

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, December 30, 2019

DRAFT! — December 30, 2019 — DRAFT!

Supervisor: Professor Hong-Linh Truong
Advisor: Professor Hong-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:	December 30, 2019	Pages:	68
Major:	Computer Science	Code:	SCI3042
Supervisor:	Professor Hong-Linh Truong		
Advisor:	Professor Hong-Linh Truong		
<p>Stock markets have been a point of interest for researchers for a long time. 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 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 Unfortunately, the information about these tools in the context of stock data is scattered, somewhat outdated and hard to find, which can drive away potential valuable innovations. The goal of this thesis is to provide information that can make this analysis more accessible to novice 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 proposes a method that could help novice data scientists to build their own pipeline. Our study finds that using this method, the resulting pipeline can be built cost-efficiently while keeping the pipeline sustainable and easy to reproduce. However, the price of this is that the pipeline still runs on singular hardware without promises of scaling with the data. The sample size of our study is also small to make more viable claims, because of time constraints and more research would be needed on the subject.</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:	30. joulukuuta 2019	Sivumäärä:	68
Pääaine:	Tietotekniikka	Koodi:	SCI3042
Valvoja:	Professori Hong-Linh Truong		
Ohjaaja:	Professori Hong-Linh Truong		
<p>Maailman talous pyörii osakemarkkinoiden ympärillä. Informaatio, jota osake-markkina tarjoaa ei pelkästään kerro talouden nykytilasta vaan se myös heijas-taa kaikkia osakemarkkinoihin vaikuttavia ilmiöitä. Tämän takia uusia tutkimuk-sia tehdään jatkuvasti jotta sekä taloutta että näitä ilmiöitä voitaisiin ymmärtää paremmin.</p> <p>Osakemarkkinoiden tuottaman datan määrä vuodessa on teratavujen mitta-luokassa. Jotta tätä dataa voitaisiin siis tutkia, tarvitaan työkaluja jotka pys-tyvät toimimaan tämän mittaluokan kanssa. Valitettavasti informaatio näistä työkaluista osake datan kontekstissa on hajanaista, jossain määrin vanhentu-nutta ja vaikeasti löydettävissä, mikä jossain tapauksissa estää mahdollisesti tärkeän analyysiin tekemisen. Tämän diplomityön tarkoituksena on tarjota tie-toa ja työkaluja, jotta tämän tyyppinen tutkimus olisi mahdollisimman help-popääsyistä aloitteleville 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 ehdotamme menetelmää, jota alottelevat tutkijat voivat hyödyntää aloittaessaan rakentamaan omia järjestelmiään. Huomaamme, että tällä menetelmällä on mahdollista toteuttaa, ilman suuria kuluja, kestävä ja helposti toistettava järjestelmä osakedatan analysointiin, mutta tämän haitta-puolena on, että järjestelmä ajetaan vielä yhden fyysisen koneen päällä antamatta mitään lupauksia skaalautuvuudesta. Tutkimuksemme otos on myös valitettavan pieni, jotta kunnon johtopäätöksiä voitaisiin vetää, johtuen työlle annetuista ra-joitteista ja tämän takia tarvitsisikin enemmän tutkimusta aiheesta.</p>			
Asiasanat:	osakkeet, big data, pilvilaskenta		
Kieli:	Englanti		

Espoo, December 30, 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
LSTM	Long Short-Term Memory
MFF	Multilayered Feed Forward Network
ELK	Elasticsearch Logstash Kibana
RNN	Recurrent Neural Network
CNN	Convolutional Neural Network
HDFS	Hadoop Distributed File System
QoS	Quality of Service
DL4J	Deeplearning4j Machine learning library
sbt	Interactive build tool for Scala
CSV	Comma-Separated Values
JSON	JavaScript Object Notation
CAP theorem	Theorem that states that no system can guarantee consistency, availability and partition tolerance at the same time

Contents

Abbreviations and Acronyms	5
1 Introduction	8
1.1 Application scenario	10
1.2 Audience and Research Questions	12
1.3 Structure of the Thesis	13
2 Background	15
2.1 Stock market analysis	15
2.1.1 Stock price prediction	16
2.1.1.1 Statistical methods	16
2.1.1.2 Machine learning methods	17
2.1.2 Existing pipelines	19
2.1.2.1 Data sources	20
2.1.2.2 Used technologies	20
2.1.2.3 Monitoring	21
2.2 Possible Big Data Technologies	22
2.2.1 Data ingestion	22
2.2.1.1 Apache Flume	23
2.2.1.2 Apache Kafka	24
2.2.1.3 Apache NiFi	25
2.2.2 Data storage	25
2.2.2.1 Hadoop distributed file system	25
2.2.2.2 Apache HBase	26
2.2.2.3 Apache Cassandra	27
2.2.3 Machine learning frameworks	27
2.2.3.1 Apache Spark	28
2.2.3.2 Deeplearning4J	29
2.2.3.3 Apache SystemML	29
2.2.4 Monitoring	30
2.2.4.1 Log monitoring	30

2.2.4.2	Machine learning monitoring	31
3	Requirements	32
3.1	Novice data scientist	32
3.2	Inaccessibility of big data stock analysis	33
3.2.1	Obstacles in stock data analysis	34
3.2.2	Obstacles in big data frameworks	35
3.2.3	Requirements for the methology	36
4	Proposed methology	38
4.1	Overview	38
4.1.1	Data source	39
4.1.2	Main programming language	41
4.1.3	Technology selection	41
4.1.4	Implement and share	44
4.2	Use case	45
4.2.1	Description	45
4.2.2	Analysis of results	46
5	Use case	48
5.1	Final Architecture	48
5.2	Usage	51
5.3	Adversities during development	52
5.3.1	Kafka-HDFS integration	53
5.3.2	Apache Zeppelin dependency management	54
5.3.3	DL4J sbt project	54
5.3.4	Running the application in Spark	55
6	Discussion	57
6.1	Method evaluation	57
6.1.1	Advantages	58
6.1.2	Disadvantages	59
6.2	Future improvements	59
7	Conclusions	60

Chapter 1

Introduction

Stock markets have been a point of interest for researchers for a long time. 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 [42].

When talking about stock analysis, people usually connect this with making profit the most optimal way. However, stock analysis can be much more than this. Because of the nature of stock market, where buyers define the actual price of the individual stock, stock market reflects a lot more information than one would think. This hidden information can help us understand a lot of other phenomenons that can influence stock market and this understanding can help us develop new innovations.

This reflected information also works the other way. Information from other sources such as news and social media affect the stock market and both can be used to analyse it. This is partly why the role of big data in stock data analysis has grown tremendously in the previous years.

Unfortunately, because of the monetary gains one can get from analysing stock data, publicly available information of the state of art stock analysis is quite hard to get. This kind of information is usually kept private by companies that gain from them. Innovating new ways to analyse stock is

becoming harder and harder as breaking into the field needs more and more resources.

In this thesis, we try to alleviate the effort needed for novices start their learning process in this field and try to enable more smaller scale individuals to get in this field. But first, in order to get a better understanding of how complex these systems can be, we are going to look stock analysis through a fictional but realistic application scenario. We would like to note that, this thesis will be focusing only on the part of the depicted application, but the idea of the following section is to give the reader better overall view of the underlying domain.

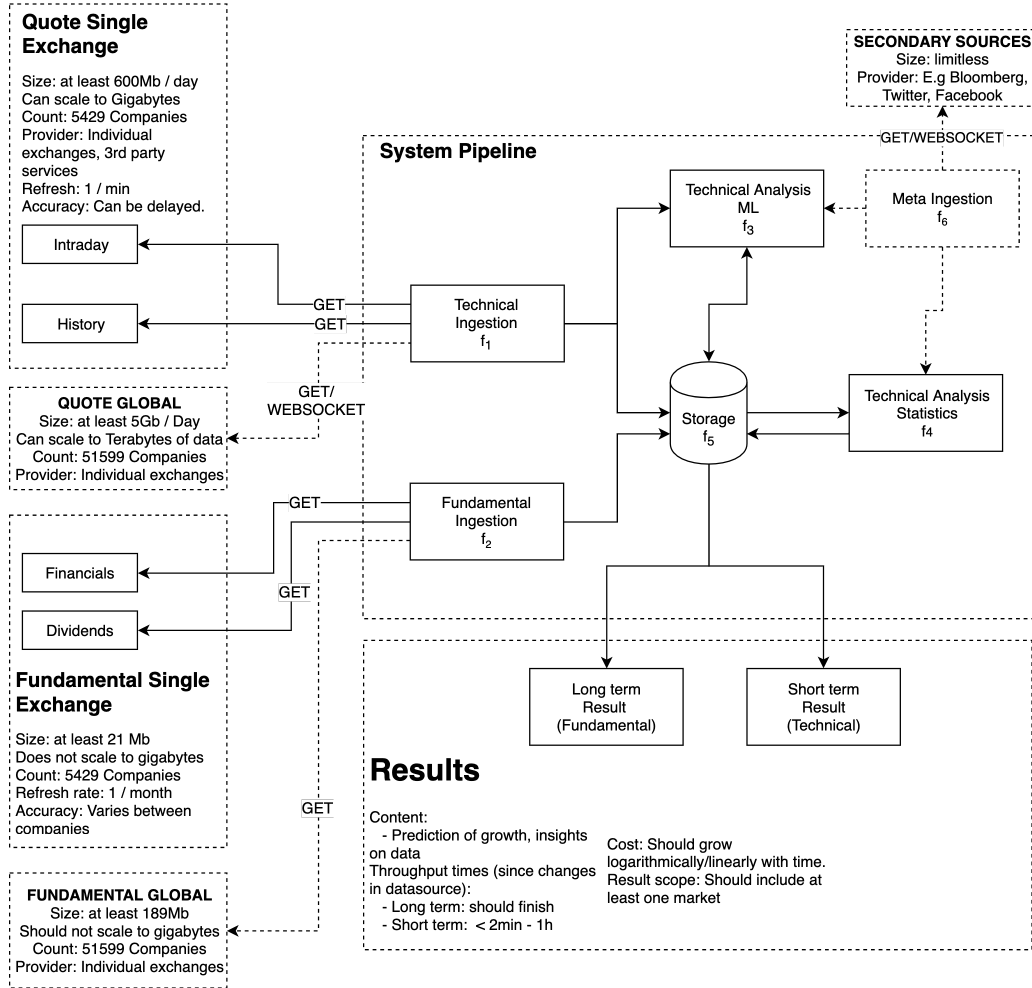


Figure 1.1: Example of a stock data pipeline

1.1 Application scenario

In order to understand the system depicted above, we need to understand couple of key terms. Stock analysis can be divided into two overlapping categories: fundamental analysis and technical analysis. Fundamental analysis is based on the idea that each stock has an intrinsic value that can be larger than the actual price of the stock in the market and buying these will eventually lead to profits [24]. Fundamental analysis is usually based on financial metrics and from a computational perspective the amount of data needed is not that large. The technical analysis, on the other hand, focuses on the real-time market values and needs new data constantly. This the kind of analysis where the amount of data can grow tremendously. The technical analysis focuses on finding recognizable patterns through this data [51].

Somewhat normal system for analysing both fundamental and technical analysis is presented in Figure 1.1. 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 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 1.1, 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 Figure 1.1, 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 the maximum size of this extendable data would be 5Gb per day. This 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 [70]. 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 [73], [62] and [50]. 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 1.1) 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 systems, 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 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, such as [34], [38] and [69], 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 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. This definition is not based on any existing literature and is only for clarity when we examine our target group. This thesis could be useful for the reader 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 a place to find knowledge easily on how to get started.

As the system can be very complex as seen in the previous section, we will be focusing on 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 and the data can come from multiple different sources. 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 more accessible 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 and an example case with solutions that can be used to navigate through these different problems.

The research questions that this thesis tries to answer are: What is the current state of art in big data stock market analysis?, What makes stock analysis possibly inaccessible for novice data scientists? and How to build a pipeline that novice data scientist can use or benefit from?

For the first research question, we will be conducting a literary research on both academic and industrial stock market analysis using big data methods. With the second research question, we want to understand what are the problems with the current environment for big data analytics on stock data for novice data scientists. For this, we will be conducting literary study on the subject and we will analyse the results of this.

For the third and final research questions we will be building method that novice data scientists can use to build their own pipelines and we will analyze how this method works through a specific example case. We will be focusing on the way the pipeline is build and what are the challenges that we face during it that could gravely affect novice data scientists.

1.3 Structure of the Thesis

We will start this thesis by examining the necessary background information in order to understand better the underlying domain and answer our first research question. 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.

After this we move on to solve the second research question "What makes stock analysis possibly inaccessible?" and derive the requirements for our proposed method through this question. We will define more precisely to whom this thesis is aimed at and examine what kind of challenges these people will face when starting big data analytics on stock market data. Using this and the knowledge from the previous background chapter, we will derive the requirements for our method.

We will then define our key contribution that we think that would fulfill these requirements and present reasoning behind the choices we have made. With this we will conclude the theoretical part of the thesis and start the empirical part. In chapter 4.2, we will start by defining a special use case which we will be using to test this method in practice. In the following chapters, we will report the process and end result of using this method.

After this we will move to analyse this test on what are the pros and cons on using this method in practice. We will try our best to validate the claims that we make about the example case and try to 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, to finish the thesis, we will have a conclusion chapter that summarizes the results of the thesis.

Chapter 2

Background

In this first chapter, we are going to look at the current stock market data analysis in big data environments. The goal of this chapter is to alleviate the problem with the scattered data of stock data analysis by aggregating information about the state of art stock market analysis. The other goal of this chapter is to introduce background information needed in the following chapters where we will be introducing the empirical part of this thesis to help reason the choices we make later. If the reader is familiar with these subjects we recommend moving to the next chapter as this chapter mostly aggregates existing information.

This chapter is divided into two 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 [41].

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 [28].

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 [35]. 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 [52]. 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 [35].

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 [65]. 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 [29]. 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 predict prices [60].

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. Whereas 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 indicator's bias when using multiple different indicators [29].

Moving away from the momentum indicators, Autoregressive integrated moving average (ARIMA) models have also been a successful way of analysing time series stock data and they have been also used with big data [68]. However, with the new advances in machine learning, there seems to be better performance to be achieved with machine learning methods [43]. 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 a couple of recent examples of the usage of this kind of data: Day et al. [35] tried to link oil prices to stock prices using global financial data streams with stochastic oscillator techniques. Kyo [45] combined technical and fundamental analysis by using a regression 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 methodology that has gained a lot of interest. Both supervised [60] and unsupervised learning techniques [34] have 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 [58].

Although MFF is seemingly the most popular, better results in price predictions have been able to achieve with convolutional neural networks (CNN) and long short-term memory (LSTM) [58]. 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 [58]. 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 [58]. 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 [58].

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 whatever the price is changing more or less than median price [36]. 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 [60].

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 [47]. 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 24 different stock data analysis pipelines that handle or are able to handle big data. [68] [35] [52] [23] [41] [46] [48] [64] [30] [56] [39] [47] [33] [34] [60] [61] [75] [69] [72] [38] [32] [55] [63] [31] 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 consists 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. [41] [34] [46] [60] 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 [49]. 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 of the ones that did report their ingestion was found to be custom scripts. [52] [46] [48] [30] [34] [60] [69] [38] 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 [56] and

[33]. Another framework that is commonly reported in the ingestion step is the Apache Kafka streaming platform [47] [31]. 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 [55].

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.

No pipeline reported a usage of mixing cloud platform products with open-source products. From the cloud platform products, there are information on usage of Amazons S3 and Googles BigTable with BigQuery [63] [55]. In the open-source side HBase [39] 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 [41] [30] [47] [34]. Where Apache Hive and Apache Pig were previously most used technologies, now Spark seems to be taking their place. [63]

In the industry side, there is indications of Apache Storm being used with real-time classification [31]. 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 [44].

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 data with machine learning methods, the quality of data can be somewhat important [59]. 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 section, 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 Possible 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 [40]. 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 [40].

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 [40].

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. [40]

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 [57]. After being released from incubation in 2012 many other big industrial companies such as eBay and Uber [74] 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 [57].

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 [26].

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 [37].

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 [37].

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 [37].

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) [37].

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 [71].

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 [53].

Unlike HBase, Cassandra does not guarantee the consistency of data due to CAP theorem but instead guarantees the availability at all times. Cassandra implements "tunably consistent" paradigm where user can define for each write/read to be consistent with the cost of availability. This high availability, makes Cassandra better choise for example web application where the data must be available at all times, but the data doesn't necessary have to be up to date [53].

2.2.3 Machine learning frameworks

When we researched existing pipelines we noticed that machine learning libraries such as Tensorflow and Theano are still quite used even when data can scale to enormous amounts [41] [48] [46]. 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 accross 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 [22].

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

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 preprocessing 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 Elasticsearch Logstash Kibana (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 [12].

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 [12].

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, a pure ELK stack is still very popular option to handle log data [12].

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 [66].

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 [17].

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 [17].

That concludes the needed background information needed to build a pipeline for historical stock market analysis using big data. In the next chapter, we will start the empirical part of this thesis by examining more closely who this thesis is for and defining the requirements for our implementation part based on the results from this and the following chapter.

Chapter 3

Requirements

In this chapter we examine the problems that a novice data scientist faces when starting their process to implement big data stock market pipeline. In this chapter we will expand the problem introduced in the introduction chapter and build the requirements for the methodology we will propose in the following chapters. We will start by defining the target group that these problems affect the most. These problems are not universal and to understand better to whom this thesis is benefitting the most we will be defining the target audience and the reasoning why would they be interested in the subject of this thesis. After this we will be defining the problem itself by dividing it to smaller parts and examining these parts individually turning them to the needed requirements.

3.1 Novice data scientist

As stated in the introduction, this thesis is meant for novice data scientists that have little to no experience on developing big data systems, but have general information on data analysis and want to use state of art methods to analyse stock market data. We use the term data scientist as it is usually used when talked about a person that does data analysis with enormous amount of data. This is in contrary to data analyst, which is in the same data analysis domain, but works with smaller amounts of data where one does not need to worry about optimization of calculations when the data is processed [67]. Because of this optimization aspect, data scientist must have knowledge not only about complex data analysis methods but also knowledge about developing and programming systems that scale. They usually also do not have much experience in "DevOps". The term "DevOps" can be defined as "...a set of practices intended to reduce the time between committing

a change to a system and the change being placed into normal production, while ensuring high quality.” [25]. This is usually not a point of consideration when developing ad-hoc data analysis systems on a small scale but when the system is used for a longer time in much bigger scale this becomes important [25].

There is still a lot of hidden potential when it comes to stock market. With the arise of the amount of data from social medias and other new sources, there are new ways to analyse the stock data and at the same time there are new uses for the stock market data itself to understand these new sources of data. There is a lot of potential for new innovations when it comes to this data. Companies who want to benefit from these kind of innovations have their own highly professional teams that have resources to analyze this kind of data in large scale. However, the situation is different for smaller startups, which might have great ideas concerning the field but no resources to implement them. This is where the need for novice data scientists is and this is where one can find persons who benefit from this thesis the most.

As there exists a lot companies that have products to analyze stock data why would the reader then be building their own pipeline instead of using some ready-made product. Stock market analysis is a subject that has been researched for decades now and because of this there exists already a lot ready made tools for it which those who have enough money can use [16] [20]. However, with the rise of modern machine learning as the de facto way to conduct stock market analysis, these tools have not grown with this and are usually only made for statisticians or have very limiting capabilities concerning machine learning. This makes them not suitable for any data scientist to use who wants to try the latest methods and experiment ideas that have not been done before.

Other good question is that why would the reader then be interested in using methods to analyse the big data instead of using traditional methods that the data analyst would use. As the stock market data has been researched for decades, there is already a lot of knowledge about the stock analysis using the timeseries data with statistics. This is why the current interest is in using big data with machine learning to learn new things about this field which has been researched a lot.

3.2 Inaccessibility of big data stock analysis

Now that we have expanded our typical novice data scientist term and examined the reasons why would they be interested in the field of stock data in big data context, next we are going to examine what are the obstacles

these types of people face when trying to start data analysis on stock market data using state of art methods. These are not problems that everyone who works in the field faces or do they mean that there is something majorly wrong with the current methods. These are more of a problems that can make it harder for new poeple in the field to get their ideas heard. We will use these problems to define the requirements for our proposed methology in the following chapters.

The problems can be seen to stem from three main factors; the information needed is scattered, the information needed is outdated and the needed information is lacking all together. We start by examining the obstacles in basic stock data analysis and then move to the technical side and examine obstacles in big data technologies to implement this analysis. Finally we end this chapter with a section that using this knowledge aggregates requirements for the system.

3.2.1 Obstacles in stock data analysis

As stated before stock data analysis is a subject that has been researched a lot and there exists a lot of materials on the subject that anybody can obtain. However, the material that is publicly available can be quite outdated compared to the state of art research. Much of the research is can be assumed to be carried by private companies which keep their findings as market secrets which is understandable taking into account their monetary value, but this makes it hard to not invent the wheel over and over again.

The papers published by researchers give us a better look at how, for example, machine learning is used in the stock market analysis. The problem with these papers is usually that they report their findings on a very high level usually focusing on theoretical side and keep their actual implementations private [46] [34] [41]. So the papers do provide valuable information on the subject but usually leave out important information how to reproduce the models which is vital information especially for the target audience of this thesis.

The problem is somewhat similar in the industrial side, although where in academic papers the theoretical side of algorithms is usually well explained, in industrial public information this side is usually left out. The information about industrial pipelines is very limited and usually the information available focuses on promoting some product that the pipeline uses instead of the pipeline itself [55] [63]. For novice data scientist this gives a climpse on how the industrial pipelines are implemented, but leaves a lot of details out on what are the algorithms these pipelines implement on data and how they are implemented in practice.

Availability of the actual stock data is also a somewhat of a problem. During the rise of internet, stock market data has been publicly available through services such as Yahoo Finance and Google Finance. However, as the price of the stock data has grown over the years both of these services have been shut down [49]. There exists some services which offer the same types of data, but there are still many papers quote the Yahoo Finance as their data source although it has been 2 years since it was shut down [61] [46]. So the needed stock data can be also hard to obtain.

3.2.2 Obstacles in big data frameworks

The situation of available information is much better on the big data tooling side. There are a lot of open-source solutions for many different use cases and these are usually very tested as they are. The problems arise when we have to filter the right technologies for the current domain as the amount of information is large. If one finally finds the right technologies to their use case, they are faced with a challenge to integrate and run these tools. In this section we are going to examine these problems.

For a novice data scientist the number of possible technologies you can use for stock data analysis can be overwhelming. The reason for this is that as there are a lot of technologies which all have different methods of solving their problems, none of these really stand out as the de facto solution for the stock market domain. So to choose the perfect solution for each use case can be a tremendous job as there are a lot of different factors to weight in.

This itself is not hard to overcome but what makes this really hard is that although the advantages and disadvantages of a particular technology for one specific purpose is well documented, the integration of one technology to a plethora of others is usually not. This is due to the constant development of each technology and the vast amount of different possible integrations. So it is really hard for a novice data scientist to pick technologies for their pipelines that seem to be the best solution individually while not leading to a deadend because the chose technologies do not work together.

Once the technologies, that seem to work with the problem instance you have, have been chosen the next step is to run them. This is usually documented very well for quickly starting the development on a single machine [3] [2], but these are usually instructions that are only meant for testing a single instance setup and are far from production ready ways to run the program or even develop it efficiently. To put this in other terms, the methods described are not sustainable in the long term. What this usually means is that the novice user either naively starts to develop their program on top these instructions that need a lot of revisions in order to run in production

or have to spend a lot of time learning the framework and its nuances in order to build themselves a future proof development setup.

In this chapter, we examined the obstacles a lot of novice data scientists can face when starting their journey on big data stock market analysis. In the following chapters we will try to tear down these obstacles, while confirming that these problems do exist. Our goal is not to solve every problem that is listed here, but instead alleviate some of the burdens that these can introduce to a starting novice data scientist.

3.2.3 Requirements for the methodology

As stated in the introduction, the goal of the methodology is to make big data analysis for stock market data more accessible to novice data scientists and as such we want it to fulfil two key goals:

- Make stock big data pipeline building easier for novice data scientist using it
- Make stock big data pipeline building easier for novice data scientist that read about/use the pipelines built using this method

The problems we are trying to solve here are not possible to solve with one person alone. This is why we want the methodology itself to be able to alleviate these problems if it is used enough and this is why we define the second goal here as it can help us define more meaningful requirements.

As we saw in the previous two sections there are a lot problems that can hinder the development process of a novice data scientist. However, we are only going to focus only on a subset of these as we do not have the ability to solve everything although as we will see the other problems do affect our lives. The specific problems we want to alleviate and the requirements that are derived from these can be seen in table 3.2. To make referencing these requirements in the following parts of the thesis easier, we have labeled these and following requirements with labels R1 - R8.

As the mapping from novice data scientist problems to requirements is quite trivial and we have already examined the problems in previous chapters we then move to the requirements that the underlying domain imposes. We have gathered these requirements in the following list based on the literary research we did on the stock market analysis in the previous chapter:

- Ability to add and ingest multiple data sources (e.g social media feeds) (R5)

Table 3.1: Novice data scientist problems and requirements they impose

Problem	Requirement
Lack of information/Outdated information on used data sources	A way to find relevant data sources (R1)
Vast amount of possible technologies	A way to filter relevant technologies (R2)
Documentation lack info on building sustainable pipelines	Provide a way to build pipelines that are sustainable (R3)
Papers and industrial presentation lack practical information to rebuild the pipelines	Provide a way to build pipelines that are easy to reproduce (R4)

- Ability to efficiently store possible big data (R6)
- Ability to efficiently process possible big data (R7)
- Ability to support deep learning models (R8)

The three first requirements are somewhat trivial for a normal big data system, but the last one is the one might propose some restrictions to our choices. These are purely selected on based on the state of art research and the thought process on how to enable these methods in the developed pipeline. This concludes our chapter for requirements based on the problems novice data scientists face and couple of requirements that the underlying domain imposes.

Chapter 4

Proposed methodology

In this chapter we will finally propose a methodology to try to solve the problems introduced in the previous chapter. 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.

As stated in the chapter about requirements, there are multiple points that can hinder the development of big data pipeline for stock data analysis for novice data scientist. In this chapter we will introduce our method which can help novice data scientist to build more efficiently big data pipelines for stock analysis.

We will start by introducing the overview of the method and the steps that it consist of. We then examine each step giving our reasoning why to do it like depicted in the step. Finally, we end the chapter by defining how are we going to test this method in the following chapters.

4.1 Overview

Our method to build novice friendly big data pipelines for stock market analysis is divided into 5 steps, which are mostly focused on planning the pipeline. The steps are the following:

1. Pick a data source for stock data
2. Pick either Scala or Python as main programming language
3. Pick technology for each part of pipeline
4. Implement the architecture using Docker and Docker Compose

5. Share the pipeline (if possible) and start developing

The process is quite simple and each of the steps are explained more thoroughly in the following sections with some examples, but here we examine the basic reasoning behind each step. Every process starts with picking a data source for stock data as the format of the data specifies what are the requirements for ingestion of the system and as such is a pretty trivial step. The main problem with data sources is the low amount of relevant data sources which is why we have listed some in the following section.

We recommend picking either Scala or Python as the main programming language and using it as much as possible in the pipeline. Both have the needed libraries that can be used to analyse stock data with machine learning methods while taking into account the big data aspect. Choosing one reduces the complexity of the pipeline and because both have a very versatile supply of libraries choosing one over the other does not reduce the options of the user. Choosing Scala, however, might restrict the amount of available information online as Python has very active community [54].

For selecting technologies we unfortunately cannot give any universal instructions. As said in the previous section, there are too many integrations to consider on how each component works with another to ensure that the system is optimal for each use case. However, we have provided one example on the technology selection. The last two steps are more thoroughly reasoned in the Section 4.1.4.

4.1.1 Data source

Our first step in the methodology has to do with picking up the data source and is directly linked to our requirement R1. Of course, if the user of this method already has a data source for stock data then we definitely recommend skipping this step and moving to the next one. As most of the data sources concerning stock data are closed and heavily priced we have gathered here the most promising data sources that novice data scientist could use as their starting point. These are picked based on their pricing and the amount of public information these services provide about themselves. A list of this kind of services is presented in Table 4.1.

We recommend the user of this methodology to pick from one of these, because of their low pricing and quite available data. Which one the user picks from these is case specific and is mostly determined by how much historical data the user needs and do they need also current intraday data for more up to date predictions. As there are multiple data sources that can

Table 4.1: Stock data sources

	Type of data	Minimum 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 [21]	Intraday and Historical (from 1980)	0\$/month	JSON/CSV	5 symbols/request, 250 request/day and 25 intraday requests/day
Intrinio [15]	Intraday	52\$/month	JSON / CSV / Excel	120 request/min
Quandl [18]	Historical <i>EOD</i> ¹ (from 1996)	15\$/month ²	JSON / Excel	none
Barchart [7]	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

provide quite similar amounts of data the next step is to check the pricing and restrictions and choose from these based on the users needs.

To better grasp Table 4.1 we have gathered here some notes that can help when choosing one of the sources. For acquiring free data test data fast, the IEX provides the best option as its limits are not restricted to time intervals, 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 the user 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).

4.1.2 Main programming language

The next step is to pick one programming language for the pipeline. The key thing here is to pick only one language and try stick with it through the development process. Using multiple languages can be useful because of the supply of different libraries can vary language basis, but we recommend sticking to one as each language introduces complexity to the system. Using multiple languages also decreases the sustainability of the system as the following persons who develop onto the system have to have knowledge on same set of languages.

In this step we have also gone a step further and limited this choice to only two different languages: Python and Scala. This is because at the moment these are the de facto languages that are currently used in this field and they both are taught introductory programming courses (at least in Aalto University). Java is also an option as most of the tooling still runs on top of JVM, but as Scala starts to be able to do everything it can and still provide a better support for functional programming our recommendation would be to use Scala instead of Java.

4.1.3 Technology selection

Next, we move on to choosing technologies for the pipeline. As we are working with the historical data part of the pipeline (f_1 , f_4 , f_5 from Section 1.1) we focus on three major parts of the pipeline: ingestion, storage and analysis. We have depicted the technology selection process in the Figure 4.1. The usage of figure is quite self-explanatory: User starts from the top and chooses at least one technology from each box following arrows based on their choices.

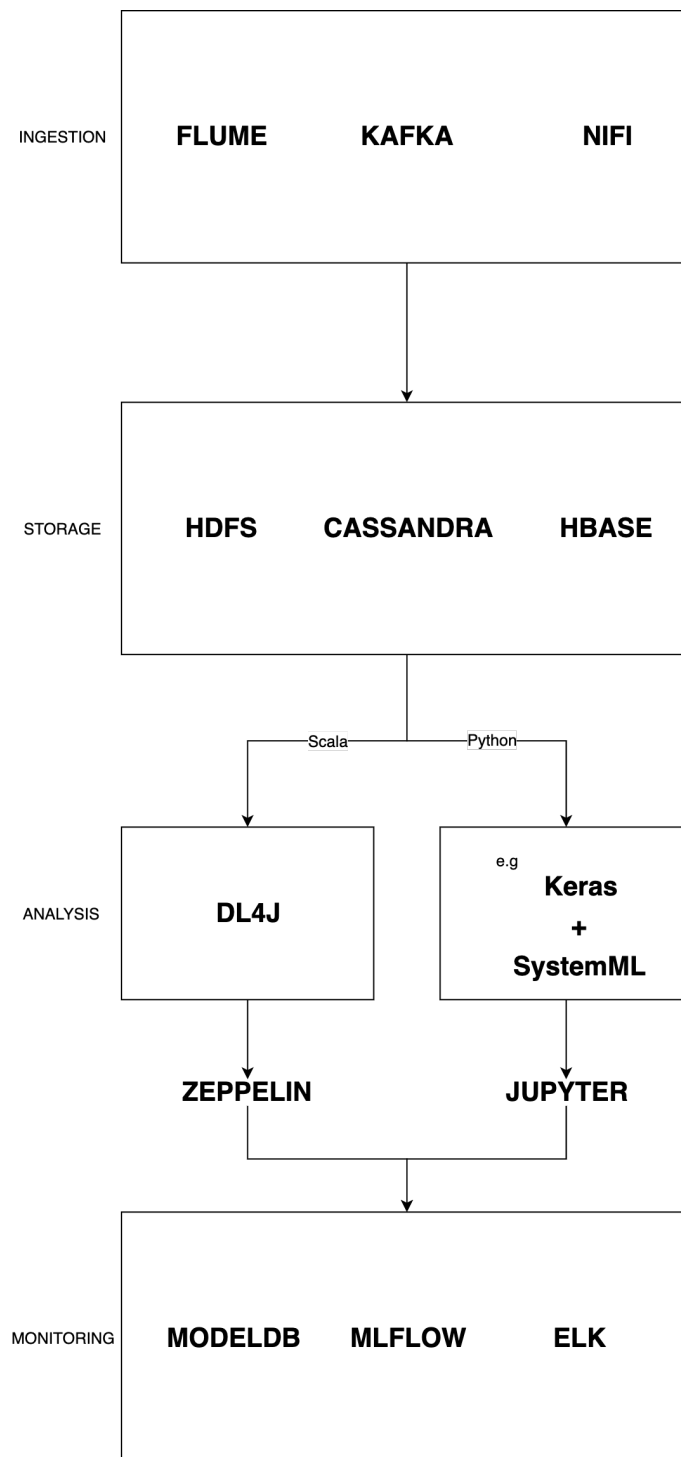


Figure 4.1: Example technology selection flow

We have narrowed down the number of possible technologies a lot to fulfill our requirement R2. The recommendations here are not absolute and the user is encouraged to use the technologies they are comfortable with if they know technologies that would fit their pipeline better. The goal of this step is to make the technology selection easier for novices by recommending technologies that have been seen to work in this domain.

The technologies that have been chosen here are the ones that we examined in the background chapter. These choices were made mostly based on technologies that are already in use and have been through this confirmed to work in stock data environment. We have added couple of new technologies based on their reputation in this field and how they at least theoretically fulfill the requirements that we defined in the previous chapter. From the requirements perspective all the ingestion technologies should in some sense fulfill the requirement R5, all the storage should fulfill R6 and all the analysis part should fulfill R7-R8.

The ingestion step is the most complicated one as the format of ingested data can be almost anything and the technologies we have listed can be used as they are or with conjunction with each other. Custom ingestion choices are also common in this field. Our recommendation with the data sources introduced in previous section is to use custom ingestion with Apache Kafka as most of the technologies listed here have somewhat little support for HTTP based sources. Flume works well if the data source is file or log format (e.g the sources providing excel formats) and NiFi can be used instead of Kafka if the user wants different kind of data routing management.

The storage selection is more straightforward. HDFS works with almost any case as it has a rich development environment, but the downside of it is that it lack proper query languages at it is. Cassandra and HBase do provide these but they both have their own use cases based on the CAP (Consistency, Availability and Partition tolerance) theorem which states that no system can guarantee consistency, availability and partition tolerance at the same time [27]. Cassandra provides availability and is recommended to use in application where the data is fetched directly to the user and the inconsistency of data is not a problem. HBase is the vice versa of this by promising consistency of data instead of availability. Which one the user chooses depends on the users use case.

For the analysis, as we saw in the background chapter, Spark is the de facto environment to run distributed processing. However, our requirement R8 limits quite a bit our options. If the user has chosen Scala as the main programming language we recommend using DL4J for the analysis at it provides support for distributed training in Spark and deep learning models. In the Python side the situation can be challenging in the sense that user

might have to combine couple of libraries to make the models train in truly distributed manner. Some of the de facto deep learning libraries in Python, such as Pytorch, have been designed mainly to be run on one machine. Users can go around these limitations if the library supports exporting the models into some universal model language that can be run in distributed environment using other specific libraries for this. We have included one example of this with Keras and SystemML which allows running Keras models in Spark, but as the number of possible combinations is quite large we recommend the user choosing combinations that they are most comfortable with.

After this recommend adding a language specific notebook library into the pipeline. Zeppelin if the language of the pipeline is Scala and Jupyter if the language is Python. This is to make the pipeline more reproducible (R4) and sustainable (R3) as the user has a way to share their experiments in more compact form instead of sharing the whole pipeline for each experiment.

Finally we have the monitoring which is somewhat optional field. We have collected here choices that could be most useful for novice data scientist and should not require much work to integrate into pipeline, but user should pick ones based on their own needs. Adding each should increase sustainability of the pipeline (R3) as they make easier to monitor the experiments in the pipeline and what could go wrong with these experiments.

4.1.4 Implement and share

Final two steps are implementation and sharing. For novice data scientists, we recommend building the architecture/configuration using Docker and Docker Compose container technologies [11]. These allow the user to develop multi-host systems in their local machine in a realistic virtualized environment. We have chosen this approach, because of multiple reasons. First, because Docker allows user to develop multi-host system in a single machine, the user does not have to rent or own multiple machines to develop the system. This makes the development of configuration very cost-efficient as there is no costs from testing it in a distributed environment. Second, because of the nature of docker Containers, the system becomes very easy to share and run on new machines. This directly linked to the requirement R4.

Once the user has managed to produce the architecture and configuration in docker environment we recommend the user to share this docker project at this point. At this point of the development the configuration probably includes already a lot of valuable information for new novice data scientists that can use it as a reference. As the development is not very far at this point there is no harm sharing this as the pipeline probably does not contain much classified information. The only way to fight the lack of information is

by producing more of it and this would be a low effort way to do it.

This concludes the examination of our proposed methodology. Next we move on to test it in practice via specific example case.

4.2 Use case

In the implementation part of this thesis we will be trying this method in a use case that reflects one possible situation that this method could be used. Ideally, we would want multiple use cases, but because of the time constraints that this thesis has we have settle for one. In this and the following section we will describe this use case and define how we are going to valuate the results from this use case.

4.2.1 Description

We will try the method in a situation where we have a novice data scientist that has the following background knowledge and needs. In this fictional situation, we will call this novice data scientist as Novice A for more compact text.

Novice A has computer science background and has studied machine learning in the master's degree level. Novice A, however does not know anything about big data pipelines and little about DevOps. Now Novice A is faced with a problem where he is tasked in his new job at small startup that was just established to develop a system that predicts stock market prices using state of art methods. The startup wants to use this system in the future to predict social media trends so the system must be able to handle big data from social media in the future and possibly big data from stock market if the idea works.

As the Novice A does not know anything about the stock market analysis and big data pipelines, he uses this method with the following choices.

1. Novice A chooses IEX as their data source from the list of data sources as it fulfills the data needs
2. Novice A chooses Scala as the main programming language as this is the language Novice A is most familiar with due to education
3. Novice A chooses the following technologies:
 - (a) Custom ingestion with Kafka for ingestion because of the data source they have chosen

- (b) HDFS for storage because of its mature developing environment
- (c) DL4J for deep learning because no other choices really exist for Scala with Spark
- (d) Zeppelin for notebooks because no other choices exist for Scala and Novice A wants a way to share their results easily
- (e) MLFlow for ML monitoring because there probably will be a lot of iterations to develop a working model because of the age of the startup

In the next chapter we will examine more thoroughly the implementation step of these technologies and see what kind of challenges this selection brings.

4.2.2 Analysis of results

Now that we have seen our specific use case, the next question is how are we going to test whether the methodology actually helped? We will be using qualitative analysis focusing on the requirements that we have set to this method. So we will be examining the sustainability of the pipeline, the reproducibility of the pipeline and novice friendliness of the pipeline, but we will also discuss about traits that big data pipeline should have e.g scalability.

We will be comparing the pipeline with two different pipelines using these qualitative metrics. Because of the scarcity of possible points of reference the pipelines have a bit different purposes, but we are more interested how do these pipelines compare with our metrics. We have picked one pipeline from the academia and one pipeline from the industrial side. The pipelines are depicted in very high level which can make it hard to make any reasonable deductions, but given the scarcity of data available it is the best we can do. We are more interested in how they are run than the actual components that they are made of but to give a better picture we also introduce the architectures.

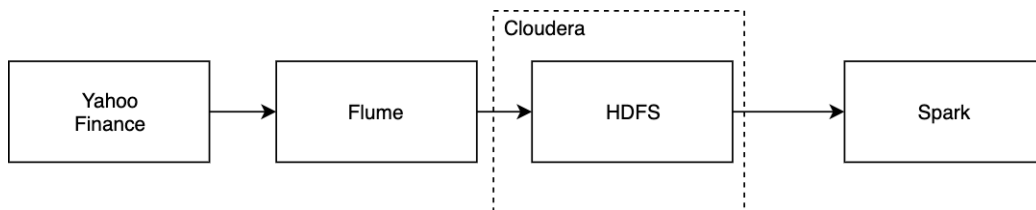


Figure 4.2: Simplified architecture of the academic pipeline

The pipeline we have picked from academia is the pipeline proposed in [56]. The reported architecture in the paper is recaptured in Figure 4.1.

From the paper we know that the system was ran locally using local Cloudera instance. Python was used with PySpark as the programming language and the pipeline ran normal machine learning models. Although this information is not much, we can use it to make some observations on our system.

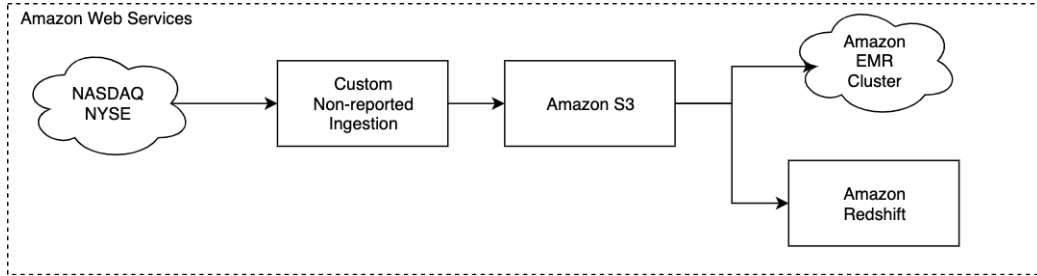


Figure 4.3: Simplified architecture of the industrial pipeline

The industrial pipeline that resembles the most of our pipeline here is the pipeline recaptured in [63]. The high level architecture shown is presented in Figure 4.2. What can be seen in the architecture is that we do not know much about the parts of programs that analyse the data. Amazon EMR can be used to run almost many different big data frameworks but given that they report using it to analyse data it is highly likely that they are running multiple Spark clusters there. But what we do know for sure is that the system is run on Amazon Web Services in large scale and we can use this to analyse our own system from the perspectives that we listed above.

We will also be quickly reflecting on the reported challenges. We will question whether these are significant in the sense that they could happen to any novice data scientist and is the cause of these challenges in this methodology or could they have been avoided using other methods.

This concludes our chapter on the proposed methodology and its planned usage in the empirical part of this thesis. In the next chapter we will report the process of the implementation this in the depicted special case and see what kind of challenges were faced during it.

Chapter 5

Use case

In this chapter we will report the implementation for the specific use case of the methodology we described in the previous chapter. This chapter can be 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. We will also explain how the end result can be run in this part. The code for the whole system is available at <https://github.com/Shipuli/masters-thesis-pipeline> per our final step of the methodology. Then as there were a lot of challenges during the implementation, these challenges are explained in the latter part of this chapter and we continue with them in the following discussion chapter.

5.1 Final Architecture

The final architecture of the pipeline that was planned using our method can be seen in 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.

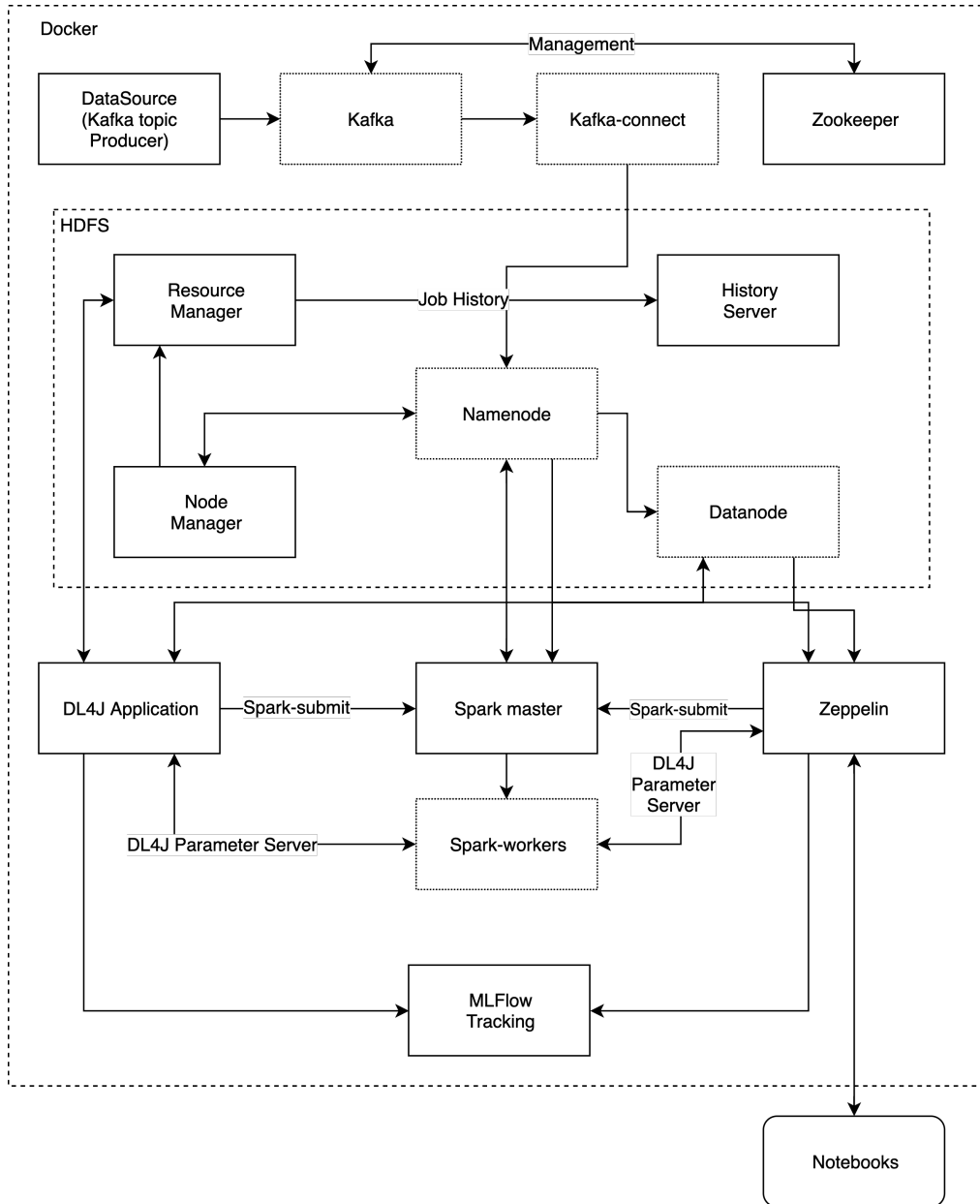


Figure 5.1: Final Architecture of the implementation section

It reads the data from a source, in this example case from JSON files that were exported from IEX, and produces the data into a Kafka topic. It also creates the needed Kafka topics on startup using *Kafka 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 producer and would need excessive scripting in order to automate this process which did not seem as maintainable as the Admin API approach. This is why admin API was used in the data source server.

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 pipeline 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 [9]. 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 software [8]. 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 5.1: 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 5.2: 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 the code runs correctly.

5.2 Usage

We made an initialization script for this case in case we had time to test multiple pipelines. 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. 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.

5.3 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. The purpose of this section is to highlight the problems that do exist and show that our methodology cannot really directly solve these because of their complexity, but can alleviate them indirectly if information is shared more.

5.3.1 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 the developer can only see if the whole system works or if it does not from HDFS management console. If only some configuration was not right, Connect would only log this as a minor error in the stderr and continue running like nothing had happened although no data was flowing through it. This made debugging the connection surprisingly difficult, but after all the configurations were in place the component worked as it should throughout the rest of the development.

5.3.2 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 have allowed better dependency management with maven.

5.3.3 DL4J sbt project

The DL4J application instance without Zeppelin is build with sbt. sbt is a build tool developed specifically for Scala applications [19] (although it supports Java also) 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.

5.3.4 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. However, it does not log 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 a 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 methodology we used to plan this software. We will also discuss more about what went well, what went wrong and what could have been improved.

Chapter 6

Discussion

In this final chapter before conclusion, we will be evaluating our methodology through the pipeline that was build in the previous chapter. We will be using the qualitative metrics that we defined in the planning chapter before implementation. Finally, we end this chapter by contemplating on how this pipeline and process could improved in the future.

6.1 Method evaluation

We have divided this evaluation strictly to advantages and disadvantages, and we will be using the metrics that we depicted in the methodology chapter. We have summarized these in the figure 6.1 and we will be opening up these in the following sections. We will be comparing our way of building the pipeline with the two somewhat similar existing pipelines to see which are the advantages and disadvantages of using one for novice data scientist.

Table 6.1: Advatanges and disadvantages of the methodology

Advantages	Disadvantages
Cost efficient	No promises on scaling/Runs only on one hardware
Sustainable	Docker introduces more complexity
Easy to reproduce	

6.1.1 Advantages

Let us start with the most obvious advantage of using this method: it is very cost efficient when it comes to running the pipeline. No cloud subscription or local cluster is needed to run and test the pipeline. As we saw in the previous chapter, this methodology does not erase the problems that arise integrating the technologies, but using it does erase the cost of debugging this against cloud platform where every run of the program costs money. We cannot make direct comparisons with academic pipeline as we do not have exact information on if the pipeline is run on a local cluster but we can say that the pipeline is much more cheaper than the industrial pipeline.

The sustainability of the pipeline is also an advantage in this case. Compared to the local academic pipeline, where each components version has to be manually confirmed when it is installed and updating components means that everyone who wants to run newer versions of the pipeline has to do the same, our version does not suffer from these problems as these are defined in the docker configuration files. Compared to the industrial pipeline, however, the process is probably taken into step further as Amazon can keep the components better up to date by automatically updating minor versions, which our current methods do not necessarily do. Of course, this depends a lot on the docker image provider and their methodology on versioning, but in our example case this seemed to provide efficient way to handle versioning.

The reproducibility of the pipeline is an another advantage. Using docker and docker compose to build the pipeline we were able to limit the number of external dependencies to only these two when a new user wants to run the pipeline and with only couple minor exceptions we were able to automate the local setting up process. Also the usage of Zeppelin notebooks allows the user to share their models with anybody who has the same pipeline or other means to run them. When we compare this to the academic pipeline which had all of its components as local dependencies, the amount of possible work is greatly smaller as the user does not have to install each component to their local machine and worry about their versions. As what comes to the automation of the setting up the pipeline, we have no information on this to make any reasonable claims. When compared to the industrial pipeline, although not reported, we can pretty much assume that the industrial pipeline is highly reproducible as there are probably ten to hundred employees working with the pipeline and in this sense is probably a lot more reproducible than the pipeline presented here. This could be due to inherent nature of containerized applications in AWS and the tooling that continuous integration tools and Amazon provide.

6.1.2 Disadvantages

For the disadvantages, clearly, our pipeline has only run on one machine and we still have no proof that it would scale to big data order of magnitude. Compared to the industrial pipeline we can clearly state that our pipeline does not yield same kind of efficiency as a pipeline run on cloud. This is a clear disadvantage of this methodology.

Related to this is the inherent nature of running the pipeline only on one machine. The virtual distributed environment does allow us to simulate networks that are needed to configure the system, but the underlying hardware does not allow us test it with an amount of data that is too much for one machine.

Adding docker to the mix of technologies does bring a bit more complexity to the system. As we saw in the previous chapter this can introduce quite challenging problems when trying to integrate these technologies together.

6.2 Future improvements

The lack of reference points made it hard to compare this way of building a pipeline with other possible ways. Hopefully in the future, the information about these systems would become more open so that we could better understand them, but until then we have to work with what we have. The use of only one example is also a clear point of improvement as it does not prove much about the functionality of our methodology. More specific cases would be needed to make more believable deductions.

For the methodology itself, better instructions on how to deploy the pipeline to the cloud would have probably solved a lot of disadvantages that we saw in the previous section. Also a better instructions on running the pipeline in distributed GPU environment would have probably greatly improved the efficiency of model training. Unfortunately due to lack of time this was not possible, but would be a great way to improve the information here as there is not that many examples online of this kind of process.

Chapter 7

Conclusions

In this thesis, we examined big data pipelines that analyze stock market data from the point of view of a novice data scientist. We examined what are the challenges that novice data scientists face when starting their development of these kind of pipelines and this was examined from both the stock analysis side and the big data technology side. In the stock analysis side we saw that the available practical information on stock analysis is very limited, whereas in the big data technology side the problems were mostly due to the lack of intermediate level information or the fact that it is very scattered.

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 big 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.

Based on the requirements set by our target group and the current big data stock analysis enviroment, we proposed a novel five step method to help novice data scientists build big data pipelines for stock analysis. Our goal was to bring down the gap between beginner level documentation and corporate level complex systems and highlight the problems that novice data scientists could face while developing such a system.

We provided an experimental use case and using our method we were

able to build such a pipeline and report multiple significant challenges while developing the pipeline which could affect our target group. We saw from this use case that while our method makes the building of a pipeline more cost efficient and the end result seems to be quite sustainable and easy to produce, it also might add challenges and complexity to the system because of our technology choices. However, as our sample size was only one in this case, further research would be needed to make viable claims.

In the future, hopefully, more information can come available if and if 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] Barchart free market apis. <https://www.barchart.com/ondemand/free-market-data-api> Accessed 23.09.19.
- [8] Big data europe project. <https://www.big-data-europe.eu/> Accessed 25.09.19.
- [9] Big data europe project docker images. <https://hub.docker.com/u/bde2020> Accessed 30.12.19.
- [10] Deeplearning4j documentation. <https://deeplearning4j.org/docs/latest/> Accessed: 30.06.19.
- [11] Docker documentation. <https://docs.docker.com/> Accessed 30.12.19.
- [12] Elastic stack documentation. <https://www.elastic.co/elk-stack> Accessed: 30.06.19.

- [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] Metastock homepage. <https://www.metastock.com/> Accessed 01.12.19.
- [17] Mlflow documentation. <https://mlflow.org/docs/latest/index.html> Ac-
cessed: 30.06.19.
- [18] Quandl us eod data documentation.
<https://www.quandl.com/data/EOD-End-of-Day-US-Stock-Prices>
Accessed 23.09.19.
- [19] sbt documentation. scala-sbt.org/1.x/docs/index.html Accessed
30.12.19.
- [20] T2000 homepage. <https://www.worden.com/> Accessed 01.12.19.
- [21] World trading data documentation. <https://www.worldtradingdata.com/>
Accessed 23.09.19.
- [22] AMIRGHODSI, S., ALLA, S., KARIM, M. R., AND KIENZLER, R. *Apache Spark 2: Data Processing and Real-Time Analytics*. Packt Pub-
lishing, 2018.
- [23] BALAM, S., CHANDRAKUMAR, T. A., AND CHAKRABORTY, S. A time
series analysis of the it stock market during the 2007 – 2009 recession.
2018 IEEE International Conference on Big Data (Big Data) (2018).
- [24] BARTRAM, S. M., AND GRINBLATT, M. Agnostic Fundamental Anal-
ysis Works. *Journal of Financial Economics (JRE)* (June 20 2017).
- [25] BASS, L., WEBER, I., AND ZHU, L. *DevOps: A Software Architect's*
Perspective. Addison-Wesley Professional, 2015.
- [26] BRIDGWATER, A. Nsa 'nifi' big data automation project out in the open.
[https://www.forbes.com/sites/adrianbridgwater/2015/07/21/nsa-](https://www.forbes.com/sites/adrianbridgwater/2015/07/21/nsa-nifi-big-data-automation-project-out-in-the-open/#79a9319655d6)
[nifi-big-data-automation-project-out-in-the-open/#79a9319655d6](https://www.forbes.com/sites/adrianbridgwater/2015/07/21/nsa-nifi-big-data-automation-project-out-in-the-open/#79a9319655d6)
Accessed: 20.06.19.
- [27] BUTTERFIELD, A., NGONDI, G. E., AND KERR, A. *A Dictionary of*
Computer Science (7 ed). Oxford University PressPrint, 2016.

- [28] BZDOK, D., ALTMAN, N., AND KRZYWINSKI, M. Statistics versus machine learning. *Nature Methods* (2018). <https://www.nature.com/articles/nmeth.4642> Accessed: 30.12.19.
- [29] CHEN, J. *Essentials of Technical Analysis for Financial Markets*. John Wiley & Sons Inc, 2010.
- [30] 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.
- [31] 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.
- [32] CURTIS, P. Streaming stock market data with apache spark and kafka. Available at: <https://www.youtube.com/watch?v=0tSZo8I2924> Accessed: 08.05.19.
- [33] 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).
- [34] 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).
- [35] 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).
- [36] FISCHER, T., AND KRAUSS, C. Deep learning with long short-term memory networks for financial market predictions. *European Journal of Operational Research* (2017).
- [37] GEORGE, L. *HBase: The Definitive Guide, 2nd Edition*. Packt Publishing, 2018.
- [38] GOSSWAMI, T., SAHA, S. K., AND HASAN, M. Stock market data analysis and future stock prediction using neural network. *International Journal of Computer Science and Information Security (IJCSIS)* (2018).

- [39] 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).
- [40] HOFFMAN, S. *Apache Flume: Distributed Log Collection for Hadoop - Second Edition*. Packt Publishing, 2015.
- [41] 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).
- [42] JOHN L. PERSON. *Mastering the Stock Market: High Probability Market Timing and Stock Selection Tools*. John Wiley & Sons, Incorporated, 2013.
- [43] 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).
- [44] 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.
- [45] 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).
- [46] 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).
- [47] 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).
- [48] LI, Y., JIN, T., XI, M., LIU, S., AND LUO, Z. Massive text mining for abnormal market trend detection. *2018 IEEE International Conference on Big Data (Big Data)* (2018).
- [49] 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.
- [50] MUN, F. W. Big data, small pickings: Predicting the stock market with google trends. *The Journal of Index Investing* 7 (2017).

- [51] MURPHY, J. J. *Technical analysis financial markets: A comprehensive guide to trading methods and applications*. New York Institute of Finance, 1999.
- [52] 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).
- [53] NEERAJ, N., MALEPATI, T., AND PLOETZ, A. *Mastering Apache Cassandra 3.x - Third Edition*. Packt Publishing, 2018.
- [54] NGUYEN, G., DLUGOLINSKY, S., BOBAK, M., TRAN, V., GARCÌA, A. L., HEREDIA, I., MALIK, P., AND HLUCHY, L. Machine learning and deep learning frameworks and libraries for large-scale data mining: a survey.
- [55] 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.
- [56] PENG, Z. Stocks analysis and prediction using big data analytics. *International Conference on Intelligent Transportation, Big Data & Smart City* (2019).
- [57] 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).
- [58] 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.
- [59] SESSIONS, V., AND VALTORTA, M. The effects of data quality on machine learning algorithms.
- [60] SEZER, O. B., OZBAYOGLU, A. M., AND DOGDU, E. An artificial neural network-based stock trading system using technical analysis and big data framework.
- [61] SEZER, O. B., OZBAYOGLU, A. M., AND DOGDU, E. A deep neural-network based stock trading system based on evolutionary optimized technical analysis parameters. *Procedia Computer Science* 114 (2018).

- [62] 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).
- [63] SNIVELY, B. AWS re:Invent 2018: Big Data Analytics Architectural Patterns & Best Practices. Available at: <https://youtu.be/ovPheIbY7U8> Accessed: 08.05.19.
- [64] UMADEVI.K.S, GAONKA, A., KULKARNI, R., AND KANNAN, R. J. Analysis of stock market using streaming data framework.
- [65] UTTHAMMAJAI, K., AND LEESUTTHIPORNCHAI, P. Association mining on stock index indicators. *International Journal of Computer and Communication Engineering* 4 (2015).
- [66] VARTAK, M., SUBRAMANYAM, H., LEE, W.-E., VISWANATHAN, S., HUSNOO, S., MADDEN, S., AND ZAHARIA, M. Modeldb: A system for machine learning model management.
- [67] VOULGARIS, Z. *Data Scientist: The Definitive Guide to Becoming a Data Scientist*. Technics Publications, 2015.
- [68] WANG, W. A big data framework for stock price forecasting using fuzzy time series. *Multimedia Tools and Applications* 77 (2018).
- [69] WENGA, B., LUA, L., WANG, X., MEGAHEDB, F. M., AND MARTINEZ, W. Predicting short-term stock prices using ensemble methods and online data sources. *Expert Systems With Applications* 112 (2018).
- [70] WORLD FEDERATION OF EXCHANGES. Monthly report january 2019. <https://www.world-exchanges.org/our-work/statistics> Accessed 15.04.19.
- [71] YARABARLA, S. *Learning Apache Cassandra - Second Edition*. Packt Publishing, 2017.
- [72] YU, F., LI CAO, S., CHENG HUANG, S., AND FEI GAN, P. Stock transaction analysis system based on hadoop and capital flow. *Sixth International Conference on Advanced Cloud and Big Data* (2018).
- [73] 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).

- [74] YUAN, D. Stream processing in uber. <https://www.infoq.com/presentations/uber-stream-processing/> Accessed: 19.06.19.
- [75] ZHANGA, J., CUIA, S., XUA, Y., LIA, Q., AND LIB, T. A novel data-driven stock price trend prediction system. *Expert Systems With Applications* 97 (2018).