

Presentation of the Foreign Exchange (FX) Market

1 Introduction

The Foreign Exchange (FX) market, also known as the forex market, is the largest and most liquid financial market in the world, with a daily trading volume exceeding \$6 trillion. It facilitates the exchange of currencies between participants globally. The FX market operates over-the-counter (OTC), meaning that trades are conducted directly between parties, typically via electronic trading platforms, without a centralized exchange.

2 Participants in the FX Market

2.1 Central Banks

Central banks, such as the Federal Reserve (USA), European Central Bank (ECB), and the Bank of Japan (BoJ), are key players in the FX market. They influence currency prices through monetary policy, interest rate decisions, and currency interventions.

2.2 Commercial Banks

Commercial and investment banks are major liquidity providers, facilitating currency transactions for clients and other market participants. They also engage in speculative trading and risk management for their own accounts.

2.3 Corporations

Large multinational corporations participate in the FX market to hedge currency risk arising from international trade and cross-border transactions. For instance, a U.S.-based company may need to exchange U.S. dollars (USD) for euros (EUR) to pay a European supplier.

2.4 Hedge Funds and Speculators

Hedge funds and speculative traders engage in FX trading with the goal of profiting from short-term movements in exchange rates. They add liquidity to the market and contribute to price discovery but are primarily focused on generating returns.

2.5 Retail Traders

Individual retail traders participate in the FX market through brokerage platforms. Although they account for a smaller portion of the overall trading volume, retail trading has grown due to the availability of online platforms and leverage.

2.6 Market Makers

FX market makers are financial institutions or firms that provide liquidity by offering buy (bid) and sell (ask) quotes for currency pairs. They enable continuous trading by being ready to either buy or sell at publicly quoted prices.

2.7 Governments and Sovereign Wealth Funds

Governments and sovereign wealth funds engage in FX transactions for reserves management and international investments.

3 The Role of an FX Market Maker

A market maker in the FX market is responsible for providing liquidity by continuously quoting buy and sell prices for currency pairs. Market makers operate under two main objectives:

3.1 Liquidity Provision

They ensure that there is always a market for currencies by standing ready to buy or sell, regardless of the market's immediate direction. This helps keep the market functioning smoothly.

3.2 Spreads and Profit Generation

Market makers profit by capturing the spread between the bid (buy) and ask (sell) prices. They manage risk by carefully balancing their books, hedging positions, and analyzing market conditions.

4 Importance of Modeling Client Flow

For a market maker, it is critical to understand and model client order flow. Client order flow refers to the volume and types of transactions made by clients, including both speculative trades and hedging activities. Here's why it matters:

4.1 Market Impact and Pricing

The volume and direction of client flows can affect exchange rates. For instance, if there is a large demand to buy EUR/USD, this can push the price of EUR/USD higher. By modeling client flows, market makers can better anticipate these moves and adjust their pricing to remain competitive while managing their risk.

4.2 Risk Management

Knowing the size and timing of client trades helps market makers manage their own risk exposure. By analyzing order flows, they can hedge positions more effectively and reduce the impact of large trades on their books.

4.3 Profitability

By predicting the flow of client orders, market makers can adjust their spreads and liquidity provision to optimize profit. For instance, if a market maker expects a surge in demand for a particular currency, they might widen their spreads or position themselves in advance to capitalize on the move.

4.4 Behavioral Analysis

Analyzing client flows also provides insights into market sentiment and the potential behavior of different participants. For example, if retail clients are aggressively buying a currency, this might signal bullish sentiment, but if institutional clients start selling, it could indicate a shift in expectations. Understanding these patterns can help market makers adjust their strategies accordingly.

5 Why It's Important to Predict the Market Impact of Client Trades

5.1 Adverse Selection Risk

Market makers face the risk of adverse selection, where clients with better information take positions that leave the market maker exposed to losses. Predicting the impact of trades can help market makers mitigate this risk by adjusting quotes based on expected market reactions.

5.2 Liquidity Management

Large client trades can significantly affect market liquidity. By predicting the impact of these trades, market makers can manage liquidity provision more efficiently, ensuring that they can absorb large trades without suffering major slippage or price disruption.

5.3 Reducing Slippage

Slippage refers to the difference between the expected price of a trade and the actual execution price. Large client orders can cause slippage, especially in less liquid currency pairs. Predicting the impact of client trades allows market makers to adjust prices proactively, reducing the likelihood of slippage and protecting both their clients and themselves.

5.4 Improving Client Relationships

A market maker that understands and predicts the market impact of client trades can offer better pricing and execution. This enhances trust and helps in maintaining long-

term relationships with key clients, such as institutional investors and hedge funds.

6 Conclusion

In summary, FX market makers play a critical role in maintaining liquidity and enabling efficient trading in the forex market. By modeling and analyzing client flows, market makers can manage risk, improve profitability, and anticipate the market impact of trades, which is crucial for staying competitive in a highly dynamic and fast-paced environment.

I am writing my internship report on client flow modelling and order flow prediction on FX market. Here is a part of my report. Write this in latex, make it more compact with paragraph Studies on Order Flow in Asset Classes like Equities Order flow modeling in financial markets has been extensively studied, particularly in the context of equities, where the order book structure plays a central role. In equity markets, order books provide detailed information on the buy and sell orders placed by market participants. These studies focus on:

Price Discovery: Order flow plays a crucial role in the price discovery process. By analyzing the size, frequency, and timing of buy and sell orders, researchers can infer how new information is incorporated into stock prices.

Market Microstructure: Equity markets are structured around centralized exchanges with transparent order books, which allows for detailed analysis of liquidity, volatility, and market impact. Models like the Easley and O'Hara model and the Kyle model analyze how informed and uninformed traders impact price and liquidity.

Market Impact of Trades: Studies have shown that large trades impact prices, especially when liquidity is thin. Predicting the market impact of large orders is essential for minimizing transaction costs. In equities, researchers have developed models to forecast the impact of block trades on prices and liquidity.

Non-Random Behavior of Market Participants: The order flow in equity markets is not random. Different market participants, such as retail traders, institutional investors, and market makers, have different objectives and strategies. For example:

Institutional investors typically break up large orders to avoid signaling their intentions to the market. Retail traders might act based on momentum or sentiment. Market makers provide liquidity and attempt to maintain a balanced order book while profiting from the bid-ask spread. These non-random behaviors have been well-documented, leading to a better understanding of how different players contribute to price dynamics and market efficiency in equity markets.

Why Few Studies Have Focused on the FX Market In contrast to the equity market, fewer studies have focused on order flow and client behavior in the FX market. Several reasons explain this gap:

Decentralized Nature of the FX Market: Unlike equities, the FX market is largely decentralized. It operates over-the-counter (OTC), with no centralized exchange or public order book. Instead, trading happens through a network of banks, financial institutions, and electronic platforms. This lack of transparency in order books makes it difficult to gather granular data, which is essential for in-depth order flow analysis.

Request-for-Quote (RFQ) System: A significant portion of FX trading happens through the Request-for-Quote (RFQ) system, where clients request quotes from dealers before executing trades. Unlike exchange-traded assets with a continuous auction system, the RFQ system is more opaque, limiting the availability of real-time order book data for researchers.

Complexity of Client Behavior: FX participants range from central banks and corporations to retail traders, each with different motives, strategies, and access to information. This diversity makes modeling order flow more challenging in FX than in equities, where participants' actions may be more homogeneous.

Liquidity and Market Depth: The FX market is often considered more liquid than equities, especially in major currency pairs. This deep liquidity can mute the impact of individual trades, making it more challenging to isolate and predict market reactions based on order flow.

Problematic of Client Flow Modelization and Order Flow Prediction in FX The challenges of client flow modelization and order flow prediction in FX markets raise several key questions and problems that need to be addressed:

Market Impact and Liquidity Management

How do large client trades affect FX prices in real-time?: While the FX market is deep and liquid, large institutional trades can still impact prices. Predicting this impact is crucial for market makers to avoid losses and manage liquidity effectively. What factors drive price movements in illiquid or exotic currency pairs?: While major currency pairs (like EUR/USD) are liquid, emerging market or exotic currencies are not. Predicting market impact in less liquid markets requires a different approach. Predicting Order Flow Behavior

How do client types differ in their trading behavior?: Institutional clients, central banks, hedge funds, and retail traders have different objectives and strategies. Can these behaviors be systematically modeled to predict future order flow? How does the FX market respond to macroeconomic news and geopolitical events?: In FX, external factors like central bank policies, interest rate changes, and geopolitical risks can trigger large flows. How do these exogenous events shape client behavior and subsequent price movements? Adverse Selection and Information Asymmetry

How can market makers detect informed trading?: Adverse selection arises when market makers cannot distinguish between informed and uninformed clients. This problem is exacerbated in FX, where some clients (e.g., hedge funds) may have access to superior information. How can market makers use order flow data to mitigate this risk? How can we model the flow of orders in the presence of information asymmetry?: Predicting market impact becomes difficult when market participants have different levels of information about future currency movements. Risk Management and Hedging

How can market makers effectively hedge their positions based on client order flow?: FX market makers need to continuously hedge their risk exposure to avoid significant losses from adverse price movements. How can they predict client order flow to improve their risk management strategies? How do market makers adjust their spreads and liquidity provision in response to order flow?: Predicting when large flows will come into the market is crucial for setting appropriate spreads. How do market makers adjust pricing

strategies based on client flows? Impact of Technological Advancements

How do algorithmic and high-frequency traders (HFT) influence order flow dynamics?: In recent years, the rise of algorithmic and high-frequency trading in FX has changed the nature of market participation. How do these traders affect order flow predictability and market impact? Can machine learning and AI improve order flow prediction models?: The increasing availability of data and computational power allows for more advanced modeling techniques. How can market makers leverage these technologies to improve client flow predictions? RFQ Systems and Data Limitations

How can market makers optimize their response to Request-for-Quote (RFQ) systems?: RFQ-based trading provides limited visibility into order flow, but it is still essential to predict market impact from quotes. What data can be extracted from RFQ interactions to improve order flow models? Can we develop new ways to track client behavior without a centralized order book?: Given the decentralized nature of the FX market, market makers must rely on fragmented data from different trading venues. How can they use this partial information to better model and predict client behavior? In conclusion, while much of the literature on order flow and market impact is focused on equities, the FX market presents unique challenges due to its decentralized structure, diversity of participants, and RFQ trading mechanisms. Understanding and modeling client order flow is critical for market makers in managing risk, optimizing pricing strategies, and improving liquidity provision, but it raises a host of questions about how to deal with the complexities inherent in this global and multi-faceted market.

Method and Analysis Introduction The Foreign Exchange (FX) market is one of the largest and most liquid financial markets globally, with daily trading volumes exceeding 6trillion. *Thisimmensescalepresentsuniquechallengesandopportunitiesforanalyzingclientimpactson activities[Reference : StudyonClientCo-activities], weaimtorefineourunderstandingofclientbehavi*

Quantifying Client Activity: The Imbalance Ratio Defining the Imbalance Ratio A key aspect of our analysis involves defining a quantitative measure to capture client activity. The imbalance ratio serves as a critical indicator of trading behavior. It is calculated as follows:

Imbalance Ratio = $\frac{\text{Buy Orders} - \text{Sell Orders}}{\text{Buy Orders} + \text{Sell Orders}}$ Imbalance Ratio = $\frac{\text{Buy Orders} - \text{Sell Orders}}{\text{Buy Orders} + \text{Sell Orders}}$

The imbalance ratio provides a measure of the net trading position of a client, reflecting whether their trades are predominantly buy or sell orders. High positive values indicate a buying bias, while high negative values indicate a selling bias. This metric is valuable in understanding the intensity and directionality of a client's trading activity [Reference: Research on Imbalance Ratio in High-Frequency Trading].

Practical Implications The imbalance ratio's utility extends beyond mere quantification. For example, in the context of market making, understanding the imbalance ratio can help in predicting price movements and managing risk. Clients with high imbalance ratios may exert greater influence on market prices, making them crucial for developing effective trading strategies.

Clustering Client Activities: Methodological Challenges Traditional Clustering Techniques Once client activity is quantified, the next step is to group clients based on their trading behaviors. Traditional clustering methods, such as k-means, are commonly used

for this purpose. K-means clustering relies on minimizing the within-cluster variance, measured using the Euclidean distance between data points. However, this method can struggle with capturing the complex patterns of client activities, especially when these activities exhibit synchronicity or asynchronicity [Reference: Limitations of k-means in Financial Clustering].

Limitations of Euclidean Distance Euclidean distance assumes that clusters are spherical and of similar size, which may not align with the diverse patterns observed in FX trading. For instance, clients with highly correlated trading activities may form distinct clusters that are not easily captured by k-means. This limitation necessitates the use of alternative clustering approaches that can better accommodate the nuances of client behavior.

Advanced Clustering with Statistically Validated Networks (SVNs) Overview of SVN To overcome the limitations of traditional clustering methods, we employ Statistically Validated Networks (SVNs). SVN approaches leverage statistical validation techniques to identify meaningful clusters of clients based on their trading activities. Unlike k-means, SVN does not rely on rigid assumptions about cluster shapes or distances. Instead, it uses statistical measures to assess the validity of clusters and their stability over time [Reference: SVN Methodology and Applications].

Methodological Advantages SVN offers several advantages over traditional methods:

Flexibility: SVN can accommodate various cluster shapes and sizes, making it more adaptable to the complex patterns in trading data. **Robustness:** By using statistical validation, SVN reduces the risk of overfitting and provides more reliable cluster assignments. **Insightfulness:** SVN can reveal underlying structures in client activities that may not be apparent with simpler methods. **Analyzing Market Impact: Challenges and Considerations** The Market Impact Phenomenon Market impact refers to the effect that a trade has on the market price of a security. In the FX market, understanding market impact involves analyzing how client trades influence price movements. While extensive research has been conducted on market impact in equities, studies focusing on FX markets are relatively scarce [Reference: Market Impact Studies in Equities vs. FX].

Liquidity and Trade Volume One critical factor affecting market impact is liquidity. Highly liquid currencies, such as EUR/USD, tend to exhibit smaller price impacts per trade due to the large volume of transactions. In contrast, less liquid currencies, such as USD/JPY or GBP/USD, may experience more significant price changes from individual trades. The availability of trade data and the volume of trades within the dataset play a crucial role in determining the accuracy of market impact predictions [Reference: Liquidity and Market Impact in FX].

To illustrate the influence of liquidity, consider BNP Paribas, a major European FX trader. BNP Paribas's significant trading volume in EUR/USD provides a rich dataset for analyzing market impact. In contrast, trading in less liquid currencies may result in less precise predictions due to the lower volume of data. This example highlights the importance of considering liquidity and trade volume when evaluating market impact.

Predictive Modeling: Tree-Based Models Advantages of Tree-Based Models For predicting market impact and price changes, we utilize Tree-Based Models, specifically LightGBM (LGBM). Tree-Based Models are well-suited for capturing complex, non-linear relationships in data. Unlike linear models, which may struggle to model intricate interactions,

Tree-Based Models can effectively handle the non-linearities inherent in FX trading data [Reference: Advantages of Tree-Based Models in Financial Prediction].

LightGBM (LGBM) LGBM is chosen for its efficiency and scalability in large datasets. It employs gradient boosting techniques to build an ensemble of decision trees, optimizing performance through advanced algorithms. LGBM's ability to handle large volumes of data and its robustness in modeling non-linear relationships make it an ideal choice for predicting market impacts

Implementation and Results **Model Calibration** The LGBM model is calibrated using historical trading data, incorporating features such as trade size, timing, client behavior, and liquidity conditions. The model's performance is evaluated based on its accuracy in predicting price changes and its sensitivity to variations in liquidity.

Performance Evaluation Preliminary results indicate that the LGBM model provides a robust framework for analyzing market impact. Performance metrics demonstrate its effectiveness in capturing the complex dynamics of the FX market. The model's ability to incorporate non-linearities and interactions among features enhances its predictive accuracy, offering valuable insights for market participants.

Conclusion The methodology presented in this report provides a comprehensive approach to understanding client impacts in the FX market. By defining client activity through the imbalance ratio, employing SVN for advanced clustering, and utilizing Tree-Based Models for prediction, we gain deeper insights into market dynamics. This approach addresses the limitations of traditional methods and leverages advanced statistical and machine learning techniques to deliver actionable insights for market analysis and strategy development.

7 Dataset Overview

Our analysis leverages an extensive dataset of trades executed by NP, encompassing the second semester of 2022. This dataset covers major global financial markets, including the US, European, and Asian markets, thus providing a comprehensive view of trading activities across various regions.

The dataset represents a rich source of information about trading behavior in the Foreign Exchange (FX) market. It captures the high-volume nature and diversity of FX trading, reflecting numerous aspects of market liquidity and client activity.

8 Trade Volume and Liquidity

8.1 Volume Distribution

The dataset includes a substantial number of trades across a wide array of currency pairs. Given NP's extensive trading activities, the dataset offers insights into different liquidity conditions for various currencies. For example, it is crucial to compare high liquidity pairs such as EUR/USD with lower liquidity pairs like USD/JPY.

To visualize trade volume distribution, we utilize histograms and kernel density plots. These graphical representations illustrate the spread of trade volumes across currency

pairs, enabling us to assess the relative trading activity and liquidity of each pair.

8.2 Liquidity Comparison

Liquidity comparison is critical for understanding how trade volumes relate to the typical values observed in the order book. By comparing trade volumes and liquidity metrics across currency pairs, we can uncover patterns related to market depth and price stability.

We expect higher trade volumes and liquidity for major pairs like EUR/USD compared to more exotic ones. A graph comparing average trade volumes and liquidity metrics for various currency pairs provides insights into how liquidity impacts market behavior.

9 Trade Characteristics and Data Quality

9.1 Attributes of Trades

The dataset captures several attributes for each trade, including:

- **Timestamp:** The exact time when the trade occurred.
- **Currency Pair:** The specific FX pair involved in the trade.
- **Trade Size:** The volume of the trade.
- **Trade Direction:** Whether the trade was a buy or sell action.

These attributes provide a detailed view of trading activities. However, it is important to address data quality issues, such as the presence of negligible trades. These trades, which have minimal impact or very small volumes, constitute approximately X% of the total trades. For a more focused analysis, we will exclude these negligible trades.

9.2 Quantile Filtering

To ensure our analysis concentrates on significant trades, we filter the dataset to include only those trades above the 1% quantile. This approach helps in focusing on trades that have a meaningful impact on market dynamics.

Visualizing the distribution of trades and applying the 1% quantile threshold can be depicted using cumulative distribution functions (CDFs) or threshold-based histograms. This visualization highlights the concentration of impactful trades while filtering out less significant data.

10 Client Activity and Trade Distribution

10.1 Client Trade Distribution

The dataset reveals a non-uniform distribution of trades among clients. Certain clients exhibit notably higher trading volumes and market impact compared to others. This variance in client activity suggests that some clients play a more substantial role in influencing market dynamics.

To illustrate this, we can use bar charts or pie charts to display the trade volume distribution among the top clients. For example, a chart showing the top 10 clients by trade volume will highlight which clients are most influential.

10.2 Focus on Top Clients

To streamline the analysis, we will focus on the top 100 clients based on their trade volume. This selection allows for a more in-depth examination of the most influential market participants and their trading behaviors.

Providing a detailed list or chart of the top 100 clients, sorted by trade volume, will offer insights into the key contributors to market activity. This focused approach ensures that the analysis is both manageable and insightful.

11 Conclusion

The dataset provides a robust foundation for analyzing client behavior, liquidity dynamics, and trading impact in the FX market. By focusing on substantial trades and key market players, our analysis will uncover patterns and insights into market activities. The next phase will involve clustering these top clients to reveal deeper insights into their trading behavior and its impact on market dynamics.

12 Introduction to Clustering Methods

Clustering is a fundamental technique in data analysis used to group similar entities based on their attributes. In the context of order flow modeling, clustering plays a crucial role in understanding and grouping clients based on their trading activities. This section explores two classic clustering methods: K-Means and Spectral Clustering, which are commonly employed in such analyses.

13 K-Means Clustering

13.1 Overview

K-Means clustering is a widely used method that partitions data into k clusters, with each data point assigned to the cluster with the nearest mean. The algorithm aims to minimize the within-cluster variance, also known as the Mean Squared Error (MSE), making it effective for well-separated clusters in a Euclidean space.

13.2 Algorithm Steps

The K-Means algorithm operates through the following steps:

1. **Initialize k centroids randomly:** Select k initial points as the starting centroids for the clusters.
2. **Assign each data point to the nearest centroid:** Each data point is assigned to the cluster with the closest centroid based on Euclidean distance.

3. **Update the centroids:** Recalculate the centroid of each cluster by computing the mean of all data points assigned to that cluster.
4. **Repeat steps 2 and 3 until convergence:** Continue assigning points and updating centroids until the centroids no longer change significantly or a maximum number of iterations is reached.

13.3 Mathematical Formulation

The objective function for K-Means is to minimize the following cost function:

$$J = \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2 \quad (1)$$

where:

- C_i represents the set of points in cluster i ,
- μ_i is the mean of points in cluster i ,
- $\|\cdot\|^2$ denotes the squared Euclidean distance.

13.4 Graphical Representation

A typical visualization of K-Means clustering shows data points partitioned into clusters with centroids marked. However, K-Means struggles with non-spherical clusters and varying densities, which can be problematic in high-dimensional spaces or when clusters have different shapes.

13.5 Example

Consider clustering clients based on their trading volumes and imbalance ratios. While K-Means may effectively group clients with similar trading volumes, it might not capture the complex, non-linear relationships in trading patterns, especially if the data is not uniformly distributed.

14 Spectral Clustering

14.1 Overview

Spectral Clustering is a more advanced technique that leverages graph theory to perform clustering. It constructs a similarity graph of the data and then applies dimensionality reduction to this graph before performing clustering. This method is particularly useful for detecting clusters that are not necessarily convex or spherical.

14.2 Algorithm Steps

The Spectral Clustering algorithm follows these steps:

1. **Construct a similarity graph G :** Build a graph where nodes represent data points and edges represent similarities between them.
2. **Compute the Laplacian matrix L :** Calculate the Laplacian matrix L from the adjacency matrix A of the graph, where $L = D - A$ and D is the degree matrix.
3. **Compute eigenvectors of L :** Find the eigenvectors of the Laplacian matrix and form a new matrix with these eigenvectors.
4. **Apply K-Means clustering:** Use K-Means clustering on the rows of the matrix formed from the eigenvectors to identify clusters.

14.3 Mathematical Formulation

Spectral Clustering uses the Laplacian matrix L , which is defined as:

$$L = D - A \tag{2}$$

where D is the degree matrix (a diagonal matrix where each element D_{ii} represents the sum of the weights of edges connected to vertex i), and A is the adjacency matrix of the graph.

14.4 Graphical Representation

Spectral Clustering can capture complex cluster structures and is particularly effective for non-convex shapes. Visualizations often reveal data points clustered in a manner that reflects underlying structures, which may not be apparent with K-Means.

14.5 Example

For clustering clients based on trading patterns, Spectral Clustering may uncover intricate structures in the data, such as distinct trading strategies or patterns that are not easily separable with K-Means.

15 Challenges in Order Flow Modeling

In the context of order flow modeling, traditional clustering methods encounter several challenges:

15.1 Metric Selection

Both K-Means and Spectral Clustering are sensitive to the choice of distance metrics. K-Means uses Euclidean distance, which may not capture the full complexity of trading activities, especially when patterns are non-linear or involve varying scales. Spectral Clustering relies on the similarity graph, which requires careful construction and choice of similarity measure.

15.2 Data Distribution

The distribution of trading data can be highly non-uniform and exhibit varying densities. K-Means assumes spherical clusters and equal cluster sizes, which may not align with the real-world distribution of client activities. This assumption can lead to suboptimal clustering results when dealing with diverse trading behaviors.

15.3 Temporal Dynamics

Trading behaviors are dynamic and can change over time. Static clustering methods may fail to capture these temporal dynamics, potentially missing significant trends or shifts in client activities. Clustering methods that do not account for temporal changes might not reflect the current state of market dynamics.

16 Mathematical Definition of Client Activity

To address these challenges, we define the activity of a client i as a function of their imbalance ratio, which is calculated for each trade and aggregated over a specific time period.

16.1 Signed USD Amount Calculation

The signed USD amount for a trade is given by:

$$\text{signedUSDAmount} = \text{Buy} \times \text{USDAmount} - \text{Sell} \times \text{USDAmount} \quad (3)$$

where:

- Buy and Sell indicate the direction of the trade,
- USDAmount is the amount of USD traded.

16.2 State Definition

For each time slice $I = [t, t + S]$, the client state a_i is defined based on the signed USD amount:

$$a_i = \begin{cases} 1 & \text{if signedUSDAmount} \geq p \\ 0 & \text{if signedUSDAmount} < p \end{cases} \quad (4)$$

where p is a threshold value determining the cutoff for significant trading activity.

16.3 Graphical Representation

A plot of client states over time can illustrate how different clients' activities evolve. This visualization can highlight periods of high or low trading intensity, providing insights into trading behavior patterns.

17 Conclusion

In summary, while K-Means and Spectral Clustering provide foundational approaches to clustering, they face limitations in the context of order flow modeling. Challenges related to metric selection, data distribution, and temporal dynamics underscore the need for more sophisticated methods. By defining client activity through the imbalance ratio and state definitions, we establish a robust framework for understanding and clustering client activities. Future work may involve developing or adapting clustering methods to better handle the complexities of trading data and dynamic market conditions.

Statistically Validated Network (SVN): A Comprehensive Approach

18 Introduction to Statistically Validated Networks

In financial data analysis, particularly for complex datasets such as those involved in order flow modeling, traditional clustering techniques like K-Means and Spectral Clustering often fall short due to their inherent limitations and assumptions. To address these challenges, the concept of Statistically Validated Networks (SVNs) has emerged as a robust method for clustering and network analysis. SVN leverages statistical validation techniques to assess the reliability and significance of the detected structures within a network, making it particularly valuable in analyzing high-dimensional and temporal data typical in financial markets.

19 The Concept of Statistically Validated Networks

Statistically Validated Networks (SVNs) are designed to address the need for rigorous statistical validation in network analysis. This approach is particularly useful in financial contexts where understanding the relationships between entities, such as traders or clients, is crucial. SVN involves the following key components:

19.1 Network Construction

The first step in SVN involves constructing a network where nodes represent entities (e.g., clients or traders) and edges represent the relationships or interactions between these entities. In the context of order flow, an edge might represent a correlation or similarity in trading behavior between two clients.

19.1.1 Metrics for Network Construction

The construction of the network relies on various metrics that quantify the relationships between entities. Key metrics include:

- **Imbalance Ratio:** This measures the difference between buy and sell orders, providing insight into the trading behavior of each client. It is calculated as:

$$\text{Imbalance Ratio} = \frac{\text{Buy Volume} - \text{Sell Volume}}{\text{Buy Volume} + \text{Sell Volume}} \quad (5)$$

- **Correlation Measures:** These quantify how similarly two clients trade over time. For instance, the correlation coefficient between the trading activities of two clients i and j can be calculated as:

$$w_{ij} = \text{corr}(x_i, x_j) \quad (6)$$

where $\text{corr}(x_i, x_j)$ denotes the Pearson correlation coefficient between the trading activities of clients i and j .

19.2 Statistical Validation

Unlike traditional clustering methods that rely solely on distance metrics, SVN incorporates statistical techniques to validate the significance of the detected clusters or network structures. This involves comparing the observed network structures against null models to determine if the detected patterns are statistically significant.

19.2.1 Null Models

Null models are used to generate random networks with similar properties to the observed network. Common null models include:

- **Erdos-Renyi Model:** This model generates random graphs where each edge is included with a fixed probability.
- **Configuration Model:** This model preserves the degree distribution of the observed network while randomly connecting nodes.

By comparing the observed network to these null models, researchers can assess whether the observed structures are due to chance or represent significant patterns.

19.2.2 Permutation Tests

Permutation tests involve permuting the data and recalculating network metrics to determine if the observed structures are significantly different from what would be expected under random conditions. For example, if G_{obs} represents the observed network and G_{null} represents the null network, a permutation test might involve calculating the distribution of a metric (e.g., modularity) across multiple null networks and comparing it to the observed metric.

19.2.3 Statistical Significance

Measures such as p-values can be used to assess the significance of the observed network structures. For instance, if a particular cluster appears significantly more cohesive in the observed network compared to null networks, this cluster is considered statistically significant. Modularity Q is a common measure of cluster significance, defined as:

$$Q = \frac{1}{2m} \sum_{i,j} \left[A_{ij} - \frac{k_i k_j}{2m} \right] \delta(c_i, c_j) \quad (7)$$

where:

- A_{ij} is the adjacency matrix element representing the weight of the edge between nodes i and j ,
- k_i and k_j are the degrees of nodes i and j ,
- m is the total number of edges,
- $\delta(c_i, c_j)$ is 1 if nodes i and j are in the same cluster and 0 otherwise.

19.3 Stability Analysis

Stability analysis is crucial for assessing the reliability of the detected network structures over time. In financial markets, where trading patterns can evolve rapidly, it is essential to ensure that the identified clusters or network structures are not transient artifacts but reflect persistent underlying patterns.

19.3.1 Temporal Stability

Temporal stability involves examining how network structures change over different time periods. For example, one might compare network structures in different months or quarters to determine if the detected clusters remain consistent. Techniques such as dynamic community detection can be employed to track how clusters evolve over time.

19.3.2 Cross-Validation

Cross-validation techniques involve dividing the dataset into multiple folds and performing clustering on each fold. The stability of clusters across different folds can provide insights into their robustness. This approach helps to verify if the detected clusters are consistent and reliable across various subsets of data.

20 Comparative Advantages of SVN over K-Means and Spectral Clustering

Statistically Validated Networks (SVNs) offer several advantages over traditional clustering methods such as K-Means and Spectral Clustering, particularly in the context of financial data analysis and order flow modeling. Here, we compare SVN with these methods across several dimensions.

20.1 Handling Complex Data Structures

20.1.1 K-Means Clustering

K-Means clustering is effective for partitioning data into well-separated, spherical clusters. However, it struggles with complex data structures due to its reliance on Euclidean distance and the assumption of equal cluster sizes and shapes. This can be problematic in financial data, where clusters may not be spherical or may vary in density.

20.1.2 Spectral Clustering

Spectral Clustering is more flexible than K-Means and can handle non-convex clusters by using graph-based approaches. It excels at capturing complex structures but may still face challenges with high-dimensional data and the choice of similarity measures. Additionally, it requires the computation of the Laplacian matrix, which can be computationally intensive.

20.1.3 Statistically Validated Networks (SVNs)

SVNs are designed to handle complex data structures by incorporating statistical validation techniques. SVN does not rely solely on distance metrics but instead assesses the significance of network structures through comparison with null models and permutation tests. This allows SVN to detect and validate complex, non-spherical clusters and relationships that may not be captured by K-Means or Spectral Clustering.

20.2 Validation and Robustness

20.2.1 K-Means Clustering

K-Means clustering does not inherently include a validation mechanism for the detected clusters. The results are sensitive to initial centroid positions and the choice of k , which can lead to variability and potential overfitting.

20.2.2 Spectral Clustering

Spectral Clustering can provide some validation through eigenvalue analysis, but it lacks a rigorous statistical framework for assessing the significance of detected clusters. It also requires careful selection of parameters and similarity measures.

20.2.3 Statistically Validated Networks (SVNs)

SVNs incorporate statistical validation methods to assess the reliability and significance of network structures. By using null models and permutation tests, SVN provides a robust framework for validating clusters and network patterns. This statistical rigor helps ensure that detected structures are not artifacts of random fluctuations but reflect meaningful relationships in the data.

20.3 Temporal Dynamics and Stability

20.3.1 K-Means Clustering

K-Means clustering is typically static and does not account for temporal dynamics. It may not effectively capture changes in trading patterns over time.

20.3.2 Spectral Clustering

Spectral Clustering is also static and may not address temporal dynamics unless specifically adapted for time-series analysis. It may require additional techniques to track changes in clusters over time.

20.3.3 Statistically Validated Networks (SVNs)

SVNs explicitly address temporal dynamics through stability analysis. By examining network structures over different time periods and using dynamic community detection, SVN can track changes and ensure that detected clusters reflect persistent patterns. This capability is crucial for analyzing evolving trading behaviors in financial markets.

21 Conclusion

Statistically Validated Networks (SVNs) provide a comprehensive and statistically rigorous approach to clustering and network analysis, addressing many of the limitations of traditional methods such as K-Means and Spectral Clustering. SVN's emphasis on statistical validation, network construction based on meaningful metrics, and stability analysis makes it particularly well-suited for complex and dynamic datasets in financial markets.

While K-Means and Spectral Clustering offer valuable insights into data structures, SVN's ability to handle complex, non-spherical clusters, validate detected patterns statistically, and analyze temporal dynamics provides a more robust framework for understanding and clustering client activities in order flow modeling. Future research may further refine SVN techniques and explore their application to other domains where high-dimensional and temporal data pose significant challenges.

22 Introduction

Statistically Validated Networks (SVNs) offer a powerful framework for analyzing imbalance ratio activity in client trading. This approach transforms trade data into a format amenable to network analysis, constructs networks based on this transformed data, and statistically validates the resulting network links to ensure their significance. This methodology provides a comprehensive way to uncover meaningful patterns and relationships within financial trading activities.

23 1. Data Preprocessing and Network Generation

The initial step in applying SVN involves preprocessing the trading data to define the imbalance ratio activity for each client. This process includes several critical steps:

23.1 1.1 Calculating Signed USD Amount

For each trade, we compute the signed USD amount to quantify the net buying or selling activity of a client. This is given by:

$$\text{signedUSDAmount} = (\text{Buy} \times \text{USDAmount}) - (\text{Sell} \times \text{USDAmount}), \quad (8)$$

where:

- **Buy** and **Sell** are binary indicators for whether the transaction was a buy or sell,
- **USDAmount** represents the dollar amount of the transaction.

23.2 1.2 Defining Client States

For each time slice $I = [t, t + S]$, where S is the sampling period, we categorize each client's trading behavior into different states based on their signed USD amount. The client's state a_i is defined as follows:

$$a_i = \begin{cases} 0 & \text{if signedUSDAmount} < p, \\ 1 & \text{if signedUSDAmount} \geq p, \end{cases} \quad (9)$$

where p is a threshold value that separates high trading activity from low activity.

23.3 1.3 Constructing the Bipartite Network

The transformation allows us to define each client's state over the given period. We use this information to construct a bipartite network $G = (V, E)$, where:

- One set of nodes V_A represents the clients,
- The other set V_B represents the trading states,
- An edge between a client $v_i \in V_A$ and a state $v_j \in V_B$ signifies that the client v_i was in state v_j during the corresponding period.

24 2. Projection and Link Validation

To analyze client interactions more directly, we project the bipartite network onto V_A . This projected network G' contains edges between clients who share common trading states, reflecting their co-occurrence in specific trading activities.

24.1 2.1 Statistical Validation of Links

To validate the links in the projected network, we use the following statistical framework:

24.1.1 2.1.1 Null Hypothesis

We assume that the connections between clients are random, preserving the degree distribution of the clients. The probability p of two clients i and j sharing X common trading states, given their individual degrees N_i and N_j , follows the hypergeometric distribution:

$$H(X | N_B, N_i, N_j) = \frac{\binom{N_i}{X} \binom{N_B - N_i}{N_j - X}}{\binom{N_B}{N_j}}, \quad (10)$$

where:

- N_B is the total number of states in V_B ,
- $\binom{\cdot}{\cdot}$ denotes the binomial coefficient.

24.1.2 2.1.2 p-Value Calculation

The p-value for observing $N_{i,j}$ or more common states between clients i and j is computed as:

$$p(N_{i,j}) = 1 - \sum_{X=0}^{N_{i,j}-1} H(X \mid N_B, N_i, N_j). \quad (11)$$

This p-value indicates the likelihood of observing such a number of common states by chance.

25 3. Multiple Hypothesis Testing Correction

Given the large number of potential client pairs, multiple hypothesis testing corrections are necessary to control for false positives. Two commonly used methods are:

25.1 3.1 Bonferroni Correction

The Bonferroni correction adjusts the significance level α by dividing it by the number of tests m . If the original significance level is α , the Bonferroni-corrected threshold is:

$$\text{Bonferroni threshold} = \frac{\alpha}{m}. \quad (12)$$

This conservative method reduces the risk of Type I errors but can increase Type II errors, making it less effective in settings with many tests.

25.2 3.2 False Discovery Rate (FDR)

The FDR approach, such as the Benjamini-Hochberg procedure, controls the expected proportion of false positives among the declared significant results. The procedure involves:

1. **Sorting p-values:** Arrange all p-values in ascending order.
2. **Calculating the FDR threshold:** Determine the threshold for significance using:

$$\text{FDR threshold} = \frac{i}{m} \times Q, \quad (13)$$

where i is the rank of the p-value, m is the total number of tests, and Q is the desired FDR level.

3. **Selecting Significant Edges:** Retain edges with p-values below the FDR threshold.

The FDR method is less conservative than Bonferroni and offers a better balance between false positives and false negatives.

26 4. Detection of Under-Expressed Links

In addition to over-expressed links, where the observed number of common states exceeds the expected number under the null hypothesis, SVNs can identify under-expressed links.

This involves testing whether the number of common states is significantly lower than expected. The p-value for under-expression is calculated as:

$$p(N_{i,j}) = \sum_{X=N_{i,j}}^{N_B} H(X \mid N_B, N_i, N_j). \quad (14)$$

This two-tailed approach reveals less frequent or avoided relationships, offering additional insights into the dynamics of client interactions.

27 Conclusion

By applying Statistically Validated Networks to imbalance ratio activity, we can generate and validate a network that highlights both significant interactions and under-expressed relationships. This comprehensive methodology allows for a nuanced understanding of trading behaviors, facilitating more informed decision-making and strategy development.

Clustering Algorithms for Network Analysis: Infomap and Louvain Author Name September 16, 2024

28 Introduction

When analyzing networks, identifying clusters or communities within the network is crucial for understanding the underlying structure and dynamics. Two prominent algorithms for community detection are the Infomap algorithm and the Louvain algorithm. Both methods offer robust techniques for partitioning networks into clusters based on different principles. This section provides a comprehensive overview of these algorithms, detailing their approaches, mathematical foundations, and applications.

29 Infomap Algorithm

The Infomap algorithm is a method based on information theory, particularly focusing on compressing the description of the network's structure. The core idea of Infomap is to partition the network such that the flow of information (or the movement of a random walker) within the network is maximally compressed.

29.1 Principle of Infomap

Infomap employs the concept of information entropy to detect communities. The algorithm treats the network as a system where a random walker moves between nodes. The goal is to minimize the description length of the random walker's trajectory, which can be achieved by partitioning the network into communities where the flow within communities is high, and the flow between communities is low.

29.2 Mathematical Formulation

To formalize the Infomap approach, consider a network with nodes V and edges E . The network is represented by its adjacency matrix A , where A_{ij} denotes the weight of the edge between nodes i and j .

The algorithm uses a stochastic process to model the random walk. Define the transition probability matrix P , where P_{ij} represents the probability of moving from node i to node j . For a network with edge weights, P_{ij} can be computed as:

$$P_{ij} = \frac{A_{ij}}{\sum_k A_{ik}}. \quad (15)$$

Here, P_{ij} represents the probability of transitioning from node i to node j , normalized by the total degree of node i .

The Infomap algorithm aims to minimize the map equation, which is a measure of the expected description length of the random walker's trajectory. The map equation is defined as:

$$L_{\text{map}} = \frac{1}{N} \left[\sum_{i \in C} \sum_{j \in C} P_{ij} \log \frac{P_{ij}}{q_{ij}} + \sum_{i \notin C} \sum_{j \in C} P_{ij} \log \frac{P_{ij}}{q_{ij}} \right], \quad (16)$$

where q_{ij} is the probability of transitioning between communities C and \bar{C} . Minimizing this equation partitions the network into communities that optimize the balance between internal and external information flow.

29.3 Algorithm Steps

1. **Initialization:** Start with an initial partition of the network, where each node is its own community.
2. **Refinement:** Iteratively merge communities to minimize the map equation. This involves moving nodes between communities to find the configuration that minimizes the description length.
3. **Optimization:** Use a hierarchical approach to refine partitions, often employing techniques like simulated annealing to avoid local minima and achieve a globally optimal solution.

Infomap is highly effective for detecting hierarchical structures in networks and is known for its accuracy in identifying meaningful communities.

30 Louvain Algorithm

The Louvain algorithm is a widely used method for community detection based on modularity optimization. Modularity measures the strength of the division of a network into communities, quantifying the difference between the observed fraction of edges within communities and the expected fraction if edges were distributed randomly.

30.1 Principle of Louvain

The Louvain algorithm aims to maximize the modularity of a network partition. Modularity Q is defined as:

$$Q = \frac{1}{2m} \sum_{i,j} \left[A_{ij} - \frac{k_i k_j}{2m} \right] \delta(c_i, c_j), \quad (17)$$

where:

- m is the total number of edges in the network.
- A_{ij} is the weight of the edge between nodes i and j .
- k_i and k_j are the degrees of nodes i and j , respectively.
- $\delta(c_i, c_j)$ is 1 if nodes i and j are in the same community and 0 otherwise.

30.2 Mathematical Formulation

The modularity function Q measures the density of edges within communities compared to a null model where edges are placed randomly. The goal is to find a partition that maximizes Q .

30.2.1 Algorithm Steps

1. **Initial Partitioning:** Start with a partition where each node is its own community.
2. **Community Aggregation:** Group nodes into communities based on modularity optimization. Nodes in the same community form a new meta-node in a reduced network.
3. **Iterative Optimization:** Repeat the community aggregation process on the new meta-network to further refine the partition.

Phase 1 involves calculating the gain in modularity if each node were moved to the community of each of its neighbors. Nodes are assigned to the community that maximizes this gain. Phase 2 constructs a new network where each community is represented as a single node, and the process is applied recursively.

31 Comparative Analysis

Both Infomap and Louvain algorithms offer robust methods for community detection but differ in their approaches:

- **Infomap** focuses on information theory and random walk-based compression of network descriptions. It is well-suited for networks with hierarchical structures and complex community arrangements.
- **Louvain** is based on modularity optimization and is highly effective for large networks. It provides a hierarchical view of communities through iterative refinement and aggregation.

In summary, both algorithms provide powerful tools for understanding network structures and can be selected based on the specific characteristics of the network and the goals of the analysis.

32 References

1. Rosvall, M., Bergstrom, C. T. (2008). Maps of random walks on complex networks reveal community structure. *Proceedings of the National Academy of Sciences*, 105(4), 1118-1123.

2. Blondel, V. D., Guillaume, J. L., Lambiotte, R., Lefebvre, E. (2008). Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2008(10), P10008.

To assess the effectiveness of your clustering method for client interaction studies, two key approaches can be used: UMAP for visualizing the clusters and Adjusted Rand Index (ARI) for numerical evaluation. In this section, we will detail how each of these methods is applied, provide the relevant mathematical equations, and illustrate their practical application—especially in the context of comparing the output of your Statistically Validated Network (SVN) pipeline to alternative clusterings, such as client category clustering, client location clustering, and client business type clustering.

33 UMAP for Visualization

UMAP (Uniform Manifold Approximation and Projection) is an advanced technique for visualizing high-dimensional data in lower dimensions. It helps to visually evaluate the quality of clustering results by projecting them into a 2D or 3D space, allowing you to see how well the clients are grouped based on their interaction patterns.

33.1 Mathematical Foundation of UMAP

UMAP works by constructing a graph in high-dimensional space, where each data point (client) is connected to its nearest neighbors. The primary steps in UMAP are:

33.1.1 Constructing the High-Dimensional Graph

UMAP uses a k-nearest neighbors algorithm to build a graph that connects each data point to its closest neighbors in the original feature space. The distance between each pair of points is used to assign weights to the edges of the graph. For any two points i and j , the weight w_{ij} is computed as:

$$w_{ij} = \exp\left(-\frac{d(x_i, x_j)^2}{2\sigma_i^2}\right)$$

where $d(x_i, x_j)$ is the distance between points x_i and x_j , and σ_i is a parameter controlling the spread of the neighborhood for point i .

33.1.2 Optimization in Low-Dimensional Space

UMAP then finds a lower-dimensional embedding (typically 2D or 3D) by minimizing the discrepancy between distances in the high-dimensional and low-dimensional space. The goal is to preserve the local structure of the data. The loss function minimized by UMAP is a form of cross-entropy:

$$C = \sum_{(i,j)} w_{ij} \log\left(\frac{w_{ij}}{w'_{ij}}\right)$$

where w'_{ij} are the edge weights in the low-dimensional projection.

33.2 Using UMAP for Client Interaction Clustering

When applied to client interaction data, UMAP allows you to project complex client relationships (such as those derived from the SVN pipeline) into 2D or 3D space. For example, after generating the client clusters using your SVN pipeline, UMAP can be used to visualize how well-separated the clusters are, and whether there are clear groupings of clients based on their trading behavior or interaction patterns.

34 Adjusted Rand Index (ARI) for Numerical Evaluation

While UMAP gives an intuitive visual representation of the clusters, Adjusted Rand Index (ARI) provides a robust numerical method for comparing different clusterings. Specifically, you can use ARI to compare the output of your SVN pipeline to client category clustering, client location clustering, and client business type clustering.

34.1 Mathematical Foundation of ARI

The Adjusted Rand Index (ARI) measures the similarity between two clusterings while adjusting for the possibility of chance groupings. ARI is a modification of the Rand Index, which counts how often pairs of points (clients) are consistently assigned to the same or different clusters in two clustering solutions.

Let C be the ground truth clustering (e.g., client category, location, or business type) and C' be the clustering obtained from your SVN pipeline. ARI is computed as follows:

Define:

- a : the number of pairs of points that are in the same cluster in both C and C' .
- b : the number of pairs of points that are in different clusters in both C and C' .
- c : the number of pairs of points that are in the same cluster in C but in different clusters in C' .
- d : the number of pairs of points that are in different clusters in C but in the same cluster in C' .

The Rand Index (RI) is given by:

$$RI = \frac{a + b}{a + b + c + d}$$

Since random clusters could result in some agreement, ARI adjusts for this by subtracting the expected agreement by chance:

$$ARI = \frac{RI - E[RI]}{\max(RI) - E[RI]}$$

Here, $E[RI]$ is the expected Rand Index for random clustering, and $\max(RI)$ is the maximum possible Rand Index. This ensures that ARI values are normalized between -1 and 1:

- $ARI = 1$ indicates perfect agreement between the two clusterings.
- $ARI = 0$ indicates that the clustering is no better than random.
- $ARI < 0$ suggests a worse-than-random clustering, which is rare in practice.

34.2 Example of ARI Calculation

Let us assume you have clustered 100 clients using the SVN pipeline and you want to compare the results to a known client location clustering. If your SVN pipeline successfully groups clients based on geographic proximity, the ARI should be high, indicating that the SVN clustering aligns with the location-based clustering.

34.3 Using ARI to Compare SVN Pipeline with Client Categories

You can use ARI to compare the results of the SVN pipeline with client category clustering, client location clustering, and client business type clustering. If the ARI is high, it indicates strong agreement between the SVN clusters and the predefined categories.

35 Conclusion

In summary, the combination of UMAP for visualization and Adjusted Rand Index (ARI) for numerical evaluation offers a powerful approach to assess the effectiveness of clustering in client interaction studies. UMAP provides a visual insight into the grouping of clients, while ARI gives a rigorous, quantitative comparison between the clustering output of the SVN pipeline and alternative groupings based on client category, location, or business type.

By applying ARI to compare the SVN clustering with these alternative clusterings, you can better understand the underlying client interactions and refine your clustering approach to capture meaningful relationships in the data.

36 Introduction

In financial markets, the dynamics of price variation are closely linked to the order flow—the net volume of buy and sell orders executed by clients. Understanding this relationship is crucial for market participants, as it helps in predicting price movements based on the imbalance between supply and demand. This paper explores the connection between client order flow and price variation, utilizing a LightGBM (LGBM) regression model.

37 Price Variation as a Percentage Change

To quantify price movements, we define the price variation as the percentage change between two time periods. Mathematically, the price variation from time t to $t + 1$ is given by:

$$\text{Price Variation}_{t+1} = \frac{p_{t+1} - p_t}{p_t}, \quad (18)$$

where:

- p_t is the asset price at time t ,
- p_{t+1} is the asset price at time $t + 1$.

This equation captures the relative change in price over time, allowing us to monitor upward or downward movements in percentage terms.

38 Order Flow as a Predictor of Price Variation

Order flow represents the net buying or selling volume from clients and is a key driver of price evolution. The total order flow at time t , denoted as order_flow_t , is calculated as:

$$\text{order_flow}_t = \sum_{i=1}^N q_{i,t}, \quad (19)$$

where:

- N is the number of clients,
- $q_{i,t}$ is the trade volume of client i at time t (positive for buy orders, negative for sell orders).

38.1 Hypothesized Relationship: Order Flow and Price Variation

The hypothesis is that the order flow at time t influences the price variation at time $t + 1$, as an imbalance in buy/sell orders leads to price adjustments. We can express this relationship as:

$$\text{order_flow}_t \sim \text{price_variation}_{t+1}. \quad (20)$$

This implies that the price variation at time $t + 1$ is at least partially driven by the order flow at time t .

39 LGBM Regression to Model Price Variation

We use a LightGBM regression model to capture the relationship between client order flow and price variation. LightGBM is a gradient-boosting framework that iteratively builds decision trees, focusing on correcting prediction errors.

39.1 Objective Function: Mean Squared Error

The objective function for this regression model is the Mean Squared Error (MSE), which measures the average squared differences between predicted and actual values. The MSE is expressed as:

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2, \quad (21)$$

where:

- y_i is the actual price variation at time $t + 1$,
- \hat{y}_i is the predicted price variation,
- N is the total number of time steps.

The LightGBM model minimizes MSE by iteratively correcting residuals from previous trees.

39.2 Model Evaluation Metrics

39.2.1 R-Squared: Coefficient of Determination

We use the R^2 metric to evaluate the model's performance. It measures the proportion of variance in the price variation explained by the model. R^2 is given by:

$$R^2 = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2}, \quad (22)$$

where \bar{y} is the mean of the actual values y_i . An R^2 value close to 1 indicates a high explanatory power.

39.2.2 Pearson Correlation

The Pearson correlation coefficient ρ measures the linear relationship between the predicted and actual price variation:

$$\rho_{y,\hat{y}} = \frac{\text{Cov}(y, \hat{y})}{\sigma_y \sigma_{\hat{y}}}, \quad (23)$$

where:

- $\text{Cov}(y, \hat{y})$ is the covariance between the actual and predicted values,

- σ_y and $\sigma_{\hat{y}}$ are the standard deviations of the actual and predicted values, respectively.

A Pearson correlation close to 1 suggests a strong positive linear relationship.

40 Order Flow as Input Features for LGBM

The input features for the LGBM model are the order flows of different client clusters. The order flow for cluster i at time t is defined as:

$$\text{order_flow}_{t,i} = \sum_{j=1}^{N_i} q_{j,t}, \quad (24)$$

where N_i is the number of clients in cluster i , and $q_{j,t}$ is the trade volume of client j . The feature vector at time t is:

$$\mathbf{x}_t = (\text{order_flow}_{t,1}, \text{order_flow}_{t,2}, \dots, \text{order_flow}_{t,K}), \quad (25)$$

where K is the total number of clusters, and the target variable is $\text{price_variation}_{t+1}$.

41 Handling Differences in Distribution

Order flow often exhibits a skewed distribution, while price variation tends to follow a more continuous distribution. The LGBM model is capable of learning non-linear relationships between these differently distributed variables.

42 Choice of LGBM over Other Models

While linear models assume a constant linear relationship between input features and the target, they are insufficient for capturing the complexities of financial markets. Neural networks, though flexible, often underperform on tabular data and may overfit. LGBM strikes a balance by handling non-linearity efficiently and performing well on tabular data without the risks associated with neural networks.

43 Conclusion

The LGBM model provides a robust framework for modeling the relationship between client order flow and price variation. It captures non-linearities, handles differences in distribution, and performs efficiently on large datasets. Metrics such as R^2 , MSE, and Pearson correlation allow for thorough model evaluation, ensuring that the predictions provide meaningful insights into price movements.

44 Analysis of Client-Specific Trade Data for Price Variation Prediction

44.1 Introduction

In financial markets, data scarcity can significantly hinder the performance of predictive models. This issue is particularly acute when modeling price variation based on client order flow. To overcome this limitation, we propose a novel model that treats each client's trade as an individual data point. This approach substantially increases the sample size, thereby enhancing the model's ability to predict future price variations. By expanding the dataset through individual client trades, we aim to improve the accuracy and robustness of our predictions.

44.2 Problem Statement

The primary challenge of our analysis is the limited amount of data, which constrains the model's performance. Traditional models might struggle to generalize due to insufficient data points. To address this, we propose using each client's trade as a separate data point in our model. This approach allows us to effectively multiply the number of samples by the number of clients, thereby improving the model's robustness and predictive accuracy.

44.3 Model Framework

The core idea of our model is to predict future close price variation based on trades made by clients during a given sampling period. The model operates as follows:

- **Client Trade Data Collection:** Each trade executed by a client during a sampling period is treated as an individual data point.
- **Target Variable:** For each trade, the target variable is the price variation observed at the end of the next sampling period.
- **Simultaneous Predictions:** We predict the future price variation simultaneously for all clients, using trades from the same sampling period to inform these predictions.

Mathematically, this can be formulated as follows:

Let $q_{i,t}$ represent the trade volume of client i at time t , and p_{t+1} denote the asset price at the end of the next sampling period. The price variation at time t is given by:

$$\text{price_variation}_{t+1} = \frac{p_{t+1} - p_t}{p_t}$$

where p_t is the price of the asset at time t , and p_{t+1} is the price at the end of the next sampling period.

For each client i , we have a set of trades $\{q_{i,t}\}$, and the corresponding target variable for each trade is the price variation $\text{price_variation}_{t+1}$.

44.4 Assumptions and Rationale

- **Similarity of Trade-Price Variation Function:** We assume that the relationship between trade activity and price variation is similar across clients. This assumption is supported by empirical studies suggesting that the impact of order flow on price changes is consistent across different market participants.
- **Leveraging Other Clients' Information:** By considering all clients' trades simultaneously, we leverage information from other clients' activities. This is particularly valuable when specific clients are not active during a sampling period. The collective trading behavior helps fill in gaps and enhances the model's predictive performance.
- **Client Influence:** We recognize that some clients may possess more significant market influence than others. Consequently, our model benefits from incorporating the activity of influential clients to enhance prediction accuracy. This assumption also enables us to develop a client-following model, where we can track and predict based on influential clients' trading movements.

44.5 Model Implementation

To implement this model, we need to construct the dataset such that each client's trade can be used as a data point. The following steps outline the process:

- **Data Aggregation:** Collect trade data for each client during the sampling period. Each trade will be linked to the price variation observed in the subsequent period.
- **Feature Engineering:** Define features for each trade. These features might include:
 - Trade volume $q_{i,t}$
 - Time of trade
 - Client's historical trading behavior
 - Market conditions at the time of trade
- **Model Training:** Use machine learning algorithms, such as LightGBM (LGBM), to train the model. The input features will be derived from the trades, and the target variable will be the future price variation.
- **Evaluation Metrics:** Assess the model using performance metrics such as Mean Squared Error (MSE), R^2 (coefficient of determination), and Pearson correlation coefficient. These metrics will help quantify the model's effectiveness in explaining price variation based on client trades.

44.6 Mathematical Formulation

- **Objective Function:** The objective function for training the model is to minimize the Mean Squared Error (MSE):

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

where y_i is the actual price variation for trade i , and \hat{y}_i is the predicted price variation.

- **Coefficient of Determination (R^2):** The R^2 score is given by:

$$R^2 = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$$

where \bar{y} is the mean of the actual price variations.

- **Pearson Correlation Coefficient:** The Pearson correlation coefficient between the predicted and actual values is:

$$\rho_{y,\hat{y}} = \frac{\text{Cov}(y, \hat{y})}{\sigma_y \sigma_{\hat{y}}}$$

where $\text{Cov}(y, \hat{y})$ is the covariance between actual and predicted values, and σ_y and $\sigma_{\hat{y}}$ are the standard deviations of actual and predicted values, respectively.

44.7 Model Advantages

- **Increased Sample Size:** Treating each client's trade as an individual data point effectively increases the number of samples, which improves model performance.
- **Enhanced Predictive Accuracy:** By incorporating the collective information from all clients, the model can provide more accurate predictions, even when individual clients are inactive.
- **Identification of Influential Clients:** The model can highlight clients with significant market influence, allowing for targeted analysis and strategy development.

44.8 Conclusion

The proposed model addresses the data limitation issue by treating each client's trade as a separate data point, thereby increasing the sample size and improving predictive accuracy. By leveraging the collective trading behavior of all clients and recognizing the varying influence of individual clients, the model offers a robust framework for analyzing and predicting price variation in financial markets.