

Aalto University  
School of Science  
Master's Programme in Computer, Communication and Information Sciences

Ville Vainio

# Engineering analytics of big data pipelines for stock market data

Master's Thesis  
Espoo, May 1, 2019

**DRAFT! — June 30, 2019 — DRAFT!**

Supervisor:	Professor Linh Truong
Advisor:	Professor Linh Truong

Aalto University  
 School of Science

Master's Programme in Computer, Communication and  
 Information Sciences

ABSTRACT OF  
 MASTER'S THESIS

<b>Author:</b>	Ville Vainio		
<b>Title:</b>	Engineering analytics of big data pipelines for stock market data		
<b>Date:</b>	May 1, 2019	<b>Pages:</b>	42
<b>Major:</b>	Computer Science	<b>Code:</b>	SCI3042
<b>Supervisor:</b>	Professor Linh Truong		
<b>Advisor:</b>	Professor Linh Truong		
The abstract provides goal, motivation, background, and conclusions of the work. It has to fit to one page together with the bibliographical information. !FIXME <b>Add abstract.</b> FIXME!			
<b>Keywords:</b>	stock market, big data, cloud computing, stream processing		
<b>Language:</b>	English		

Aalto-yliopisto

Perustieteiden korkeakoulu

Tieto-, tietoliikenne- ja informaatiotekniikan maisteriohjelma

DIPLOMITYÖN

TIIVISTELMÄ

<b>Tekijä:</b>	Ville Vainio		
<b>Työn nimi:</b>	Osakemarkkina dataa käsittelevien big data järjestelmien tekninen analyysi		
<b>Päiväys:</b>	1. toukokuuta 2019	<b>Sivumäärä:</b>	42
<b>Pääaine:</b>	Tietotekniikka	<b>Koodi:</b>	SCI3042
<b>Valvoja:</b>	Professori Linh Truong		
<b>Ohjaaja:</b>	Professori Linh Truong		
!FIXME <b>Lisää tiivistelmä</b> FIXME!			
<b>Asiasanat:</b>	osakkeet, big data, pilvilaskenta		
<b>Kieli:</b>	Englanti		

Espoo, May 1, 2019

Ville Vainio

# Abbreviations and Acronyms

EMH	Efficient Market Hypothesis
RSI	Relative Strength Index
MACD	Moving Average Convergence Divergence
TF-IDF	Term Frequency–Inverse Document Frequency
HDFS	Hadoop Distributed File System
QoS	Quality of Service
SLA	Service Level Agreement

# Contents

<b>Abbreviations and Acronyms</b>	<b>5</b>
<b>1 Introduction</b>	<b>8</b>
1.1 Application scenario . . . . .	9
1.2 Research Questions . . . . .	11
1.3 Expected Outcome . . . . .	12
1.4 Structure of the Thesis . . . . .	13
<b>2 Stock market analysis</b>	<b>16</b>
2.1 Stock price prediction . . . . .	16
2.1.1 Statistical methods . . . . .	17
2.1.2 Machine learning methods . . . . .	18
2.2 Existing pipelines . . . . .	19
2.2.1 Data sources . . . . .	20
2.2.2 Used technologies . . . . .	21
2.2.3 Characteristics . . . . .	22
2.2.3.1 Monitoring . . . . .	22
<b>3 Big Data Technologies</b>	<b>23</b>
3.1 Data ingestion . . . . .	23
3.1.1 Apache Flume . . . . .	24
3.1.2 Apache Kafka . . . . .	25
3.1.3 Apache NiFi . . . . .	25
3.2 Data storage . . . . .	26
3.2.1 Hadoop distributed file system . . . . .	26
3.2.2 Apache HBase . . . . .	27
3.2.3 Apache Cassandra . . . . .	28
3.3 Machine learning frameworks . . . . .	28
3.3.1 Apache Spark . . . . .	29
3.3.2 Deeplearning4J . . . . .	29
3.3.3 Apache SystemML . . . . .	30

3.4	Monitoring . . . . .	30
3.4.1	Log monitoring . . . . .	31
3.4.2	Machine learning monitoring . . . . .	31
<b>4</b>	<b>Methods</b>	<b>33</b>
<b>5</b>	<b>Implementation</b>	<b>34</b>
<b>6</b>	<b>Evaluation</b>	<b>35</b>
<b>7</b>	<b>Discussion</b>	<b>36</b>
<b>8</b>	<b>Conclusions</b>	<b>37</b>

# Chapter 1

## Introduction

The modern economy revolves around stock market. Stock market is a way for companies to obtain capital which they can invest into their own business. In exchange, the person who invests into the companies stocks technically owns a piece of the company which can return profit to the investor two different ways. The stock can grow in value, which allows the investor to sell the stock in higher price or the company itself can pay dividends to investors based on the number of stocks the investor owns from the company.

The price of the stock is simply determined by the law of supply and demand. If somebody is willing to pay a higher price for the stock then the price of the stock can grow. Because of this the stock market is in continuous fluctuation where people are selling and buying the stocks with the price they think the stock is worth using stockbrokers as the middleman. [26] All of this has lead to the question, how can we invest most optimally into stocks? This is where the following computational methods come in.

There are many strategies on how to invest into these stocks which depend on multiple factors such as; how much do you expect to profit with your investment, how much are you willing to take risk, do you want to make money by selling the stocks or by receiving dividends and so on. The underlying principle with every strategy is to minimize the risk you need to take in order to gain as much as profit as possible. Some of the strategies are based on subjective evaluation of the companies, but more technical strategies use metrics that are calculated from the financial statistics or the real-time market values. Strategies using the former data are called fundamental analysis and the strategies using latter data technical analysis. Neither of these approaches can predict the future of the market, but can statistically decrease the probability of larger losses in the market for the investor although the probability of large losses is still not zero with these methods. [21]

Fundamental analysis is based on the idea that each stock has a intrinsic



value that can be larger than the actual price of the stock in the market and buying these will eventually lead to profits.[12] The fundamental analysis focuses on the financial metrics that consist of companys overall statistics. These are for example how much the company has made profit, how much the company has paid dividends and what is companys cash flow. These tell a lot about the growth of the company and how the future of the company looks like. These metrics are usually published quarterly four times a year and present more long-term statistics about the company. Because of this, the amount of data these values present is quite small in terms of space.

The technical analysis that focuses on the real-time market values, on the other hand, needs new data almost daily. Stock exchanges are usually open from morning, opening around 8 to 10am, until evening, closing around 5 to 7pm on weekdays. Before and after this there are more limited pre- and after-hours trading which lasts usually around 1 to 2 hours depending on the exchange in which more limited stock trades can be made. During these hours multiple values are recorded on the prices of the stock from which the most important ones being: the highest price the stock was sold, the highest price the stock was sold and the number of stocks traded during the time interval. The technical analysis focuses on finding recognizable patterns through this data. [35] Where the data used by the fundamental analysis was relatively small, these values can generate gigabytes of raw data in a week.

Developing a system that can do both of these analyses automatically would mean that the system should be planned from the start to be able to handle large amounts of these data as time progresses. As such task is not trivial, the goal of this thesis is to provide developers, who want to analyse this data efficiently, basic knowledge on what are the best current solutions on handling this data. With this knowledge these developers can save considerable amount of time without the need of trial and error when developing this kind of system from the ground up.

## 1.1 Application scenario

To have a better picture of how a normal system that handles stock data operates, we are next going to introduce an application scenario which is based on real-life system. From this real-life example we are going to apply the used data sources, the requirements of the possible extensions and the requirements of the results to create a more realistic scenario. This scenario will be then used throughout the thesis to give reader more practical and meaningful results.

This application scenario is presented in figure 1. In the figure, the com-

ponents marked with solid lines represent the core system functions, and the dashed lines represent add-on functionalities that should be possible to extend into the system in the future. The requirements of the system should be as follows; The system should produce long ( $r_1$ ) and short ( $r_2$ ) term predictions of stock prices. The computation of long term predictions can take time because of the nature of these predictions but the short term predictions should be between two minutes to one hour available. Long term predictions are the result of fundamental analysis ( $f_2$ ) and historical technical analysis ( $f_4$ ) whereas the short term prediction should come mostly from quick technical analysis ( $f_3$ ). The cost of the system should grow logarithmically/linear with time meaning that the cost of processing and storing data should not exponentially increase over time. Finally, the core system should be able to fulfill these requirements for at least the +5000 companies in the major U.S stock markets.

The data to the application is ingested from two main types of data sources quote and fundamental. These sources consist of values that were briefly described previously in technical and fundamental analysis respectively. The quote data is usually updated with minute intervals depending on the provider whereas the fundamental data does not change so often. Theoretically, the fundamental data can change anytime, because of dividends which can be paid whenever the companies want but this does not happen often and these values are mostly used in the long term fundamental analysis so longer update intervals are acceptable. In the figure, these are separated into the U.S market ones ( $d_1$  and  $d_2$ ) and the other sources ( $d_3$  and  $d_4$ ) that provide the same data on other global markets. The extendable global data sources are grouped into one box but in reality this data would be ingested from numberable different providers as there is no single entity at the time of writing this that provides all of this data. Theoretically the maximum size of this extendable data would be 5Gb per day which is extrapolated in from the U.S market data based on statistics that in January 2019 there were globally 51 599 companies listed in the stock markets [48]. This amount can and will fluctuate as companies enter and exit the markets but it gives us the scale of data we are working with.

Today, stock market analysis has also a large focus on predicting stock prices using secondary data sources that can have reflect and affect the prices of stocks. These secondary data sources can be anything but at the moment one of the most researched sources are traditional media and social media data. Examples of using this kind of data to predict predict stocks can be found in [51], [43] and [34]. This is why the system should have the ability to extend to ingest data from arbitrary secondary sources ( $f_6$  and  $d_5$  in the figure) to provide more versatile predictions about the stocks. The amount

of this data can be unlimited but is restricted to relevant sources.

Data is ingested from these data sources mainly using HTTP-protocol as this is the main method that these services ( $d_1 - d_4$ ) provide. Other possible methods that are usually available are Excel sheets and sometimes websockets, of which the websockets can be actually useful in cloud system, but as the HTTP-methods are currently the most used technology, this thesis is also going to focus on these. Here we have separated the main ingestion functions into two main types of functions  $f_1$  and  $f_2$ .  $f_1$  is constantly polling and processing data whereas  $f_2$  handles batch processing. Both of these function store their raw output into the storage, but  $f_1$  passes this also to the immediate technical analysis.

The system has two technical analysis functions. For methods that allow streaming updates there is  $f_3$ , which can for example be cumulative/reinforced ML models and for methods that need historical data in order to calculate the prediction, there is  $f_4$ . For fundamental analysis, there is no a specific function as the introduced data sources usually provide these values pre-calculated and these values are usually easy to calculate dynamically with little to none amount of processing.

Finally, at the center of the system is the storage which is used to store the calculated predictions as well as the raw data from the data sources for later analysis. For historical technical analysis,  $f_4$ , the storage should provide reasonable range query times when quering historical data and for the results the storage should provide efficient point queries for the results ( $< 1s$ ).

## 1.2 Research Questions

The focus of this thesis is going to be the ingestion of stock data ( $f_1$  and  $f_2$ ) using big data technologies. As the field of possible technologies is large, the main focus of this thesis would be the comparison of current relevant technologies in the context of stock data to help developers to decide what technologies to use in their own projects. The main result of this thesis would be information and possible tools that developers could use to develop this kind of system more quickly. Developers here can be people from companies that either do stock analysis as their main business or just want to analyze stock data efficiently. These results could also be used in research to implement analyzing pipelines more efficiently so that the research group can focus on the analyzing of the stock data instead of worrying with getting the data to these parts.

Analyzing the stock data is also what most of the research today is focused on as this is the part that can actually produce profit. This means that most

papers ignore the steps of ingesting and storing this data. Examples of this kind of papers are [49], [28] and [51]. So one of the goals of this thesis would be to bring also some practical knowledge on subject that has not been researched that much before.

This thesis will approach this subject from three different perspectives that cumulate on one another. These perspectives are represented by the following three research questions:

What are the needs and requirements of this kind of system data-wise? Which is important knowledge when starting to design or implement any kind of big data system. This thesis plans to provide the information on these so that when a reader is developing their own system they have a point of reference which they can use to evaluate their data sources.

What are technological options to implement this in practice? As there are enormous amount of different technological frameworks, this thesis plans to provide the reader information on the most relevant ones currently. This way the reader does not need to go through and learn variety of technologies in order to build their system.

How do the most viable big data ingestion options compare to one another in the context of stock data performance- and development-wise? Finally the thesis plans to provide comparison and prototype implementations of stock data ingestion using the technologies that seem most prominent. The metrics used to compare these systems would be the response time, possible scalability and the time it takes to develop these systems. This is to save possible time of a person developing these when the trial and error is done beforehand and if the prototype system fits the developed architecture it can also be used as a basis for further development by the reader.

### 1.3 Expected Outcome

For the first research question "What are the needs and requirements of this kind of system data-wise?" this thesis plans to provide an analysis of the necessary stock market data and its usages. From this analysis, the thesis would derive the main requirements for the system to fulfill in order to satisfy the needs of the possible subsequent analysis stage. This result could then be used in the future if one would want to build their own ingestion system from the ground up as a base.

For the second question "What are technological options to implement this in practice?" the thesis would perform an analysis on the current trends in data ingestion solutions. The thesis would provide information on the latest open-source technologies that could be used to implement this kind

of system, how these would fulfill the requirements introduced by the first research question with application scenario and conclude this with a comparison of these technologies on what are the advantages of using one over another. The result of this part could be used to decide what seems to be the most suitable technology to use to implement the application scenario technically.

For the last question "How do the most viable big data ingestion options compare to one another in the context of stock data performance-wise?" the thesis would implement open-source prototype solutions based on the results of the second research question. This prototype could be used by anybody (company or individual) as it is or as a base to build a more complex system on top of it. The system would be targeted to practically implement the described application scenario, which could be used by companies with similar products allowing them to write more complex analysis algorithms based on the larger amount of data.

## 1.4 Structure of the Thesis

In the first chapter of the thesis we will be focusing on solving the first research question "What are the needs and requirements of this kind of system data-wise?". The thesis would start by going through scientific papers about stock markets and stock analysis. There would be first text generally about the characteristics of the stock markets, what are they based on, what affects them, can they actually be predicted (random walk hypothesis) and so on. Then the thesis would go through the main directions of analysing the stock markets (fundamental analysis, technical analysis) and explain briefly some of the methods (about three methods per direction) that characterize these directions (Gordon model, Magic formula, LSTM etc.) focusing on the data that these methods need in order to calculate their predictions. This part would be concluded by deriving the requirements for the system based partly on these analysis methods and their data needs, and partly on the definitions that make up a scalable, secure and stable cloud system. The next step would be solving the second research questions "What are technological options to implement this in practice?". This step would perform literary research on what is currently used to perform big-data ingestion and storing, selecting from the list of technologies mostly those that seem to fulfill the requirements derived in the step 1 and would fit in the application scenario introduced in figure 1, specifically in  $f_1$  and  $f_2$ . For data ingestion, these technologies would probably be Apache NIFI, Apache Flume, Fluentd etc. For data storage, these could be HDFS, HDFS, Apache HBase, Apache Cassandra etc.

This step would consist of first introducing all of the selected technologies and going through how do they work, what are they supposed to solve and what are the advantages and disadvantages of using one. After this, the section would do a comparison of these technologies in the context of stock data and conclude with analysis on which of the technologies would be the most prominent ones to solve this problem. After this would start the experimental part of the thesis and the rest of thesis would be focusing on solving the final research question "How do the most viable big data ingestion options compare to one another in the context of stock data performance-wise?". Based on the results of the step 2, I would implement couple of the most prominent solutions as a prototype systems. This step would include the actual implementations and reporting of these implementations. The report would consist of technical details; what parts does each of the system consists of, what versions were used, where the system was run etc. And subjective remarks; was it easy to implement, was there parts that did not fit together etc. The reporting part would also explain the metrics that would be measured for the subsequent analysis. For these metrics, the dataset used to test the system would be the open-source REST API from IEX [9], which is open for consumers until first of June, 2019. These metrics could be, for example, the time it takes to batch process, processor usage, memory usage, database fetching times etc. As these systems would be run inside Docker containers, the tools used to measure these metrics would most probably be programming language specific functions and Docker specific statistics tools. In the final step, the thesis would first inspect the results from the step 3 and based on these make remarks on what could be the best potentially the best implementation in this context. After this there would be an wrap-up on each of the previous sections concluding in retrospective what could've been possibly done better and what could be done in the future, concluding in recommendation what could be based on this thesis the best technical solution to implement system described in the application scenario.

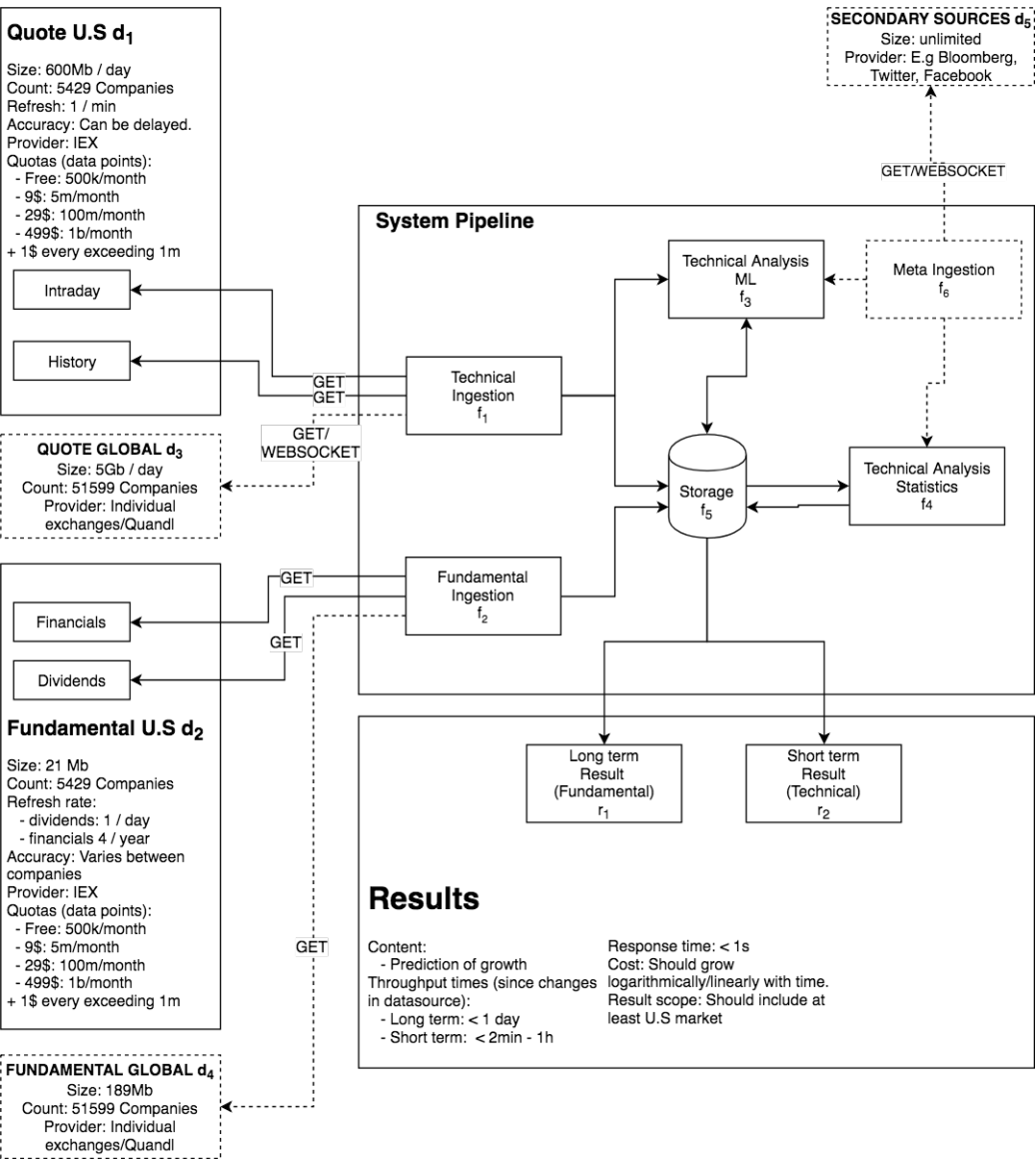


Figure 1.1: Application scenario for the thesis

## Chapter 2

# Stock market analysis

In order to build practical stock analysis systems, we start by looking at the methods of stock data analysis to understand the underlying domain. Stock market analysis is an enormous domain by itself and there are multiple ways to approach this data. For example possible problems to solve are finding anomalies in the data, predicting the future price and analysing possible causalities. In this thesis we will be approaching this from the price prediction perspective which can also be used in other problems such as anomaly detection [25].

In this chapter we examine what are the state of art methods that researchers use to analyse stock market data and we are especially focusing on methods that in some way use big data for this analysis. We will be also looking at the implementation side of things and examine some of the existing big data pipelines for stock market analysis and what are the technologies used to build these. As we will see, big data can be part of the pipeline in the form of complete stock data in which case the main method of analysing this is technical analysis. Or the big data can aspect can come from other sources of information such as news and social media in which case the main method of analysis is fundamental analysis.

## 2.1 Stock price prediction

The current methods on stock price prediction can be divided roughly into two different categories; statistical methods and machine learning methods. Both of these categories can be further divided into fundamental and technical analysis categories depending on what data they use as a source of analysis, but these categories overlap a little because of new ensembled models. With statistical methods we mean here the traditional mathematical models



that need quite a lot of understanding of the domain in order to derive them. With machine learning methods, we mean algorithms and statistical models that derive underlying principles from data using patterns and inference.

Both of these methods are trying to beat the efficient market hypothesis (EMH). EMH states that all the current public information available should already be seen in the price of the stock. In other words, the only thing that can affect the price of the stock would be unknown new information and the randomness of the system, which leads to a claim that stock prices can not be predicted using historical data. However, there have been multiple studies shown that this hypothesis could possibly be beaten using big data. [36] This hypothesis has also been challenged with overreaction hypothesis that states that the market overreacts to new information making it possible to somewhat predict the market before the prices change [19].

### 2.1.1 Statistical methods

The driving force of stock analysis in the past have been the statistical methods and there are still a lot of new papers analysing stocks using these methods. There are hundreds of ways to approach this problem from statistician point of view so here we have listed only some of those that have some relevance to the subject of this thesis.

From the technical analysis perspective, where only the stock data is analysed, we are going to be looking at momentum indicators. [45] Momentum indicators, as the implies, measure the speed of change in the stock prices and investors make decisions based on the thresholds of these values whenever a stock is being overbought or oversold. [14] When searching research on big data usage in stock analysis there are some key indicators that show up frequently. These are relative strength index (RSI), moving average convergence divergence MACD and Williams %R. These indicators can be used as they are but these have been also used as features that machine learning algorithms use to predict prices. [42]

All of these indicators measure the momentum of the price, but they all do it differently. RSI and Williams %R are both called oscillators because their values oscillate between maximum and minimum values they can get. Where as MACD compares long-term and short-term trends to predict if there is currently a notable trend. Because of the differences in the formulas and the factors that these values are measuring, the results from these formulas can be conflicting. [14]

Autoregressive integrated moving average (ARIMA) models have also been successful way of analysing time series stock data and they have been also used with big data. [47] However, with the new advances in machine

learning, there seems to better performance to be achieved with machine learning methods. [27] Due to this and the complexity of ARIMA models, we are going to only briefly make a note of them here and focus more on the machine learning methods in the next section.

In the fundamental analysis side, most of the recent developments have happened with textual data from news and social media using machine learning methods meaning that traditional statistical methods have not gained that much interest recently. However, as there is a lot of relevant financial big data that can be used with traditional methods, here are couple of recent examples of the usage of this kind data: Day et al. [19] tried to link oil prices to stock prices using global financial data streams with stochastic oscillator techniques. Kyo [30] combined technical and fundamental analysis by using regressive model that takes into account business cycles in Japan's stock market.

### 2.1.2 Machine learning methods

As stated before, machine learning is the current trending stock analysis methodology that has gained a lot of interest. Both supervised and unsupervised learning techniques has been tried to predict the prices but recently neural networks especially have gained a lot of interest due to their adaptability to any non-linear data. Neural networks have the advantage that they usually do not need any prior knowledge about the problem and this is the case when we are looking at seemingly random stock market data.

Let us start by looking at neural networks used to predict the market using only the market data (technical analysis data). There are multiple different neural network architectures to choose from and the most popular one currently seems to be a basic multilayered feed forward network (MFF) when analysing only the stock market data. There is some results that have shown that increasing the size of MFF, by adding layers and neurons, can produce better results. This however, increases the amount of needed computation and will probably have some upper bound before it starts to overfit the training data. [41]

Although MFF is seemingly the most popular, better results have been able to achieve with convolutional neural networks (CNN) and long short-term memory (LSTM). With convolutional neural networks, the upper hand to regular MFF seems to be the ability to express same amount of complexity with smaller amount neurons, although no direct comparisons have been made with these two. With LSTM however, it has been directly shown that it performs better than other memory-free machine learning methods including MFF. As stock data is literally a time series, the LSTM's ability to remember

previous states seems to separate it with other methods. There currently seems to be no papers on how the LSTM performs compared to CNN so we can only assume that the performance of these two is pretty close when it comes to the complexity of the network. After these standalone architectures the next step to seem to be hybrid architectures combining more than one model into one and this has already achieved seemingly better results than these standalone models. [41]

In order to train a neural network to predict stock markets we need features and classes to label these features correctly. Because there is such a huge amount of data in stock transactions, manual labeling is not an option. The simplest way automatically generate features is to take calculate change in price in each time step. This way the network can for example output simple binary classification telling whether the price is changing more or less than median price. [20] When we take this to a bit more complex level, similar features can be made from RSI, MACD and Williams %R values which we introduced in the previous section. [42]

When we move on to the fundamental analysis, similar kind of networks are used to with financial news and social media data. Again LSTM and ensemble models are used to connect this data from outside with the price trend of stock market. There are research using just general social media data that concerns the whole market but there are also usages of stock specific posts. Fundamental data that is in textual form, term frequency-inverse document frequency (TF-IDF) is a common choice to present features [32]. This data is then classified based on the general sentiment that they seem to represent. Positive sentiment toward company would lead to the rise of stock price and negative sentiment vice versa.

## 2.2 Existing pipelines

Now that we have some kind of understanding what the current stock analysis is, let's move on to examine the actual implemetations that do this analysis. For this section, we examined 19 different stock data analysis pipeline that handle or are able to handle big data in some sense. What we mean here by "in some sense" is that the big data might not have been the actual stock data but for example social media data that was used to analyse the actual stock data. Information about these pipelines were all publicly available, although most of the pipelines were not open-sourced. Fourteen of these pipelines were reported in academic publications and the rest five are, or at least have been, in industrial use. There were also multiple open-source pipelines that have been developed for big data stock analysis, but information about the actual

usage of these pipelines were not found. This is why we excluded these from this study in order to get results that might have more relevance to pipelines that are actually in use.

We have divided this section into subsections based on different parts from the pipeline starting from the furthest away of the possible clients and moving our way up towards the analysis phase. The biggest problem here is that the companies that work in the forefront of stock analysis, seldom share their software architectures for business reasons. Nevertheless, with the public information available we can get an excellent view of the state of art pipelines in academic world and a general idea how the pipelines are build in the industry side.

It is also not that easy to compare technologies used in academia and industry. Both have different needs and goals that they wish to achieve. In academia, the pipeline is usually made in ad-hoc manner trying to minimize the time used for development and maximise the time needed for testing. This is possible, because researchers do not have same client side restrictions that industry might have. In industry side, it is usually important that the pipeline stays operational during the whole live cycle of the system. Where researchers only need these pipelines in order to conduct a couple of experiments, these industry pipelines have to survive possibly multiple of years of usage. These industry pipelines also serve the results to multiple clients that all expect small latencies from the system in order to use the system efficiently whereas in academia, this is not a requirement but helps the research to be made efficiently. So bearing these domain-specific requirements in mind, different technologies can be used to achieve optimal solution in different situations.

### 2.2.1 Data sources

Stock market data is not cheap. This is mostly because the exchanges that run the stock markets, usually make most of their profit with trade data and thus do not want to give this data freely away. There are however services who do provide part of this data with a much lower cost available to companies and researches which cannot possibly afford the complete real-time data. This partial data usually consist of historical end of day data, which is the aggregated statistics of the stocks after the market has closed for the day. Although this data is complete in the sense that it tells all the necessary information about stocks price evolution on day level, we will be referring this data as the aggregated data during this thesis and reserve the term complete stock market data for the data that contains all of the transactions. What this means for the systems is that the stock data can

exponentially smaller at the beginning of a financial companys life when they possibly have the ability to access the complete stock market data, but system must be able to scale to this complete data set size.

In academic papers, the Yahoo Finance API has been the de facto service used as the source of this partial stock data. [25] [18] [31] [42] Unfortunately, although this service have been used in papers published in 2019, this API was shutdown by Yahoo in 2017. Today, there are no service that provides the same amount of data that this API did, but some substituting services do exist.[33]

### 2.2.2 Used technologies

Next let's look at the actual technologies and frameworks that are used to implement this pipeline by starting from the furthest away from the client; ingestion layer. What we mean by ingestion here is the fetching and preprocessing of the data. This step with stock market data varies a lot depending who is fetching the data and for what purpose. Most common method for stock market data ingestion is just using custom scripts. This is usually enough with the aggregated stock market data but it does not scale.

Common big data technologies usually come in to play when the system has to ingest for example huge amounts of textual data such as news and social media feeds. For these purposes, the most common technology in scientific papers is Apache Flume which for example used in [39] and [17]. Another framework that is commonly reported in the ingestion step is the Apache Kafka streaming platform. [32] [16] There are also indications of usage of Google Dataflow in industry side. [38]

As we have now this ingested data, we need to store it somewhere. Stock market data is quite structured but the data usually used to enrich this data is mostly unstructured. This puts a restriction on what are the possible storage options available. Academic papers usually do not have to worry about this as the systems just have support ad hoc calculations but in the company context this becomes a crucial part of the system as it is the component that enables low latency responses without having to fetch data all the way from the original data source. With big data systems there is usually two options which are either cloud platform specific products or open-source HDFS-based systems and this is also the case with stock analysis systems. From the cloud platform products, there are information on usage of Amazons S3 and Googles BigTable with BigQuery. [44] [38] In the open-source side HBase [23] and Cassandra are the two used database solutions.

Then as we finally have the data in control, we can apply some analysis to it. In the analysis step, there is not that much variety in used technolo-

gies. Apache Spark with its MLlib library is dominating this field with its capabilities to process data efficiently. [25] [15] [32] [18] Where Apache Hive and Apache Pig were previously most used technologies, now Spark seems to be taking their place [44].

In the industry side, there is indications of Apache Storm being used with real-time classification [16]. The strong point of Spark is that same models can be used with Spark Streaming framework, but better latencies can be achieved with Storm making it possibly better choice for companies which have resources to implement two separate systems [29].

## 2.2.3 Characteristics

### 2.2.3.1 Monitoring

In most of the pipelines, the system is produced so that it can scale with the input allowing no real performance issues to rise. None of the cases however, report how they monitor the system while it is running. In cases where for example Spark is used, we can assume that the process is monitored using Spark's own monitoring tools which measure load and resource usages. These tools are valuable to measure the performance of the application, but they do not always tell the whole truth about the applications state.

Let's look at an example case, where we have machine learning model training pipeline implemented in Spark which reads stock data from the database and trains a model based on this. What would happen if this data corrupted on the way to the training algorithm. In some cases the corruption would make the training algorithm possibly crash which would be noticed with performance monitoring tools. However, in some cases the corruption could only make the value of the data change ever so slightly that it would pass as a input to the algorithm. In this case, this could affect the final trained model making the model performance decrease tremendously. Which could then lead to even false conclusions on the model.

So in order to make right decisions concerning models and save time while training, it is crucial to identify these kinds of problems early as possible. This is where good monitoring comes into play. In the next chapters, we are first going to take a look at different monitoring solutions and then implement these on the most common stock data machine learning pipeline in order to solve this problem.

## Chapter 3

# Big Data Technologies

In this chapter, we will examine more carefully the technologies that can be used to implement a stock data pipeline in the big data context and try to answer the second research question. We have gathered here technologies that we saw used in existing big data stock pipelines in the previous chapter. We also added couple of promising frameworks that could be used in the pipelines but there is no public information about companies using them yet. To keep this section compact and more useful for majority of the developers that do not possibly have access to company level resources, we have included only the open-source solutions here leaving outside the services that cloud providers offer such as Amazon S3 for storage or Google Dataflow for ingestion.

### 3.1 Data ingestion

Here we are going to be using the same definition for ingestion as before; data fetching and preprocessing which includes data validation and feature extraction. This is one of the crucial parts of the pipeline as it is responsible of turning arbitrary data into facts that the rest of the pipeline can use and rely on. This also makes it the hardest part to develop as the developer must understand the data well enough that these parts can work efficiently and prevent erroneous states in the later stages of the pipeline.

We have gathered here three main technologies that are usually mentioned when academics are talking about data ingestion, but in reality all of these frameworks have a bit different tasks they try to fulfill. This makes comparing these technologies harder and it usually just means that the better technology here is usually the one that is better suited for the current problem, which does not make the other ones any worse than the selected one. On top

of Apache Flume and Apache Kafka which were already used in existing pipelines, we have also included here Apache NiFi which promising features that could also be used for this task in this kind of system.

### 3.1.1 Apache Flume

Apache Flume is one of the Apache Software Foundation (ASF) projects that we will be seeing a lot especially in the context of ingestion. From the projects official website we have definition that Flume is "distributed, reliable, and available service for efficiently collecting, aggregating, and moving large amounts of log data" [1]. So the main focus of Flume is to manage log data, but as we have already seen, this is not the only use case for this tool.

Flume was first introduced by Cloudera in 2011. In same year, it was officially moved under ASF and came out of the incubation phase the following year. During this incubation, developers had already started to refactor the Flume and the result of this is the 1.X lineage of Flume which is still going to this day. [24] At the time of writing this the latest production-ready version is 1.9.0.

Flume's high level architecture can be divided into 3 different parts: sources, channels and sinks which all run inside an agent which is an abstraction for one Flume process. Data is inputted to the system from the sources, then it goes through channels and is finally written into sinks. While going through channels the data can be processed using functions that Flume calls Interceptors. [24]

The main unit used of data in Flume is called Event which is structured like most of the other message formats you can find in network protocols. Event has a header, which is key-value pairs of meta data and a body which usually is the actual data. Overall, Flume architecture is message driven which allows it efficiently to multiplex data across multiple computing instances levitating the work load between these machines. [24]

When processing large amounts of data, it is important to avoid bottlenecks that can form in the pipeline and this is why knowing the amount of data flowing through Flume is extremely important. Bottlenecks that can form in Flume are for example cases where data is coming from sources faster than the Flume is able to write to sinks. These kind of situations can be avoided with the knowledge of the data domain when configuring Flume instances but also by monitoring Flume processes in order to respond to these kind of situations. [24]



### 3.1.2 Apache Kafka

Apache Kafka describes itself to be a "distributed publish-subscribe messaging system" and it can be used for three different purposes: as a messaging system, as a storage system and as a stream processor. [2] Because we are examining Kafka in the context of data ingestion, we are mostly interested in its messaging and stream processing capabilities, but it's good to keep in mind that it is possible to store data with Kafka for a longer time if necessary.

Before being a ASF project, it was first developed by LinkedIn to gather user activity data in 2010. [40] After being released from incubation in 2012 many other big industrial companies such as eBay and Uber [52] have taken kafka into use to manage their own enormous data systems. Today, Kafka is already in its version 2 and can be integrated with almost any modern big data framework. [40]

Kafka operates on publish-subscribe architecture where producers input data into the system by producing data and consumers subscribe to the data which they want to consume/receive. To make this work with big data, Kafka has a core abstraction called topic which has multiple immutable queues called partitions. When producers produce new data, this data is appended to a partition queue where where Kafka keeps track which items consumers have already consumed by tracking the offset of a consumer in a queue. With this kind of architecture Kafka makes promises that the data retains its order throughout the system and does not arrive to consumers out of order.

As for the scalability of this kind of system, a single partition must fit onto a single server, but topic which has multiple partitions does not have this limitation and this allows topics to scale over the Kafka cluster. Fault-tolerancy can be achieved by just replicating these partitions over the cluster. Multiple consumers can form consumer groups that consume some topic simultaneously from multiple partitions meaning that in addition of Kafka being itself scalable cluster, it allows its consumers to be also scalable cluster and without any additional complexity. [2]

### 3.1.3 Apache NiFi

Apache NiFi was made to dataflow management and its design is inspired by flow-based programming. It was originally developed by National Security Agency (NSA) as a system called "Niagara Files" and was moved under ASF in 2014 making it the newest addition under ASF out of these three introduced ingestion technologies. [13]

As we have seen already with the last two ingestion frameworks, all these frameworks have their own abstractions for data and the processes that han-

dle this data. NiFi's core abstraction is FlowFile which represents the data that flows through the system. These are processed with FlowFile Processors that are connected to each other by Connections and these can be grouped into Process Groups. These processors and their connections are then managed by a Flow Controller which acts as the brain of each node in a NiFi cluster. [3]

NiFi cluster follows zero-master clustering paradigm meaning that the cluster does not have clear master nodes and vice versa no slave nodes. Every node in NiFi cluster processes data the same way and the data is divided and distributed to as many nodes as needed. Apache ZooKeeper is used to handle failure of nodes in the cluster.[3]

## 3.2 Data storage

Data storage forms the center of the pipeline and is one of the major places where bottlenecks can form. Possible bottlenecks are the query latency and the writing speed which can slow the whole application down if not implemented efficiently enough. If this non-trivial task was not enough, the data storage must also be able to resist

### 3.2.1 Hadoop distributed file system

Hadoop distributed file system (HDFS), which is the core part of Apache Hadoop ecosystem, is one of the current defacto ways to store big data. It has gained a lot of popularity with the map-reduce programming paradigm. HDFS build so that it does not have to be run on high quality hardware, normal commodity hardware is more than enough to run the system.[8]

HDFS offers all the same functionalities that a traditional file system would offer. Clients can create, edit and delete files which can be stored into directories that form hierarchies. What makes HDFS differ from a normal file system is its ability to handle a lot larger data sets and at the same time have a better fault-tolerance than a traditional file system.[8]

HDFS is based on a master-slave architecture where the master is called NameNode and slaves are called DataNode. The NameNode stores and manages the state of the whole system and the actual client data never flows through this node. The data is stored into the DataNodes which manage this data based on the instructions that the NameNode gives them.[8]

For reliability, data is replicated between DataNodes so that the outage of individual DataNodes does not affect the overall performance of the system. In order to identify and react to these kind of failures, each DataNode

send heartbeat-like messages to the NameNode, which then conducts needed actions if a heartbeat is missing. The NameNode itself, on the other hand, is a single point of failure. It does record the changes and the state of the system into files called EditLog and FsImage which can be used to replay the changes in the system if NameNode goes down temporarily and because of these files are crucial to get the system back up, usually multiple copies of these are stored into disk. [8]

### 3.2.2 Apache HBase

HDFS by itself is only made to store very large files, so the methods it provides are quite limited. This is where the Apache HBase comes in which is implemented based on Google's BigTable framework. Where BigTable uses Google's own file system, GFS, HBase is built on top of HDFS. HBase is a NoSQL database although it does not natively have many features that a normal database would have. Notably, it does not have its own advanced query language.[22]

HBase's record structure is somewhat similar to its relational counterparts. Data is stored into rows that consist of columns, that are identified by a unique row key. Rows form tables and tables can be further grouped into namespaces. The difference to typical relational data model comes after this. Columns may have multiple versions and timestamp information about the column and its version are stored into separate entity called cell. The columns do not form table like structures with rows, and instead act like key-value pairs which can be grouped into column families. This way values that would be for example 'null' in normal relational database, do not take any room in HBase as such key-value pairs can be omitted. [22]

HBase's core abstraction of scalability is called Region. Region is a set of continuous rows that are split when one Region grows too large. Each region is served by a single RegionServer and each RegionServer can serve multiple Regions. So Hbase cluster is scaled by adding these RegionServers which can serve more Regions to clients. [22]

HBase does not instantly store changes into disk. Changes are first recorded into a log called write-ahead log (WAL) and after this the change is stored into a memstore which is in memory. After the changes expire in the memory, the changes in memory are flushed into the disk as they are, into a file called HFile. In order to avoid large amounts of small HFiles in the disk, HBase periodically merges these files using a process called compaction. These are done in file scale (minor compaction), but also in larger scale where all the files inside one region are merged into one and data market for deletion can be cleaned in the process (major compaction). [22]

### 3.2.3 Apache Cassandra

Apache Cassandra is a NoSQL database that is build on peer to peer architecture. Cassandra provides its users same tools that a normal relational database would. Its record structure is the same with rows, columns and tables such as with typical relational database and it provides a SQL-like query language. What differs Cassandra from a normal relational database is its scalable architecture which we are going to look next. [50]

Cassandra's main scalable units are called Nodes that are a single instance of Cassandra running in a single machine. These Nodes are grouped together based on the data they serve and these form Rings. Data distributed and replicated between the nodes based on the hash of data's partition key. With consistent hashing algorithm, every Node has the same amount of data. [37]

Unlike HBase, Cassandra does not guarantee the consistency of data due to CAP theorem but instead guarantees the availability at all times. Cassandra implements "tunably consistent" paradigm where user can define for each write/read to be consistent with the cost of availability. This high availability, makes Cassandra better choise for example web application where the data must be available at all times, but the data doesn't necessary have to be up to date. [37]

## 3.3 Machine learning frameworks

When it comes machine learning in big data context, libraries such as Tensorflow and PyTorch are still quite used because of their support of processing data with GPU. This allows them to perform really well in single machine instances while needing minimal time to develop. But when these are run in big data context, development becomes a bit more harder as this is not the environment these are designed to run at.

As we saw in the previous chapter, the Apache Spark framework with its Spark ML library is currently the most used big data machine learning platform. Spark, however, natively supports only classical machine learning models and currently has no ready-made deep learning solutions, which are the state of art methods we are interested in. We want to still be able to do the deep learning training on Spark environment because of its good integrations with the frameworks we already introduced and its mature ecosystem. That is why we are next going to briefly look at how the Spark ecosystem works and then move on to see tools which we can use to train deep learning models in Spark ecosystem without writing implementations from scratch.

### 3.3.1 Apache Spark

Apache Spark is a distributed in-memory data processing system that provides tools for scalable data processing. Spark has a master-slave architecture, where a driver node acts as a master and executes Spark program through Spark Context. The driver node passes tasks to worker nodes that the Spark Context together with Cluster Manager manages. The most popular cluster manager currently seems to be YARN but Mesos or Spark's own standalone could also be used for this job. Each worker node then has its own executor process which runs the given tasks in multiple threads when necessary.[4]

Until version 2.0, Spark's main programming interface was a data structure called Resilient Distributed Dataset (RDD). RDD is a collection of elements that can be partitioned throughout the cluster allowing them to be processed in parallel across multiple servers. RDD's are still in use for backward compatibility reasons, but from version 2.0 onward Spark's main abstraction has become structure called DataSet. DataSet is similar to RDD in high level, but it provides richer optimizations under the hood and higher level methods to transform the data. With the DataSets, a new abstraction called DataFrame was introduced which is can be compared to a table in relational database. DataFrame is a DataSet of Rows that can be manipulated with similar methods that you could do for DataSet.[4]

Spark has divided its functionalities into sub-modules that are specified into specific tasks such as Spark Streaming for stream processing and SQL for processing data in tabular form. We are mostly interested here in Spark ML module. Spark also has a module called Spark MLlib which some of the research still use. The main difference between Spark ML and MLlib is that Spark ML provides newer DataFrame API whereas MLlib uses older RDD datastructures.[11]

Spark ML provides whole set of classical machine learning algorithms such as linear regression, naive bayesian and random forests. It also has a concept of pipelines which consist of different transformers and estimators, which are made to make model developing easier. However, the main problem with Spark ML is that it does not have native tools to implement deep learning models on Spark.

### 3.3.2 Deeplearning4J

Deeplearning4J (DL4J) is distributed framework of deep learning algorithms that work on top Spark and Hadoop ecosystems. It provides all the most popular deep learning models such as multilayered networks, convolutional

networks, recurrent networks. DL4J also provides a lot of tools for preprocessing the data before fed for training in the form of sub-projects. [6]

In academia, DL4J is not that used as its main language is Java, but this choice of language makes it really suitable for industrial use. However, DL4J supports Keras model import which allows user to use other languages than the JVM based alternatives. This makes python based prototype systems be able to run in industrial Spark cluster. [6]

### 3.3.3 Apache SystemML

Apache SystemML is a bit different machine learning system when compared to Spark ML and DL4J. Similarly to DL4J, SystemML also runs top of Spark, but unlike DL4J it is not a straight forward programming library. SystemML has its own R- and Python-like declarative machine learning languages (DML), which can be used to define machine learning algorithms that run on the Spark. [5]

Currently, these languages cover about the same use-cases that Spark ML does. What makes SystemML differ from Spark ML is that deep learning models developed in Keras or Caffe can be converted into DML. So theoretically SystemML supports deep learning through these libraries although its core methods do not have these methods. This allows it to be run on classical command-line interface, but also from jupyter notebooks which are currently vastly in use in academia. [5]

## 3.4 Monitoring

We have already seen quite a bit of framework specific monitoring solutions, which cover monitoring individual components in the pipeline. Next we are going to take this monitoring a step further and look at monitoring solutions that work outside of these components and can be used to monitor the global state of the pipeline. This can make monitoring easier when all the information is available in one place and it also helps to infer cause-effect relations.

We have gathered here two different monitoring solutions for two different needs. The ELK stack for monitoring the data flow inside the system and individual components and the ModelDB to specifically monitor the Machine learning models that are produced by the pipeline.

### 3.4.1 Log monitoring

The most common solution for monitoring logs in bigger system is the ELK stack. ELK stack which is an acronym for Elasticsearch, Logstash and Kibana stack, is a common stack used to implement collection of logs and their visualization in Big Data environment. The need for such as elaborate stack for just collecting logs, comes from the fact that when you have a big data system, the logs of such a system form a big data problem of their own, so in order to solve this problem this stack was developed. It provides scalable data ingestion system, with a distributed storage that can be accessed and visualized in efficient manner. [7]

In this stack, Logstash, an open-source data ingestion pipeline tool, is used to ingest the log data. This ingested data is the stored into Elasticsearch which is a distributed search and analytics engine, which provides REST interface. Finally, the data stored into Elasticsearch can be visualized and monitored with Kibana, which is a tool developed to do just this. [7]

Because of different user needs the stack has evolved into stack that the company behind these technologies calls Elastic Stack. This stack really only differs from ELK-stack by a component called Beats, which can be used to build more lightweight stack for simpler needs. However, pure ELK stack is still very popular option to handle log data. [7]

### 3.4.2 Machine learning monitoring

Open-source monitoring tools that are specifically designed for machine learning are not that common, but couple of tools do exist. These tools do not only help monitor the models performance, but also provide tools for deploying and versioning these models just like developer would use git to manage their code base.

ModelDB is a system developed in Massachusetts Institute of Technology (MIT) for ML model management. It provides tools that can be used to monitor the performance of models and compare these results with each other. It also allows logical versioning of models and helps to reproduce the models this way. ModelDB, however, does not support models from many different ML libraries and currently the only supported libraries are Spark ML and scikit-learn. The library is said to have second version coming up, but at the time of writing this the library has not had major update for an year. [46]

Another, newer option is an open-source project called MLFlow. It does all the same things that the ModelDB does, but markets itself as a more end-to-end solution. On top of tools for the managing explained in ModelDB

paragraph, MLFlow puts more emphasis on packaging the models and the deployment of these models into actual use. [10]

Unlike ModelDB, MLFlow does not care about what is the library that generates models. It does this by providing CLI and REST API's that can be integrated with the system ignoring the underlying technologies. For convenience, it also provides APIs for Python, R and Java which can be used for tighter integration. MLFlow is as project quite young having its first full version released in June 2019, but it has a healthy development community and is actively developed. [10]



## Chapter 4

# Methods

## Chapter 5

# Implementation

## Chapter 6

# Evaluation

You have done your work, but that's<sup>1</sup> not enough.

You also need to evaluate how well your implementation works. The nature of the evaluation depends on your problem, your method, and your implementation that are all described in the thesis before this chapter. If you have created a program for exact-text matching, then you measure how long it takes for your implementation to search for different patterns, and compare it against the implementation that was used before. If you have designed a process for managing software projects, you perhaps interview people working with a waterfall-style management process, have them adapt your management process, and interview them again after they have worked with your process for some time. See what's changed.

The important thing is that you can evaluate your success somehow. Remember that you do not have to succeed in making something spectacular; a total implementation failure may still give grounds for a very good master's thesis—if you can analyze what went wrong and what should have been done.

---

<sup>1</sup>By the way, do *not* use shorthands like this in your text! It is not professional! Always write out all the words: “that is”.

## Chapter 7

# Discussion

## Chapter 8

# Conclusions

# Bibliography

- [1] Apache flume documentation. <https://flume.apache.org/> Accessed: 16.06.19.
- [2] Apache kafka documentation. <https://kafka.apache.org/> Accessed: 16.06.19.
- [3] Apache nifi documentation. <https://nifi.apache.org/docs.html> Accessed: 20.06.19.
- [4] Apache spark documentation. <https://spark.apache.org/docs/latest/index.html> Accessed: 26.06.19.
- [5] Apache systemml documentation. <http://apache.github.io/systemml/index.html> Accessed 30.06.19.
- [6] Deeplearning4j documentation. <https://deeplearning4j.org/docs/latest/> Accessed: 30.06.19.
- [7] Elastic stack documentation. <https://www.elastic.co/elk-stack> Accessed: 30.06.19.
- [8] Hdfs documentation. [https://hadoop.apache.org/docs/r1.2.1/hdfs\\_design.html](https://hadoop.apache.org/docs/r1.2.1/hdfs_design.html) Accessed: 22.06.19.
- [9] Iex trading documentation. <https://iextrading.com/> Accessed 15.04.19.
- [10] Mlflow documentation. <https://mlflow.org/docs/latest/index.html> Accessed: 30.06.19.
- [11] AMIRGHODSI, S., ALLA, S., KARIM, M. R., AND KIENZLER, R. *Apache Spark 2: Data Processing and Real-Time Analytics*. Packt Publishing, 2018.

- [12] BARTRAM, S. M., AND GRINBLATT, M. Agnostic Fundamental Analysis Works. *Journal of Financial Economics (JRE)* (June 20 2017).
- [13] BRIDGWATER, A. Nsa 'nifi' big data automation project out in the open. <https://www.forbes.com/sites/adrianbridgwater/2015/07/21/nsa-nifi-big-data-automation-project-out-in-the-open/#79a9319655d6> Accessed: 20.06.19.
- [14] CHEN, J. *Essentials of Technical Analysis for Financial Markets*. John Wiley & Sons Inc, 2010.
- [15] CHEN, L., AND YANG, C.-Y. Stock price prediction via financial news sentiment analysis. <https://github.com/Finance-And-ML/US-Stock-Prediction-Using-ML-And-Spark> Accessed: 08.05.19.
- [16] CHENG, J. Real time machine learning architecture and sentiment analysis applied to finance. Slide set available at: <https://www.slideshare.net/Quantopian/real-time-machine-learning-architecture-and-sentiment-analysis-applied-to-finance> Accessed: 08.05.19.
- [17] DAS, S., BEHERA, R. K., KUMAR, M., AND RATH, S. K. Real-time sentiment analysis of twitter streaming data for stock prediction. *International Conference on Computational Intelligence and Data Science* (2018).
- [18] DAVID ANDREŠIĆ AND PETR ŠALOUN AND IOANNIS ANAGNOSTOPOULOS. Efficient big data analysis on a single machine using apache spark and self-organizing map libraries. *12th International Workshop on Semantic and Social Media Adaptation and Personalization* (2017).
- [19] DAY, M.-Y., NI, Y., AND HUANG, P. Trading as sharp movements in oil prices and technical trading signals emitted with big data concerns. *Physica A: Statistical Mechanics and its Applications* 525 (2019).
- [20] FISCHER, T., AND KRAUSS, C. Deep learning with long short-term memory networks for financial market predictions. *European Journal of Operational Research* (2017).
- [21] FOX, J. *Myth of the Rational Market*. Harper Business, 2009.
- [22] GEORGE, L. *HBase: The Definitive Guide, 2nd Edition*. Packt Publishing, 2018.

- [23] GU, R., ZHOU, Y., WANG, Z., YUAN, C., AND HUANG, Y. Penguin: Efficient query-based framework for replaying large scale historical data. *IEEE Transactions on parallel and distributed systems* 29 (2018).
- [24] HOFFMAN, S. *Apache Flume: Distributed Log Collection for Hadoop - Second Edition*. Packt Publishing, 2015.
- [25] ISLAM, S. R., GHAFOOR, S. K., AND EBERLE, W. Mining illegal insider trading of stocks: A proactive approach. *IEEE International Conference on Big Data* (2018).
- [26] JOHN L. PERSON. *Mastering the Stock Market: High Probability Market Timing and Stock Selection Tools*. John Wiley & Sons, Incorporated, 2013.
- [27] KHASHEI, M., AND HAJIRAHIMI, Z. A comparative study of series arima/mlp hybrid models for stock price forecasting. *Communications in Statistics Simulation and Computation* (2018).
- [28] KIARASH, A., AND ABBAS, K. A New Approach to Predict Stock Big Data by combination of Neural Networks and Harmony Search Algorithm. *International Journal of Computer Science and Information Security* 14 (2016).
- [29] KUMAR, H., AND KUMAR, M. P. Apache storm vs spark streaming. <https://www.ericsson.com/en/blog/2015/7/apache-storm-vs-spark-streaming> Accessed: 19.05.19.
- [30] KYO, K. Big data analysis of the dynamic effects of business cycles on stock prices in japan. *15th International Symposium on Pervasive Systems, Algorithms and Networks (I-SPAN)* (2018).
- [31] LE, L., AND XIE, Y. Recurrent embedding kernel for predicting stock daily direction. *IEEE/ACM 5th International Conference on Big Data Computing Applications and Technologies* (2018).
- [32] LEE, C., AND PAIK, I. Stock market analysis from twitter and news based on streaming big data infrastructure. *IEEE 8th International Conference on Awareness Science and Technology* (2017).
- [33] LOTTER, J. C. Bye yahoo, and thanks for all the fish. <https://financial-hacker.com/bye-yahoo-and-thank-you-for-the-fish/> Accessed: 19.05.19.



- [34] MUN, F. W. Big data, small pickings: Predicting the stock market with google trends. *The Journal of Index Investing* 7 (2017).
- [35] MURPHY, J. J. *Technical analysis financial markets: A comprehensive guide to trading methods and applications*. New York Institute of Finance, 1999.
- [36] NAM, K., AND SEONG, N. Financial news-based stock movement prediction using causality analysis of influence in the korean stock market. *Decision Support Systems* 117 (2019).
- [37] NEERAJ, N., MALEPATI, T., AND PLOETZ, A. *Mastering Apache Cassandra 3.x - Third Edition*. Packt Publishing, 2018.
- [38] PALMER, N., RICKER, T., AND PAGE, C. DATA & ANALYTICS: Analyzing 25 billion stock market events in an hour with NoOps on GCP. Available at: <https://www.youtube.com/watch?v=fq0paCS117Q> Accessed: 08.05.19.
- [39] PENG, Z. Stocks analysis and prediction using big data analytics. *International Conference on Intelligent Transportation, Big Data & Smart City* (2019).
- [40] RAO, T. R., MITRA, P., BHATT, R., AND GOSWAMI, A. The big data system, components, tools, and technologies: a survey. *Knowledge and Information Systems* (2018).
- [41] SENGUPTAA, S., BASAKA, S., SAIKIAB, P., PAULC, S., TSALAVOUTISD, V., ATIAHE, F., RAVIF, V., AND PETERS, A. A review of deep learning with special emphasis on architectures, applications and recent trends.
- [42] SEZER, O. B., OZBAYOGLU, A. M., AND DOGDU, E. An artificial neural network-based stock trading system using technical analysis and big data framework.
- [43] SKUZA MICHAL AND ROMANOWSKI ANDRZEJ. Sentiment analysis of twitter data within big data distributed environment for stock prediction. *2015 Federated Conference on Computer Science and Information Systems (FedCSIS)* (9 2015).
- [44] SNIVELY, B. AWS re:Invent 2018: Big Data Analytics Architectural Patterns & Best Practices. Available at: <https://youtu.be/ovPheIbY7U8> Accessed: 08.05.19.

- [45] UTTHAMMAJAI, K., AND LEESUTTHIPORNCHAI, P. Association mining on stock index indicators. *International Journal of Computer and Communication Engineering* 4 (2015).
- [46] VARTAK, M., SUBRAMANYAM, H., LEE, W.-E., VISWANATHAN, S., HUSNOO, S., MADDEN, S., AND ZAHARIA, M. Modeldb: A system for machine learning model management.
- [47] WANG, W. A big data framework for stock price forecasting using fuzzy time series. *Multimedia Tools and Applications* 77 (2018).
- [48] WORLD FEDERATION OF EXCHANGES. Monthly report january 2019. <https://www.world-exchanges.org/our-work/statistics> Accessed 15.04.19.
- [49] YANBIN, W., GUO YIQIANG, L. L., NI, H., AND LI, W. Trend analysis of variations in carbon stock using stock big data. *Cluster Computing* 20 (2017).
- [50] YARABARLA, S. *Learning Apache Cassandra - Second Edition*. Packt Publishing, 2017.
- [51] YU-CHENG, K., JONCHI, S., AND JIM-YUH, H. eWOM for Stock Market by Big Data Methods. *Journal of Accounting, Finance & Management Strategy* 10 (2015).
- [52] YUAN, D. Stream processing in uber. <https://www.infoq.com/presentations/uber-stream-processing/> Accessed: 19.06.19.