

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, October 9, 2019

DRAFT! — November 27, 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

Author:	Ville Vainio		
Title:	Engineering analytics of big data pipelines for stock market data		
Date:	October 9, 2019	Pages:	62
Major:	Computer Science	Code:	SCI3042
Supervisor:	Professor Linh Truong		
Advisor:	Professor Linh Truong		
<p>The world economics revolve around stock markets. The information in the stock market does not only represent the current state of economy but it also reflects other phenomenons that affect the prices of the stocks. This is why there are constantly new academic studies so that we could understand both global economy and these phenomenons better.</p> <p>The size of the data that the stock market alone produces in a year is in the scale of terabytes. In order to analyze this data, tools that can handle this much infomation are necessary. Unfortunately, the information about these tools in the context of stock data is scattered, somewhat outdated and hard to find, which can drive away potential valuable research. The goal of this thesis is to provide information and tools that can make this analysis more accessible to data analysist.</p> <p>This thesis performs a literary research on what is the current state of stock data analysis in big data environment focusing on historical data analysis pipelines and based on this research implements an open-source pipeline that anyone can use. The literary study shows that the current trend in stock data analysis is deep learning models giving the pipelines requirement to support these methods. In the implementation part we see that running deep learning libraries in distributed enviroment can be quite challenging and provide information about these challenges and how to solve them. The resulting pipeline does fullfil most of its requirements, but the results of the study are a bit incomplete without better comparative study.</p>			
Keywords:	stock market, big data, cloud computing, data analysis		
Language:	English		

Aalto-yliopisto

Perustieteiden korkeakoulu

Tieto-, tietoliikenne- ja informaatiotekniikan maisteriohjelma

DIPLOMITYÖN

TIIVISTELMÄ

Tekijä:	Ville Vainio		
Työn nimi:	Osakemarkkina dataa käsittelevien big data järjestelmien tekninen analyysi		
Päiväys:	9. lokakuuta 2019	Sivumäärä:	62
Pääaine:	Tietotekniikka	Koodi:	SCI3042
Valvoja:	Professori Linh Truong		
Ohjaaja:	Professori Linh Truong		
<p>Maailman talous pyörii osakemarkkinoiden ympärillä. Informaatio, jota osakemarkkina tarjoaa ei pelkästään kerro talouden nykytilasta vaan se myös heijastaa kaikkia osakemarkkinoihin vaikuttavia ilmiöitä. Tämän takia uusia tutkimuksia tehdään jatkuvasti jotta sekä taloutta että näitä ilmiöitä voitaisiin ymmärtää paremmin.</p> <p>Osakemarkkinoiden tuottaman datan määrä vuodessa on teratavujen mittaluokassa. Jotta tätä dataa voitaisiin siis tutkia, tarvitaan työkaluja jotka pysyvät toimimaan tämän mittaluokan kanssa. Valitettavasti informaatio näistä työkaluista osake datan kontekstissa on hajanaista, jossain määrin vanhentunut ja vaikeasti löydettävissä, mikä jossain tapauksissa estää mahdollisesti tärkeän analyysiin tekemisen. Tämän diplomityön tarkoituksena on tarjota tietoa ja työkaluja, jotta tämän tyyppinen tutkimus olisi mahdollisimman helpopääsyistä tutkijoille, jotka eivät ole tämän alan ammattilaisia.</p> <p>Tässä diplomityössä suoritamme kirjallisuuskatsauksen tämän hetken osake datan tutkimukseen big data ympäristössä. Keskitymme historialista dataa käsitteleviin järjestelmiin ja toteutamme avoimen lähdekoodin järjestelmän, jota kuka tahansa voi hyödyntää. Kirjallisuuskatsaus osoittaa että tämän hetken kehityssuunta osake datan tutkimisessa on deep learning menetelmät, jonka vuoksi järjestelmien on tuettava näitä metodeja. Työn kokeellisessa osuudessa huomaamme että deep learning ohjelmistokirjastojen käyttö hajautetussa ympäristössä voi olla hyvin hankalaa ja tarjoamme tietoa siitä mitä kaikkea voi mennä vikaan ja miten korjata nämä. Työn tuloksena saamme kehitettyä järjestelmän joka täyttää suurimman osan sille annetuista vaatimuksista, mutta tulokset jäävät hieman vajaavaisiksi, koska vertailevaa testausta ei ehditä toteuttaa.</p>			
Asiasanat:	osakkeet, big data, pilvilaskenta		
Kieli:	Englanti		

Espoo, October 9, 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
DL4J	Deeplearning4j Machine learning library

Contents

Abbreviations and Acronyms	5
1 Introduction	8
1.1 Application scenario	9
1.2 Audience and Research Questions	12
1.3 Expected Outcome	14
1.4 Structure of the Thesis	14
2 Background	16
2.1 Stock market analysis	16
2.1.1 Stock price prediction	17
2.1.1.1 Statistical methods	17
2.1.1.2 Machine learning methods	18
2.1.2 Existing pipelines	20
2.1.2.1 Data sources	20
2.1.2.2 Used technologies	21
2.1.2.3 Monitoring	22
2.2 Big Data Technologies	23
2.2.1 Data ingestion	23
2.2.1.1 Apache Flume	24
2.2.1.2 Apache Kafka	25
2.2.1.3 Apache NiFi	25
2.2.2 Data storage	26
2.2.2.1 Hadoop distributed file system	26
2.2.2.2 Apache HBase	27
2.2.2.3 Apache Cassandra	28
2.2.3 Machine learning frameworks	28
2.2.3.1 Apache Spark	29
2.2.3.2 Deeplearning4J	30
2.2.3.3 Apache SystemML	30
2.2.4 Monitoring	31

2.2.4.1	Log monitoring	31
2.2.4.2	Machine learning monitoring	32
3	Planning	33
3.1	Building a developer friendly pipeline	33
3.2	Data source	34
3.3	Technology selection and requirements	36
3.4	Pipeline evaluation	37
4	Implementation	39
4.1	Overview	39
4.1.1	Dead-ends at the start of development	39
4.1.2	Final Architecture	40
4.1.3	Usage	43
4.2	Adversities during development	44
4.2.1	Kafka-HDFS integration	45
4.2.2	Apache Zeppelin dependency management	46
4.2.3	DL4J sbt project	46
4.2.4	Running the application in Spark	47
5	Evaluating the results	49
5.1	Perfomance metrics	49
5.2	Pipeline management	51
5.3	Flaws in the end product	52
5.4	Reflection	53
5.5	Potential improvements	53
6	Conclusions	55

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. [35] 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. [30]

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.[20] 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. [43] 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 be used to conduct technical analysis means that the system should be planned 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 and researchers, who want to analyse this data efficiently, basic knowledge and tools on what are the best current solutions on handling this data. With this knowledge these data scientists 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 what are the challenges of developing a system for technical analysis are, we will next look at a example case. This is to give the reader more complete view of a system that handles stock market data analysis. After this, for the rest of the thesis, we are going to focus on a part of this system, but the purpose of this example is to give a better view of the system as a whole in order to understand some of the requirements that these other components give to the focus of this thesis.

Somewhat normal system for analysing both fundamental and technical

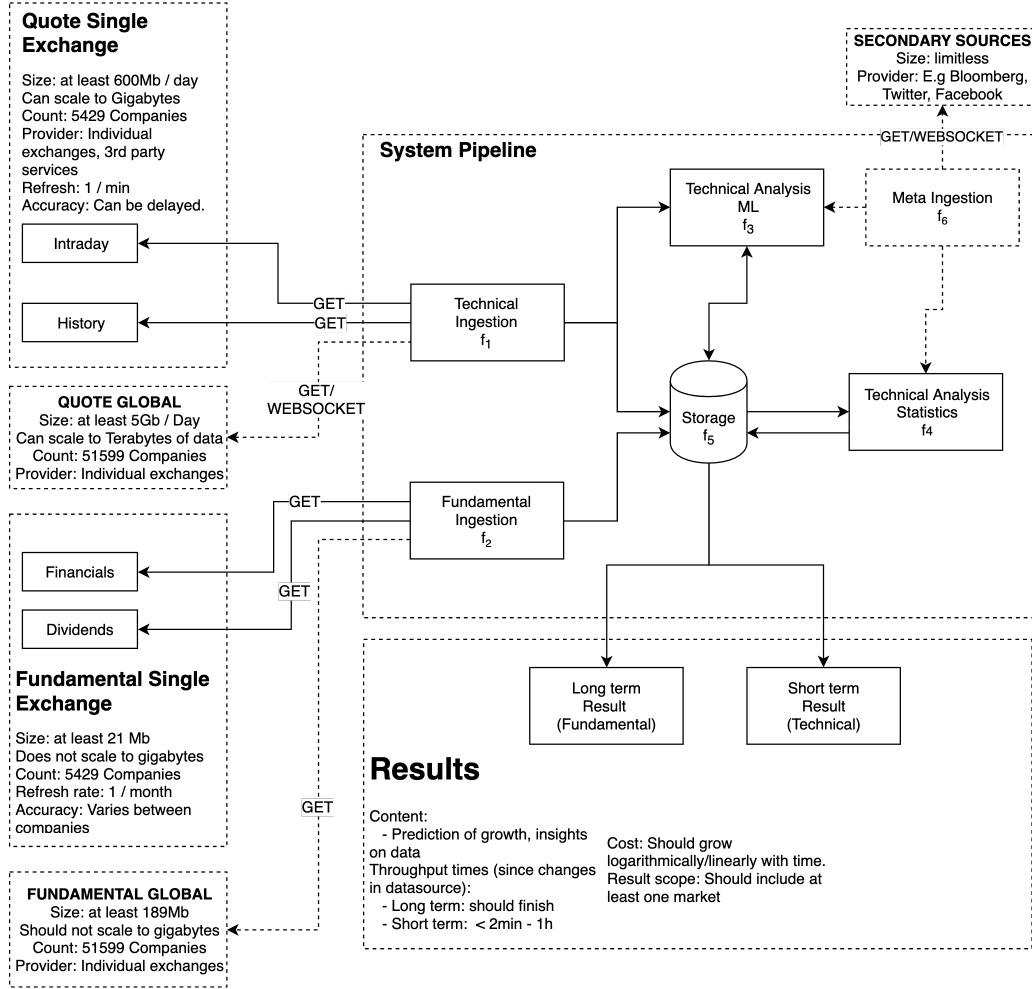


Figure 1.1: Example of a stock data pipeline

analysis is presented in figure 1. In the figure, the components 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 are usually 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 which is the most accessible and individually significant market.

The data to the application can be 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 single market ones (d_1 and d_2), which represent the minimum of at least including the U.S stock market and the other sources (d_3 and d_4) that provide the same data on other global markets.

In the figure, 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 minimum of a 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 [57]. 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 [59], [52] and [42]. 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 usually 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 for example Excel sheets and sometimes websockets, of which the websockets can be actually useful in cloud system, but as the amount of data grows more custom file transfer methods are used. However, as the HTTP-methods are currently the most publicly available technology, this thesis is going to focus on these. The main ingestion functions are usually separated 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 usually store their raw output into the storage, but f_1 passes this also to the immediate technical analysis.

The system should usually have two types of methods on analysing data. 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, there 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 querying historical data and for the results the storage should provide efficient point queries for the results ($< 1s$).

1.2 Audience and Research Questions

We saw the possible complexity of a stock data analysis pipeline in the previous section. There are a lot of articles about how to build a stock analysis pipeline but these usually are focused on the analysis of small amount of data in order to test some hypothesis. Building a large-scale stock data analysis pipeline is a vastly complex task and to alleviate this process, this thesis plans to provide tools and information on how to build one.

This thesis is aimed at novice data scientists that work or want to work with stock market data, but do not have any experience in big data technologies. We define the term "novice" here to mean persons who have knowledge about analysing data using modern data analysis methods and know how to conduct this onto small datasets, but their proficiency lies strictly on the data analysis side. So the reader can be very good at analysing data, but lacks proficiency at the engineering side of big data technologies. The reason why the reader would want to read this could be because they have interest in the subject or are, for example, working for a startup which does not have resources to hire multiple persons with high degree on the subject. This thesis is meant for this kind of people, as the place to find easily knowledge on how to get started.

As the system can be very complex as seen in the previous section, the focus of this thesis is going to be the part of the pipeline that calculates

models based on historical data. This means the ingestion (f_1), storage (f_5) and analysis (f_4) of the historical data in a environment where the time restrictions are not that tight but the amount of data is enormous. We want the pipeline to be able to have the possibility to ingest meta data from third-party sources as seen on the example case in the previous section. We have chosen this part of the pipeline as this is the part that has the highest probability to scale to big data first and is the part of the pipeline that is usually referred when talked about data analysis on stock market data.

The main goal of this thesis is going to be to make the data analysis in this big data context as accessible as possible within the time constraints and resources that this thesis has. What we mean by the word "accessible" here is that we want to minimize the time needed for the novice data scientist to start experimenting with their analysis, while making sure that developing models can continue smoothly in the future. In this sense, things that make the process of analysing stock data in big data context "inaccessible" are things such as scattered information about options, outdated information online and lack of information on how to integrate different parts together. This thesis plans to alleviate these aspects by providing aggregated information, testing and software that can be used to navigate through these different problems.

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 kind of pipelines are currently used in stock market data analysis in big data context? What are the methods and tools currently used to conduct technical analysis and possibly what kind of data they use. This thesis plans to aggregate this scattered information into more available format so that the reader can more easily read about the subject.

How do these pipelines compare in the perspective of novice data scientist? As there are enormous amount of different technological frameworks, this thesis plans to provide the reader information on the most promising ones currently and do a comparison between them so that user can have some point of reference when choosing technologies on their own.

Which one is the best for novice data scientist as a starting point? Finally the thesis plans to provide a prototype implementations of stock data pipelines using the technologies that seem most prominent. These are compared with one another using relevant metrics that reflect how manageable and easy it is to develop on the implemented pipeline. If the prototype system fits readers developed architecture it can be used as a basis for further development by the reader This is to save time of a person developing these when the trial and error has been done beforehand.

1.3 Expected Outcome

For the first research question "What kind of pipelines are currently used?" this thesis plans to provide a study on currently used stock market analysis methods and used pipelines. The outcome will be a literature study on pipelines that are used in academia and industry based on public information.

For the second question "How do these pipelines compare in the perspective of novice data scientist" the thesis would perform an analysis on the most promising big data pipeline technologies. This will be done in a form of comparison of different types of solutions that are either used in practice based on the results in the literature study for the first question or just seem promising for the use case. The result of this part could be used to decide what seems to be the most suitable technology to use to implement the stock analysis pipeline technically.

For the last question "Which one is the best for novice data scientist as a starting point?" the thesis would implement open-source prototype solutions based on the results of the second research question. The expected outcome is to find a pipeline that a novice data scientist could develop onto which at the same time fulfills the aforementioned requirements. These prototypes could be used by anyone (company or individual) as it is or as a base to build a more complex system on top of them or use them as it is to perform data analysis.

1.4 Structure of the Thesis

In the first chapter of the thesis we will be focusing on solving the first research question "What kind of pipelines are currently used?" by offering aggregated background knowledge on the subject. The thesis will start by going through scientific papers about stock markets and stock analysis. We will examine trends and state-of-art methods on stock analysis that have been used in recent scientific research. Then we move on to examine what kind of pipelines exists in real-life. We will be examining both commercial and research pipelines in order to give a more complete picture on options that exist. We will also perform a literary research on technologies that are currently used to perform big-data ingestion, storing and analysis, selecting from the list of technologies mostly those that have been seen to be used in practice in this context while introducing couple of promising ones.

In the next chapter we will be focusing on the second research question "How do these pipelines compare in the perspective of novice data scientist?" and perform an comparison of the technologies introduced in the previous

chapter. This will be used also as a preparation for the experimental part of the thesis by choosing the technologies that we are going to try to use in practice. After we have seen how these technologies compare to one another we start the experimental part by defining what are we going to be building and how we are going to evaluate it. All this in preparation to answer the final research question "Which one is the best for novice data scientist as a starting point?".

We will then start the actual implementation part. In the following chapter, we will describe what was done and what challenges we faced while implementing the systems. The chapters from this on will be focused on the technical and practicalities that appeared during development.

After this chapter we move on to evaluate the results that we got from the implementation chapter. We will try our best to evaluate the quality of solutions and bring up fair criticism that the end result could possibly have. We then finish this chapter with a discussion on what could have been done better and how the systems developed could be improved in the future. Then finishing the thesis we will have a conclusion chapter that summarizes the results of the thesis. However, next we start the thesis by examining the field of stock data analysis.

Chapter 2

Background

In this first chapter, we are going to look at the current stock market data analysis in big data environments. This chapter is divided into two distinctive parts. In the first part we will examine the current state of stock market data analysis in big data environment. This includes the methods and pipelines used to conduct this in real life. In the second part we will use the result of the this first section and study the most promising big data technologies that could be used to implement pipelines that analyze the stock data.

2.1 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 the possible causalities. In this thesis we will be approaching this from the price prediction perspective which can also be used as a tool in other problems such as anomaly detection [34].

In this chapter we examine what are the state-of-the-art methods that researchers use to analyse stock market data. We are especially focusing on methods that in some way use big data for this analysis. We will be also inspecting real-life implementations of stock data analysis and examine some of the existing big data pipelines for stock market analysis focusing on what are the technologies used to build these.

2.1.1 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 the 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. [44] 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 [27].

2.1.1.1 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 countless 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 data related to the price of a stock is analysed, we will be looking at momentum indicators. [54] 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. [22] When searching research on big data usage in technical analysis there are some key indicators that show up frequently. These are relative strength index (RSI), moving average convergence divergence MACD and Williams %R which are all momentum indicators. These indicators can be used as they are but these have been also used as features that machine learning algorithms use to predics prices. [51]

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 which can help to recognize one indicators bias when using multiple different indicators. [22]

Moving away from the momentum indicators, 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. [56] However, with the new advances in machine learning, there seems to better performance to be achieved with machine learning methods. [36] Due to this and the complexity of ARIMA models, we will 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. [27] tried to link oil prices to stock prices using global financial data streams with stochastic oscillator techniques. Kyo [38] combined technical and fundamental analysis by using regressive model that takes into account business cycles in Japan's stock market.

2.1.1.2 Machine learning methods

As stated before, machine learning is the current trending stock analysis methology 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.

We 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. [50]

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. [50]

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. [29] 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. [51]

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 [40]. 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.

2.1.2 Existing pipelines

We then 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. The big data might not be 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 more relevant results.

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 achitectures for business reasons. Nevertheless, with the public information available we can get an excellent view of the state-of-the-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 lifecycle 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.

2.1.2.1 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 a 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 pipelines lifecycle when there possibly is no way to access the complete stock market data, but system must be able to scale to this complete data set size when time progresses.

In academic papers, the Yahoo Finance API has been the de facto service used as the source of this partial stock data. [34] [26] [39] [51] 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.[41] These services will be introduced more completely in chapter 4 where we will evaluate which one to use in the implementation chapter.

2.1.2.2 Used technologies

We then move on to examine the actual technologies and frameworks that are in use by starting from the furthest away from the client, the ingestion layer. With ingestion we mean the part of the pipeline that handles fetching and initial processing of the data. This step with stock market data varies a lot depending who is fetching the data and for what purpose. The most common method for ingesting stock market data was found to be custom scripts. This is because the format of ingested data can vary greatly and scripts are relatively easy to write. This is usually enough with the aggregated stock market data, however, it does not scale.

Common big data technologies are usually introduced when the system has to ingest for example great 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 is used, for example, in [48] and [25]. Another framework that is commonly reported in the ingestion step is the Apache Kafka streaming platform. [40] [24] Both of these are open-source frameworks which makes their usage relatively cost-free. In the industry side the usage of ready-made services from cloud providers are common. These services are for example the Google Dataflow as in [47].

After the ingestion, the data has to be stored. Stock market data is quite structured by itself but the data from third-party sources used to enrich stock data is usually not. 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. 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. [53] [47] In the open-source side HBase [32] and Cassandra seem to be the two most openly used database solutions.

Then as we finally have the data in control, we can apply some analysis to it. In the analysis phase, there is not that much variety in used technologies. Apache Spark with its MLlib library is dominating this field with its capabilities to process data efficiently. [34] [23] [40] [26] Where Apache Hive and Apache Pig were previously most used technologies, now Spark seems to be taking their place [53].

In the industry side, there is indications of Apache Storm being used with real-time classification [24]. 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 [37].

2.1.2.3 Monitoring

Finally we are going to examine one specific characteristic which is usually ignored when pipelines are reported. Good monitoring is essential to ensure that the pipeline is running correctly and optimally. In the next few paragraphs we explain a bit more about the importance of monitoring as this is one of the areas that we will focus while trying to create a stock data pipeline.

None of the public sources on pipelines that were examined for this chapter reported exactly how they monitored their pipelines in practice. In cases where for example Spark is used, we can assume that the process is monitored using Sparks 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.

With pipelines that analyze massive amounts of data, the quality of data

is very important. This is not only when choosing what data source to use but the data must be kept from corrupting during the whole pipeline and to ensure this good monitoring is essential. Performance monitoring can pick up these corruptions if the corruption is large enough to make the program crash. However, in most cases this corruption can be only a change in values that would pass as an input. This would then lead to erroneous results which could take quite a bit of resources to calculate.

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 and this is why we need good monitoring. In the next chapter, we will be looking more closely on the technologies that could be used to build a pipeline that can conduct modern stock analysis. We will also examine couple of different ways to monitor the system in order to prevent problems that were raised in this section.

2.2 Big Data Technologies

In this section, we will examine more carefully the technologies that can be used to implement a stock data pipelines in the big data context and try to answer the second research question in the process. We have gathered here technologies that we saw used in existing big data stock pipelines in the previous section. 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 novice data scientists that do not possibly have access to costly 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.

2.2.1 Data ingestion

Here we will be using the same definition for ingestion as before; data fetching and initial processing which includes data validation. 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 developers are talking about open-source 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 has promising features that could be suitable for our use case.

2.2.1.1 Apache Flume

Apache Flume is one of the Apache Software Foundation (ASF) projects that is quite popular 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" [2]. 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. [33] 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. [33]

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. [33]

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. [33]

2.2.1.2 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. [3] 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. [49] After being released from incubation in 2012 many other big industrial companies such as eBay and Uber [60] 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. [49]

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. [3]

2.2.1.3 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. [21]

As we have seen already with the last two ingestion frameworks, all these frameworks have their own abstractions for data and the processes that handle 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. [4]

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.[4]

2.2.2 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 error states that could occur when the hardware malfunctions for example with power outages or network outages. We have gathered three major big data storing technologies that were used in existing pipelines: HDFS, HBase and Cassandra. We start with the HDFS and move towards more industrially used systems.

2.2.2.1 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.[13]

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.[13]

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.[13]

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. [13]

2.2.2.2 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.[31]

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. [31]

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. [31]

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). [31]

2.2.2.3 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 will be looking next. [58]

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. [45]

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 choice 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. [45]

2.2.3 Machine learning frameworks

When we researched existing pipelines we noticed that machine learning libraries such as Tensorflow and PyTorch are still quite used even when data can scale to enormous amounts. This is partly due to their support for 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 with data which size is scaled to terabytes, 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 native deep learning solutions, which are the state-of-the-art methods we are interested in. Spark is still a highly efficient

processing framework for big data and the de facto technology in the market with 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.

2.2.3.1 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.[5]

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.[5]

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.[19]

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 natively tools to implement deep learning

models on Spark.

2.2.3.2 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 pre-processing the data before fed for training in the form of sub-projects. [10]

In distributed environment DL4J trains its models using Stochastic Gradient Descent (SGD). This is implemented in parallel on each node of the cluster using either of the two methods that the DL4J offers, gradient sharing or parameter averaging. From these two, the first one has become the preferred way to implement SGD in DL4J starting from its latest release and this is why we are going to ignore the technicalities of the latter one for now.

In the gradient sharing approach, the gradient is calculated asynchronously on each node in the cluster. The main idea behind DL4Js implementation of this, is that not every update on the gradient is sent to every other node. Each node has a threshold which defines when a change is large enough to be shared with the global gradient. This combined with its own heartbeat mechanism provides efficient and fault tolerant way of training a model in distributed environment. [10]

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. [10]

2.2.3.3 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. [6]

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. [6]

2.2.4 Monitoring

We have already seen quite a bit of framework specific monitoring solutions, which cover monitoring individual components in the pipeline. Next we will 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.

2.2.4.1 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. [11]

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. [11]

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. [11]

2.2.4.2 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. [55]

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. [16]

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. [16]

Chapter 3

Planning

Next we will look at how do we can use the technologies in the last chapter in order to build a pipeline that can be used to examine stock data when the amount of data might be too much for one computer to handle. We will be reviewing, what we will be doing in the implementation part of this thesis and analyze the information that we have seen to this point. Finally, we will look how we will evaluate the end-product and how is it better or worse than other alternatives.

3.1 Building a developer friendly pipeline

The main goal of the following implementation section is to provide information on developing these pipelines. As we said in the introduction we want to build the pipelines so that the reader can either use them as they are or as a basis for more complex system. This is why we want to develop a system that is easy to develop onto. We will be focusing making the system run in local environment for the analysts who do not possibly have the resources to invest much time in order to build a analysis pipeline that can scale to big data. This will be the base for most of the decisions that are being made for some of the used technologies.

The technical requirements for this is that the pipeline can be run almost any computer with enough computing resources. This means that the parent operating system should not require any excessive operations from the developer in order to start developing on their machine. This is why in this thesis all the parts of the pipeline will be build on top of Docker and Docker Compose container technologies. These are chosen as the container technology here as they are easy to develop onto and cloud providers such as Google Cloud Platfrom and Amazon Web Services both support them. [7]

[12] This way we can create a developer ready pipelines that work flexibly in their local machines and can be scaled to the cloud without too many excessive measures.

Most of the open-source Big Data products are build on top of JVM (Java Virtual Machine) so in order to avoid excessive complexity that can come from using multiple programming languages we will use single programming language that runs above JVM throughout the whole pipeline where it is possible. For this thesis that language is Scala, because of its good interoperability with Apache Spark and its programming interfaces for both object oriented and functional paradigms. On top this, Scala offers static typing which is extremely helpful to prevent errors that can occur with long computations saving developers time and resources. [46]

3.2 Data source

As stated before, most of the data sources concerning stock data are closed and heavily priced. However, there exists services that offer data for smaller scale development and analysis. A list of this kind of services is presented in the table 4.1. In the table the cheapest possible plan for each service is presented to give a better understanding what kind of possibilities there are for acquiring stock data for analysis purposes and what kind of data formats these offer.

From this table we can make couple of notes. For acquiring free data test data fast, the IEX provides the best option as its limits are not restricted to time, but the 500k datapoint restriction does not provide nearly enough data for any serious application. For academic use, Quandl provides at the moment of writing this, the most affordable API for historical data analysis. Historical data is usually cheaper because with historical data alone you usually cannot make money as stock market revolves around the most recent data, but for training models and other data analysis it is ideal. Other thing to note is that the amount of historical data can vary greatly between services from 6 months (Barchart) to 39 years (World Trading Data).

For this thesis we have chosen to use data from IEX API which was an open API until 30.09.2019 when the company decided to close this in order to capitalize with the closed API which free plan is in the table. This data is from the 5 year interval between 2014 and 2019 and the relevant values that it contains are opening prices, closing prices, highest prices, lowest prices and volumes. In order to test freely with this data, we will implement our own server that serves this data into our pipeline.

Table 3.1: Stock data sources

	Type of data	Price	Data Structure	Restrictions
IEX [14]	Intraday and Historical (15 years)	0\$/month	JSON	500k data-points/month
Alpha Vantage [1]	Intraday and Historical (20+ years)	0\$/month	JSON/CSV	5 requests/min and 500/day
World Trading Data [18]	Intraday and Historical (from 1980)	0\$/month	JSON/CSV	5 symbols/request, 250 request/day and 25 intraday requests/day
Intrino [15]	Intraday	52\$/month	JSON / CSV / Excel	120 request/min
Quandl [17]	Historical <i>EOD</i> ¹ (from 1996)	15\$/month ²	JSON / Excel	none
Barchart [8]	Intraday and Historical (6 months)	0\$/month	JSON / CSV / XML	25 symbols /request and 400 requests/day

¹ EOD = End of Day² For academic use

3.3 Technology selection and requirements

Next we will explain what technologies we are trying to use to build each part of the pipeline. The choices are highly based on the technologies that have been seen already in use, but there are also choices that seem promising but do not have real-life examples yet. This is hopefully to provide new meaningful knowledge about the possibilities of these technologies.

The requirements for ingestion can be derived from the chapters 1 and 2. We want the system to be able to ingest normal structured stock data, but have the ability to extend to third-party metadata that can have any form and is usually unstructured. Because stock data by itself can already scale to gigabytes per day we want the ingestion have the ability to scale with input.

For the ingestion we will try the technologies that we introduced in the previous chapter: Apache Flume, Apache Kafka and Apache NiFi. All of these are scalable data ingestion frameworks that have different paradigms of handling data as seen before. For only local development these products are a bit heavy weighted, but for scalability these are necessary in order to handle massive amounts of ingested data from varying sources.

From the storage we want the ability to scale with the possible terabytes of data. This means that even when amount of data is enormous, we want the queries to be executed in somewhat manageable time. We also want the storage to fault tolerant in a distributed environment in order to ensure that the data waiting to be analyzed retains its quality throughout the wait.

To fulfill these requirements, we will be testing Apache Cassandra, plain HDFS and, if there is time, Apache HBase. All of these storage formats have been developed to be used in big data environment and from these HDFS and HBase were both used in some existing pipeline. HDFS does not have as good as query capabilities that could be hoped for but offers easy integrations to other technologies and a mature development environment which is why we also included it. As for Cassandra, we are trying to see whether it is easy to integrate with the other technologies and can it bring anything to the pipeline with its high availability features.

With analysis part of the pipeline, we want the pipeline to be able to pre-process the data for training algorithms that use it and again the amount of data can be from gigabytes to terabytes. As we saw in the chapter 1, the current cutting edge methods for stock data analysis are deep learning methods. This is why we want the analysis frameworks have the ability to build and train these models without having to oneself implement them from scratch.

For analysis, there does not exist that many options that work well with

our requirements. The main difficulty here is to keep our pipeline using mainly Scala for programming. There are only two machine learning libraries that has Scala support and can be used in big data context; Apache Spark ML and Deeplearning4j libraries which were introduced in the previous chapters. Because deep learning, and specifically LSTM networks, are the current trend in stock analysis, the library needs to have support for these. Unfortunately, Spark ML does not have these natively as the writing of this, so that leaves us with only Deeplearning4j.

Because building and demonstrating data analysis applications can be a time consuming and complex job, data science community has adopted the usage of Jupyter notebooks when implementing data analysis in Python. For the same reasons, we will be integrating Apache Zeppelin notebooks into this pipeline as these have a support for Scala programming language and allow submitting applications to Spark cluster. This way we can offer more streamlined model implementation and testing.

Finally, we have the monitoring of the pipeline. The requirements here is to allow the novice data scientist to have a simple and intuitive way to manage and monitor the the state of the whole pipeline. The main goal is to give the data scientist an ability to notice errors in the pipeline as soon as possible this way preventing the errors to possibly escalate to the latter parts of the pipeline. Other requirements include good tracking on progress of machine learning model development in the analysis stage to give the analyst tools to track the results of training.

To monitor the machine learning model development we will integrate the MLFlow tracking into the analysis phase. This is done by implementing a separate tracking server to allow this process scale separately from the actual analysis. Thus also enabling responsive tracking UI usage. We will try to also develop a global ELK stack which was described in the previous chapter to give global monitoring on the entire pipeline.

3.4 Pipeline evaluation

After the we have the implemented pipeline, we need some ways to measure whatever our system fulfills the requirements it has been given and how does it compare with other software that might be as good as it. In this section we will introduce the measurements that we will be using in the chapter 6 after the implementation. Some of these are quantative such as the performance metrics, but these usually give very shallow view on the system. Because of this we will be using more qualitative methods in order to grasp the actual benefits of the pipeline.

Our focus with the qualitative metrics is going to be on how easy the pipeline is to develop onto and how well the user can monitor changes in the pipeline. These are usually important factors when continuously developing a system and can affect greatly when new developers are choosing what pipelines to use. This also includes how automated the processes in the pipeline are and does developing and deploying the software require excessive amount of manual work.

Some quantitative metrics are taken in order to asses the fitness of the pipeline for any kind of development where it is important that iterations of the software can be swiftly build and that the whole service can run comfortably on the developers machine. Namely for these reasons we will examine the memory consumption of these services on a normal developing machine. Other quantitative statistics that we can examine are processor usage, network usage and disk usage which, although not as restricting as the memory usage, can tell a lot about the performance of the system. For these metrics, we will be using Dockers native tools, and especially the command 'docker stats', to measure the containers usage of resources.

As there will be a bit different pipelines with different technologies, we can get results that can be used to compare components with one another. This would give the reader better information and alternatives on what technologies should they could choose for their own project. Now we move on to the implementation chapter where we examine more closely what was done during this thesis.

Chapter 4

Implementation

In this chapter we will introduce the software created for this thesis and the challenges faced while developing it. This chapter is divided into two parts. In the first part we will examine technically what was the end result while quickly commenting on some of the choices that can be seen in the final architecture. Because of time constraints and challenges during development we were able to only produce one working pipeline. These challenges are explained in the latter part of this chapter and we continue with them in the following evaluation chapter.

4.1 Overview

In this section we will explain the practical application that was developed. We start by addressing some of the deadends that lead to the final pipeline to give reader a clearer view what changed in practice after the previous chapter. The final architecture is the end product of this thesis. Some of the reasoning why some of the technologies were prioritized leading to not being to test every technology is presented here but more in-depth look at for specific integrations and problems will be examined in the next section. Finally in the section 5.1.3, we introduce how the pipeline can be run with any novice data scientists local machine.

4.1.1 Dead-ends at the start of development

Two technologies that were mentioned in the previous chapter were eliminated at the start of the development of each. These were Apache Flume and Apache Cassandra. As both seemed to lead to a dead-end in the development and there existed seemingly better alternatives than these, focus was moved

to other technologies.

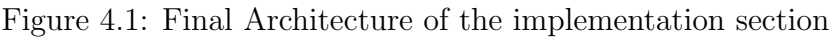
Apache Flume was supposed to be used as a component in the ingestion step, but due to its characteristic of being mostly aimed at data which could be fetched in log form, it could not be used with the REST APIs that most of the data services use. Flume also did not seem to have any native ways to connect to storages such as Cassandra. There exists some ready-made solutions for this, but these were unmaintained and seemed possible dead-ends. This is why Flume was ruled out.

Apache Cassandra was dropped from the development mostly because integrating it with other services would have required a lot of extra work. Further examination of Cassandra also showed that Cassandra is not very good choice for this specific purpose that we want to use it for. Cassandra is known for its ability to write enormous volumes of data, but when it is time to read from the storage the task becomes non-trivial. Reading from Cassandra is based on the Cassandras query language, CQL, which allows somewhat efficient queries but you have to define these high performance queries at the time you create tables for data. [28] Because our pipeline should be used mostly by analysts who can have varying queries into the storage this raises a challenge for future development. This is why we decided to discontinue with Cassandra.

4.1.2 Final Architecture

The final architecture can be seen in the figure 5.1. The figure presents the main dependencies between the components each arrow usually presenting the flow of data in the system with labels sometimes added to clarify better the relation. In the figure each solid rectangle represents one actual server with its own process with the exception of rounded rectangle to represent notebooks that are stored into the local file system. In the case of this thesis these servers were virtual machines but theoretically the system could be run with right configuration with multiple physical machines that are in the same network. Rectangles with small dashed lines such as with the spark workers represent that the server can have easily more than one instance allowing horizontal scaling. The longer dashed lines here represent bigger entities such as the HDFS cluster to make the figure more readable.

The pipeline starts from the data source that acts as the kafka producer. It reads the data from a source, in this example case from JSON files, and produces the data into a kafka topic. It also creates the needed kafka topics on startup using Kafkas admin API. Although most of the kafka documentation recommends using the command line client, this did not seem production-ready way to do this as this decouples the topic creation logic from the actual



Then we have Kafka which is integrated into HDFS cluster with Kafka Connect HDFS 2 Sink connector which is a component made by Confluent.

This acts as a kafka consumer and listens to the topic that the producer defines. Connect can be scaled and contains some basic configurations that can be used for example to define the format in which the data is stored into HDFS. We have chosen JSON as the format which the data is stored but Apache Avro is also an option. Other notable configurations that can be used to tweak the pipeline are the rate in which the topic is consumed into the HDFS and schema validation.

In the middle of the pipelin is the HDFS cluster. In this cluster there is the normal HDFS cluster with namenodes and datanodes which can be accessed to directly read the data. There is also a YARN resource management on top of this which can be used to for better manage IO routines in the cluster. This consist of the resource manager, node manager and history server nodes, that can be used as an alternative route to access data.

The HDFS cluster docker containers, as well as the containers for Spark cluster later, are a work of Big Data Europe project. The Big Data Europe is project funded by European union which one of the goals is to produce open-source tools for big data development without the need to use closed softwares.[9] As these containers were the most maintained hadoop containers at the moment of writing this, they are were the best option for this case, but they brought a couple of problems which we will examine later.

Table 4.1: Storage and ingestion components versions

	Hadoop	Kafka	Kafka Connect	Zookeeper
Version	3.1.1	2.3.0	5.2.1	3.5.5

In table 5.1 we have gathered the versions of the technologies used in the first half of the pipeline. Most are the newest versions of each software at the writing of this. The version of hadoop is specially tricky as its clients are used in multiple parts of the pipeline coupled with other software libraries which have support for only some of the older versions but do not complain if used with the newer one. This is why there can be other versions than this in the pipeline e.g in the spark cluster, but as they do not currently cause any visible errors and due to the time constraints that this project has, the versions can mismatch for now.

After the HDFS cluster comes the analysis part of the pipeline which has Spark cluster in the middle and two options to run analysis code on it. The application part allows writing production grade scala applications that can be run like any normal scala spark application. The Zeppelin is a Apache Zeppelin instance which allows writing and running notebooks that can be saved into local file system for distribution. Both approaches

submit the spark application to spark using spark-submit script and the DL4J communicates with itself with its own parameter server.

In default case, the spark application first preprocesses the json formatted data and saves this back to HDFS as csv files. In this process, it normalizes the data and appends labels to it that in our example case is just values 1 and 0 whenever the value of stock grew in n-days after the datapoint or not respectively. The data is stored back as a CSV file because DL4J is quite picky about the data format that it accepts and does not have simple default way of transforming spark DataFrames with sequential data to the Dataset format that it internally uses. After this the data goes through the training and evaluation pipeline which logs its parameters and results into MLFlow server where user can monitor the process of their different experiments.

Table 4.2: Analysis Software versions

	Spark	Scala	Java	DL4J	sbt	Zeppelin	MLFlow
Version	2.4.3	8.x	2.11.12	1.0.0-beta5	1.26	0.8.1	1.2.0

In the table 5.2 we have collected the versions of different components in this part of pipeline. Almost all of them are the latest releases of each technology at the time of implementing with the exception which is the version of Scala. Currently, the latest version of Scala is 2.13. Spark mainly supports scala 2.12, because as of 2.4.1 version of spark, the version 2.11 of scala is deprecated and there is still no support for 2.13. However, Zeppelin does not support Scala version 2.12 so the only version of Scala that still works with all of the latest releases of these two technologies is 2.11 which is why it had to be chose for the version. Java is still version 8 which is already at its end of life stage, but it was used as it is currently the safest way to ensure that your code runs correctly.

4.1.3 Usage

The initial idea of the project was to provide modular design for the pipeline where the user can cut and paste the technologies to the pipeline making it easy to produce multiple different pipelines for comparison. That is why the there is a initialization script called pipeline.sh. This is a command-line interface that asks user what technologies they want and builds a docker-compose file into build folder with all the necessary files. Sample usage is depicted in the figure 4.2. Only thing user has to do before running the script is to download manually the Kafka Connect component from confluent website and copy its content to a location noted in the project readme. This

step is mandatory because unfortunately Kafka Connect did not have public distribution that would have allowed access without authentication. If this component is missing during the build, the initialization script notifies user about this before continuing further.

After the initialization, using the software is the same as with any other normal docker-compose project. Users can copy the contents of the build folder as the base of their own project or use it as it is. The project can then be build with "docker-compose build" command and run using the command "docker-compose up". After this as there will be over 10 containers running simultaneously, tools such as lazydocker are recommended for proper management.

When the containers have finally finished the startup procedures, every services user interface can be found in their default ports in localhost. The ports needed to be open to this happen, can be found in the docker-compose.yml file. To start analysing data, the user has to only open localhost:8090 where Zeppelin notebooks are running. In the project there exists an example notebook that implements a simple LSTM model. To run this example on spark cluster user has to first manually set two configurations in the spark interpreter menu that can be found under interpreter settings in the user interface. The value "master" must be set to "spark://spark-master:7077" in order to let the zeppelin use the cluster instead of local spark client. Also in the dependencies has to be added "org.apache.commons:commons-lang3:3.6" in order to have common versions in the zeppelin machine and in the cluster. The reason for the manuality of these tasks will be discussed in following sections. After this the user has to only run the cells in the notebook and the result will be saved into HDFS.

To add third party dependencies to the notebooks such as the DL4J library, the "spark.deb" syntax in notebooks is encouraged as this is the best way to preserve the dependencies for possible distribution of the notebook. This way the notebook is as independent of its environment as it can be. It is currently also the only way to preserve these if the container happens to be destroyed. Reasons for this will again be discussed in following sections.

4.2 Adversities during development

Most of the time during development went to trying to integrate these services with each other. The reason for this was mostly features that were not documented well and other mishaps such as bugs that happen when integrating two services that are not that common together. There was a lot of unique problems that did not have any documented solutions on the

internet which made solving them very time consuming. In this section we will be going through the ones that took most time to solve and present the solutions at high level to these problems. More precise technical problems with their solutions can be found in this thesis' git repository that contains all the practical results.

4.2.1 Kafka-HDFS integration

Integrating Kafka topics to HDFS proved to be a non-trivial task. Unlike products like Flume, Kafka does not have any build-in sinks that would allow easy integration of different services. So the developer is left with the task of filling this hole.

The requirements for this component are also non-trivial, as it should be able to scale with the kafka and the HDFS cluster. So in order to not introduce a bottleneck to the pipeline and to do this with time limits that this thesis has, a ready-made solution was necessary. Previously, a good solution to this has been LinkedIn's Camus API, but this was phased to Apache Gobblin in 2015. The problem with Apache Gobblin, however, is that it is not the lightest tool and it can move data only to Hadoop system. That is why if we were to change the HDFS storage to try some other storage option, the Gobblin should be entirely changed to something else.

Technology that solves these problems is Kafka Connect. Kafka Connect acts like a sink plugin to Kafka. Unlike with Gobblin, developers can download only the sink they actually need instead of sinks for every possible case. Because Kafka Connect has sinks for most major big data storages in the market, user is not constrained as much to Hadoop ecosystem like with the Gobblin.

The downside is that Kafka Connect is clearly Confluent product although it is open-sourced and free. Confluent has its own platform that it promotes as open-source event-streaming platform that can be used as a enterprise solution. This is why most of the Kafka Connect documentation is only about how to use it on top of this Confluent platform, although it has capabilities to run independently. This lack of good documentation introduced some obstacles for integrating this component between Kafka and HDFS especially when it came to installing and running the component.

Setting up the Kafka Connect was not that easy as it does not offer much monitoring in itself. This means that as developer you could only see if the whole system worked or did not from HDFS management console. If only some configuration was not right, Connect would only log this as an minor error in the stderr and continue running like nothing had happened although no data was flowing through it. This made debugging the connection surpris-

ingly difficult, but after all the configurations were in place the component worked as it should throughout the rest of the development.

4.2.2 Apache Zeppelin dependency management

Major problems came when the development reached to the point of adding DL4J dependencies to the Zeppelin instance. Zeppelin has multiple ways of adding dependencies and the main way that is defined in its documentation is by its dependency UI that is located in its interpreter settings panel. This again saves these settings into `interpreter.json` file which is a normal JSON configuration file.

At the start of this development in order to remove the manual step of adding dependencies and other settings, we tried to save this interpreter configuration file between container restarts. This worked with every other setting except dependencies. Dependencies would conflict for unknown reasons and the error message would change per restart. Sometimes by mere luck the dependencies would build correctly for a certain amount of time but then fail again. This would be easy to debug if Zeppelin would be a normal maven project as it is. But with ready-made containers there did not exist such an option.

Apache Zeppelin has official container which were used, but this had to be extended in order to have custom spark-submit client because the versions on spark cluster would conflict. When thought afterwards, the whole zeppelin container probably should have been build from scratch as this would've allowed better dependency management with maven.

4.2.3 DL4J sbt project

The DL4J application instance without Zeppelin is build with sbt. sbt is a build tool developed specifically for Scala applications and this why it was also chosen for this project. DL4J's main build tool that the documentation is build around is maven because it is a Java library, but it has its own sbt example in their project repository. This, however, turned out to be highly outdated and caused developing problems.

To have third-party dependencies added to a spark project there are two possible ways to implement this. First one is using the spark-submit scripts flags to define maven repositories. This was first tried but it only lead to dependency conflicts and other bugs with ivy dependency management that the script uses.

The alternative and more manageable way is to build the whole application into so called fat jar, which contains all the code to run the application

including the third-party dependencies. This is what the outdated sbt example did so we opted to update this example to work with newer versions of sbt and sbt-assembly which is a plugin needed for building fat jars on sbt. To make the example project work, we migrated it from sbt version 0.X to sbt version 1.X syntax. The migration of general syntax was relatively easy, but with the newer version of DL4J library, a new merge strategy had to be written for sbt-assembly. Merge strategy is what sbt-assembly uses when it encounters duplicate files during the building process. Writing a updated merge strategy for DL4J was not a trivial task as the strategy that was usually suggested did not work in this case and lead to errors in runtime. After quite a bit of trial and error, we thankfully managed to create a merge strategy that did not throw away necessary files needed at the runtime.

4.2.4 Running the application in Spark

Now we move on to the problems that were really hard to debug when the problems in the previous two subsections were persisting at the same time. When running a Zeppelin task on Spark and the program receives error, Zeppelin correctly logs the stack trace from Spark workers stderr. What it does not do, is logging the the whole stack trace but only the two to three first errors. This made it seem like the error was a some upper level error when instead the actual error that happened in lower part of the stack. This was resolved by noticing the full stack when the application was run on application server.

After this, the next problem was that Intels mkl-dnn library did not seem to be in the classpath of DL4J subdependency library javacpp. This is the code that is responsible of allowing the java code to use high efficiency feature of C++ language. This of course did not have any clear solutions online except some vague conversation logs about possibility of glibc missing in the parent OS and by accident this turned out to be the exact problem in this case too. The Spark worker docker container, which was a result of the Big Data Europe project, was build on top of alpine linux image which has phased out to use musl libc library instead of glibc. There exists tools that can be used to install glibc to alpine instance and we tried this, but after multiple errors in multiple different attempts to install the library, we opted out to fork the docker image to use debian based image. This resolved all the problems that had been until this point and the application finally compiled in worker machine.

The problems after this did not cause the application process on worker to crash but were logged into worker machine while the process continued to run. These were somewhat minor but laborous problems such as allocating

enough disk for `/dev/smh` shared memory implementation and opening and defining right ports for DL4Js parameter server that handled the parameter update logic. Once all of these were fixed, the next problem was that the training process in spark did not finish at all and seemed to go into a state of infinite loop without giving any errors on the driver logs. This seemed like a dead-end until we noticed that a spark job before this had silently failed. It seems that the DL4Js training code was beforehand trying to temporarily save the data into hdfs and this failed as something went wrong. This step in the process could be skipped with configurations and after this the model training process finally succeeded.

The final problem we faced was a randomized index out bounds exception. Raising the number of epochs we noticed that this occurred at random, and unlike the previous errors, this seemed like an error in the library. The code that raises this error handles parallel code so if we were to make conjectures from this it would seem like there would be an error with shared memory access in some part of the asynchronous code. However at this point, the time allocated to this project was already running out so we could not confirm this definitely.

Now that we have seen what we were able to implement during this thesis we move on to evaluate the quality of this software. We will also discuss more about what went well, what went wrong and what could have been done better.

Chapter 5

Evaluating the results

In this final chapter before conclusion, we will be examining the pipeline which we build in the previous chapter. We start with basic performance metrics that were briefly introduced in the section 4.4. Then we move on to qualitative metrics and questions that we also defined in the section 4.4 and discuss about the pipeline using them. We then move to examine general flaws in this study to give more complete image of the results and follow this by reflection on what could have been possibly done better and why some detrimental choices were made. Finally we close this chapter with a section where we briefly list possible improvements that could improve the current pipeline, but could not be implemented with the resource and time restrictions that this project had.

5.1 Performance metrics

We start by defining the environment where the following results were derived. The following metrics were taken in a machine that has a Windows 10 as parent operating system, Intel Core i5 2500K processor and 16GB of DDR3 memory. The docker version that the code was run against was 19.03.2 and it had 3 CPU cores and 10GB of memory as its use. Not all of these resources were needed, but this was done to in order to remove the factors of running out of memory that could cause inexplicable errors. These were tested only in one machine which would usually not be enough to make any kind of actual inferences but because the code is run in a containerized environment the role of underlying hardware becomes less important.

The resource usage in idle state after the initial data has gone through the pipeline and stored into HDFS can be seen in the figure 6.1. The figure is the output of the "docker stats" -command. While inspecting the names

CONTAINER_ID	NAME	CPU %	MEM USAGE / LIMIT	MEM %	NET I/O	BLOCK I/O	PIDS
5d8bf0dc7f3f	mlflow	0.02%	413.6MiB / 9.744GiB	4.15%	2.93kB / 0B	0B / 0B	7
28680599d9bf	spark-master	0.05%	177.9MiB / 9.744GiB	1.78%	18.5kB / 3.51kB	0B / 32.8kB	30
d9efa04f8d31	zeppelin	0.05%	506.2MiB / 9.744GiB	5.07%	173kB / 6.9kB	524kB / 0B	35
002c027f5a0e	analyzer	0.00%	560KiB / 9.744GiB	0.01%	2.43kB / 0B	0B / 0B	2
3ebee8023724	build_zookeeper_1	0.07%	71.9MiB / 9.744GiB	0.72%	89.5kB / 103kB	28.7kB / 324kB	38
fb5cf08f5d4c	build_nodemanager_1	0.27%	240.2MiB / 9.744GiB	2.41%	12.8kB / 31.7kB	13.2MB / 184kB	82
b96cdf8b2780	build_datanode_1	0.17%	207.4MiB / 9.744GiB	2.08%	13.1MB / 156kB	9.1MB / 1.45MB	59
88e52733c29d	build_resourcemanager_1	0.62%	279.8MiB / 9.744GiB	2.88%	43.8kB / 26.5kB	13.8MB / 184kB	228
0b39c0f8e6e7	build_historyserver_1	0.06%	180.4MiB / 9.744GiB	1.81%	3.67kB / 2.09kB	979kB / 1.18MB	47
5f92fab3069	build_namenode_1	0.34%	232.3MiB / 9.744GiB	2.33%	168kB / 85kB	74.5MB / 2.93MB	57
4c91c2f1e3a7	spark-worker-1	0.06%	172.9MiB / 9.744GiB	1.73%	5.62kB / 16.1kB	0B / 32.8kB	29
c9805ee983eb	build_kafka_1	3.37%	319.7MiB / 9.744GiB	3.20%	13.1MB / 15.3MB	602kB / 295kB	70
41a3ad4c2ac4	build_connect_1	0.30%	1010MiB / 9.744GiB	10.12%	15.3MB / 13.3MB	173MB / 0B	40

Figure 5.1: Resource usage in idle state

of each container the reader should be able to map these into corresponding components in the architecture figure in section 5.2, but to clarify the most ambiguous namings, the "analyzer" container refers to the standalone DL4J application and build_connect_1 refers to the Kafka Connect instance.

The total memory usage here was around 3.8GiB. This seems to match with the fact that most of the components ran on top JVM and we did not have the time to optimize each components garbage collection. This is partly because components such as the Kafka Connect, which had a notably large memory usage, would silently fail if not enough memory were offered. Additionally, this would have probably lead to premature optimization. For future improvements on the accessibility of the pipeline, this would be a great place to start.

Processor usage in idle state was minimal and did not have any notable jumps except a bit higher usage for Kafka. The other statistics did not have any notable deviations that would affect greatly on the usage of the pipeline. Minor thing to note is that the DL4J application instance is not running actively during these results as user usually are not using both it and the zeppelin instances simultaneously.

CONTAINER_ID	NAME	CPU %	MEM USAGE / LIMIT	MEM %	NET I/O	BLOCK I/O	PIDS
5d8bf0dc7f3f	mlflow	0.01%	411.7MiB / 9.744GiB	4.13%	3.51kB / 0B	0B / 0B	7
28680599d9bf	spark-master	0.04%	186.2MiB / 9.744GiB	1.87%	76.5kB / 15.3kB	332kB / 32.8kB	32
d9efa04f8d31	zeppelin	114.77%	2.049GiB / 9.744GiB	21.03%	963MB / 1.03GB	125MB / 1.56MB	442
002c027f5a0e	analyzer	0.00%	540KiB / 9.744GiB	0.01%	3kB / 0B	0B / 0B	2
3ebee8023724	build_zookeeper_1	0.12%	75.2MiB / 9.744GiB	0.75%	132kB / 128kB	86kB / 324kB	38
fb5cf08f5d4c	build_nodemanager_1	0.27%	242.9MiB / 9.744GiB	2.43%	76kB / 241kB	13.2MB / 184kB	82
b96cdf8b2780	build_datanode_1	1.58%	228.9MiB / 9.744GiB	2.29%	16.6MB / 61MB	10.7MB / 1.45MB	66
88e52733c29d	build_resourcemanager_1	0.49%	281.6MiB / 9.744GiB	2.82%	254kB / 89.1kB	13.8MB / 184kB	227
0b39c0f8e6e7	build_historyserver_1	0.09%	181MiB / 9.744GiB	1.81%	4.24kB / 2.09kB	1.09MB / 1.18MB	47
5f92fab3069	build_namenode_1	0.92%	237MB / 9.744GiB	2.38%	538kB / 381kB	74.7MB / 3.9MB	57
4c91c2f1e3a7	spark-worker-1	53.06%	899.9MiB / 9.744GiB	9.02%	1.06GB / 7.98MB	27.2MB / 1.36GB	164
c9805ee983eb	build_kafka_1	0.62%	323MiB / 9.744GiB	3.24%	13.4MB / 15.6MB	602kB / 856kB	70
41a3ad4c2ac4	build_connect_1	0.19%	1016MiB / 9.744GiB	10.18%	15.6MB / 13.7MB	173MB / 0B	39

Figure 5.2: Resource usage in preprocessing state

In figure 6.2 we have the same statistics while the zeppelin code is trying to preprocess the data in order to feed it into LSTM network. The pipeline has at this point run for some time and the overall memory load has increased to 6.1GiB. This is mainly due to Zeppelin and Spark worker using a lot more memory. The CPU usage for both has raised, and one can see that at the moment of taking these stats the notebook code is actually using more resources than the actual worker.

As there was not enough time to test other technologies in pipeline these results by themselves do not give us much information about the quality of the software. They do give us some information about the requirements of the machine that is needed to run locally a pipeline that is this complex. Over 8GiB of memory is almost a must to run it in active use where the data is flowing continously. The need for memory can be alleviated with proper garbage collection, but peak memory usage can reach to these numbers.

5.2 Pipeline management

The main benefits of this pipeline do not come from its performance. It is highly containerized meaning it can run with little to none re-configuration on any machine. During the development of this pipeline the pipeline was developed first on Mac OS and then because of memory requirements moved to Windows machine. The porting of code did not require any extra steps, other than making sure that the formatting of line endings stayed the same with script files that are run in containerized linux environment.

Docker containers also allow the development in isolated networked environment that is very similar to one which would be used in production. This means that most of the problems that could occur when rolling locally developed code into a production environment where services lie on different machine interconnected by a network, are already solved in here before the user has even started development. Only thing that the user has to do is to change the addresses in this code from docker based DNS to the addresses used in their own cloud environment.

As explained earlier, the technologies such as Kafka Connect have been chosen so that changing one technology does not have large impacts on the overall architecture. If the user would want to change some component in the pipeline, the integrations between components have to be done again but this has been made as easy as possible.

The pipeline is also mostly automated. No manual input is required from the user to setup the basic pipeline which makes adapting the pipeline quite trivial. Only in the analysis phase, users manual input is needed to control the pipeline.

Machine Learning pipeline monitoring is made easy with MLFlow. Although logging new experiments needs a bit of manual work for user because of the nature of MLFlow API, but nevertheless it does make managing the machine learning models easier. The only downside here is that by default MLFlow does not yet accept DL4J models to be saved in its containerized packages which would make model export easier.

For general monitoring the pipeline is still lacking because we did not have time to implement the ELK stack. Almost each component has its own monitoring UI but these are scattered making it hard to follow the bigger picture of the pipeline which was one of the problems we introduced in the chapter 2. The exception of not having an own monitoring UI is the Kafka Connect. It offers REST API to monitor its metrics which would need additional integration with some log monitoring solution such as ELK.

5.3 Flaws in the end product

In this section we examine general flaws in the pipeline that we are aware of. We try to objectively examine these and try to give some reasoning why these flaws did happen and how they could be fixed. We will leave out here the flaws in monitoring as these were mostly examined in the previous section.

Currently, the main method of training machine learning models is to use computers graphics card (GPU) which is multiple times faster than using CPUs. This makes the current pipeline vastly underperforming when it comes to training machine learning models and this is a clear downside for this system. It needs quite a bit of CPUs and good networking before it even reaches the performance of a single high-end GPU. The DL4J, does however, support GPU training on spark cluster meaning that adding GPU support to this pipeline with small amount of work is possible. However, we did not have time to implement it as it needs some work for docker to access CUDA resources on parent OS.

Unfortunately, because we were only able to try the HDFS as the storage method, the storage is currently lacking a proper query properties meaning that the user has to pull the entire timeseries for each stock. This in big data context is not desirable behaviour, but integrating a technology like HBase should not need much additional work. However, for us it is already out of the scope of this project.

With the resources this project had, the pipeline could not be tested with actual complete stock market data, which leaves means that we have no prove that the pipeline actually can scale with the data or if this scaling is even possible with the pipeline. We can only speculate this based on the each of technologies abilities but the possibility for horizontal scalability should be there theoretically.

5.4 Reflection

In this section we will reflect back on the implementation part and discuss what possibly went wrong on the development of the pipeline. We will try our best to explain the choices taken while being aware that there could have been better alternatives.

Choosing DL4J as the machine learning library turned out to be a time consuming option. After the implementation part, it was clear why data scientist usually use python libraries such as Keras and PyTorch, the amount of resources is substantially larger. Many of the problems faced during development had quite simple solutions to them, but they were really hard to solve without any reference and empty search results on Google. This combined with the out-of-date examples lead to problems seen in the previous chapter.

One can question whether we should have opted out from DL4J during development and changed it to some other more compatible library. We did not do this because there was not that many sensible choices and the documentation and resources online of those that seemed promising were even worse than with the DL4J. But although the original plan to build to multiple pipelines was not possible partly because of this, we think that the information produced in the implementation chapter is quite valuable on its own and can help a lot of novice data scientists who have restrictions on what technologies they can use.

Of course, the scope of the implementation part could have been narrowed down as it was quite ambiguous. Especially, during development things that seemed simple turned out to be quite more complex than initially thought. The inexperience about some of the technologies used did affect the results, but it also indicates that the documentation of some technologies could have been better. Quite a bit of things were assumed in the documentations, which lead to errors that probably could not have happened with right documentation.

5.5 Potential improvements

There are quite a bit of improvements already mentioned at the previous sections as these usually are part of the problems, but we have gathered them and a couple others to this section to give a clearer view on what could be improved. These are quite straightforward to implement but implementing them can take some time.

As stated before, the performance of the pipeline could be improved. This could be done by adding support for GPU usage in Spark cluster which

would require the configuration of docker to give the process access to these. Memory usage could be probably optimized by configuring better garbage collection on components that take most of the memory.

To make the pipeline more automated and improve its mobility, urls for servers could be made dynamic with for example environmental variables that could be configured by the user. Other such configurations could be made in order to make porting the code to a different environment easier.

The final giveaway in this chapter is that there was quite a bit of things that could have been done better. This is partly because of the ambiguous scope of the project, inexperience of the author and lack of correct information online. Despite these, we think that the complications during the development and the final product do offer valuable insights on developing a pipeline for stock data analysis using these technologies and can give somebody a basis to start developing their own pipeline. Now we move on to the final chapter where we conclude and summarize everything that we have learned during this thesis.

Chapter 6

Conclusions

In this thesis, we examined big data pipelines to analyze stock market data. We conducted literary research on current trends on methods that are used to analyze stock market data. From this we saw that deep learning models are currently the driving force in stock market analysis. Other statistic methods are also used with conjunction of data from other varying sources. These methods are not only used to predict prices in hopes for profit, but can also be used to analyze e.g causalities in other phenomenons.

We researched what are the technologies currently used in pipelines that analyze stock market data covering both academic and industrial use and saw that public information about these is quite limited. With the information publicly available, we saw that usually only the analysis phase was reported and other aspects of the systems were dismissed. These other aspects were most of the time monitoring of the system, but also ingestion and storage were not reported in many cases.

We did a literary study on technologies that were used in the existing pipelines and technologies that seemed promising to be used in these pipelines. These included technologies for ingestion, storage, monitoring and analysis parts of the pipeline and we saw that the analysis part is dominated by Spark solutions at the moment of writing. In ingestion and storage there exists a lot of alternatives but none of them dominated field because of different use cases. For our pipeline, where we wanted the pipeline to be as flexible as possible we opted to Kafka for ingestion and HDFS for storage.

We planned and implemented a basic stock data analysis pipeline based on the results from the previous literary studies. Our goal was to build a pipeline that would make analysis of stock market data in big data environment more accessible to data scientist who do not have that much knowledge on big data system development. The pipeline we built is open-source that anyone can use as a base or use it as it is.

We were not able to build multiple pipelines for comparative results because of ambiguous scope and problems with the technologies used, but we managed to produce a lot of valuable practical information about developing these systems. Our pipeline did manage to fulfil most of the requirements it was given at the start of the development, but fields such as good monitoring are still lacking.

We can say that building a machine learning pipeline in a big data environment is not trivial yet. A lot of information is scattered, missing or going out of date as newer versions and technologies are introduced. Building a big data pipeline requires integration of multiple technologies and because of the amount of the possible integrations, there exists pairs of technologies that are not that compatible yet. This thesis was able to test only a couple of these possible integrations but the ones that were tested did not integrate with each other without challenges.

In the future, hopefully, more information can come available if and when these technologies gain popularity and the role of big data grows. Until then we have to cope with the information available and just try to produce more information to make the field more accessible for new novice data scientists. This thesis hopefully alleviates someones process, in the field of stock data analysis.

Bibliography

- [1] Alpha vantage documentation. <https://www.alphavantage.co/> Accessed: 23.09.19.
- [2] Apache flume documentation. <https://flume.apache.org/> Accessed: 16.06.19.
- [3] Apache kafka documentation. <https://kafka.apache.org/> Accessed: 16.06.19.
- [4] Apache nifi documentation. <https://nifi.apache.org/docs.html> Accessed: 20.06.19.
- [5] Apache spark documentation. <https://spark.apache.org/docs/latest/index.html> Accessed: 26.06.19.
- [6] Apache systemml documentation. <http://apache.github.io/systemml/index.html> Accessed 30.06.19.
- [7] Aws: How to deploy docker containers. <https://aws.amazon.com/getting-started/tutorials/deploy-docker-containers/> Accessed: 23.09.19.
- [8] Barchart free market apis. <https://www.barchart.com/ondemand/free-market-data-api> Accessed 23.09.19.
- [9] Big data europe project. <https://www.big-data-europe.eu/> Accessed 25.09.19.
- [10] Deeplearning4j documentation. <https://deeplearning4j.org/docs/latest/> Accessed: 30.06.19.
- [11] Elastic stack documentation. <https://www.elastic.co/elk-stack> Accessed: 30.06.19.

- [12] Gcp: Quickstart for docker. <https://cloud.google.com/cloud-build/docs/quickstart-docker> Accessed: 23.09.19.
- [13] Hdfs documentation. https://hadoop.apache.org/docs/r1.2.1/hdfs_design.html Accessed: 22.06.19.
- [14] Iex trading documentation. <https://iextrading.com/> Accessed 15.04.19.
- [15] Intrinio documentation. <https://intrinio.com/> Accessed 23.09.19.
- [16] Mlflow documentation. <https://mlflow.org/docs/latest/index.html> Accessed: 30.06.19.
- [17] Quandl us eod data documentation. <https://www.quandl.com/data/EOD-End-of-Day-US-Stock-Prices> Accessed 23.09.19.
- [18] World trading data documentation. <https://www.worldtradingdata.com/> Accessed 23.09.19.
- [19] AMIRGHODSI, S., ALLA, S., KARIM, M. R., AND KIENZLER, R. *Apache Spark 2: Data Processing and Real-Time Analytics*. Packt Publishing, 2018.
- [20] BARTRAM, S. M., AND GRINBLATT, M. Agnostic Fundamental Analysis Works. *Journal of Financial Economics (JRE)* (June 20 2017).
- [21] 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.
- [22] CHEN, J. *Essentials of Technical Analysis for Financial Markets*. John Wiley & Sons Inc, 2010.
- [23] 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.
- [24] 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.

- [25] 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).
- [26] 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).
- [27] 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).
- [28] DRONAVALLI, A. Why apache cassandra is not good for timeseries data & analytics. <https://medium.com/@abhidrona/why-apache-cassandra-is-not-good-for-timeseries-data-analytics-a1d65a369048> Accessed 24.09.19.
- [29] FISCHER, T., AND KRAUSS, C. Deep learning with long short-term memory networks for financial market predictions. *European Journal of Operational Research* (2017).
- [30] FOX, J. *Myth of the Rational Market*. Harper Business, 2009.
- [31] GEORGE, L. *HBase: The Definitive Guide, 2nd Edition*. Packt Publishing, 2018.
- [32] 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).
- [33] HOFFMAN, S. *Apache Flume: Distributed Log Collection for Hadoop - Second Edition*. Packt Publishing, 2015.
- [34] 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).
- [35] JOHN L. PERSON. *Mastering the Stock Market: High Probability Market Timing and Stock Selection Tools*. John Wiley & Sons, Incorporated, 2013.

- [36] 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).
- [37] 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.
- [38] 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).
- [39] 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).
- [40] 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).
- [41] 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.
- [42] MUN, F. W. Big data, small pickings: Predicting the stock market with google trends. *The Journal of Index Investing* 7 (2017).
- [43] MURPHY, J. J. *Technical analysis financial markets: A comprehensive guide to trading methods and applications*. New York Institute of Finance, 1999.
- [44] 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).
- [45] NEERAJ, N., MALEPATI, T., AND PLOETZ, A. *Mastering Apache Cassandra 3.x - Third Edition*. Packt Publishing, 2018.
- [46] NICOLAS, P. R., MANIVANNAN, A., AND BUGNION, P. *Scala: Guide for data science professionals*, 2017.
- [47] 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.

- [48] PENG, Z. Stocks analysis and prediction using big data analytics. *International Conference on Intelligent Transportation, Big Data & Smart City* (2019).
- [49] 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).
- [50] 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.
- [51] SEZER, O. B., OZBAYOGLU, A. M., AND DOGDU, E. An artificial neural network-based stock trading system using technical analysis and big data framework.
- [52] 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).
- [53] SNIVELY, B. AWS re:Invent 2018: Big Data Analytics Architectural Patterns & Best Practices. Available at: <https://youtu.be/ovPheIbY7U8> Accessed: 08.05.19.
- [54] UTTHAMMAJAI, K., AND LEESUTTHIPORNCHAI, P. Association mining on stock index indicators. *International Journal of Computer and Communication Engineering* 4 (2015).
- [55] VARTAK, M., SUBRAMANYAM, H., LEE, W.-E., VISWANATHAN, S., HUSNOO, S., MADDEN, S., AND ZAHARIA, M. Modeldb: A system for machine learning model management.
- [56] WANG, W. A big data framework for stock price forecasting using fuzzy time series. *Multimedia Tools and Applications* 77 (2018).
- [57] WORLD FEDERATION OF EXCHANGES. Monthly report january 2019. <https://www.world-exchanges.org/our-work/statistics> Accessed 15.04.19.
- [58] YARABARLA, S. *Learning Apache Cassandra - Second Edition*. Packt Publishing, 2017.

- [59] 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).
- [60] YUAN, D. Stream processing in uber. <https://www.infoq.com/presentations/uber-stream-processing/> Accessed: 19.06.19.