



UNIVERSIDAD NACIONAL DE COLOMBIA

# Order book's microstructure visualization: the Colombian case

**Andrea Marcela Cruz Moreno**

Universidad Nacional de Colombia  
Engineering School  
Bogotá, Colombia  
2015



# Order book's microstructure visualization: the Colombian case

**Andrea Marcela Cruz Moreno**

Submitted to the Engineering School of the Universidad Nacional de Colombia, in partial  
fulfillment of the requirements for the degree of:

**Master of Science in Systems and Computer Engineering**

Under the guidance of:  
Germán Jairo Hernández Pérez, Ph.D.  
Associate Professor  
Engineering School

ALGOS Research Group  
Universidad Nacional de Colombia  
Engineering School  
Bogotá, Colombia  
2015



To 61.



# Acknowledgements

This thesis would not have been possible without the help, support, guidance and trust of my thesis advisor, Professor Germán Hernández.

Also, I am grateful to my amazing parents and siblings: Laura, Tom, and Diego who have always given me their unconditional love, support, encouragement and advice. Thanks to my wonderful husband, Michael.





# Abstract

Order book data provide a valuable source of information in financial markets, for this reason, it is an excellent candidate in the construction of new trading tools and models. Order book representation is an still active study branch in quantitative finance.

This work addresses the problem of information visualization of financial data from the Colombian Forex Market using two approaches: a heatmap representation, and a wavelet based representation in order to filter high frequency noise. To do so, is imperative to deal with a massive amount of data coming from the Colombian Forex Market Limit Order Book, a register with all the buy and sell intentions of the market's participants.

The experimental evaluation shows that the proposed strategies are able to identify frequent patterns within the presented visualizations tools. Furthermore, and more important, it is possible to associate some of those frequent patterns with a trend with a probability greater than 0.5. This result is useful in order to generate buy and sell signals for a trader.

**Keywords:** Order book, Wavelet Coefficients, Scientific Visualization, Financial Engineering, Machine Learning, Heatmap, Information Representation.

# Content

<b>Acknowledgements</b>	<b>VII</b>
<b>Abstract</b>	<b>IX</b>
<b>1. Introduction</b>	<b>2</b>
1.1. Thesis goals . . . . .	2
1.2. Main contributions . . . . .	3
1.2.1. Colombian Forex Market Order Book Visualization . . . . .	3
1.2.2. Frequent Patterns Exploration in a reasonable amount of time . . . . .	3
1.2.3. Association among some frequent patterns and a specific trend . . . . .	3
1.3. Thesis organization . . . . .	4
<b>2. Problem Statement</b>	<b>5</b>
2.1. Scope . . . . .	5
2.2. Previous work . . . . .	5
2.2.1. Order book . . . . .	6
2.2.2. Strategies based on the order book analysis . . . . .	10
2.2.3. Discussion . . . . .	11
2.2.4. Remarks . . . . .	15
<b>3. Basic Notions and Definitions</b>	<b>16</b>
3.1. Basic Financial Concepts . . . . .	16
3.1.1. Forex Markets . . . . .	16
3.1.2. Order book . . . . .	16
3.1.3. Spread . . . . .	17
3.1.4. Trader . . . . .	17
3.2. Scientific Visualization . . . . .	17
3.3. Heatmap . . . . .	18
3.4. Bag of Words . . . . .	18
3.5. Wavelets . . . . .	20
3.5.1. Wavelet bases . . . . .	20
3.5.2. Quick glossary . . . . .	21
3.5.3. Filter Bank . . . . .	21

3.5.4. Wavelets for images . . . . .	21
3.5.5. Sampling . . . . .	22
3.5.6. Compression . . . . .	22
3.5.7. Denoising . . . . .	22
3.5.8. Time Frequency Dictionaries . . . . .	23
3.5.9. Heisenberg Uncertainty . . . . .	23
3.5.10. Windowed Fourier Transform . . . . .	23
<b>4. Microstructure Visualization Tools</b>	<b>25</b>
4.1. Heatmap based approach . . . . .	25
4.1.1. First approach . . . . .	25
4.1.2. Towards a better information understanding . . . . .	26
4.2. Use of Wavelets in the representation . . . . .	26
<b>5. Experimental Setup</b>	<b>31</b>
5.1. Methodology . . . . .	31
5.1.1. Dataset description . . . . .	32
5.1.2. Pattern exploration . . . . .	34
5.1.3. Performance measurement . . . . .	37
5.2. Results and Discussion . . . . .	38
5.2.1. Market Trend Visual Bag of Words Informative Patterns in Limit Order Books . . . . .	38
5.3. Conclusions . . . . .	51
<b>6. Adaptive Method for Market Microstructure Exploration</b>	<b>53</b>
6.1. Method's description . . . . .	53
<b>7. Conclusions and Future work</b>	<b>56</b>
7.1. Conclusions . . . . .	56
7.2. Future work . . . . .	56
<b>A. Appendix: Heatmap approach results</b>	<b>57</b>
<b>B. Appendix: Wavelet based approach results</b>	<b>59</b>
B.1. Accuracy for Wavelet transform approach using only differences, one minute time slot. . . . .	59
B.2. Accuracy for Wavelet transform approach using only averages, one minute time slot. . . . .	60
B.3. Accuracy for Wavelet transform approach using only differences, ten minutes time slot. . . . .	62

---

B.4. Accuracy for Wavelet transform approach using only averages, ten minutes time slot. . . . .	63
<b>References</b>	<b>64</b>

# List of Figures

<b>2-1.</b>	References distribution by category. . . . .	12
<b>2-2.</b>	Main techniques used for exploring each category. . . . .	13
<b>2-3.</b>	Publications' timeline for each category: line 1 presents information about LOB's Information content, line 2 presents information about LOB's Dynamics, line 3 presents information about LOB's Order placement, line 4 presents information about LOB's Representation and modeling, line 5 presents information about LOB's Trading strategies, line 6 presents information about LOB's Consequences and forecasting and line 7 presents information about LOB's Dynamics and Representation and modeling simultaneously. . . . .	14
<b>3-1.</b>	Example of Heat Map of Order Book Depth, Todd et al. [45]) . . . . .	19
<b>3-2.</b>	Time-frequency boxes representing the energy spread of two windowed Fourier atoms [31]. . . . .	23
<b>4-1.</b>	First approach to the order book heatmap visualization. . . . .	26
<b>4-2.</b>	Different gray level thresholds for the previous image. . . . .	27
<b>4-3.</b>	Example of position of frequent patterns in a heatmap order book image for 2 days of trading. . . . .	27
<b>4-4.</b>	Example of image produced by four levels of compression using haar wavelet transform. . . . .	28
<b>4-5.</b>	Example of order book visualization using wavelets based approach (filtering). . . . .	29
<b>4-6.</b>	Example of order book visualization using wavelets based approach (filtering). . . . .	30
<b>4-7.</b>	Three months of trading using one minute resolution as basis for the image construction. . . . .	30
<b>5-1.</b>	Visual representation of the described methodology . . . . .	32
<b>5-2.</b>	Informative Region . . . . .	33
<b>5-3.</b>	Limit Order Book example [34]. . . . .	40
<b>5-4.</b>	Order execution example [54]. . . . .	41
<b>5-5.</b>	Bid-ask distribution example [?]. . . . .	41
<b>5-6.</b>	Operations in the order book, modification of [?] for illustration purposes. . . . .	42
<b>5-7.</b>	Example of the proposed order book visualization. . . . .	43
<b>5-8.</b>	Image Representation through Bag-of Visual-Ngrams, extracted from [30] . . . . .	44

---

<b>5-9.</b> Predictor's accuracy using different 1 minute pattern sizes. . . . .	47
<b>5-10.</b> Predictor's accuracy using different 10 minutes pattern sizes. . . . .	48
<b>5-11.</b> Example of patches associated with clusters. . . . .	49
<b>5-12.</b> Clusters matrix and centroids. . . . .	50
<b>5-13.</b> Frequency at which patterns are associated with a certain cluster. . . . .	50
<b>5-14.</b> Clusters and their performance. . . . .	51
<b>5-15.</b> Cumulated retrurn usd cop in six months. . . . .	52
<b>5-16.</b> Patterns accuracy within and outside a global trend. . . . .	52
 <b>6-1.</b> Convergence behavior type 1. . . . .	 54
<b>6-2.</b> Convergence behavior type 2. . . . .	54
<b>6-3.</b> Convergence behavior type 3. . . . .	55

# List of Tables

<b>2-1.</b> References classification by subject. . . . .	12
<b>2-2.</b> Markets classificated by geographic region and associated with the order book subjects explored on them (IC, D, OP, RM, TS, CF stands for information content, dynamics, order placement, representation and modeling, trading strategies, and consequences and forecasting, respectively). . . . .	15
<b>3-1.</b> Todd et al. [45] definitions of data, information and knowledge in computational space. . . . .	18
<b>A-1.</b> Pattern Sizes . . . . .	57
<b>A-2.</b> Experimental Setup 1 (Raw data) 1 minute . . . . .	57
<b>A-3.</b> Experimental Setup 1 (Raw data) 5 minutes . . . . .	57
<b>A-4.</b> Experimental Setup 1 (Raw data) 10 minutes . . . . .	58
<b>B-1.</b> Accuracy for Wavelet transform approach using only differences, one minute time slot, first iteration. . . . .	59
<b>B-2.</b> Accuracy for Wavelet transform approach using only differences, one minute time slot, second iteration. . . . .	59
<b>B-3.</b> Accuracy for Wavelet transform approach using only differences, one minute time slot, third iteration. . . . .	60
<b>B-4.</b> Accuracy for Wavelet transform approach using only differences, one minute time slot, fourth iteration. . . . .	60
<b>B-5.</b> Accuracy for Wavelet transform approach using only averages, one minute time slot, first iteration. . . . .	60
<b>B-6.</b> Accuracy for Wavelet transform approach using only averages, one minute time slot, second iteration. . . . .	61
<b>B-7.</b> Accuracy for Wavelet transform approach using only averages, one minute time slot, third iteration. . . . .	61
<b>B-8.</b> Accuracy for Wavelet transform approach using only averages, one minute time slot, fourth iteration. . . . .	61
<b>B-9.</b> Accuracy for Wavelet transform approach using only differences, ten minutes time slot, first iteration. . . . .	62
<b>B-10</b> Add caption . . . . .	62

---

<b>B-11</b> Accuracy for Wavelet transform approach using only differences, ten minutes time slot, second iteration. . . . .	62
<b>B-12</b> Accuracy for Wavelet transform approach using only averages, ten minutes time slot, first iteration. . . . .	63
<b>B-13</b> Accuracy for Wavelet transform approach using only averages, ten minutes time slot, second iteration. . . . .	63



# 1. Introduction

Limit Order Book has become a valuable source of knowledge for traders, taking an important role in research and as support tool for making financial decisions [1] [3] [6] [14] [53].

One of the main issues in research is the size of the order book, due to the number of orders placed and cancelled every minute. For example, for London Exchange Market, for a single stock, several orders were placed within one millisecond in July 2009 [18]. This provides vast repositories of valuable data, which is hard to process and manage duly.

Due to this situation, the development of strategies for extracting information from these data is required, producing an increasing interest in this research field lately (See Chapter 2). Likewise, the demand of such strategies has increased in recent years.

Foreign exchange markets are essential for the correct operation of the world economy. By the time of writing of this thesis (See chapter 2), there is no evidence of publications of work related to the Colombian Forex Market Limit Order Book, a fact that motivates this work.

This work addresses the problem of information extraction from financial datasets. The main goal of this research was to study high frequency trading strategies using order book information from the Colombian Forex Market and its potential in the construction of predicting models. This work additionally presents a visualization tool in order to facilitate the trader's understanding of large amounts of Limit Order Book data.

## 1.1. Thesis goals

The main goal of this research was to study high frequency trading strategies using order book information from the Colombian Forex Market and its potential in the construction of predicting models. The following is the description of this research specific targets:

- Providing a survey of the methods published to date for detecting trading strategies using order book information, via a systematic literature review.
- Selecting or designing a methodology able to represent in a summarized and efficient way the order book information.

- Establishing a time window in order to preserve relevant order book information.
- Selecting or designing a methodology that allows representing properly the Colombian Forex Market Order Book information dynamics.
- Selecting or designing a trading strategies detection system for the Colombian Forex Market using Order Book information.
- Evaluating the performance and feasibility of the proposed system, in supporting the financial decision making process in the Colombian Forex Market.

## 1.2. Main contributions

The following is the summary of the main contributions of this work:

### 1.2.1. Colombian Forex Market Order Book Visualization

A Colombian Forex Market Order Book Visualization is presented. This visualization provides the trader with a framework which allows the interpretation of large sections of the limit order book at a glance. It shows relationships between price, volume and time directly.

This work was published as a contributed talk named «Order Book Microstructure Visualization: The case of Colombian High- Frequency Foreign Market. XIII Latin American Congress of Probability and Mathematical Statistics CLAPEM. September, 2014.»

### 1.2.2. Frequent Patterns Exploration in a reasonable amount of time

Algorithms for Frequent Patterns Exploration are presented. These algorithms have reduced the amount of time required for mining a dataset up to two orders of magnitude depending on the pattern size, thanks to the use of a pattern summary function. The use of a wavelets based approach, in some time windows, can reduce the initial dataset without loss of accuracy for the classifier, so it reduces the amount of non valuable information.

### 1.2.3. Association among some frequent patterns and a specific trend

Algorithms for association between frequent patterns and a specific trend are depicted. These algorithms allow calculating the probability of each pattern of being associated with a bearish trend, a bullish trend or with no trend, labeling each pattern accordingly. The use of these algorithms allowed to detect patterns seasonality in the Colombian Forex Market

Order Book.

This work will be published in the 6th Annual Stevens Conference on High Frequency Finance and Analytics (HF2015) that will be held on October 29th-31st, 2015 at Stevens Institute of Technology, Hoboken, NJ, USA. It will be published as a presentation titled "Market Trend Visual Bag of Words Informative Patterns in Limit Order Books".

### **1.3. Thesis organization**

The organization of this thesis is as follows:

- In chapter 1, an introduction and the main contributions of this work will be presented.
- In chapter 2, a brief state of the art on the limit order book will be presented and the research problem will be stated.
- Some basic notions and definitions are introduced in chapter 3.
- The proposed visualization tools for the Colombian Forex Market order book are depicted in chapter 4.
- The experimental setup is detailed in chapter 5. An adaptive method for mining patterns is published in chapter 6.
- Finally, in chapter 7, conclusions and future work are described.

## **2. Problem Statement**

### **2.1. Scope**

This research studied the problem of order book information extraction and its application in the Colombian Forex Market. The main objective of this work is to determine if the information provided by the Colombian Forex Market Order Book is enough in order to recognize propitious trading scenarios.

Two main components have been investigated in the proposed strategy: a proper visualization and the frequent patterns exploration. A proper visualization facilitates the task of user's data understanding, according to the information needs. The frequent patterns exploration deals with the processing of a huge amount of available data in a reasonable time lapse. These components will be tested with real tick data from the Colombian Forex Market of the year 2012.

### **2.2. Previous work**

This section presents a review of recent work done about the information content, dynamics, order placement and strategies in the order book during the last decade, approximately. The purpose is to provide to a non expert reader with a general understanding of the subject. This is achieved by explaining the order book operation, what is relevant in the order book information, how it operates, how is modeled and represented, and presenting some methods of using the order book to support investment strategies. Finally, an analysis and a discussion about future trends and possible enhancement methods is introduced.

The recent progress of the order book exploration in several markets worldwide and in simulated markets is presented in this chapter. It was made a classification of the literature in six categories: information content, dynamics, order placement, representation and modeling, trading strategies, and consequences and forecasting. This classification introduced a general view that serves as a basis for trend analysis.

Goods exchange has been one of the activities developed by humankind, necessary for civilizations to thrive. Nowadays, this activity has evolved towards automatic boosting profit

methodologies in sophisticated markets, such as Forex market.

In order to perform a profitable investment, two questions should be solved: where to invest and when to trade. In this paper, the order book, a tool that provides information which supports the decision making process necessary to answer the second question, is explained.

The order book was employed under the assumption that there is market information which allows the discovery of price behavior pattern association with a future market trend, that would provide some probably optimal trading points in time which would support the traders' decision making task.

The main aim of this work is to provide elements to understand the informational potential of the order book by means of a systematic literature review, analyzing trends, defining study categories and determining the study state of the markets surveyed.

The rest of the chapter is organized as follows: subsection 2 presents a brief review of the main concepts and definitions necessary to understand the order book dynamics and, connects High Frequency Trading with Forex markets and order book analysis; subsection 3 introduces a discussion about the strategies based on the order book analysis; lastly, subsection 4 presents the conclusions and some suggestions for improving trading strategies.

### **2.2.1. Order book**

With the purpose of registering the prices at which traders would buy or sell a financial instrument and at which volume, the order book was created. This tool allows accessing information in order to characterize the asset behavior, and based on this model, being able to generate trading strategies that would increase the investor wealth.

#### **Order book information content**

For the sake of understanding the order book information content impact, the way in which the order book is presented has been studied by [53] in order to determine whether it influences the order placement strategies. Yu shows that the limit order book information does have a different impact in the order submission in the bid side than in the ask side, using a generalization of a linear regression model which assumes a discrete dependent variable and a probit model.

Another study concerning the order book content information analysis [22], proposes a strategy formed on a combination of dynamic focus and naive price adjustment (NPA). The results presented by Jiaqi et al. outperform the use of NPA exclusively, in a higher level if the asset presents more liquidity. Order book dynamic volume adjustment points are shown

to be useful for controlling an adverse price selection and information covering. It is showed that price and volume information in the book are complementary and essential to select informational features.

Limit order book were increased from three to five, the top price levels displayed in Chinese stock market. Li et al. [29], study the effects in price discovery. The results show that there exists significant difference between the pre and post transparency rise. The new quotes added have little information content, but is useful for traders to improve the price discovery process.

Fletcher et al. [11] use SVM's in junction with multiple kernel learning (MLK) for analyzing LOB information. MLK was used to mitigate the error over wrong kernel selection with mixed results. Temporal window lengths were defined after experimentation. Price movement of the EUR/USD currency was tested under this framework but profitability was not achieved. MKL and simple kernel methods did not show the expected performance difference but MKL was able to identify the most informative feature subsets. Future work can be interesting using reduced feature spaces.

### **Order book dynamics**

In the study of the order book dynamics, two aspects are considered: its depth and to which extent is informative. The helpfulness of the whole content of the book has been studied in [37], distinguishing between two kind of data: those belonging to the best quotes, and further. The conclusion was that the whole book helps to determine if the trader provides or consumes liquidity and that more patient traders use deeper order book information. The study presented by [17] models the order book high frequency dynamics in the London Stock Exchange using second level data (best quotes and volume at different prices), this model allows to capture the arrival of orders of different sizes.

Changes in order book have been studied in [36], where forecasting price change and the direction of such change are problems addressed separately. However, this approximation did not produce statistically significant results. Rinaldo [39] also studies changes in order book shape, analyzing transaction aggressiveness and the order book flow using Swiss Stock Exchange data. This study shows how traders using market orders and those using limit orders react in opposite ways to market changes. Also, market equilibrium is associated with weak aggressiveness in trading and imbalance between ask and bid is associated with higher aggressiveness, this provides information about how order book shape changes before any submission operation, forming up a thinner shape in that side.

Kercheval et al. [23] defined a set of classes (dynamic metrics) and feature spaces (raw data from books, «economical set» or refined raw data after entropy reduction). Support vector machines were used in order to predict the market based on his strength of optimal classification over linear separable samples. Those metrics were created focused on the profitability and the search of a possible competitive advantage (for example, a future bid versus actual ask). Multi-class SVM and kernel transformations were taken into account to improve reliability and over-adjustment. Results from both feature spaces did not show significant differences. On the other hand, SVM's as classification engine proved to be reliable (over 98.5 % on precision and recall) for the less risky scenarios (profit sense).

### Orders placement

Finding an optimal point to buy or sell is another widely discussed problem in finance. Using the order book, in [26], price change point quickest detection is pursued employing a social learning model and a protocol for quickest detection. It is evidenced that an optimal decision policy has several thresholds.

The order book time-varying composition expressed by each time stamp, is studied in [21]. Jiang et al. aim to model the limit price distribution in the ask and the bid side at each instant, using a gamma distribution to analyze its impact in volume and the existence of seasonal patterns. The results present strong parameters distribution seasonality, providing a model of how markets evolve in time.

Malik et al. [32], selected LOB data over bid-ask spread seeking a much deeper connexion between the data. Liquidity is defined as trading opportunity over large volumes and curve fitting inside a large time window is the most radical approach. Besides no profitability is found in the long run, the market has shown nonlinear, time-varying characteristics. Trade scheduling algorithms can be improved with this framework.

### Order book representations and modeling

With the aim of harnessing order book information content, finding a representation which allows handling its volume, and a model that reflects its changes is required. For markets operated in discrete time periods, in [47], a Markov process is modeled providing a minimum number of constant parameters, considering a one-product-market where the prices have a finite and small number of possibilities. In order to provide a shorter order book representation, in [19], Jiang et al. reduce the representation to four parameters by snapshot. A Kalman filter is used to estimate a linear dynamic system state and provide a liability measure, being used for prediction and filtering. The results, obtained from the London Stock Exchange,

show that jumps in the estimated parameters are always detected.

Another factor used to represent the order book information content is its slope [5], Cheng et al. study how order placement contributes to the price formation process, measuring how the quantity supplied changes as price does (elasticity). It is concluded that order book limits contain the current supply information and its relationship with the demand, allowing the creation of more precise investment strategies.

Ahn et al. [2], performs a quotes' price clustering over a maximum of five quotes to limit orders. Deeper quotes presented higher clustering than the best ones, pointing that the further from the best a quote is, the less information provides. This study uses data from the Hong Kong Stock Exchange.

### Consequences and forecasting

The Triennial Central Bank Survey[42] presents a foreign exchange turnover report evidencing the growth in forex exchange trading, with a turnover of USD 5.3 trillion per day in April 2013 and a growth rate of 32.5 % compared with April 2010. This survey shows that the dollar is still the dominant currency in FX deals with a 87 % of presence in April 2013. This information provides an insight on the importance of studying currencies behavior, specially the United States Dollar.

High frequency trading (HFT) is defined by the International Organization of Securities Commissions (IOSCO) as «a very quantitative trading form» with the following common features<sup>1</sup>:

- It involves the use of sophisticated technological tools for pursuing a number of different strategies, ranging from market making to arbitrage;
- It is a highly quantitative tool that employs algorithms along the whole investment chain: analysis of market data, deployment of appropriate trading strategies, minimisation of trading costs and execution of trades;
- It is characterized by a high daily portfolio turnover and order to trade ratio (i.e. a large number of orders are cancelled in comparison to trades executed);
- It usually involves flat or near flat positions at the end of the trading day, meaning that little or no risk is carried overnight, with obvious savings on the cost of capital

---

<sup>1</sup>Regulatory Issues Raised by the Impact of Technological Changes on Market Integrity and Efficiency Consultation Report. July 2010. Retrived on January 2014 from <http://www.iosco.org/library/pubdocs/pdf/IOSCOPD354.pdf>



associated with margined positions. Positions are often held for as little as seconds or even fractions of a second;

- It is mostly employed by proprietary trading firms or desks; and
- It is latency sensitive. The implementation and execution of successful HFT strategies depend crucially on the ability to be faster than competitors and to take advantage of services such as direct electronic access and co-location.

LOB analysis with reasonable execution time requires the use of HFT techniques. Consequences of the adoption of High Frequency Trading on an Electronic Market are discussed in [24], where the Flash Crash, «a brief period of extreme market volatility on May 6, 2010» is analyzed. Some of the highlights in [24] are:

- "High Frequency Traders aggressively trade in the direction of price changes. This activity comprises a large percentage of total trading volume, but does not result in a significant accumulation of inventory."
- "High Frequency Traders are not willing to accumulate large positions or absorb large losses."
- "When rebalancing their positions, High Frequency Traders may compete for liquidity and amplify price volatility."
- "HFTs did not trigger the Flash Crash, but their responses to the unusually large selling pressure on that day exacerbated market volatility"

### 2.2.2. Strategies based on the order book analysis

In an attempt to generate trading strategies with order book information, Bates et al. [3], combine the use of technical indicators with order book information, using a reinforcement learning algorithm. The results show that this combination outperforms the use of technical signals solely. Two types of indicators are employed: stop loss orders and get profit orders.

Farmer et al. [10], address the problem of constraints over intelligent agents by market institutions. It was empirically proved that even the most simple models develop a similar behavior. This work has shown circumstances where the most intelligent strategies cannot be successful.

With the objective of characterizing the behavior of traders, independently of the circumstances, Yang et al. [52], used inverse reinforcement learning to infer a solution model to the choices made by the trader, using a reward function learnt from the previous observations, producing an autonomous process. This study presents a simulation model based on agents

for the futures market. It provides evidence that is possible to identify reliably other algorithmic trading high frequency strategies, using inverse reinforcement learning. It is also shown that is possible to identify precisely a manipulative high frequency strategy among other high frequency strategies.

Informed traders adjust order submissions to the level of risk perceived, and remove mis-priced limit orders in the book. Trading strategies of the informed and the liquidity traders diverge in time. Informed traders provide liquidity through limit order submission due the dynamic order behavior[4]. Cheng et al. [6] Using simulations support the evidence that traders with analysis of LOB information produce more accurate expectation of future asset price. A market lead by LOB analysis as a strategy would show high liquidity and low volatility.

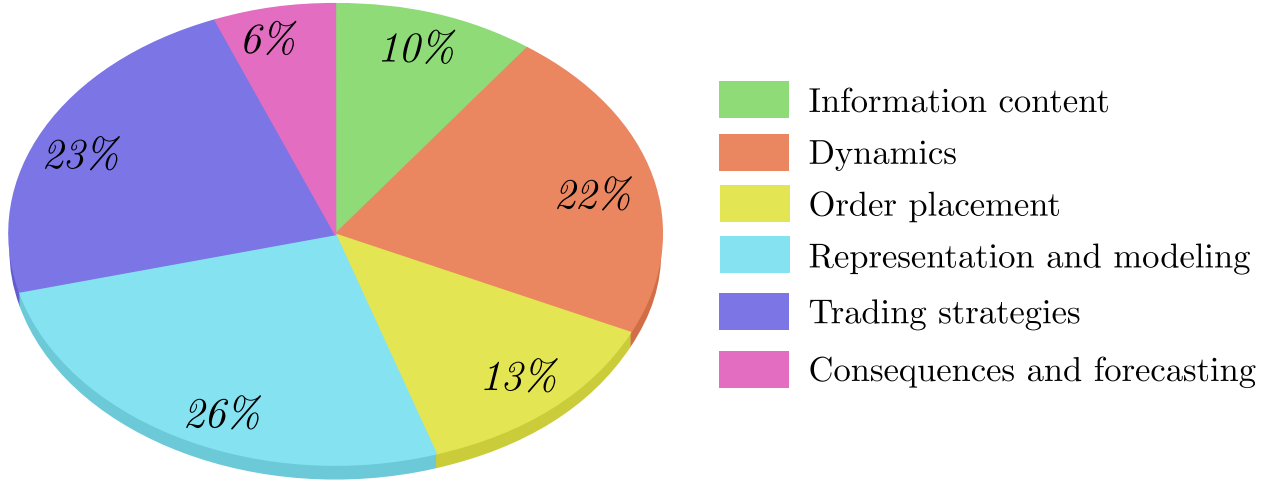
Evolving trading strategies [51] presents two main behaviors which provides a reasonable profit: buy early and hold the stock, selling solely if the price decreases and, buy if the price is lower than the average price over a certain amount of time steps.

Pascual et al.[37], provides a model for aggressiveness which outperforms the ordered probit model for forecasting. This model is a two-stage sequential ordered probit (SOP) model which separates the decision of withdraw, provide or consume liquidity, from the decision of choosing a particular type of order, and allows separating patient traders' order choices from impatient traders' order choices, using information of the whole LOB. This study also provides a confirmation of other studies:

- Asymmetry increases patient buyers operations and impatient sellers operations.
- The higher the bid-ask spread, the higher the frequency of inside-the-quotes limit orders, and increases the frequency of the most aggressive market orders.
- Patient traders increase aggressiveness if the depth at the best quote increases, impatient traders decrease aggressiveness as the thickness of the opposite quote increases.
- Patient and impatient buyers (sellers) increase aggressiveness if the book above (below) the best ask (bid) gets deeper.
- The aggressiveness of incoming impatient (patient) traders is converse to the length of the opposite (same) side of the market.

### 2.2.3. Discussion

Figure 2-1 shows that the topic that has been more widely studied from the proposed categories, is the representation and modeling of the order book, followed by trading strategies



**Figure 2-1.:** References distribution by category.

based on it. This suggests that guiding future research towards this direction could be rewarding because is a very active branch. Table 1 provides information about the references linked to the proposed categories.

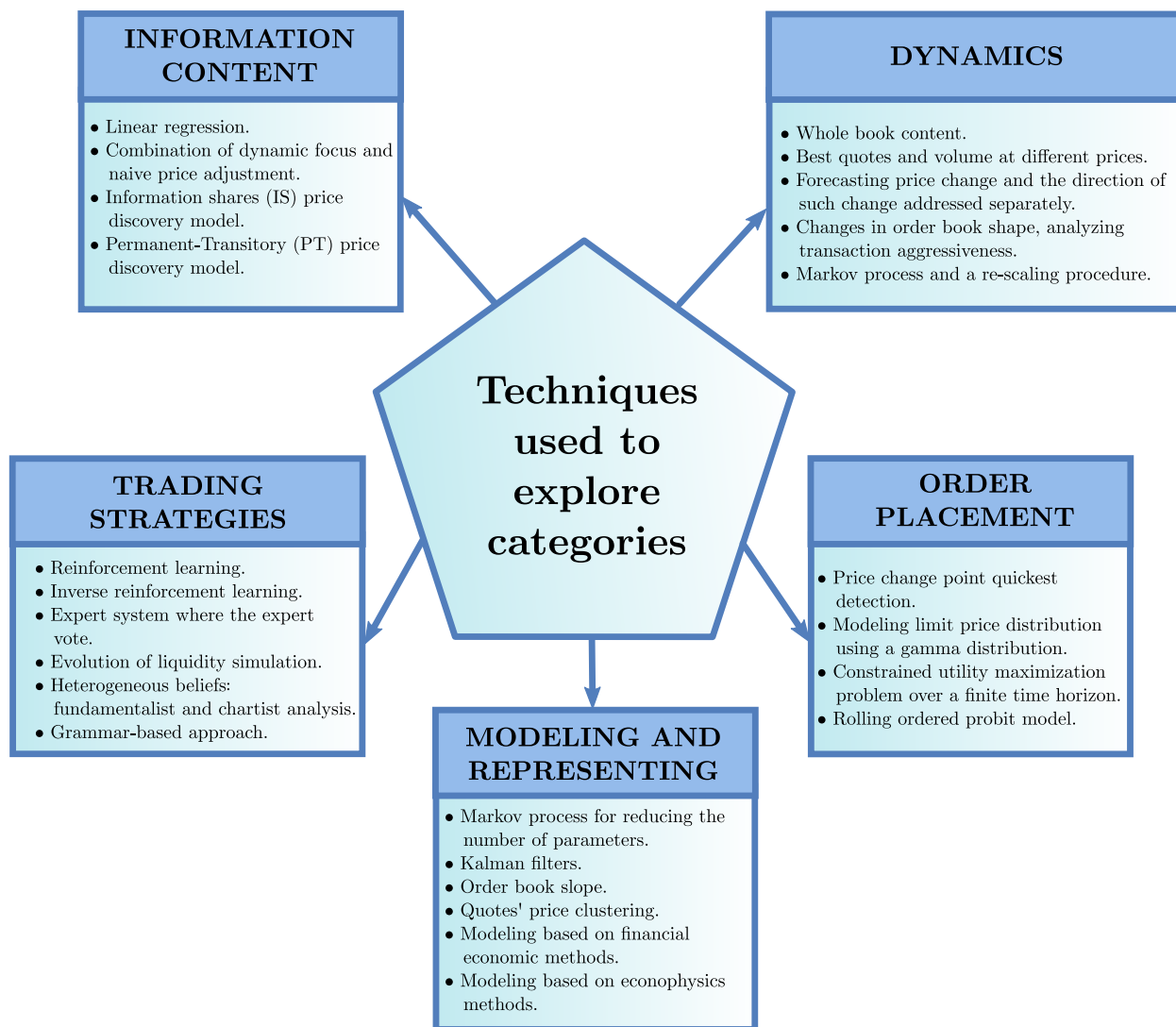
**Table 2-1.:** References classification by subject.

Order book classification	References related
Information content	[53], [22] and [29]
Dynamics	[17], [35],[36],[37], [39], [47] and [48]
Order placement	[16], [21], [26] and [44]
Representation and modeling	[2], [6],[12],[14], [17],[19],[20], [27] and [47]
Trading strategies	[1], [3], [4], [28] and [49],[51],[52]
Consequences and forecasting	[24] and [42]

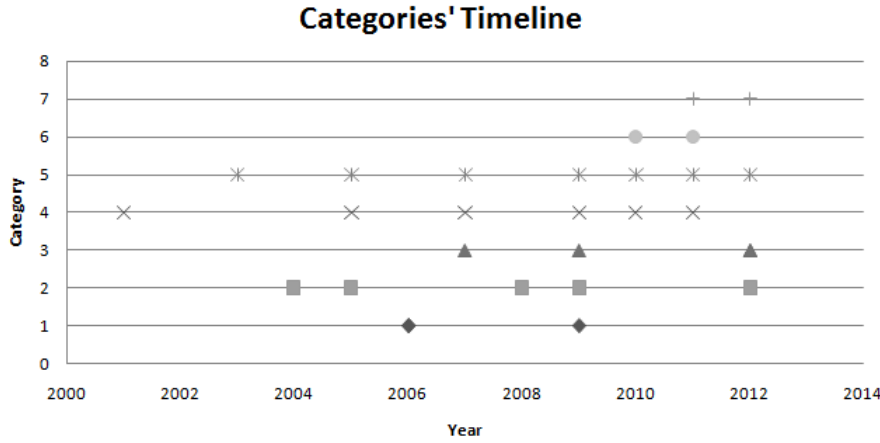
In figure 2, a graphical summary of the main techniques used in each category is provided. Techniques used for order book representation and modeling are diverse, harnessing control theory, physics, econometrics, machine learning and statistical techniques. In future research, the use of more sophisticated probabilistic graphical models (PGM) is encouraged by the author of this paper, due to the successful application of such techniques in problems which involve handling large amount of data. PGM could improve trading strategies also, given their ability to characterize hidden behavior in data.

Figure **2-2** provides a publications' timeline for each category. From the selected literature, based on the variance in the number of publications by category, it can be observed that:

- The categories that presented growing in time were LOB's Representation and modeling, LOB's Order placement, and LOB's Trading strategies.



**Figure 2-2.:** Main techniques used for exploring each category.



**Figure 2-3.:** Publications' timeline for each category: line 1 presents information about LOB's Information content, line 2 presents information about LOB's Dynamics, line 3 presents information about LOB's Order placement, line 4 presents information about LOB's Representation and modeling, line 5 presents information about LOB's Trading strategies, line 6 presents information about LOB's Consequences and forecasting and line 7 presents information about LOB's Dynamics and Representation and modeling simultaneously.

- Publications about LOB's Information content and LOB's Dynamics, showed stagnation in time.
- Articles on LOB's Consequences and forecasting, and about LOB's Dynamics and Representation and modeling simultaneously have recently emerged.

Table 2 presents a market classification by geographic region and association with the order book categories explored on them. The literature related with the Asian and the European markets is ample, but no literature describing African or south American markets was found. Studies from European markets are focused in the understanding of the order book dynamics, while studies from Asian markets are producing results in all categories.

A promising field related to trading strategies is strategy detection. Yang et al.[52] provide a methodology (inverse reinforcement learning) successfully tested which allows specific trading strategies detection. This would provide tools to design specialized strategies for getting profit from an identified strategy.

**Table 2-2.:** Markets classified by geographic region and associated with the order book subjects explored on them (IC, D, OP, RM, TS, CF stands for information content, dynamics, order placement, representation and modeling, trading strategies, and consequences and forecasting, respectively).

Market	Subjects	References	Region
Australian Stock Exchange	RM, IC	[14], [22]	Oceania
HSBC	TS	[3]	Global
Hong Kong Stock Exchange	RM	[2], [27]	Asia
Korea Stock Exchange	TS	[28]	Asia
London Stock Exchange	D,RM,OP	[17], [19], [21]	Europe
New York Stock Exchange	CF, D	[24], [36]	America
Taiwan Stock Exchange	OP,D	[16], [48]	Asia
Shanghai Stock Exchange	TS,IC	[49], [53]	Asia
Shenzhen Stock Exchange	IC	[29]	Asia
Spanish Stock Exchange	D	[37]	Europe
Switzerland Stock Exchange	D	[39]	Europe

#### 2.2.4. Remarks

High Frequency Trading does not produce periods of extreme volatility by itself, but can exacerbate market volatility [24].

FX markets experienced a growth rate of 32.5 % in the last three years, with the United States Dollar as the most traded currency[42]. This information presents the USD behavior analysis as a still interesting research area.

Based on the variance in the number of publications by category surveyed in this paper, categories that presented growing in time were LOB's Representation and modeling, LOB's Order placement, and LOB's Trading strategies. The topic that has been more widely studied from the proposed categories, is the representation and modeling of the order book, followed by trading strategies based on it. Combining growing in time and volume of publications, suggests that guiding future research towards trading strategies based on LOB could be rewarding because is a very active branch.

The relationship between the progress in the proposed categories and the geographic market location were presented in a summary table providing evidence of which markets have been more deeply characterized in the literature.

## 3. Basic Notions and Definitions

### 3.1. Basic Financial Concepts

This section was created in order to provide a guide for those readers which are not familiar with the notions related to financial markets or with concepts belonging to linear algebra and statistical learning.

#### 3.1.1. Forex Markets

*«The foreign exchange market is a global decentralized market for the trading of currencies. In terms of volume of trading, it is by far the largest market in the world.»[40]*

#### 3.1.2. Order book

Cont et al. [7] model the order book as a grid of price ticks, where:

##### Ask

The ask price is defined as [7]:

$$p_A(t) = \inf \{p = 1, \dots, n, X_p > 0\} \wedge (n + 1). \quad (3-1)$$

##### Bid

And the bid price is defined as [7]:

$$p_B(t) = \sup \{p = 1, \dots, n, X_p < 0\} \quad \vee \quad 0. \quad (3-2)$$

##### Depth

Distance from the best price. For the bid side the volume at a certain distance  $i$  is given by [7]:

$$Q_i^B(t) = \begin{cases} X_{P_A(t)-i}(t) & 0 < i < P_A(t) \\ 0 & P_A(t) \leq i < n \end{cases} \quad (3-3)$$

For the ask side:

$$Q_i^A(t) = \begin{cases} X_{P_B(t)+i}(t) & 0 < i < n - P_B(t) \\ 0 & n - P_B(t) \leq i < n \end{cases} \quad (3-4)$$

### 3.1.3. Spread

Is the difference between the best price in the ask side and the best price in the bid side [7]:

$$p_s(t) \equiv p_A(t) - p_B(t) \quad (3-5)$$

### 3.1.4. Trader

Is the agent which produce movements in the order book.

## 3.2. Scientific Visualization

*Scientific Visualization is the mapping of scientific data and information to imagery to gain understanding or insight.*<sup>1</sup>

In 1995 Moorhead and Zhu [33] encouraged scientific community to cooperate in the scope of the scientific visualization due to the massive quantity of information. They describe the visualization process of geometric objects in 3D space and point out the fact that sometimes data are directly mapped into images, skipping the geometric description. In order to represent the data, attention is drawn to the following image components: shape, color and opacity.

This thesis attempts to extract knowledge from scientific visualization according to the definition of Chen et al. [5], provided in table **3-1**.

## 3.3. Heatmap

In 2014 Todd et al. [45] depict an order book heat map to facilitate the analysis of three dimensional data. They suggest the use of heat maps in order to auditor long periods, for

<sup>1</sup>Moorhead, R.J.; Zhifan Zhu, "Signal processing aspects of scientific visualization," Signal Processing Magazine, IEEE , vol.12, no.5, pp.20,41, Sep 1995. DOI: 10.1109/79.410438



**Table 3-1.:** Todd et al. [45] definitions of data, information and knowledge in computational space.

Category	Definition
Data	Computerized representations of models and attributes of real or simulated entities.
Information	Data that represents the results of a computational process, such as statistical analysis, for assigning meanings to the data, or the transcripts of some meanings assigned by human beings.
Knowledge	Data that represents the results of a computer-simulated cognitive process, such as perception, learning, association, and reasoning, or the transcripts of some knowledge acquired by human beings.

instance, years. They use color to map buy and sell orders' depth as shown in **3-1**. They suggest using order book for market surveillance.

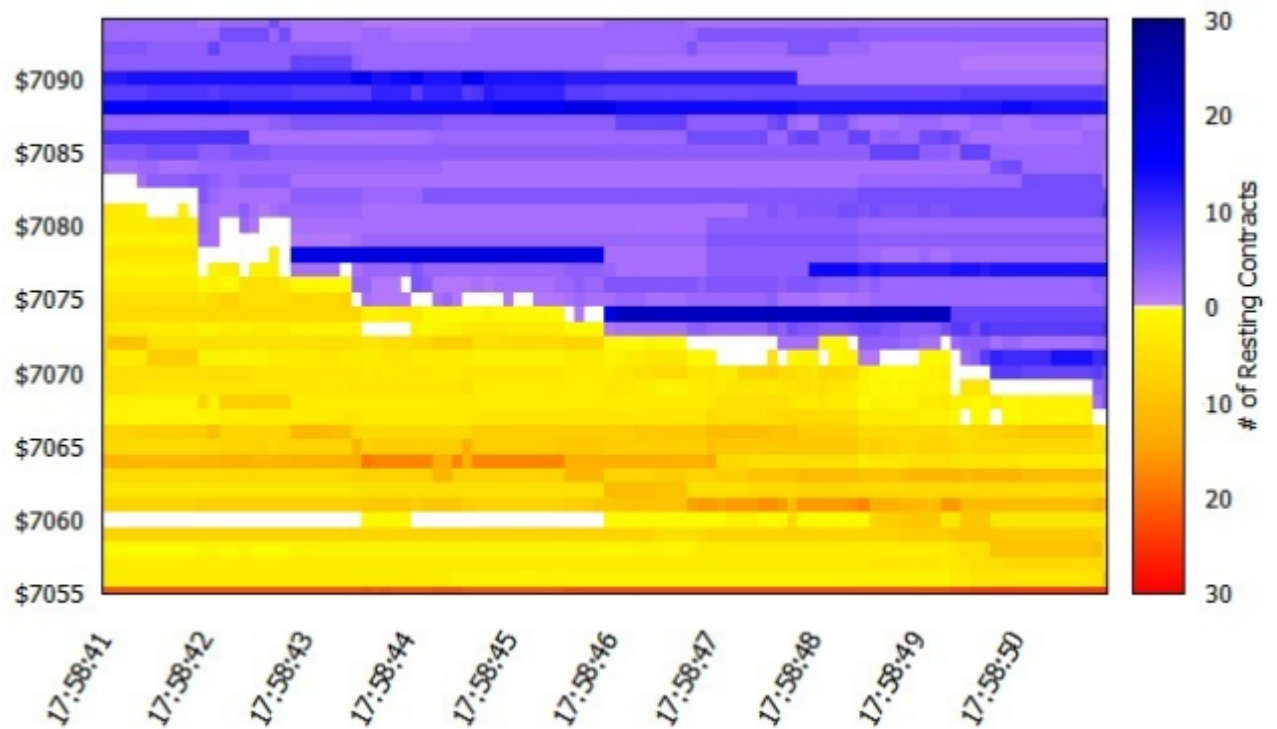
The underlying hypothesis for using this technique for the visualization of the order book is that the order placement produces an effect in the price behavior and therefore it helps in order to identify price trends.

The dataset is explored through several time window sizes in order to find a set of frequent blocks or a set of frequent groups of adjacent blocks strongly associated with a price trend.

The number of times that a pattern in the validation set matches the price trend assigned during the training will be this visualization evaluation method.

### 3.4. Bag of Words

This work is based under the assumption that frequent Price-time-volume structures within the dataset are informative. Linares, Gonzalez et Hernandez [46] indentified individual basic shapes on time series in order to build active trading strategies based on forecasting. In that vein, is reasonable to count every appearance of each pattern and store the number of times that the pattern is associated with a specific trend (in this case bullish or bearish), in order to calculate the probability of the pattern of being related with a trend. This process allows labeling frequent patterns in bullish or bearish patterns with the purpose of building a classifier.



**Figure 3-1.:** Example of Heat Map of Order Book Depth, Todd et al. [45])

The first reference to the Bag of Words method appears in [15] when Harris states that *«it is possible to define a linguistic structure solely in terms of the “distributions” (= patterns of co-occurrences) of its elements. There is no parallel meaning-structure which can aid in describing formal structure. Meaning is partly a function of distribution.»*. Later, in 2003, Sivic et al. [43] present an analogy between text and image retrieval for video retrieval where the construction of the visual vocabulary is made quantizing descriptors in clusters (using k-means) extracted from a fragment of the film. Lopez-Monroy, Gomez et al [30] show the general process for generating a bag of visual words from a set of images.

### 3.5. Wavelets

The idea behind the use of this tool for visualizing the order book is that the reduction of redundant information will make easier the patterns visual detection task.

The dataset is explored through several time window sizes in order to find a set of frequent blocks or a set of frequent groups of adjacent blocks strongly associated with a price trend, as before.

The number of times that a pattern in the validation set matches the price trend assigned during the training will be this visualization evaluation method too.

From this point this section will be an outline from the first chapter of the book «A Wavelet Tour of Signal Processing: The Sparse Way» by Stephane Mallat [31]<sup>2</sup>. Meyer and Mallat provide a systematic theory through the elaboration of multiresolution signal approximations.

#### 3.5.1. Wavelet bases

Wavelet bases reveal the signal regularity through the amplitude of coefficients, and their structure leads to a fast computational algorithm. A WB defines a sparse representation of piecewise regular signals, which may include transients and singularities. In images, large wavelet coefficients are located in the neighborhood of edges and irregular textures.

#### 3.5.2. Quick glossary

- Haar Wavelet: piecewise constant function.

---

<sup>2</sup>Mallat Stephane. A Wavelet Tour of Signal Processing: The Sparse Way. Elsevier. Third Edition. 2009. (pp. 3-16.)

- Orthonormality: two vectors in an inner product space are orthonormal if they are orthogonal and unit vectors. For Euclidean spaces two vectors are orthogonal if and only if their dot product is zero. Is an extension of the concept of perpendicularity amongst vectors.  $\ast$   $\langle, \rangle$  represents the inner product.
- Strömberg Wavelet: A piecewise linear function  $\psi$  that also generates an orthonormal basis and gives better approximations of smooth functions.
- Mayer Wavelet: A family of orthonormal wavelet bases with infinitely continuously differentiable functions.

### 3.5.3. Filter Bank

A filter bank is an array of band-pass filters that separates the input signal into multiple components, each one carrying a single frequency sub-band of the original signal.

#### Conjugate mirror filters

There is an equivalence between continuous time wavelet theory and discrete filter banks. A new interface between digital signal processing and harmonic analysis.

#### Usefulness of mixing continuous infinite analysis with discrete finite analysis

Mixing continuous infinite analysis with discrete finite analysis is useful because the asymptotic results provided the the first one is precise enough to understand the behavior of discrete algorithms but not sufficient for elaborating discrete signal-processing algorithms. The restriction of the constructions to finite discrete signals adds complexity because of the border problems, but with the understanding of the properties of the bases, this issue can be addressed.

### 3.5.4. Wavelets for images

WOB of images can be constructed from wavelet orthonormal bases of one-dimension signals. An algorithm for calculating fastly wavelet coefficients is provided in chapter 7 of [31]. «Like in one dimension, a wavelet coefficient has a small amplitude if the function which defines it is regular over the support of the mother wavelets. It has large amplitude near sharp transitions such as edges.» ( $k$  is the direction and  $2^j$  is the scale and both correspond to a subimage).

#### Approximation and processing in bases

«Sparse representations that reduce the number of parameters can be obtained by thresholding coefficients in an appropriate orthogonal basis.»

### 3.5.5. Sampling

There are two kinds of approximation errors in sampling: Linear approximation error and non linear approximation error. An approximation by thresholding is made by selecting the best vectors in the orthogonal basis of the whole analog signal space. This approximation is not linear. This is important due to the increase of the approximation resolution where the signal is irregular. Approximation support provides geometric information on the orthogonal projection, relative to dictionary, that is a wavelet basis in the given example, so the error it's smaller.

### Sparsity with Regularity

When the image is not that sharp, the non linear wavelet approximation produces small errors. As an adaptation to this issue, more representations with curvelets and bandlets can be used.

### 3.5.6. Compression

When coding a sparse representation via transform codes, the coefficients are approximated by quantized values. Mallat states that «Compression is a sparse approximation problem.». Higher coefficients are associated with geometric properties such as edges.

### 3.5.7. Denoising

Donoho and Johnstone [9] state that «Simple thresholding in sparse representations can yield nearly optimal nonlinear (noise) estimator». Bayes risk «is the expected risk calculated with respect to the prior probability distribution  $\pi$  of the random signal model  $F$ ».

Wald used deterministic models, where signals are elements of a set, without specifying their probability distribution in this set.

### Thresholding Estimators

Donoho and Johnstone proved that in an orthonormal basis, a simple thresholding of noisy coefficients provides a sparse support of the orthogonal projection from the noisy data. «Minimax risk is the lower bound computed over all operators  $D$ .»

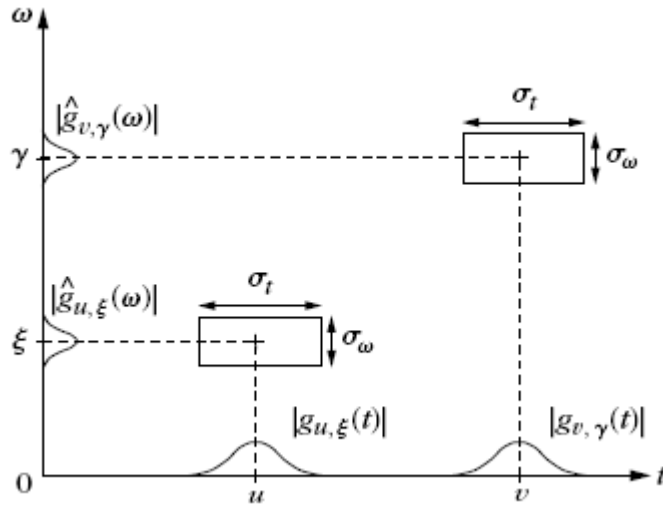
### 3.5.8. Time Frequency Dictionaries

Gabor [13] «proposed decomposing signals over dictionaries of elementary waveforms which he called time frequency atoms that have a minimal spread in a time-frequency plane.» He

also states that «The key issue is to understand how to construct dictionaries with time-frequency atoms adapted to signal properties.»

### 3.5.9. Heisenberg Uncertainty

The uncertainty principle theorem proves that this rectangle has a minimum surface that limits the joint time-frequency resolution:  $\sigma(t, \gamma), \sigma(\omega, \gamma) \geq 1/2$ . See **3-2**



**Figure 3-2.:** Time-frequency boxes representing the energy spread of two windowed Fourier atoms [31].

### 3.5.10. Windowed Fourier Transform

Mallat [31] defines wavelets for images as:

$$f(x) = f(x_1, x_2) : \left\{ \psi_{j,n}^k(x) = \frac{1}{2^j} \psi^k \left( \frac{x - 2^j n}{2^j} \right) \right\}_C \quad (3-6)$$

Where:

$$C = j \in \mathbb{Z}, n \in \mathbb{Z}^2, 1 \leq k \leq 3 \quad (3-7)$$

$$\psi(t) = \begin{cases} 1 & \text{if } 0 \leq t < \frac{1}{2}, \\ -1 & \text{if } \frac{1}{2} \leq t < 1, \\ 0 & \text{otherwise} \end{cases} \quad (3-8)$$

$\psi(x)$  denotes the mother wavelet,  $2^j$  corresponds to the scale,  $2^j n$  the translation,  $k$  to the direction,  $x_1$  and  $x_2$  are the row and the column in the matrix. Here ends the outline from Mallat's book

.

## 4. Microstructure Visualization Tools

This chapter describes in detail the information's nature, its source, the way in which visualizations are built and how to interpret each image and the meaning of the visual components.

Real tick data of foreign exchange rate USD/COP from March to May of 2012 were used in the experiments. LOB provides information about the time, the price and the volume of every request in the market; this information was summarized every minute in a price range of 120 pesos in the best quotes each 20 cents. Volumes were quantized in levels of USD 250,000, which is the minimum trading volume in this market. The maximum volume observed in an order in the analyzed period was 43500000 USD. The maximum price observed was 1862.6 COP and the minimum was 1742.2 COP.

### 4.1. Heatmap based approach

This section presents a heatmap of the discretized volumes and prices of the selected currency as a visualization tool for the market microstructure. Heatmap's construction and interpretation are explained in this section.

#### 4.1.1. First approach

The first approximation to a heatmap based representation, was gathering in a single matrix order price, volume price and its evolution over time. In the x axis, the book's depth (price) was set. In the y axis, time is represented and, finally, lightness indicates orders' cumulative volume. In this way figure 4-1 was assembled.

Thereafter, different thresholds were applied to images produced with the previously described procedure. In this way, images from figure 4.1.1 were obtained. In this figure local time evolving trends in price can be observed.





**Figure 4-1.:** First approach to the order book heatmap visualization.

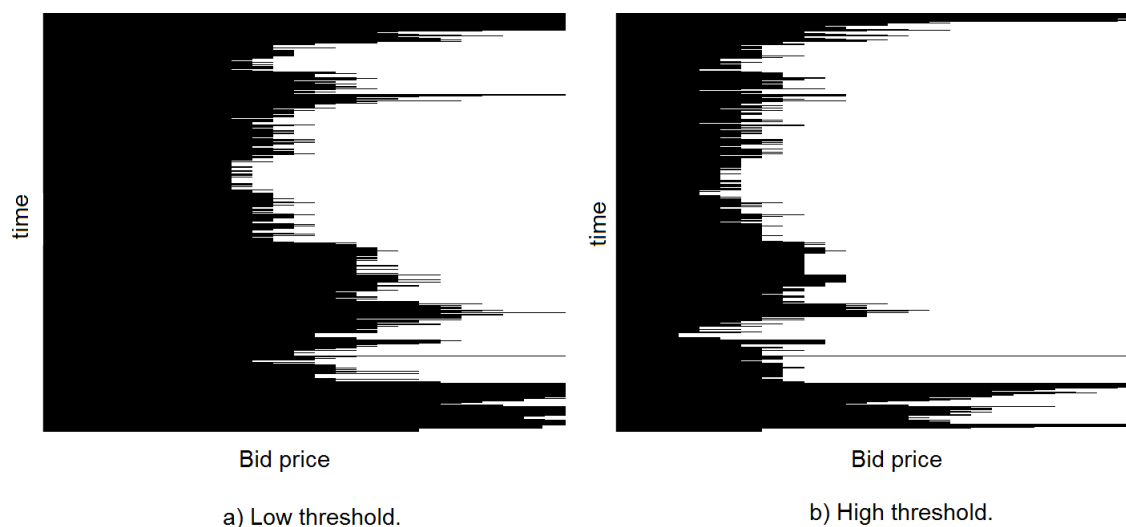
#### 4.1.2. Towards a better information understanding

In order to improve information understanding, figure 4-3 was produced. In this figure, x axis corresponds to price information, y axis come into time and gray level is used for representing orders' volume (not cumulated volume) gathered for price-time intervals. The white meander in the middle of the image corresponds to the spread, it means that in this representation ask and bid order books were combined (ask order book above the meander and bid order book below the meander).

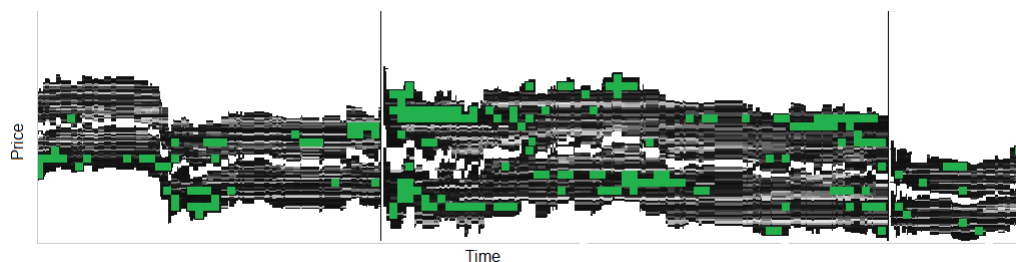
### 4.2. Use of Wavelets in the representation

Wavelets allow a multi-resolution study of the images and a frequency analysis of them. Haar wavelets were selected due to ease of computation and for the speed at which the coefficients can be calculated. Figure 4-4 shows the result of applying four times wavelet transform over an image of the order book, keeping only averages. The level of compression without shape loss can be noted in this image.

In figures 4-5 and 4-6, the effect of applying three times Wavelet Transform, keeping coefficients corresponding to the differences, can be perceived. It can be noted how the image high frequency component is filtered (the two rightmost sub images in each image were amplified in order to preserve visibility). It is expected to find patterns in the high frequency noise



**Figure 4-2.:** Different gray level thresholds for the previous image.



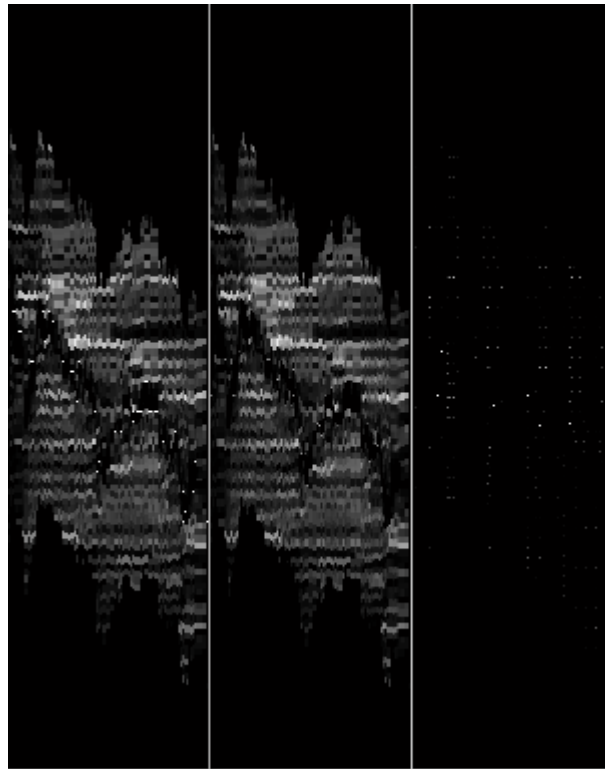
**Figure 4-3.:** Example of position of frequent patterns in a heatmap order book image for 2 days of trading.

which provide predictive information about price changes.

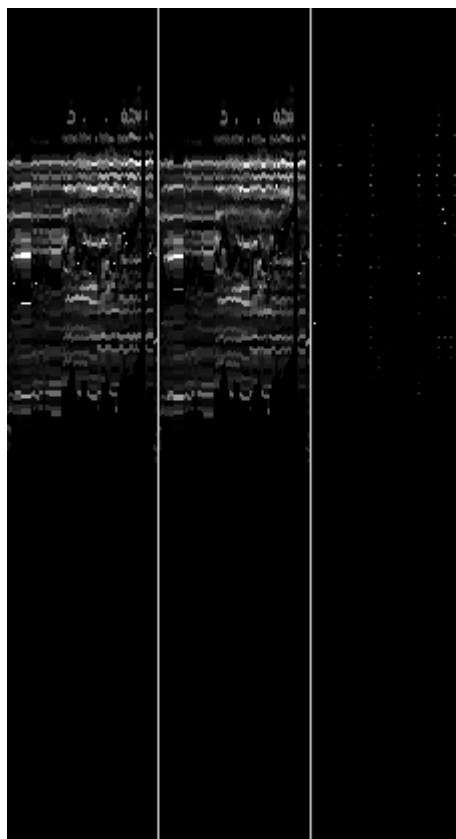
This representation allowed observing graphically long trading periods at a glance. For example, figure 4-7 provides the information of prices in time of three months of trading. The image above, in the figure, keeps high frequency information and, the image below shows the average of the orders' volumes. Recall that x axis represents time and y axis, price.



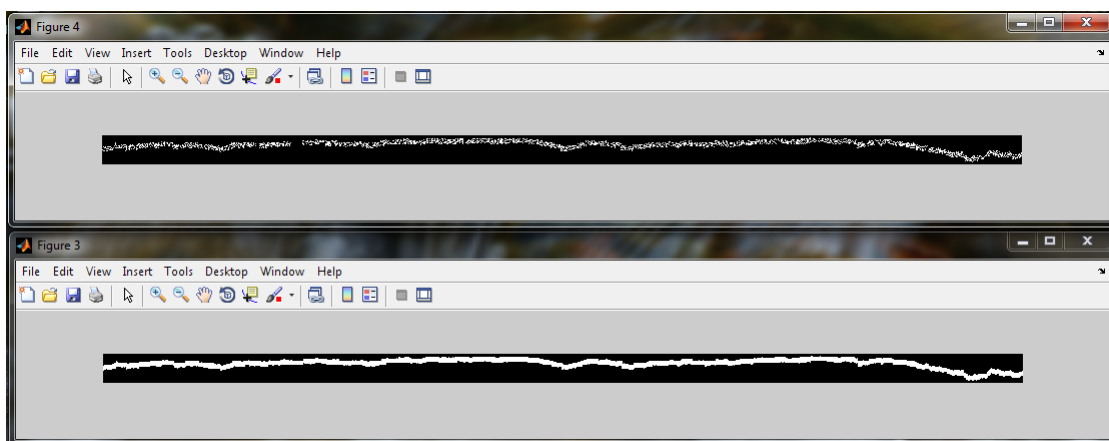
**Figure 4-4.:** Example of image produced by four levels of compression using haar wavelet transform.



**Figure 4-5.:** Example of order book visualization using wavelets based approach (filtering).



**Figure 4-6.:** Example of order book visualization using wavelets based approach (filtering).



**Figure 4-7.:** Three months of trading using one minute resolution as basis for the image construction.

## 5. Experimental Setup

This chapter presents an approximation to the extraction of informative patterns (bullish and bearish indicators) in real tick data from the Colombian Forex Order Book. This pattern exploration is performed on two scenarios: one composed of discretized events and other consisting of equally spaced samples from the discretized representation. Both scenarios are compared via accuracy in order to select the most supportive representation for investing.

Due to the massive amount of information generated in electronic markets, efficient methods and hash functions to handle and operate with this data are helpful. This work provides a methodological approach to manage this kind of data in order to extract information useful to profit generation.

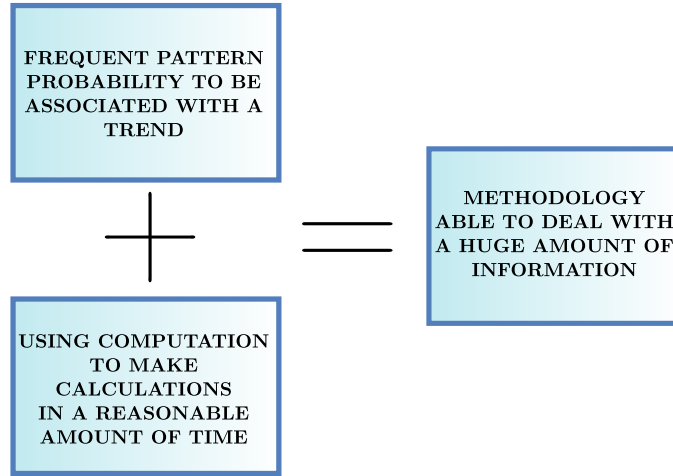
The remainder of this chapter is structured as follows: Section 2 presents the experimental setup. In section 3, results and discussion are provided. Finally, in section 4, Conclusions and future work are issued.

### 5.1. Methodology

In this section a method for handling large amounts of data in Colombian Forex Markets is proposed. Real tick data of foreign exchange rate dollar-peso from March to May of 2012 were used in the experiments. LOB provides information about the time, the price and the volume of every request in the market.

This paper is based under the assumption that frequent Price-time-volume structures within the dataset are informative. In that vein, is reasonable to count every appearance of each pattern and store the number of times that the pattern is associated with a specific trend (in this case bullish or bearish), in order to calculate the probability of the pattern of being related with a trend. This process allows labeling frequent patterns in bullish or bearish patterns with the purpose of building a classifier (See Section 3).

Based on the idea that too frequent patterns are noisy patterns and that no frequent patterns doesn't provide useful information, a threshold was defined by both sides:



**Figure 5-1.:** Visual representation of the described methodology

### 5.1.1. Dataset description

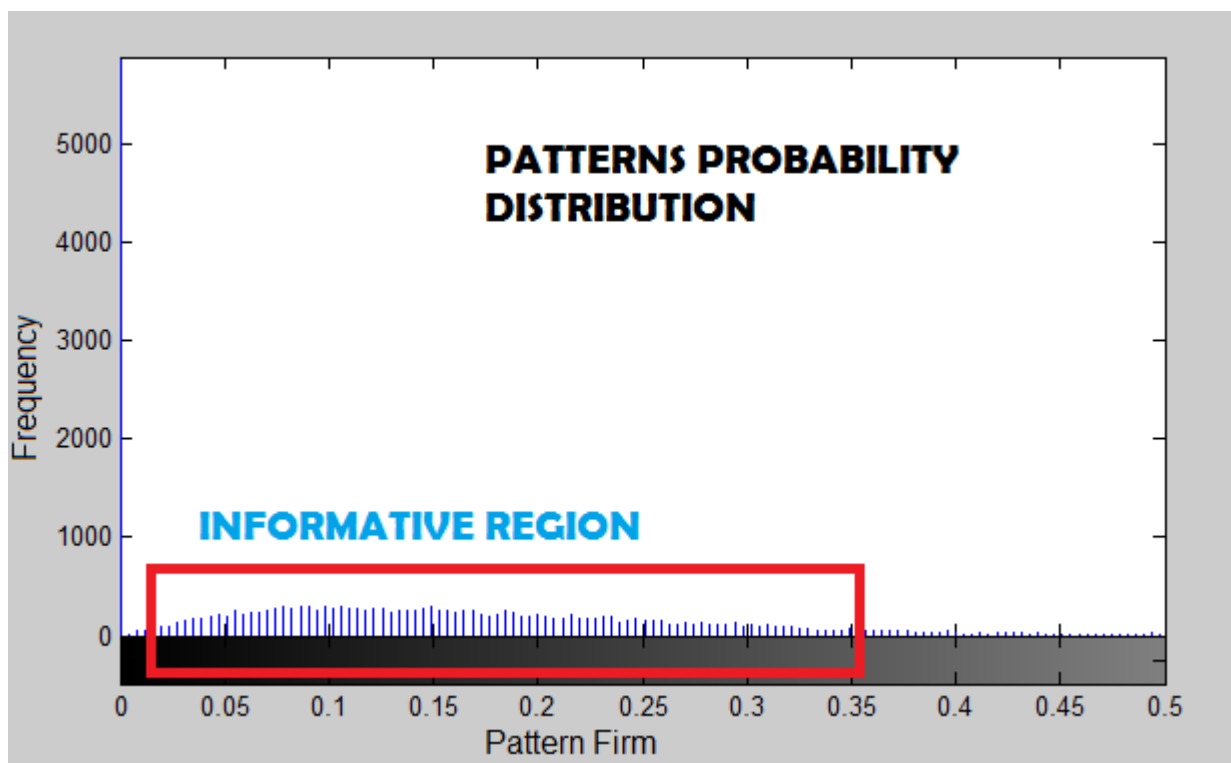
Two datasets with dollar-Colombian peso exchange rate from March to May of 2012 were available: one of them contained an orders summary minute by minute and the other one every ten minutes. For each one of those datasets, the following procedure was adopted:

- Each dataset was split into four subsets.
- Every subset was divided into two parts: one containing 30 % of the samples for training and 70 % for validation. This rare partitioning is due to the yearning of finding frequent patterns even in sets considered small. Furthermore, it works as a mechanism to avoid over fitting.
- The data was discretized in intervals of a) 1 and 10 minutes in time, b) 250000 dollars in volume and c) 5 cents in price.

The dataset arise from real data from the Colombian Forex Market, it was refined by Javier Sandoval<sup>1</sup> and it's a discretization of the order book. There are three matrices named Mv10min, Mv1min and Mv5min. These matrices were renamed as mejoresPuntas10min, mejoresPuntas1min and mejoresPuntas5min which contain the information of the best quotes for Mv10min, Mv1min and Mv5min respectively:

1. Each row in the matrix represents a level of quantization for the price.
2. Each column in the matrix represents a level of quantization for the time.
3. The value of each element represents the order's volume.

<sup>1</sup>Research Professor in Universidad Externado de Colombia, Founder and Project Manager in AlgoCodex



**Figure 5-2.:** Informative Region

The input will be a humongous matrix and it's necessary to chop it into pieces of size(5015,5015), so the function `mapPatterns(explorePatterns(choppedMatrix,width,deepness))` can find the patterns and assigns them to a certain firm (the sum).

Then, the function `findObservedTrend(choppedVectorWithBestQuotes,patternWidth)` will be applied in order to assign each pattern with a trend (bullish, bearish or no trend).

Later, the function `countPatterns` counts the number of different pattern and relates them with their most probable trend. It also saves the number of times that the pattern appears in the `choppedMatrix` (this is used in order to later determine a threshold for considering a pattern frequent). Finally, `findFrequentPatterns(frequency, threshold)` select the patterns which frequency is greater or equal to the given threshold.

Using a threshold of 10 and a probability of being associated with certain trend greater than 0.55, 28 useful patterns were found, 19 were associated with a bearish trend and 9 were associated with a bullish trend (within the next 5 minutes). So, the frequency threshold was reduced to 5.



### 5 minutes timeslots experimental setup

There are available 4524 samples corresponding to approximately 47 days. Cross validation will be made, using 4 folds:

- Each fold will have 1131 samples, 792 ( 70 %) for Training and 339( 30 %) for validation.

With the best one the rest of the dataset will be tested to measure performance.

After that, other 4 dataset will be generated with the same size randomly to compare performance. This 5 minutes experimental setup was discarded later due to its performance.

#### 5.1.2. Pattern exploration

In this section the pattern extraction basic process will be described, to do so, first a definition of pattern will be provided:

**Definition 5.1** *A pattern in this scenario is a submatrix of size  $m$  by  $n$ , where  $m$  is the number of rows and  $n$  is the number of columns.*

In order to build a pattern, the value of consecutive samples was added. Clusters of  $1 \times 10$ ,  $1 \times 20$ ,  $1 \times 40$ ,  $2 \times 10$ ,  $2 \times 20$  and  $2 \times 40$  were summarized in a hash function defined as:

$$H = \sum_{i=a}^{a+n} \sum_{j=b}^{b+m} A_{ij} \quad (5-1)$$

Where  $A$  is the matrix which contains the volumes,  $a$  and  $b$  are the initial positions of each tile and,  $n \in [1, 2]$  and  $m \in S = \{10, 20, 40\}$ .

The first step in order to find relevant patterns in the dataset is to make an exhaustive search of all the existing patterns in the dataset. Given the dataset nature, the exhaustive search could be unfeasible with the available resources (one computer with OS: Windows 7 ultimate 32 bits, processor: Intel(R) Core(TM) i3 CPU M 330 @ 2.13GHz 2.13GHZ, RAM: 4 GB (2,93 available)). Three scenarios are considered:

1. In the whole dataset, there are not repeated patterns.
2. In the dataset there are repeated patterns, but they are not associated strongly with a trend.

3. In the dataset there are repeated patterns and they are associated strongly with a defined trend.

This work is conveyed under the hypothesis that the dataset is on the third scenario. The brute force algorithm proposed for building the patterns catalogue is as follows:

The dataset should be cut off in order to eliminate those frequent patterns that are not informative (remove NaNs or zeros).

ALGORITHM 1 BasicPatternExtraction(dataMatrix,numberOfRowsPattern,  
numberOfColsPattern)

Input: Matrix with the order book information dataMatrix.

Output: Matrix with the replacement of every pattern with  
an assigned number matrixLabels.

```

1 begin
2   i=1, j=1, countOfPatterns=0;, PatternLabel=null,
   PatternsCatalogue=null.
3   for every pattern in dataMatrix
4     if(dataMatrix(i:i+(n-1),j:j+(m-1)) is not in
       PatternsCatalogue)
5       begin
6         add dataMatrix(i:i+(n-1),j:j+(m-1) in
           PatternsCatalogue
7         countOfPatterns= countOfPatterns+1
8         PatternLabel(i,j)= countOfPatterns
9       end
10    else
11      PatternLabel(i,j)= index of PatternsCatalogue
        where dataMatrix(i:i+(n-1),j:j+(m-1) appears.
12    end
13    return PatternLabel, PatternsCatalogue and countOfPatterns.
14  end

```

Another requirement in order to extract relevant patterns is to associate them with trends, so it is possible to modify the previous algorithm for calculating trends simultaneously:

ALGORITHM 2 BasicPatternExtractionWithTrends(dataMatrix,numberOfRowsPattern,  
numberOfColsPattern, tradingPricesVector)

Input: Matrix with the order book information (price, time and volume).

Output: Matrix with the replacement of every pattern with an assigned number  
matrixLabels.

```

1 begin
2   i=1, j=1, countOfPatterns=0, PatternLabel=null, PatternsCatalogue=null.
3   for every pattern in dataMatrix
4     Trend= comparison between tradingPricesVector(j:j+(m-1)) and
5       tradingPricesVector(j:j+(2m-1))
6       (Trend=-1 if tradingPricesVector(j:j+(m-1)) is greater,
7       Trend=1 if tradingPricesVector(j:j+(m-1)) is lower and
8       zero otherwise).
9
10  4   if(dataMatrix(i:i+(n-1),j:j+(m-1)) is not in PatternsCatalogue)
11    5 begin
12    6 add dataMatrix(i:i+(n-1),j:j+(m-1)) in PatternsCatalogue
13    7 countOfPatterns= countOfPatterns+1
14    8 PatternLabel(i,j)= countOfPatterns
15    9 end
16  10 else
17  11 PatternLabel(i,j)= index of PatternsCatalogue where
18    dataMatrix(i:i+(n-1),j:j+(m-1)) apperars.
19  12 end
20  13 return PatternLabel, PatternsCatalogue and countOfPatterns.
21  14 end

```

Given the amount of possible patterns (circa  $6,5 \times 10^7$ ), it is necessary to find efficient ways to convey this mining task. The idea of frequent itemset was proposed by [Agarwal et al. 1993]. In 1994 [Agrawal R, Srikant R] proposed the «A priori» method for mining patterns, using the downward closure property, which states that «A k-itemset is frequent only if all of its sub-itemsets are frequent». Several improvements and generalizations of this algorithm have been done [Savasere 1995, Toivonen 1996, Brin 1997, Cheung 96, Park 95, Agrawal schaffer 96, Geerts 2001]. Nevertheless, the A priori algorithm can generate a huge number of candidate sets or scan several times the database looking for patterns.

In 2000, [Han et al] proposed a FP- growth method which doesn't use candidate generation. The algorithm orders patterns by frequency descending order and then, reduces it into a frequent pattern tree with the association information. The suffix pattern is concatenated with the conditional FP-tree patterns, producing the trees' growth. It is used to mine long patterns because reduces the search time. This algorithm was discarded because this work is not intended to find long patterns, so the cost of building the tree is not justified.

The CLASS Transformation (Eclat) algorithm proposed by [Zaki 2000], could be useful in this work. Every time slot can be considered as a transaction. The items bought in that

transaction would be the patterns found in the column corresponding to that time slot.

When closed itemsets are used, the scalability and interpretability of the mining task is better, but is hard to verify if a pattern is closed.

Sequential pattern mining [Agrawal and Srikant (1995)] is the mining of ordered events; each itemset is a set of events which occurs in the same timeslot. Generalized Sequential Patterns [Agrawal and Srikant (1995)] include time constraints, a sliding time window, and user-defined taxonomies.

Afterwards, frequent patterns exploration was driven following two criteria: a frequency threshold and a probability threshold of being associated with a bullish or a bearish trend. Those patterns which satisfied both standards were labeled as bullish or bearish patterns, respectively. Based on those patterns, a classifier was built and its performance was measured in the remaining subset. This procedure was executed in both, 1 minute and 10 minutes subsets. This method is referred as heat map method in the results section.

With the aim of exploring different levels of resolution, the datasets were sampled using Haar Wavelet Transform four times for 1 minute patterns and two times for 10 minutes patterns. Both sets of coefficients (those from the average and those from the difference) were employed. Mallat describes this procedure as follows:

$$f(x) = f(x_1, x_2) : \left\{ \psi_{j,n}^k(x) = \frac{1}{2^j} \psi^k \left( \frac{x - 2^j n}{2^j} \right) \right\}_{j \in \mathbb{Z}, n \in \mathbb{Z}^2, 1 \leq k \leq 3} \quad (5-2)$$

Where:

$$\psi(t) = \begin{cases} 1 & \text{if } 0 \leq t < \frac{1}{2}, \\ -1 & \text{if } \frac{1}{2} \leq t < 1, \\ 0 & \text{otherwise} \end{cases} \quad (5-3)$$

$\psi(x)$  denotes the mother wavelet,  $2^j$  corresponds to the scale,  $2^j n$  the translation,  $k$  to the direction,  $x_1$  and  $x_2$  are the row and the column in the matrix.

After applying this preprocessing method, new frequent patterns were extracted and threshold criteria were applied. Finally, performance was evaluated.

### 5.1.3. Performance measurement

Performance was evaluated via accuracy. Since different frequent patterns pointing out contradictory trends could emerge several times in one sample or even, the same pattern could

arise more than one time in that sample, a voting system was introduced. In this system, every pattern instance votes for the trend to which it is associated. Lastly, the prediction for every sample is decided by the result of the sum of the votes.

The experimental setup outlined above is used to determine the optimal parametrization of the model.

## 5.2. Results and Discussion

### 5.2.1. Market Trend Visual Bag of Words Informative Patterns in Limit Order Books

Andrea Cruz

Algorithms and Combinatorics Group

ALGOS-UN

Universidad Nacional de Colombia

Bogota, Colombia

Email: amcruz@unal.edu.co

Javier Sandoval

Universidad Nacional de Colombia

Universidad Externado de Colombia

Stevens Institute of Technology NJ, USA

Email: javier.sandoval@uexternado.edu.co

Jaime Nino

Algorithms and Combinatorics Group

ALGOS-UN

Universidad Nacional de Colombia

Bogota, Colombia

Email: jhninop@unal.edu.co

German Hernandez

Algorithms and Combinatorics Group

ALGOS-UN

Universidad Nacional de Colombia

Bogota, Colombia

Email: gjhernandezp@unal.edu.co

In this paper a graphical representation that fully depicts the price-time-volume dynamics in a Limit Order Book (LOB) is presented. Based on this representation and its wavelet transform a visual Bag of Words (BoW) technique is applied to algorithmically find market trend patterns. The BoW technique is tested on information from the Colombian Foreign Exchange bulk market (USD/COP) LOB finding competitive trend prediction patterns. Market trend prediction; Limit Order Book; pattern extraction; wavelets; information extraction; financial markets;

## Introduction

The Limit Order Book (LOB) for a financial instrument gathers two lists: one of buyers orders and other for sellers orders (bids and offers) with the price and volume that they are willing to trade [53]. Figure 5-3 shows an example of a LOB for google shares. The buy orders list corresponds to the bid side and the sell orders side is equivalent to the ask side. These price and volume changing orders lists are placed by market agents who produce book

movements. The best quote is the first order in each list: this is, the maximum price for the bid side and the minimum price for the ask side.

Figure **5-3** is a several years ago snapshot of the visible segment of the order book. Four seconds before this snapshot was taken, there was a transaction at 484.9 USD. At that moment, 484.9 USD is the market price. LOB contains many more entries, but only those close to the best quote are shown. In this book representation, orders at the same price are not aggregated. Using the same color implies orders at the same price. Nevertheless, orders can be accumulated at the same price level as shown in this work.

Some markets, called matched, have an engine that pairs and executes orders according to matching criteria. Other markets, called unmatched, require the intervention of a trader to execute an order. Examples of order execution can be found in figure **5-4**. Events (transactions) define discrete time points; this is not a continuous process, in this figure, lines parallel to the vertical axis represent events.

An example of bid side distribution can be seen in figure **5-5**. The difference between the best quote of the ask side and the one in the bid side is called "spread". Within a given price, several orders from different traders can be cumulated in order to shape the volume for that quote; the priority for the order's execution is assigned according to the order placement queue.

An order is a combination of price and volume which is placed within the order book. Market actors can perform different kind of operations in the order book. Brokers can put or remove orders at the best quotes. Other orders, placed in positions other than the best quote can also be put and removed. Market participants can modify their orders changing volume or price as desired, they also have the possibility of execute an order at the best quote or subsequent orders if the volume available at the best quote is not enough for the trader's demand, see figure **5-6**

The LOB evolves by discrete time events produced by buyers and sellers which insert, modify or delete limit orders and also, by execution of orders, either by the engine on matched markets or by liquidity takers on unmatched markets [10]. An illustration of some of these events is presented in figure **5-5**.

A Limit Order Book can be considered a variable length list where transactions occur in non-uniform time intervals. At the instant when transactions occur, the spread is zero. As shown in figure **5-5**, for a given time slot, information of prices and the cumulative volume for each price are provided.



**Figure 5-3.:** Limit Order Book example [34].

Recently, there has been a strong interest in modeling and analysis of LOBs dynamics because these provide richer information about the market trends than closing prices do [17], [35], [36], [37], [39], [47] and [48]. Closing prices are in general publicly available today while LOB information is usually only available to brokers or professional investors willing to pay for this extra information.

Huang et al. [17] model order book high frequency dynamics in the London Stock Exchange using second level data (best quotes and volume at different prices) capturing the arrival of orders of different sizes. Ahn et al. [2], conclude that the further from the best a quote is, the less information it provides.

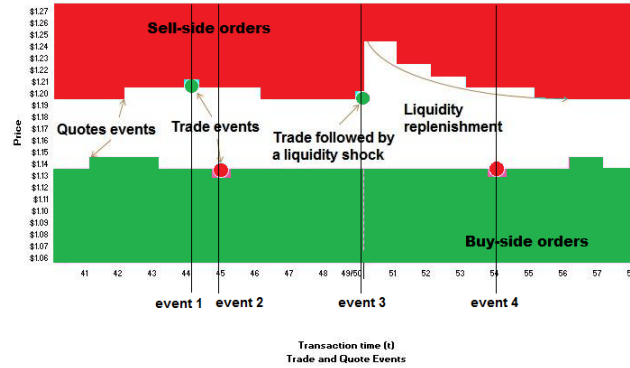


Figure 5-4.: Order execution example [54].



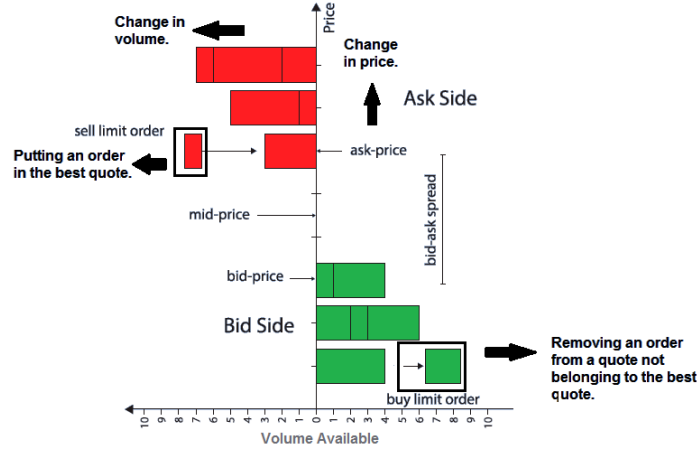
Figure 5-5.: Bid-ask distribution example [?].

Palguna et al. [36] describe changes on LOB, specifically, forecasting price changes and separately the direction of such change. This approximation does not produce statistically significant results.

There are multiple models that describe the dynamics of the order book. In terms of characterizing the LOB of multiple markets, see for example, [39]. Vvedenskaya et al. [47] have analyzed the order book shape and transactions aggressiveness of liquidity seekers and liquidity providers. Moreover, they have moved forward to describe the LOB evolution assuming markovian properties and describing the book using parametric models.

There have been found contradictory evidence about LOB information content. [37] and [17] state that the whole book (not only best quotes) helps to determine price direction.





**Figure 5-6.:** Operations in the order book, modification of [?] for illustration purposes.

However, it seems to be an indirect relationship between information content and distance from best quotes [2].

Some authors present evidence of increased price predictability when LOB information is gathered and used [41]. However, there are other authors who describe not statistically significant results [36]. It seems that final results are highly dependent on the type of market that is been analyzed.

Other studies have tried more complicated representations for the LOB. Jiang et al. [21] propose a representation of four parameters by snapshot. A Kalman filter is used to estimate a linear dynamic system state and provide a liability measure, being used for prediction and filtering. Cheng et al. [6] used order book's slope to outline its informative content.

Due to the massive amount of information generated in electronic markets [31], efficient methods and hash functions to handle and operate with this data are helpful [50]. Therefore, this paper starts providing a methodological approach to manage this kind of data in order to extract useful information to construct profitable trading strategies.

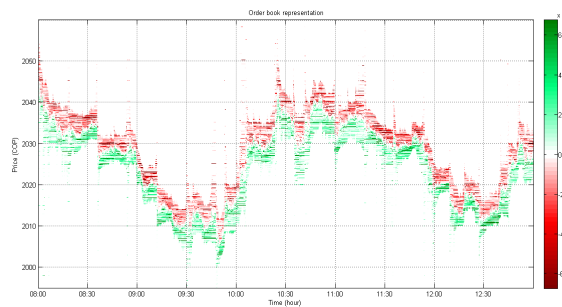
The remainder of this paper is structured as follows: Section 2 presents the experimental setup. In section 3, results and discussion are provided. Finally, in section 4, Conclusions and future work are issued.

### Experimental Setup

This section proposes a method for handling large amounts of data in Colombian Forex Markets. Experiments were conducted using real tick data of foreign exchange rate USD-COP from March to May of 2012. LOB provides information about time, price and volume for every request in the market; this information was summarized every minute, in a price range of 120 COP in the best quotes, using a 20 cents mark up. Volumes were quantized in levels of USD 250,000, which is the minimum trading volume for this market. The maximum volume observed for a particular order during the analyzed period was 43.5 USD millions. The maximum price observed was 1,862.6 COP and the minimum was 1,742.2 COP.

This paper proposes a tool to facilitate a human trader locating visual patterns. The proposed visualization is built as follows:

1. Book event's discretization: Events which happen every minute are aggregated. Increases along the x coordinate indicate progress over time.
2. Trading volume quantization: Volumes within a range of prices in an instant of time are cumulated; higher volumes are indicated with a lower level of lightness.
3. Price quantization for each time unit: An increase along the y-axis in the picture indicates a higher price.



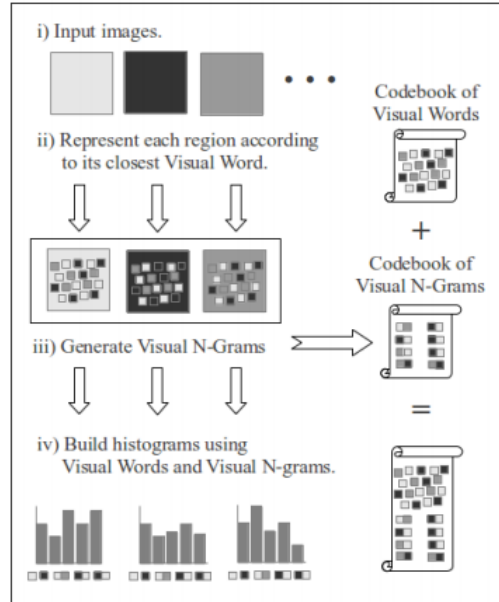
**Figure 5-7.:** Example of the proposed order book visualization.

A similar representation named *BookMap X-ray*<sup>2</sup> is provided by VeloxPro. This representation is also volume inspired, however is built in gray scale and it does not take advantage of chromatic differences or relationships.

<sup>2</sup><http://www.bookmap.com/>

This paper is based under the assumption that frequent price-time-volume structures within the dataset are informative. Linares, Gonzalez et Hernandez [46] indentified individual basic shapes in time series in order to build active trading strategies based on forecasting. In that vein, is reasonable to count every appearance of each pattern and store the number of times that the pattern is associated with a specific trend (in this case bullish or bearish), in order to calculate the probability of the pattern of being related with a trend. This process allows labeling frequent patterns in bullish or bearish formations with the purpose of building a classifier.

The first reference to the Bag of Words method appears in [15] when Harris states that *it is possible to define a linguistic structure solely in terms of the distributions (= patterns of co-occurrences) of its elements. There is no parallel meaning-structure which can aid in describing formal structure. Meaning is partly a function of distribution.* Later, in 2003, Sivic et al. [43] present an analogy between text and image retrieval for video retrieval where the construction of the visual vocabulary is made quantizing descriptors in clusters (using k-means) extracted from a fragment of the film. Lopez-Monroy, Gomez et al [30] show the general process for generating a bag of visual words from a set of images.



**Figure 5-8.:** Image Representation through Bag-of Visual-Ngrams, extracted from [30]

**Definition 2** For this context, a pattern is a matrix of market events, an aggregation of volumes spatially organized by price and date.

Market behavior

**Definition 2** A trend is the sign of the difference between the current and the previous observation.

For this case the observed phenomenon is the currency price. Two possible market behaviors are specified for this technique:

$$t(t) = \begin{cases} 1 & \text{if } p_o(t) - p_o(t-1) > 0 \\ -1 & \text{if } p_o(t) - p_o(t-1) < 0 \\ 0 & \text{otherwise} \end{cases} \quad (5-4)$$

A trend is said to be bearish when  $t(t) = -1$  and bullish if  $t(t) = 1$ .

#### Dataset description

Two datasets with USDCOP rate from March to May of 2012 were available: one of them contained an orders summary minute by minute and the other one every ten minutes. For each one of these datasets, the following procedure was adopted:

- Each dataset was split into four subsets.
- Every subset was divided into two parts: one containing 30 % of the samples for training and 70 % for validation. This rare partitioning was chosen for finding frequent patterns even in sets considered small. Furthermore, it worked as a mechanism to avoid over fitting.
- The data was discretized in intervals of a) 1 and 10 minutes in time, b) 250,000 dollars in volume and c) 20 cents in price.

#### Pattern exploration

In [38], Rajaraman et al. define a Hash function as *a function  $h$  which takes a hash-key value as an argument and produces a bucket number as a result*. In order to build a pattern, the value of consecutive samples was added. Clusters of  $1 \times 10, 1 \times 20, 1 \times 40, 2 \times 10, 2 \times 20$  and  $2 \times 40$  were summarized in a hash function  $H$ .

This function was used before the pattern extraction process to compress a whole pattern in a single value. In this way, it is possible to deal with smaller matrices.  $H$  is defined as:

$$H = \sum_{i=a}^{a+n} \sum_{j=b}^{b+m} A_{ij} \quad (5-5)$$

Where  $A$  is the matrix which contains the volumes,  $a$  and  $b$  are the initial positions of each tile and,  $n \in [1, 2]$  and  $m \in S = \{10, 20, 40\}$ .

Afterwards, frequent patterns exploration was driven following two criteria: a frequency

threshold and a probability threshold of being associated with a bullish or a bearish trend. Patterns which satisfied both standards were labeled as bullish or bearish, respectively. Based on those patterns, a classifier was built and its performance was measured in the remaining subset. This procedure was executed in both, 1 minute and 10 minutes subsets. This method is referred as heat map method in the results section.

#### Wavelet transform

With the aim of exploring different levels of resolution, the datasets were sampled using Haar Wavelet Transform four times for 1 minute patterns and two times for 10 minutes patterns. Wavelet Transform is used for sampling and compressing the dataset to work with even smaller matrices.

For ease, every matrix is treated like an image. Mallat [31] defines wavelets for images as:

$$f(x) = f(x_1, x_2) : \left\{ \psi_{j,n}^k(x) = \frac{1}{2^j} \psi^k \left( \frac{x - 2^j n}{2^j} \right) \right\}_C \quad (5-6)$$

Where:

$$C = j \in Z, n \in Z^2, 1 \leq k \leq 3 \quad (5-7)$$

$$\psi(t) = \begin{cases} 1 & \text{if } 0 \leq t < \frac{1}{2}, \\ -1 & \text{if } \frac{1}{2} \leq t < 1, \\ 0 & \text{otherwise} \end{cases} \quad (5-8)$$

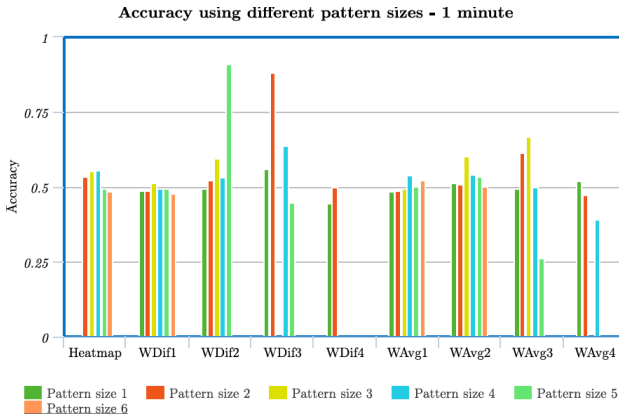
$\psi(x)$  denotes the mother wavelet,  $2^j$  corresponds to the scale,  $2^j n$  the translation,  $k$  to the direction,  $x_1$  and  $x_2$  are the row and the column in the matrix.

Both sets of coefficients (those from the average and those from the difference) were employed. After applying this pre-processing method, new frequent patterns were extracted and threshold criteria were applied. Finally, performance was evaluated.

#### Performance measurement

Performance was evaluated via accuracy. Since different frequent patterns pointing out contradictory trends could emerge several times in one sample or even, the same pattern could arise more than one time in that sample, a voting system was introduced. In this system, every pattern instance votes for the trend to which it is associated. Lastly, the prediction for every sample is decided by the result of the sum of the votes.

The experimental setup outlined above is used to determine the optimal parametrization of the model.



**Figure 5-9.:** Predictor's accuracy using different 1 minute pattern sizes.

### Cluster classification

In order to introduce the idea of similarity between patterns, k-means was applied over the patterns matrix built over the order book. The procedure made is as follows:

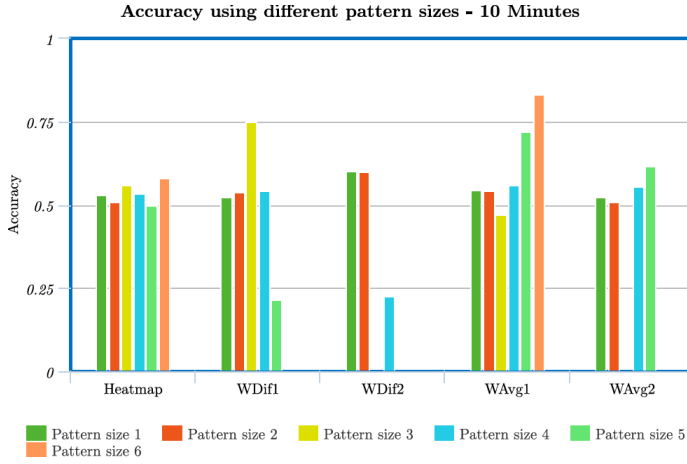
- Time-price-volume matrix was chopped and adjusted for working in regions near to the spread.
- The matrix was traversed using tiles of 30x30.
- Every pattern was assigned to a cluster using k-means.
- In a new matrix, the size of each cluster was registered in order to keep track of the frequency at which patterns were assigned to every cluster.
- Finally, each cluster was labeled with the observed trend when the probability of being associated to it was higher than 0.55.

### Results and Discussion

#### Results

The results of the parameter search for the proposed approach are shown in figures 5-9 and 5-10.

In some cases, during experimentation, no frequent patterns were found in accordance with the defined thresholds for pattern frequency and probability of being associated with a specific trend. Those patterns selected as frequent should, simultaneously, be associated with a trend with a probability higher than 0.55 and should appear more than five times in the



**Figure 5-10.:** Predictor's accuracy using different 10 minutes pattern sizes.

corresponding subset. In those cases in which no frequent patterns were found, the accuracy was infinity, hereby, those accuracy values were discarded. For that reason, some bars are missing in figures 5-9 and 5-10 to avoid misinterpretation.

The average number of frequent patterns by subset and pattern size was 70. It was found that frequent patterns lost its ability to predict trend over time, supported on this result, the future development of an adaptive method to overcome this difficulty is suggested.

The best performance for this method was provided by 1 minute frequent patterns, using as pre-processing method the application of the Haar Wavelet Transform twice, using those coefficients corresponding to the difference between consecutive values in the volumes matrix. In this case, pattern size 5 (2x20 items tile) were used. The accuracy achieved was 0.91 in the best case.

For 10 minutes frequent patterns, best performance was achieved using Haar wavelet transform once with the coefficients corresponding to the average between consecutive values in the volumes matrix and, pattern size 6 (2x40 items tile). The accuracy accomplished was 0.8333 in the best case.

Overall, the pattern size which provided the highest average accuracy was pattern size 3 (1x40 items tile) for 1 minute patterns and pattern size 6 (2x40 items tile) for 10 minutes patterns. The pre-processing method which lead to the best average accuracy was: for 1 minute patterns, three iterations of wavelet transform using difference coefficients, for this scenario the accuracy was 0.632175 at the best case. For 10 minutes frequent patterns, the best performance was obtained by using Haar wavelet transform once, selecting those coefficients corresponding to the average. This pre-processing method produced an accuracy of

0.61168 at the best case.

Figures 5-9 and 5-10 compare best performance by pattern size among the different pre-processing choices: heat map, Haar wavelet transform differences and average coefficients until 4 iterations for Figure 5-9 and until 2 iterations for Figure 5-10, every color represents a different pattern size.

#### Clusters approach performance

The uppermost picture in figure 5-12 shows the centroid patterns assigned to each cluster, and the bottommost picture presents a matrix with example patterns assigned to each cluster. A sample of these kind of patterns can be observed more closely in figure 5-11. Figure 5-13 provides the histogram that represents the size of each cluster.

Figure 5-13 introduces the performance of each cluster: in red the probability of the cluster patterns of being associated with a bearish trend, and in green, the probability of that cluster patterns of being associated with a bullish trend.



**Figure 5-11.:** Example of patches associated with clusters.

Figure 5-15 shows the cumulated return (COP amount by every invested dollar) for 6 months of trading using the cluster approach.

#### Discussion

With an optimal parametrization, the performance is better compared with a random guess (acc. 0.91 vs. 0.5). However the performance does not keep that level from one subset to another, suggesting the existence of seasonal patterns. This phenomenon is described by Jiang et al. in [21]. To test this hypothesis, two subsets with global opposite trends were identified using moving average and the performance of 400 different pattern sizes was evaluated using only the hash function. In this scenario, 35 pattern sizes achieved accuracy higher than 0.6 in the first subset, when tested in the opposite trend, only 4 kept their accuracy. Results are shown in figure 5-16.

In order to assure this method's scalability, testing with larger datasets is required. Figures 5-9 and 5-10 show the result of applying this methodology to obtain a strong classifier gathering several weak classifiers using a simple voting system.



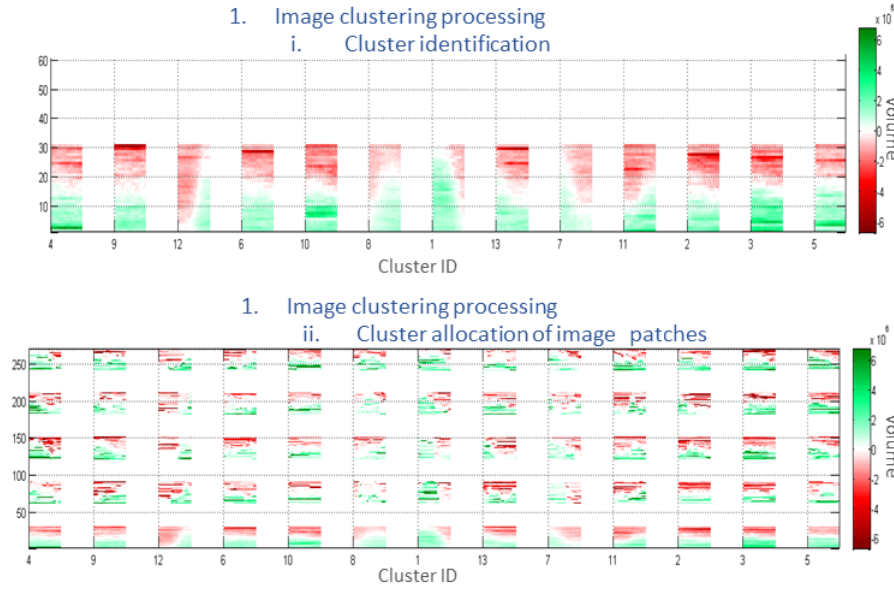


Figure 5-12.: Clusters matrix and centroids.

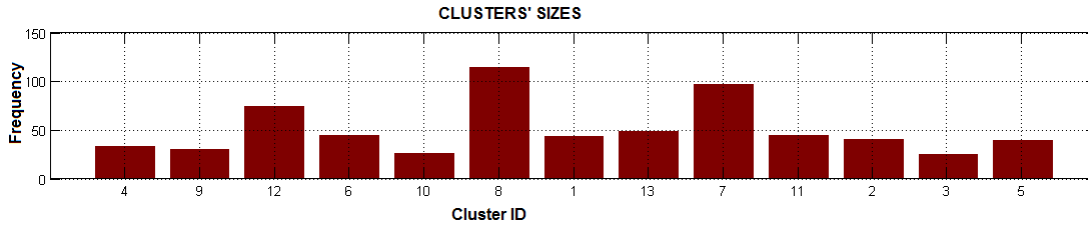


Figure 5-13.: Frequency at which patterns are associated with a certain cluster.

## Conclusions

This work introduces a method for calibrating a trend predictor for currency rates using three months of real tick data from the Colombian Forex Market (USDCOP rate). This method reduces exploration time (learning time) due to the employment of a hash function and the utilization of sampling via Haar wavelet transform. The results suggest the presence of seasonal patterns in this market.

As future work, it would be useful to introduce an adaptive method which detects when frequent patterns start to lose their ability to predict trend over time to start a new training stage for the classifier.

## Acknowledgment

The authors thank Javier Rincon from Acciones y Valores for providing the data set with real data from the Colombian Forex Market

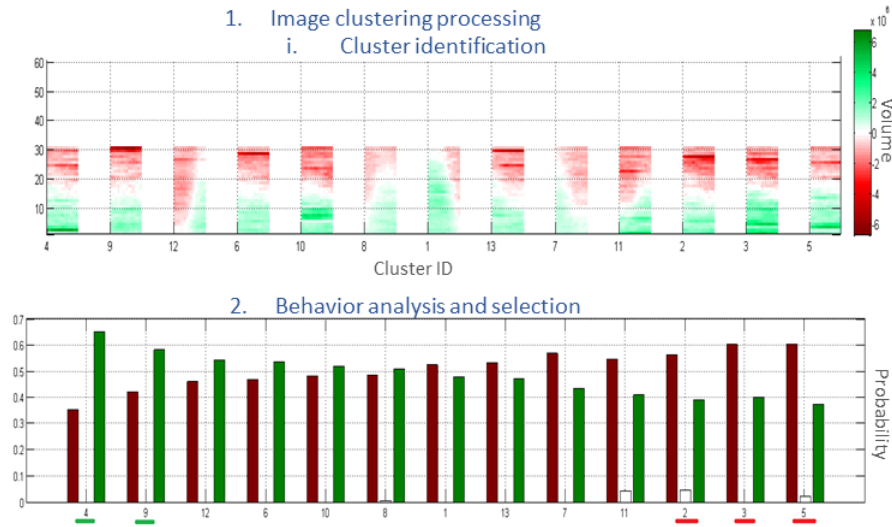


Figure 5-14.: Clusters and their performance.

### 5.3. Conclusions

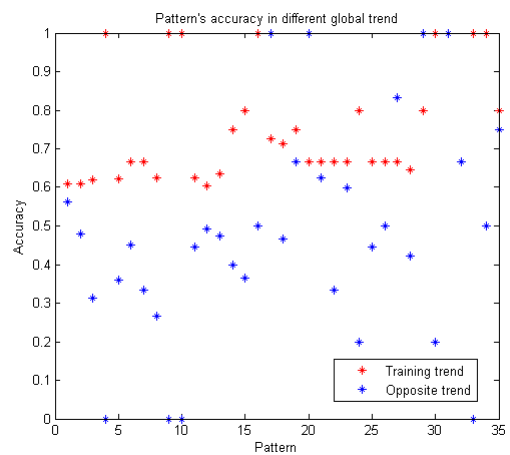
This work introduces a method for calibrating a trend predictor for currency rates using three months of real tick data from the Colombian Forex Market (USD-COP exchange rate). This method reduces exploration time (learning time) due to the employment of a hash function and the utilization of sampling via Haar wavelet transform.

Thanks to the optimal parametrization of the model, obtaining an accuracy of 0.91 was possible, far surpassing a random guess.

As future work, it would be useful to introduce an adaptive method which detects when frequent patterns start to lose their ability to predict trend over time to start a new training stage for the classifier.



**Figure 5-15.:** Cumulated retrurn usd cop in six months.



**Figure 5-16.:** Patterns accuracy within and outside a global trend.

## 6. Adaptive Method for Market Microstructure Exploration

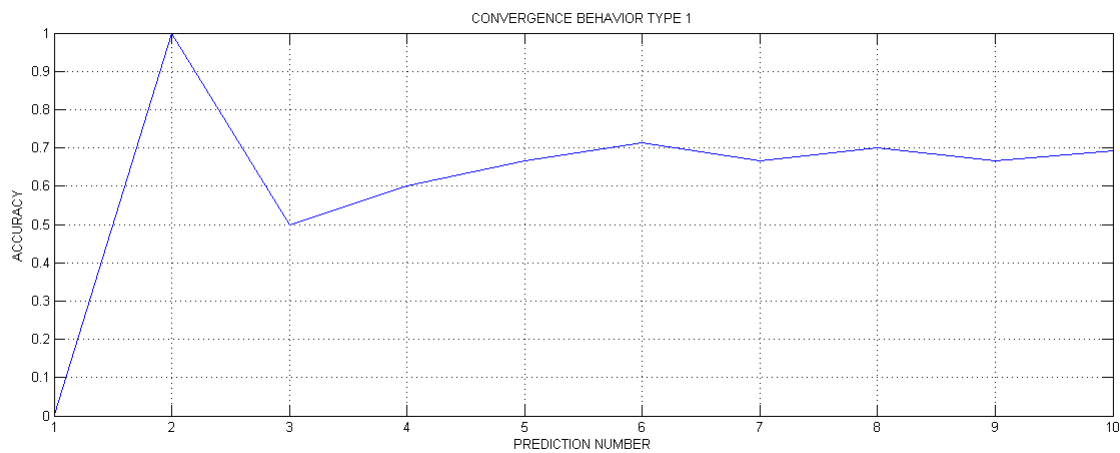
In the previous chapter was stated the need for an adaptive method which address frequent pattern selection when the ability to predict a trend is diminishing. In this chapter an online method for informative frequent patterns is presented.

### 6.1. Method's description

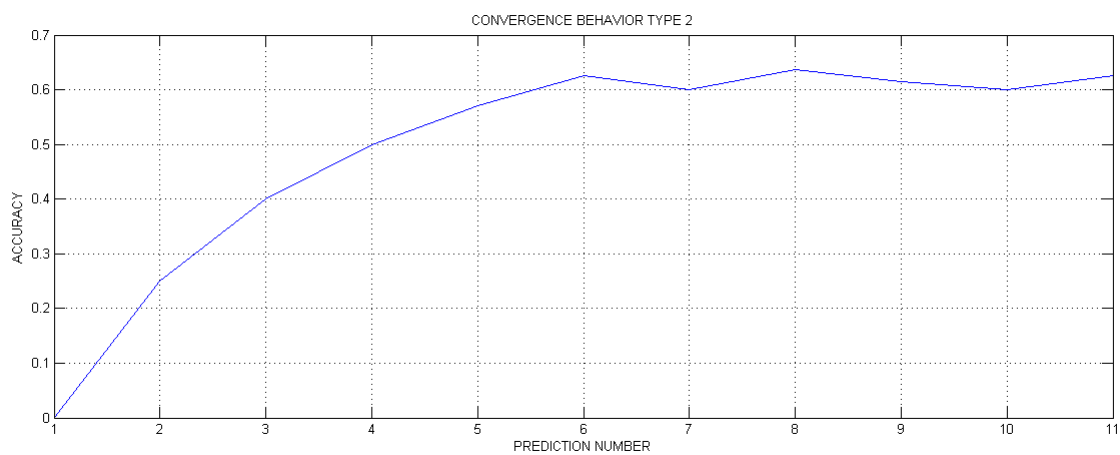
This section depicts the procedure for selecting frequent patterns in an adaptive way.

1. Make an initial training stage in which a general market trend is identified.
2. Every time that a new sample arrives for classification, recalculate the probabilities associated to each trend for the patterns found in the sample.
3. When the count for a new pattern apparitions reaches an specific amount and its probability of being associated with a determined trend surpasses a threshold, add that pattern to the informative frequent patterns' dictionary.
4. When the count of misclassifications for a pattern from the informative frequent patterns' dictionary reaches an specific number or when its probability of being associated with a determined trend falls to a determined threshold, remove that pattern from the informative frequent patterns' dictionary.

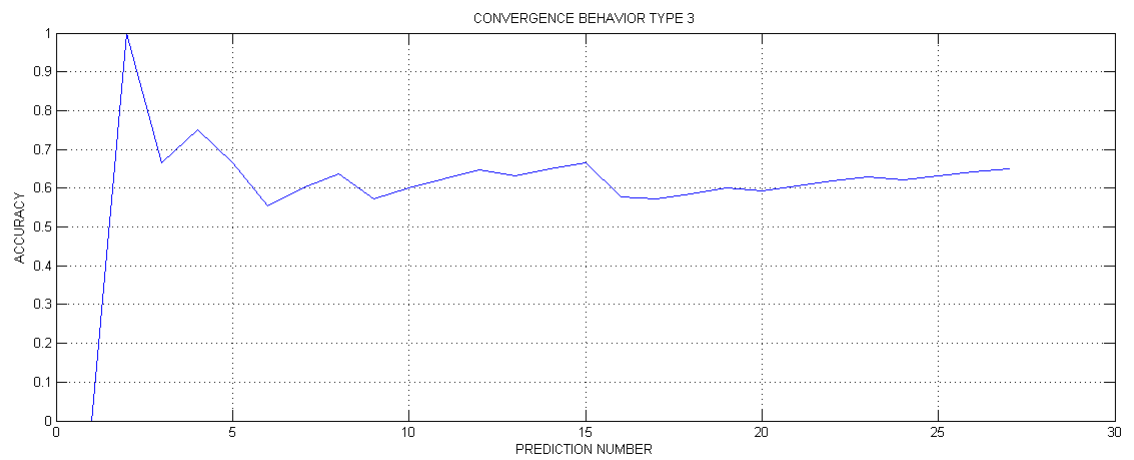
The previous procedure was tested for 300 different patterns' sizes, producing 29 configurations where patterns highly associated with a trend were found and whose accuracy was above 0.6. From those 29 configurations, 10 produce 5 or more predictions on the dataset. Three different convergence behaviors were found as it could be observed in figures **6-1,6-1,6-1**.



**Figure 6-1.:** Convergence behavior type 1.



**Figure 6-2.:** Convergence behavior type 2.



**Figure 6-3.:** Convergence behavior type 3.

## 7. Conclusions and Future work

### 7.1. Conclusions

- A systematic literature review about the Order Book was presented, there is no evidence of previous work made based on the LOB for the Colombian Forex Market until the writing date of the second chapter.
- A methodology which allows representing properly the Colombian Forex Market Order Book information dynamics is presented. The visualization tools presented in this work, can provide the user with a global understanding of a selected time interval in the Colombian Forex Market.
- Wavelet Heatmap visualization presents in a summarized and efficient way the order book information.
- A trading strategies detection system for the Colombian Forex Market using Order Book information was designed by means of a frequent patterns exploration approximation.
- Given the seasonality of the found patterns, the presented strategy was reformulated as an adaptive strategy which detects when a pattern is losing predictability in order to start a new training stage for detecting new informative patterns.
- The performance of the proposed system in supporting the financial decision making process in the Colombian Forex Market was evaluated.

### 7.2. Future work

The fact that the proposed strategies provide useful results with a relatively small dataset from one single currency, throws as a natural consequence the need of testing them in broader datasets and new assets, even for portfolio selection.

On the other hand, an important feature of the Wavelet based approach is that is highly parallelizable, allowing easy implementation in distributed systems such as GPUs. In order to reduce latency, it would be useful to implement the presented algorithms directly on hardware, for instance in a FPGA.

## A. Appendix: Heatmap approach results

Pattern Size 1	1x10
Pattern Size 2	1x20
Pattern Size 3	1x40
Pattern Size 4	2X10
Pattern Size 5	2X20
Pattern Size 6	2X40

**Table A-1.:** Pattern Sizes

DataSet	1	2	3	4
Pattern Size 1	NaN	NaN	NaN	NaN
Pattern Size 2	0.5352	NaN	0.4747	NaN
Pattern Size 3	0.4615	0.4837	NaN	0.5513
Pattern Size 4	0.4793	0.5087	0.5108	0.5553
Pattern Size 5	0.4950	0.4878	0.4807	0.4670
Pattern Size 6	0.4716	0.4855	0.4843	0.4632

**Table A-2.:** Experimental Setup 1 (Raw data) 1 minute

DataSet	1	2	3	4
Pattern Size 1	0.4700	0.4703	0.4877	0.4583
Pattern Size 2	0.4700	0.4677	0.4877	0.4570
Pattern Size 3	0.4775	0.4677	0.4876	0.4533
Pattern Size 4	0.4535	0.5081	0.4915	0.4268
Pattern Size 5	0.4649	0.5236	0.4894	0.4504
Pattern Size 6	0.4606	0.5338	0.4691	0.4139

**Table A-3.:** Experimental Setup 1 (Raw data) 5 minutes



DataSet	1	2	3	4
Pattern Size 1	0.4831	0.5126	0.4969	0.5318
Pattern Size 2	0.4517	0.4906	0.5079	0.4562
Pattern Size 3	0.4583	0.4925	0.5586	0.4586
Pattern Size 4	0.4691	0.5135	0.5325	0.4082
Pattern Size 5	0.4989	0.4640	0.4709	0.3782
Pattern Size 6	0.5807	0.5441	0.4580	0.4296

**Table A-4.:** Experimental Setup 1 (Raw data) 10 minutes

## B. Appendix: Wavelet based approach results

### B.1. Accuracy for Wavelet transform approach using only differences, one minute time slot.

DataSet	1	2	3	4
Pattern Size 1	0.4742	0.4683	0.4863	0.4748
Pattern Size 2	0.4880	0.4659	0.4870	0.4653
Pattern Size 3	0.4576	0.4646	0.5113	0.4950
Pattern Size 4	0.4733	0.4666	0.4955	0.4925
Pattern Size 5	0.4758	0.4678	0.4930	0.4417
Pattern Size 6	0.4520	NaN	0.4768	0.4776

**Table B-1.:** Accuracy for Wavelet transform approach using only differences, one minute time slot, first iteration.

DataSet	1	2	3	4
Pattern Size 1	0.4612	0.4784	0.4958	0.4569
Pattern Size 2	0.4818	0.4821	0.5225	0.5154
Pattern Size 3	0.5758	0.3125	0.5957	0.5632
Pattern Size 4	0.5046	0.4811	0.5102	0.5300
Pattern Size 5	0.9102	0.7470	0.5777	0.5476
Pattern Size 6	NaN	NaN	NaN	NaN

**Table B-2.:** Accuracy for Wavelet transform approach using only differences, one minute time slot, second iteration.

DataSet	1	2	3	4
Pattern Size 1	0.5126	0.5048	0.5414	0.5584
Pattern Size 2	0.3333	0.8824	NaN	0.5000
Pattern Size 3	NaN	NaN	NaN	NaN
Pattern Size 4	0.4187	0.4026	0.6392	0.4162
Pattern Size 5	NaN	NaN	NaN	0.4487
Pattern Size 6	NaN	NaN	NaN	NaN

**Table B-3.:** Accuracy for Wavelet transform approach using only differences, one minute time slot, third iteration.

DataSet	1	2	3	4
Pattern Size 1	0.4444	NaN	0	NaN
Pattern Size 2	0.5000	NaN	NaN	NaN
Pattern Size 3	NaN	NaN	NaN	NaN
Pattern Size 4	NaN	NaN	NaN	NaN
Pattern Size 5	NaN	NaN	NaN	NaN
Pattern Size 6	NaN	NaN	NaN	NaN

**Table B-4.:** Accuracy for Wavelet transform approach using only differences, one minute time slot, fourth iteration.

## **B.2. Accuracy for Wavelet transform approach using only averages, one minute time slot.**

DataSet	1	2	3	4
Pattern Size 1	0.4585	0.4642	0.4841	0.4617
Pattern Size 2	0.4698	0.4635	0.4882	0.4660
Pattern Size 3	0.4721	0.4714	0.4917	0.4947
Pattern Size 4	0.5094	0.5367	0.4691	0.5065
Pattern Size 5	0.4782	0.5009	0.4929	0.4741
Pattern Size 6	0.4993	0.4989	0.4702	0.5234

**Table B-5.:** Accuracy for Wavelet transform approach using only averages, one minute time slot, first iteration.

## B.2 Accuracy for Wavelet transform approach using only averages, one minute time slot

DataSet	1	2	3	4
Pattern Size 1	0.4718	0.4818	0.5125	0.4676
Pattern Size 2	0.4708	0.4653	0.5080	0.4614
Pattern Size 3	0.5013	0.4906	0.6017	0.4767
Pattern Size 4	0.5325	0.4981	0.4772	0.5398
Pattern Size 5	0.4878	0.5333	0.3618	0.5114
Pattern Size 6	NaN	NaN	NaN	0.5027

**Table B-6.:** Accuracy for Wavelet transform approach using only averages, one minute time slot, second iteration.

DataSet	1	2	3	4
Pattern Size 1	0.4652	0.4930	0.4865	0.4693
Pattern Size 2	0.5227	NaN	0.4573	0.6120
Pattern Size 3	NaN	NaN	NaN	0.6667
Pattern Size 4	0.4300	0.4126	0.4973	0.4003
Pattern Size 5	NaN	NaN	NaN	0.2629
Pattern Size 6	NaN	NaN	NaN	NaN

**Table B-7.:** Accuracy for Wavelet transform approach using only averages, one minute time slot, third iteration.

DataSet	1	2	3	4
Pattern Size 1	0.4722	NaN	NaN	0.5198
Pattern Size 2	NaN	NaN	NaN	0.4745
Pattern Size 3	NaN	NaN	NaN	NaN
Pattern Size 4	NaN	NaN	NaN	0.3892
Pattern Size 5	NaN	NaN	NaN	NaN
Pattern Size 6	NaN	NaN	NaN	NaN

**Table B-8.:** Accuracy for Wavelet transform approach using only averages, one minute time slot, fourth iteration.

### B.3. Accuracy for Wavelet transform approach using only differences, ten minutes time slot.

DataSet	1	2	3	4
Pattern Size 1	0.4914	0.5220	0.4969	0.4775
Pattern Size 2	0.4603	0.4884	0.5362	0.4231
Pattern Size 3	NaN	0.7500	NaN	NaN
Pattern Size 4	0.4628	0.5430	0.4604	0.3223
Pattern Size 5	NaN	NaN	NaN	0.2164
Pattern Size 6	NaN	NaN	NaN	NaN

**Table B-9.:** Accuracy for Wavelet transform approach using only differences, ten minutes time slot, first iteration.

**Table B-10.:** Add caption

DataSet	1	2	3	4
Pattern Size 1	0.5429	0.4407	0.6027	0.4531
Pattern Size 2	NaN	NaN	NaN	0.6000
Pattern Size 3	NaN	NaN	NaN	NaN
Pattern Size 4	NaN	NaN	NaN	0.2269
Pattern Size 5	NaN	NaN	NaN	NaN
Pattern Size 6	NaN	NaN	NaN	NaN

**Table B-11.:** Accuracy for Wavelet transform approach using only differences, ten minutes time slot, second iteration.

## B.4. Accuracy for Wavelet transform approach using only averages, ten minutes time slot.

DataSet	1	2	3	4
Pattern Size 1	0.4839	0.5026	0.5440	0.4588
Pattern Size 2	0.4917	0.4890	0.5412	0.4475
Pattern Size 3	0.4048	0.4091	0.4600	0.4706
Pattern Size 4	0.4709	0.5590	0.5067	0.3640
Pattern Size 5	0.5859	0.6090	0.7220	0.4704
Pattern Size 6	NaN	0.8333	NaN	0.3571

**Table B-12.:** Accuracy for Wavelet transform approach using only averages, ten minutes time slot, first iteration.

DataSet	1	2	3	4
Pattern Size 1	0.5000	0.5238	0.5155	0.4143
Pattern Size 2	0.5102	0.3929	0.4872	0.3636
Pattern Size 3	NaN	NaN	NaN	NaN
Pattern Size 4	0.4464	0.5383	0.5556	0.3384
Pattern Size 5	0.5097	0.2636	0.6164	0.3048
Pattern Size 6	NaN	NaN	NaN	NaN

**Table B-13.:** Accuracy for Wavelet transform approach using only averages, ten minutes time slot, second iteration.

# References

- [1] Ahmed, M., Chai, A., Ding, X., Jiang, Y., Sun, Y. (2009). Statistical Arbitrage in High Frequency Trading Based on Limit Order Book Dynamics, 1-26.
- [2] Ahn, H.-J., Cai, J., Cheung, Y. L. (2005). Price clustering on the limit-order book: Evidence from the Stock Exchange of Hong Kong. *Journal of Financial Markets*, 8(4), 421-451.
- [3] Bates, R. G., Dempster, M. A. H., Romahi, Y. S. (2003). Evolutionary reinforcement learning in FX order book and order flow analysis. In 2003 IEEE International Conference on Computational Intelligence for Financial Engineering, 2003. Proceedings. (pp. 355-362). IEEE.
- [4] Bloomfield, R., O'Hara, M., Saar, G. (2005). The «make or take» decision in an electronic market: Evidence on the evolution of liquidity. *Journal of Financial Economics*, 75(1), 165-199.
- [5] Chen, M.; Ebert, D.; Hagen, H.; Laramee, R.S.; van Liere, R.; Ma, K.-L.; Ribarsky, W.; Scheuermann, G.; Silver, D., "Data, Information, and Knowledge in Visualization, Computer Graphics and Applications, IEEE , vol.29, no.1, pp.12,19, Jan.-Feb. 2009
- [6] Cheng, W., Liu, S., Jiao, H., Qiu, W. (2009). How Does Limit Order Book Information Affect Trading Strategy and Market Quality: Simulations of an Agent-Based Stock Market. In 2009 International Conference on Management and Service Science (pp. 1-4). IEEE.
- [7] Cont, Stoikov, and Talreja: A Stochastic Model for Order Book Dynamics. *Operations Research* Vol. 58, No. 3, May - June 2010, pp. 549 - 563 issn 0030 - 364X eissn 1526 - 5463 10 5803 0549 in.
- [8] Lopez-Monroy A. P., Gomez, M. M., Escalante, H. J., Cruz-Roa, A. and Gonzalez, F. A. Bag-of-visual-ngrams for histopathology image classification, in *Proc. of SPIE* 8922, 2013, p. 89220P.
- [9] Donoho, D. and Johnstone, I. Ideal spatial adaptation via wavelet shrinkage. *Biometrika*, 81(3):425-455, 1994.

- 
- [10] Farmer, J. Doyne and Patelli, Paolo and Zovko, Ilija I., The Predictive Power of Zero Intelligence in Financial Markets (February 9, 2004). AFA 2004 San Diego Meetings.
  - [11] Fletcher, T., Hussain, Z., Shawe-Taylor, J. (2010). Multiple Kernel Learning on the Limit Order Book. In WAPA (pp. 167-174).
  - [12] Forni, M., Lippi, M. (2001). The generalized dynamic factor model: Representation theory. *ECONOMETRIC THEORY*, 17(6), 1113-1141.
  - [13] Gabor, D. Theory of communication. *J. IEE*, 93:429-457, 1946.
  - [14] Hall, A. D., Hautsch, N. (2007). Modelling the buy and sell intensity in a limit order book market. *Journal of Financial Markets*, 10(3), 249-286.
  - [15] Harris, Zellig S. Distributional structure. *Word*, Vol 10, 1954, 146-162.
  - [16] Hsin, P.-H., Wang, M.-C. (2007). Information Indicators of Limit Order Book and Optimal Dynamic Order Submission Strategy. In *Second International Conference on Innovative Computing, Information and Control (ICICIC 2007)* (pp. 197-197). IEEE.
  - [17] Huang, H., Kercheval, A. N. (2012). A generalized birth-death stochastic model for high-frequency order book dynamics. *Quantitative Finance*, 12(4), 547-557.
  - [18] Huang, R., Polak, T. (2011). LOBSTER: Limit Order Book Reconstruction System. Available at SSRN 1977207.
  - [19] Jian Jiang, Wing Lon Ng. (2010). Capturing order book dynamics with Kalman filters.
  - [20] Jiang, G., Wang, S., Dong, H. (2011). A Survey of Limit Order Book Modeling in Continuous Auction Market. In *2011 3rd International Workshop on Intelligent Systems and Applications* (pp. 1-4). IEEE.
  - [21] Jiang, J., Ng, W. L. (2009a). Revealing Intraday Market Efficiency – Estimating Diurnal Price Densities in Limit Order Books. In *2009 International Conference on Information and Financial Engineering* (pp. 8-12). IEEE.
  - [22] Jiaqi Wang, Zhang, C. (2006). Dynamic Focus Strategies for Electronic Trade Execution in Limit Order Markets. In *The 8th IEEE International Conference on E-Commerce Technology and The 3rd IEEE International Conference on Enterprise Computing, E-Commerce, and E-Services (CEC/EEE<sup>TM</sup>06)* (pp. 26-26). IEEE.
  - [23] Kercheval, Alec N. and Zhang, Yuan. Modelling high-frequency limit order book dynamics with support vector machines, *Quantitative Finance*, volume 15, number 8, pp.1315-1329. 2015.



- [24] Kirilenko, A., Kyle, A. S. (2011). The Flash Crash : The Impact of High Frequency Trading on an Electronic Market.
- [25] Koller, D. and Friedman, N. Probabilistic Graphical Models: Principles and Techniques. edited by MIT Press. (2009).
- [26] Krishnamurthy, V., Aryan, A. (2012). Quickest detection of market shocks in agent based models of the order book. In 2012 IEEE 51st IEEE Conference on Decision and Control (CDC) (pp. 1480-1485). IEEE.
- [27] Lee, S.-Y., Poon, W.-Y., Song, X.-Y. (2007). Bayesian analysis of the factor model with finance applications. QUANTITATIVE FINANCE, 7(3), 343-356.
- [28] Lee, W.-B., Choe, H. (n.d.-a). Short-term return predictability of information in the open limit order book. Asia-Pacific Journal of Financial Studies (2007) vol. 36, number 6, pp. 963-1007.
- [29] Li, Y., Zhang, X. (2009). A Comparative Study of Information Content of Limit Order Book before and after Transparency Was Increased: Evidence from Shenzhen Stock Exchange. In 2009 International Conference on Management and Service Science (pp. 1-4). IEEE.
- [30] Lopez-Monroy A. P., Gomez, M. M., Escalante, H. J., Cruz-Roa, A. and Gonzalez, F. A. Bag-of-visual-ngrams for histopathology image classification, in Proc. of SPIE 8922, 2013, p. 89220P.
- [31] Mallat Stephane. A Wavelet Tour of Signal Processing: The Sparse Way. Elsevier. Third Edition. 2009.
- [32] Malik, Azeem and Ng, Wing Lon, (2014), Intraday liquidity patterns in limit order books, Studies in Economics and Finance, 31, issue 1, p. 46-71.
- [33] Moorhead, R.J.; Zhifan Zhu, "Signal processing aspects of scientific visualization," Signal Processing Magazine, IEEE , vol.12, no.5, pp.20,41, Sep 1995. DOI: 10.1109/79.410438
- [34] Narasimhan, Priya (Carnegie Mellon University). (2006). Fault-Tolerant Distributed Systems [Course Material]. Retrieved from <https://www.ece.cmu.edu/ece749/teams-06/team3/>.
- [35] Onorato, M., Altman, E. I. (2005). An integrated pricing model for defaultable loans and bonds. EUROPEAN JOURNAL OF OPERATIONAL RESEARCH, 163(1), 65-82.
- [36] Palguna, D., Pollak, I. (2012). Non-parametric prediction of the mid-price dynamics in a limit order book. In 2012 IEEE Statistical Signal Processing Workshop (SSP) (pp. 896-899). IEEE.

- [37] Pascual, R., Veredas, D. (2009). What pieces of limit order book information matter in explaining order choice by patient and impatient traders? *Quantitative Finance*, 9(5), 527-545.
- [38] Rajaraman, Anand and Ullman, Jeffrey David. *Mining of Massive Datasets*. Cambridge University Press. New York, NY, USA. 2011
- [39] Rinaldo, A. (2004a). Order aggressiveness in limit order book markets. *Journal of Financial Markets*, 7(1), 53-74.
- [40] Record Neil, *Currency overlay*. Wiley Finance. England. 2003
- [41] Russell, Jeffrey R. and Kim, Taejin .<sup>A</sup> *New Model for Limit Order Book Dynamics, "Volatility and Time Series Econometrics : Essays in Honor of Robert F. Engle*. Oxford ; New York: Oxford University Press, 2010.
- [42] Settlements, I. (2010). Triennial Central Bank Survey Report on global foreign exchange market activity in 2010 (pp. 1-95).
- [43] Sivic, J. and Zisserman., A. (2003). Video google: A text retrieval approach to object matching in videos. In *Proceedings of the International Conference on Computer Vision, ICCV*.
- [44] Song, N., Ching, W.-K., Siu, T.-K., Yiu, C. (2012). Optimal Submission Problem in a Limit Order Boovk with VaR Constraints. In *2012 Fifth International Joint Conference on Computational Sciences and Optimization* (pp. 266-270). IEEE.
- [45] Todd, A.; Scherer, W.; Beling, P.; Paddrik, M.; Haynes, R., «Visualizations for sense-making in financial market regulation», *Big Data (Big Data)*, 2014 IEEE International Conference on , vol., no., pp.730,735, 27-30 Oct. 2014
- [46] Vasquez Linares, Mario. Gonzalez Osorio, Fabio Augusto and Hernandez Losada, Diego Fernando. *Mining Candlesticks Patterns on Stock Series: A Fuzzy Logic Approach*. *Advanced Data Mining and Applications. Lecture Notes in Computer Science*. Springer Berlin Heidelberg. 2009. pp. 661-670.
- [47] Vvedenskaya, N., Suhov, Y., Belitsky, V. (2011). A non-linear model of limit order book dynamics. In *2011 IEEE International Symposium on Information Theory Proceedings* (pp. 1260-1262). IEEE.
- [48] Wang, M.-C., Zu, L.-P., Kuo, C.-J. (2008). The state of the electronic limit order book, order aggressiveness and price formation. *Asia-Pacific Journal of Financial Studies*, 37(2).

- 
- [49] Wang Yanhong, Liu Shancun. (2011). An empirical heterogeneous trading strategy model in the Shanghai stock market of China. In MSIE 2011 (pp. 227-230). IEEE.
  - [50] Weinberger, Kilian and Dasgupta, Anirban and Langford, John and Smola, Alex and Attenberg, Josh. Feature Hashing for Large Scale Multitask Learning. Proceedings of the 26th Annual International Conference on Machine Learning. ACM. Montreal, Quebec, Canada. 2009. pp. 1113-1120.
  - [51] Whigham, P. A., Withanawasam, R., Crack, T., Premachandra, I. M. (2010). Evolving trading strategies for a limit-order book generator. In IEEE Congress on Evolutionary Computation (pp. 1-8). IEEE.
  - [52] Yang, S., Paddrik, M., Hayes, R., Todd, A., Kirilenko, A., Beling, P., Scherer, W. (2012). Behavior based learning in identifying High Frequency Trading strategies. In 2012 IEEE Conference on Computational Intelligence for Financial Engineering Economics (CIFEr) (pp. 1-8). IEEE.
  - [53] Yu, Y. (2006). The Limit Order Book Information and the Order Submission Strategy: A Model Explanation. In 2006 International Conference on Service Systems and Service Management (Vol. 1, pp. 687-691). IEEE.
  - [54] Algorithmic Trading Challenge. (2012, January 8). Retrieved from <https://www.kaggle.com/c/AlgorithmicTradingChallenge/details/Background/>.