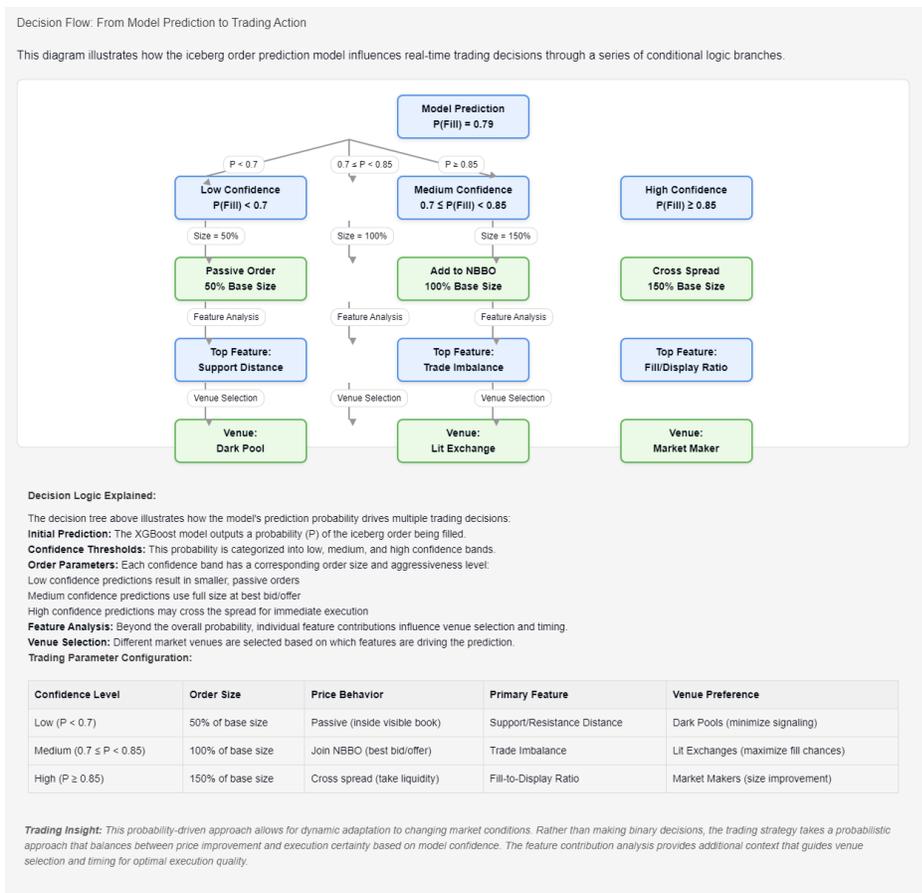
Image 3: Decision flow diagram illustrating how model predictions drive trading decisions through confidence bands and feature analysis. This probabilistic approach allows for dynamic adaptation to market conditions.

*19.a. Prediction Confidence Bands*

- **Low Confidence**: P(Fill) < 0.7 - Use smaller, passive orders
- **Medium Confidence**: 0.7 ≤ P(Fill) < 0.85 - Use full size at best bid/offer
- **High Confidence**: P(Fill) ≥ 0.85 - May cross the spread for immediate execution

*19.b. Beyond Probability: Feature Analysis*

The top features driving each prediction influence venue selection and timing:

- Support distance arrow.r Dark pool venues
- Trade imbalance arrow.r Lit exchanges
- Fill/display ratio arrow.r Market makers

This probability-driven approach allows for dynamic adaptation to changing market conditions. Rather than making binary decisions, the trading strategy takes a probabilistic approach that balances between price improvement and execution certainty.

## 20. PREDICTION FLOW

The complete prediction flow, shown below, shows the full model architecture:

## 21. Evaluation & Trading Strategy Implications

The evaluation metrics shown in Image 9 have several implications for trading strategy:

1. **XGBoost Superiority**: XGBoost consistently outperformed other models, especially taking into account the speed in which it is able to make a prediction and the flexibility it has regarding parameter optimization.
2. **Feature Transferability**: The dominance of order book position and imbalance features suggests that these signals may transfer well to other instruments beyond the ones tested.
3. **Execution Time Sensitivity**: The importance of "oneStateBeforeFill" features indicates that model prediction accuracy increases as the order approaches execution, suggesting a strategy that dynamically adjusts confidence thresholds based on order age.

## 22. Business Impact and Performance Metrics

The optimization trials demonstrate the model's effectiveness, with the best XGBoost configuration achieving a score of **0.674656** on the custom evaluation metric. The trial data reveals that shorter training windows (train_size = 2) consistently perform better across all models, indicating that recent market conditions are more predictive of execution outcomes than longer historical periods.

The custom evaluation metric (**max_precision_optimal_recall_score**) was specifically designed to balance precision (minimizing false positives) with sufficient recall (at least 50%), making it directly relevant to trading performance where both accuracy of execution signals and capturing enough opportunities are critical.

## 23. Model Persistence & Deployment Considerations

For a production trading system, model deployment requires careful handling of model artifacts:

```python
def save_and_upload_model(self, model, model_name):
    """
    Saves the given model to a file and uploads it to Neptune.
    """
    # Create a directory for saving models
    models_dir = "saved_models"
    os.makedirs(models_dir, exist_ok=True)

    # Save model for JVM export (compatible with trading systems)
    model_path = os.path.join(models_dir, f"xgbModel-
{self.current_date}.json")
    model.save_model(model_path)
```

We also export the scaling parameters to ensure consistent feature transformation in production:

```python
def write_scalers(self, model_name):
    # Extract and export scaler parameters
    scaler_params = {
        "minMaxScaler": {
            "min": self.minmax_scaler.min_.tolist(),
            "scale": self.minmax_scaler.scale_.tolist(),
            "featureNames": self.minmax_features
        },
        "standardScaler": {
            "mean": self.standard_scaler.mean_.tolist(),
            "std": self.standard_scaler.scale_.tolist(),
            "featureNames": self.non_minmax_features
        }
    }
```

This ensures that the model deployed to production uses identical scaling to what was used during training.

### 23.a. Production Integration

If implemented in a production trading system, this model would:

1. Extract real-time order book and trade imbalance features
2. Apply appropriate scaling using the persisted scaler parameters
3. Generate execution probability predictions
4. Allow trading algorithms to make more informed decisions about:

   - Order routing (to venues with higher fill probability)
   - Order sizing (increasing size when fill probability is high)
   - Aggressive vs. passive execution tactics

Additionally, the Neptune integration provides continuous monitoring capabilities to detect model drift in changing market conditions.

## 24. Potential Extensions

If I were to extend this project, I would:

1. Add market regime conditioning - adapting predictions based on volatility regimes
2. Incorporate order book depth information beyond level 1
3. Develop an adversarial model to simulate market impact of our own orders
4. Implement model confidence calibration to produce reliable probability estimates
5. Create an ensemble approach combining multiple model predictions weighted by recent performance

## 25. Conclusion: The Value of ML in Trading Systems

The complete system demonstrates how machine learning can transform trading operations by:

1. **Extracting Hidden Patterns**: The model uncovers subtle market microstructure patterns invisible to human traders
2. **Quantifying Uncertainty**: Probability-driven approaches handle the inherent uncertainty in financial markets
3. **Adaptive Decision-Making**: Dynamic parameter adjustment based on model confidence creates a more robust strategy
4. **Continuous Improvement**: Feedback loops ensure the system improves as market conditions evolve

This iceberg order prediction system represents a sophisticated approach to quantitative trading that goes beyond simple signal generation to create a complete trading ecosystem that balances predictive power with practical trading considerations.