

City University

MSc Business Systems Analysis and Design

Project Report

Can a visual data mining  
approach improve the stock  
picking process?

Sergi Vives

January 2015

Supervised by: Jason Dykes

January 2015

*By submitting this work, I declare that this work is entirely my own except those parts duly identified and referenced in my submission. It complies with any specified word limits and the requirements and regulations detailed in the assessment instructions and any other relevant programme and module documentation. In submitting this work I acknowledge that I have read and understood the regulations and code regarding academic misconduct, including that relating to plagiarism, as specified in the Programme Handbook. I also acknowledge that this work will be subject to a variety of checks for academic misconduct.*

Signed:

A handwritten signature consisting of a stylized, cursive letter 'S' followed by a straight horizontal line extending to the right.

## Abstract

This research project is the product of an internship with Thomson Reuters to work on one of the most creative activities that financial analysts perform: searching for investment opportunities. This activity consists in navigating large datasets of financial data in order to gain better insight on the data, find outliers, spot patterns, trends etc. to understand the data better and pick the most interesting equities.

If we abstract from the financial field, we are dealing with a data exploration problem. Visual data mining consists in a set of different techniques to represent the data visually in order to process it more easily.

The product of the current research is a fully functional prototype that implements visual data mining features to solve the data exploration problem.

This prototype is finally tested to answer the research question, whether visual data mining can help equity analysts to get a better insight on the data and thus improve the stock picking process.

The results obtained suggest that in fact visual data mining can help improve the stock-picking process.

## Table of Contents

<b>1. Introduction</b>	<b>9</b>
1.1. Background	9
1.1.1. Why this research?	10
1.1.2. Information visualization	11
1.2. Research question	11
1.3. Why did I choose the project	11
1.4. Beneficiaries	12
1.5. Products Of This research	12
1.6. Research Objectives	12
1.7. Scope and Definition	13
1.8. Methods	13
1.8.1. Software Development Methodology	13
1.9. Testing the objectives	13
1.10. Structure of this report	14
<b>2. Context</b>	<b>15</b>
2.1. Financial Data Visualizations	15
2.2. Visual Analytics	16
2.2.1. Information Overload	16
2.2.2. Visual Analytics	16
2.2.3. Visual Data Exploration	17
2.2.4. Visual Data Mining	17
2.3. Other Equity Screeners	18
<b>3. Methods</b>	<b>20</b>
3.1. The Design Study Methodology	20
3.1.1. Suitability	20
3.2. The 9-Stage framework	21
3.3. Precondition Phase	21
3.3.1. Learn	22
3.3.2. Winnow	22
3.3.3. Cast	23
3.4. Core Phase	24
3.4.1. Discover	24
3.4.2. Use Cases	25
3.4.3. Requirements	25
3.4.4. Design	25
3.5. The 4 levels of DESIGN	30
3.5.1. Domain situation Validation	31
3.5.2. Data/Task abstraction validation	31
3.5.3. Visual Encoding and Interaction Idiom Validation	31

3.5.4. Algorithm Validation	32
3.6. The 4 levels of design in practice	32
3.7. Implementation	34
3.8. User testing for Visual Encoding Validation	34
3.9. RE-Validation of functionalities	36
3.10. FINAL Comparison Test	37
3.10.1. DECIDE Framework	37
3.10.2. Designing the tasks	39
<b>4. Results</b>	<b>42</b>
4.1. Requirements	42
4.1.1. Use cases	42
4.1.2. Final Requirements	43
4.2. DATA	Error! Bookmark not defined.
4.2.1. TASK Abstractions	47
4.3. Design and Implementation	47
4.3.1. Current visualization	47
4.3.2. Task encoding	48
4.3.3. Interaction coding	48
4.3.4. Design Features	48
4.4. Initial Design Testing	56
4.5. Linked Views	56
4.6. Re-validation of the functionalities	57
4.7. Final Comparison Test	58
4.8. Final Testing Results	58
4.8.1. Benchmark tasks	58
4.8.2. Medium Insight Tasks	61
4.8.3. Full-insight tasks	62
<b>5. Discussion</b>	<b>63</b>
5.1. Validaty and generalisation of the results	65
5.1.1. Validity	65
5.1.2. Generalization	65
<b>6. Evaluation, Reflections and conclusions</b>	<b>67</b>
6.1. Literature review	67
6.2. Reflection TOPIC AND Objectives	67
6.3. Reflection on chosen methods	68
6.4. Reflection on PLANNING	68
6.5. What has been achieved	68
6.6. What I've learned	69
6.7. Conclusions	70
6.8. Further work	70
<b>7. Works Cited</b>	<b>72</b>

<b>8. Appendix</b>	<b>78</b>
8.1. Original Project Proposal	78
8.2. Functionality revalidation Feedback Summary	87
8.3. Abstracted Tasks	89
8.4. Questions for the comparison Test	89
8.4.1. Benchmark test	89
8.4.2. Medium Insight test	91
8.4.3. Insight exploration task	92
8.5. Final Comparison Test Results	92
8.5.1. Task Distribution	92
8.5.2. TEST A Results	93
8.5.3. TEST B Results	95
8.5.4. TEST C Results	96
8.6. Final Comparison Results Summary	97

## List of Figures

Figure 2: Technical Analysis chart.....	19
Figure 3: The 9-stages framework depicted.....	21
Figure 4: The different aspects of answering the What? question .....	26
Figure 5: The Why? Part of the triad.....	27
Figure 7: The how? Part of design.....	28
Figure 8: A slide of the mock-ups presented to the domain expert.....	29
Figure 9: The 4 levels of design.....	30
Figure 10: Validation diagram of the 4 nested levels of design .....	31
Figure 11: Overview of the entire project.....	33
Figure 12: The scatterplot on the Equity Screener .....	48
Figure 13: An example of parallel coordinates displaying 8 different dimensions.....	49
Figure 14: Brushing can be represented as drawing boxes on specific dimensions in this case highlighting data lines with 'EBIDTA Margin' >50 and 'ROE' >45.....	50
Figure 15: An example of parallel coordinates with smooth curves, note it is much easier to follow a curve rather than a straight line. ....	50
Figure 16: The scatterplot bubbles with different hue encoding the information of their corresponding category. ....	51
Figure 17: Only companies of the 'Healthcare' sector are displayed, note the median line too. ....	52
Figure 18: Depiction of the trend of two different companies, the suffixed 'E' indicates the value is estimated. ....	53
Figure 19: An example of a cluster that holds 12 equities inside.....	53
Figure 20: The equities inside the cluster are expanded in a circular shape around the cluster .....	54
Figure 21: The arcs are proportional to the number of stocks of the category.....	54
Figure 22: The QuickRank mode .....	55
Figure 23: The final scatterplot design .....	55

Figure 24: Screen capture of the final design of the prototype implementing the linked views approach..... 57

## 1. INTRODUCTION

The current project is the result of the work performed during an internship programme offered by Thomson Reuters.

Thomson Reuters (TR from now on) is leading multinational media and information firm aiming to provide professionals with quality, trustful information. TR proud themselves in delivering high quality information to professionals in a wide variety of fields ranging from finance, legal, pharmacy & life sciences to intellectual property. Since its inception in 1851 TR has been providing stock market quotations to brokers and, whilst it has widened the range of information products, it has kept clear ties with the stock market.

The current research project will be focused on producing a visualization that can help financial analysts in the process of finding interesting investment opportunities in the stock market. If successful this visualization could be a part of Eikon, the company's flagship financial product.

### 1.1. BACKGROUND

Eikon provides financial professionals with all the information required to discover investment opportunities, track assets, commodities and stocks amongst other functionalities. It is an extremely powerful tool that combines multiple information sources to offer a deep insight on market information.

The tool this research project will be focused on is a tool to screen stocks, the Equity Screener. It allows screening and analysing a vast amount of stocks being traded in the most important stock exchanges in the world. Financial analysts use it to evaluate stocks and pick the ones that are worth investing in.

The workflow is quite simple, normally equity analysts have an investment strategy that defines a series of filters, such as companies with a growth in sales higher than 5%, they filter the universe of stocks and end up with a smaller set of stocks. This operation is repeated until a manageable number of stocks are listed, and then these stocks are evaluated.

This process is not as simple as it sounds. There is no fool proof way to know what will be the price of a given stock on the future (Investopedia, 2012). Yet, by examining numerous accounting and economic factors, an investor might have a better sense of the future value of a given stock and various studies actually acknowledge this fact (Dillow, 2009) (Thomson, 2013) (Graham, 1945). The work of an equity analyst is complex and has considerable pressure associated to it as stock prices can rocket or plummet with no apparent reason (Hagstrom, 2005). The current project proposal will focus on enhancing the stock picking process enabling financial analysts to explore the data through the use of visualizations.

There are two main approaches for stock picking. In the one hand a more long-term approach that evaluates the company and decide whether to pick or not based solely on sound financial indicators. This approach is called Fundamental Analysis (FA). Its main objective is to estimate a company's *intrinsic value*<sup>1</sup> (Graham, 1945, p.27). This value is then compared to the price it's trading at. If the price is under the intrinsic value the company is undervalued so it is good to buy shares. If the intrinsic value is lower than the price (overvalued) then is better to walk away (Graham, 1945, p.38).

On the other hand, the other main approach targets solely on the price movements of the market and is known as technical analysis. Technical analysis focuses on the short-term income by trying to predict if the stock price will increase or decrease and purchase shares accordingly. This type of analysis can be automated using complex algorithms that execute trades according to a set of rules, this is known as high frequency trading (HFT). HFT is still an extremely important part of the stock market (38% of the orders in 2010, 62% in US market (CNBC, 2014)) executed all over the world it seems there has seen a decline in trading volumes (Philips, 2013). This might be due to its unintended consequences; in several occasions the stock market has plummeted without any apparent reason due to these algorithms (The Economist, 2010) . This has played against technical analysis and in favour of tools that take advantage of FA such as Eikon's Equity Screener.

Although there have been some attempts to automate FA (Islam et al., 2009) it is a difficult task as it involves evaluating both qualitative (e.g. intellectual property, brand name) and quantitative factors (e.g. growth, income, sales) that affect the company's value (Investopedia, 2013). So, FA remains a mainly manual task performed by equity analysts.

Evaluating all the indicators for a single company is complex, but evaluating these in relation to the industry, stock indexes, other companies etc. makes this task extremely challenging and time expensive (Investopedia, 2013) . Stock picking then becomes an exploration task having to navigate thought the vast amounts of data to make the decision of what stock(s) to purchase.

### 1.1.1. WHY THIS RESEARCH?

The current functionalities on Eikon's Equity Screener are quite limiting in terms of data exploration. The tool does not provide a paradigm that allows the user to explore the data, instead the exploration flow is limited by the fact that each small change on the data filters requires the data to be reloaded and re-evaluated. The current data visualization does not give the sense of exploration either as it is only capable of displaying a small part of the dataset (just 3 dimensions when there is a vast number of possible dimensions).

---

<sup>1</sup> Intrinsic Value: The actual value of a company or an asset based on an underlying perception of its true value including all aspects of the business, in terms of both tangible and intangible factors. This value may or may not be the same as the current market value.

If the current project is to be successful, it will create a visualization that can enhance the current equity screener with data exploration capabilities. This will be an example for TR on how the data they currently have can be exploited with a powerful information visualization that enables data exploration.

### 1.1.2. INFORMATION VISUALIZATION

Information Visualization (InfoVis) can be defined as "*the use of computer-supported, interactive visual representations of data to amplify cognition*" (VIAU, 2012). The last three words of their definition communicate the ultimate purpose of visualization, to amplify cognition: to help understand. InfoVis then can be regarded as a powerful tool to develop insights from data (Fekete et al., 2008). It does so by using human perception as the basis for inducing new insights from this data. Vision acts as a very fast filter making it much easier to spot anomalies and patterns on the data (Fekete et al., 2008).

Thus, InfoVis systems are best applied for exploratory tasks that involve large datasets such as the data the Equity Screener works with. InfoVis can also support cases where the expert using the InfoVis system may not have a specific goal in mind; where the user might simply be examining the data to learn more about it, to make new discoveries, or to gain greater insight (Fekete et al., 2008).

## 1.2. RESEARCH QUESTION

As we've highlighted, theory suggests that InfoVis can be used to present complex financial data in such a way that financial analysts can gain a greater insight on it. An outgrown field of InfoVis is Visual Data Mining (VDM). VDM integrates knowledge of different fields such as data mining, statistics and InfoVis itself to create a paradigm that enables data exploration. But can such approach really improve the process of stock picking? The fundamental question that the current research project aims to solve is thus:

***Does a visual data mining approach improve stock picking for equity analysts?***

## 1.3. WHY DID I CHOOSE THE PROJECT

Financial data is an interesting field of work in terms of visualization. The vast amount of information (even real-time flow information) can overwhelm analysts and hide interesting insights. From my personal perspective, I do believe information visualization techniques can be used to improve how the data is conveyed and presented to analysts to make their job easier. Not only that, by adding a considerable degree of interaction, the data not only can be perceived visually but also manipulated visually shedding new insights on the data.

Thomson Reuters is a multinational company providing vast amounts of data to different industries around the world. It is *THE* information company and supports the major financial markets in the world.

Being able to bring a possibly new perspective on data visualization into the company is an interesting challenge and an exciting opportunity for the researcher.

## 1.4. BENEFICIARIES

As previously stated, this current project is part of an internship for Thomson Reuters. Consequently, the major beneficiary of this work is Thomson Reuters as the work produced could derive a new set of features or a new visualization paradigm that could be incorporated in the Equity Screener product.

The visual data mining approach that supports this research could also be used in other domains of Thomson Reuters Eikon. If proved useful, it might prove on other fields where being able to explore the data and gather more insight is beneficial.

Additionally, not only Thomson Reuters will benefit from the current research. Researchers working on the InfoVis field for financial data could benefit from this research too. Other situations where the characteristics of data might be similar to the complex financial dataset we are working with might also benefit from the discoveries uncovered by the present this work.

Ultimately the current research project should add to the body of knowledge on visualization in general and on financial data visualization particularly.

## 1.5. PRODUCTS OF THIS RESEARCH

The product of this research will be a fully functional prototype that will feature visual data mining techniques to enable equity analysts explore this universe of data in a novel way.

## 1.6. RESEACH OBJECTIVES

To try to answer this research question a set of objectives must be set. The main objective of the current project is to design and implement a prototype that will be used to test the research question. The prototype should allow the researcher to evaluate with domain experts if a visual data mining approach improves the stock-picking process.

Several objectives can be derived from this enterprise:

- **Research Objective 1 (RO1):**  
Elicit a set of requirements with the help of domain experts for features that can help equity analysts when performing the stock-picking process.
- **RO2:**  
Design a visualization solution with visual data mining features using the DSM approach.
- **RO3:**  
Implement the prototype of a visualization solution within the project timeframe. The prototype should fulfil the requirements elicited with the domain expert.
- **RO4:**  
Evaluate the resulting prototype compared to the prior version without data-mining capabilities. This will be evaluated with domain experts. Gather enough data to test the research question.

## 1.7. SCOPE AND DEFINITION

The current project will only focus on how a visualisation or set of visualisations can be used to improve and facilitate the stock picking process. There is no intention to widen the scope of the project to evaluate how visualisation can help in other types of financial-related activates. That is, the prototype being developed will only focus on the stock picking process.

The prototype will source the information from a static source. It is not the objective of the project to focus on how to source data effectively but rather on how to visualise and interact with the data in a more effective and novel way.

## 1.8. METHODS

Due to the particular characteristics of the current research question being faced it seems a problem-driven research approach would be the best fit. A problem-driven research is such where the goal is to work with real users to solve their real-world problems (Sedlmair et al., 2012).

Design studies are a form of such problem-driven research. A design study can be defined as a project in which researchers analyse real-world problems, design a visualization that helps to solve the problem and reflect about the lessons learned to refine the field of visualization design (Sedlmair et al., 2012).

A design study approach defines a process with several phases that establishes how to conduct the research (Munzner, 2009). Although this approach is quite structured in terms of on which steps have to be taken, it is considerably loose in terms of how each step is executed. The current research will use DSM as a supporting framework iterating through the phases as needed.

### 1.8.1. SOFTWARE DEVELOPMENT METHODOLOGY

Any software development project requires a certain methodology to be followed. This project has the objective to develop iteratively a fully functional prototype. The software methodology that will then be used is Incremental Development (Larman & Basili, 2003).

The development of the prototype will be divided in small iterations. For each iteration a new set of prototypes will be developed, validated with the domain experts and only the best will be refined. This approach is known as iterative Prototyping (Mnkandla, 2005). Using multiple prototypes ensures the researcher is exploring multiple alternatives, although it requires more work it leads to better results (DOW et al., 2010).

## 1.9. TESTING THE OBJECTIVES

To answer the main research question all research objectives have to be fully met. Only by understanding the data exploration problem the researcher will be capable of deriving requirements. These requirements will be used to design and implement a prototype following the DSM methodology.

The prototype should fulfil the requirements elicited and implement a visual data exploration paradigm to be able to test the research question.

The prototype will be used to answer the research question. An evaluation study with domain experts will be performed in order to test the prototype against the current equity screener. In this test we will evaluate the visual data mining capabilities of the prototype compared against the current screener with no visual data mining capabilities.

## 1.10. STRUCTURE OF THIS REPORT

We will start giving the reader an idea about what is the problem that needs to be solved, how is being solved in the literature and academic context and understand how others have approached it.

Then we will see how the theory can be used to actually solve the current problem in hand and explain what has been the design process and what has driven the main design decisions that lead to end up with a fully working prototype. Next we will explain how we designed the user tests to evaluate the research question.

Then, we will explain the outputs of applying these methods to solve the research question.

Later, we will reflect on the results obtained by applying the chosen methods. These results will be tested against the research objectives and reflected upon.

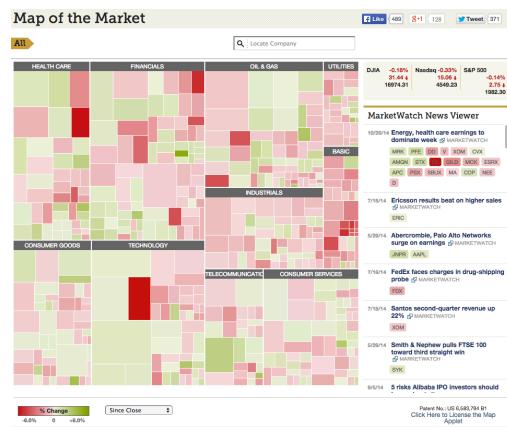
Finally, the conclusions on what has been done, how it's been done and where room for improvement has been visible will be evaluated. Additionally, what has been achieved and future work will be further elaborated.

## 2. CONTEXT

In this chapter we will evaluate how visualization is being used on the financial sector at the moment. To do so we will evaluate the current academic context around financial data visualization and highlight how visualizations can be used to convey financial data. Then we will elaborate on the theory that suggests visual data mining can effectively solve the data exploration problem through the use of visualizations.

### 2.1. FINANCIAL DATA VISUALIZATIONS

One of the most well known visualizations for financial data is the Map of the Markets (Wattenberg, 1999). It is a Tree Map representation of the whole stock market that gives an idea of the overall status of the market in a single visualization. It cleverly combines stock tendencies (price up/down into color-coded squares), market capitalization (depicted by size) and grouping by industry (by wider borders) all in a single chart.



Although this initial contribution in financial data visualization seemed promising (it was created 15 years ago), financial analysts have mainly got stuck to the traditional charts for their analysis tasks. Some studies highlight that the traditional trend line over time chart is the most commonly used visualization for financial analysts (Ziegler et al., 2010) (Schreck et al., 2007). Furthermore, very few solutions have been developed which can handle stock market data information effectively in order to gain enough insight to understand how the market works (Ziegler et al., 2010).

#### Stock data visualizations

The great majority of research focuses on representing the stock market in general to give a sense of how the market is doing. Lei and Zhang (Lei & Zhang, 2010) have conducted an interesting research on this direction. They suggested several different visualizations with features that enable a visual analytics approach that should help financial analysts. These visualizations were evaluated with stock traders with positive results on how visual analytics can improve the stock-picking process with reasonably positive results (Lei & Zhang, 2010).

Another interesting work by Schreck et al. uses different visualizations try to increase the ability of financial analysts to discover patterns on the data more easily. This research brought some positive results (Schreck et al., 2007).

Ziegler et al. use visualization not only to find patterns on the data or show the overall status of the market but both. Their tool allows comparing the features of companies, assets or even entire countries. This allows analysts to explore and discover interesting facts such as patterns that indicate that markets or countries go into turbulence (Ziegler et al., 2010).

These solutions are not based on static visualizations but on powerful interactive visualizations with the objective of enabling a deeper evaluation of the complex dataset that

is being visualized. They use techniques from what is called as Visual Analytics (VA). Some studies suggest that the field of financial data analysis can clearly benefit from information visualization and this outgrowth field of Visual Analytics (Schreck et al., 2007).

## 2.2. VISUAL ANALYTICS

### 2.2.1. INFORMATION OVERLOAD

The data the equity analyst is faced with a large and multi-dimensional dataset. This dataset needs to be broadly analysed in order to extract valuable information for the stock-picking activity. This is clearly an information overload problem as the information that needs to be processed is way too large for the analyst to digest and thus makes it difficult to make clear sensible decisions (Infogeneering, 2013).

Information Visualization techniques are gaining popularity in helping tackle this information overload problem (Ankerst et al., 1998). InfoVis techniques allow the analysis and exploration of large multidimensional data sets reducing the information overload problem by conveying sufficient information for equity analysts to make informed decisions (Lei & Zhang, 2010).

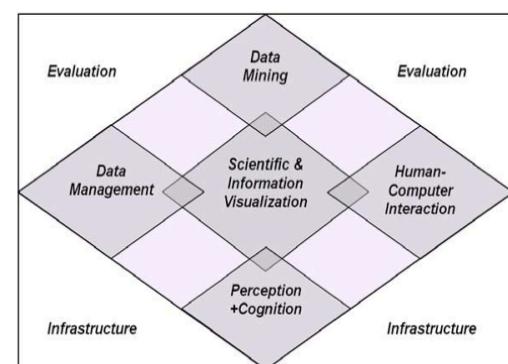
### 2.2.2. VISUAL ANALYTICS

Visual Analytics (VA) is an outgrown field of Information Visualization. It is more than a simple visualization process as it draws its techniques from several different disciplines to enable a deeper understanding of the data. Visual analytics provide solutions that combine the strengths of humans with the electronic data processing to deliver new insights on the data (Keim et al., 2008). These insights are then used to support informed decision-making. Visual Analytics then can be seen as an integral approach to decision-making based on analysing data visually with the following objectives (Keim et al., 2008):

- Synthesize information and derive insight from large datasets
- Detect the expected and discover the unexpected
- Provide timely, defensible and understandable assessments.
- Communicate assessment effectively for action

To achieve these objectives VA draws from several different disciplines (Keim et al., 2008):

- **Data Management:** Enables the management of data in an efficient way so that it can be visualized and filtered efficiently.
- **Visualization:** Presents the data in an understandable manner.
- **Perception & Cognition Principles:** To encode the information in an understandable way.
- **Data Mining:** Allows to extract patterns and structure and valuable information from the raw data).



### 2.2.3. VISUAL DATA EXPLORATION

One of the key features of a data exploration problem is that the goals set by the user when faced with a dataset are not clearly specified. That is, the user is not looking for something specific on the data; he is instead trying to get a grasp of what can be interesting. This is usually performed via an iterative hypotheses generation and verification process; the user formulates a hypothesis, tries to verify it and learns about the data throughout the process. Daniel A. Keim used the term Visual Data Exploration to name such approach, specifying that ‘visual data exploration is especially useful when little is known about the data and the exploration goals are vague’ (Keim, 2002).

Visual Data Exploration aims at making the user the center of this data exploration process, exploiting its capabilities to enable a better understanding of large data (Keim, 2002). The combination of the human’s flexibility, creativity and analytical reasoning is then put to work collaboratively with the enormous storage capacity and the computational power of the computer to enhance the exploration process (Simoff et al., 2008).

This concept of exploratory data analysis is strongly associated with visualization because visualizations can enable the user to interactively navigate the structure of the data and process it visually to gain further insights (Schulz et al., 2006).

A Visual Data Exploration approach that has proven to be a really effective method for data exploration using visualization techniques is known as Visual Data Mining (VDM).

### 2.2.4. VISUAL DATA MINING

Data mining consists on an exploratory analysis targeting to find interesting structures on the data. This requires little or no human interaction, solely relying on computer algorithms (Symanzik, 2001). Visual Data Mining is based in the same principle, but instead of drawing only from computer algorithms it is also based on human participation. This combination of automated mining algorithms together with visual analytics techniques produces great results when they are backed by a rich user interaction. The data mining algorithms can extract structures that are inherent in the data (e.g. clustering, statistical information, unusual observations) and then graphical representations present both the extracted structure as well as the real data. The human can interact with it to better understand the data and thus making it easier to highlight new knowledge (Simoff et al., 2008) (Symanzik, 2001).

The Visual Data Mining exploration process can be seen as an iterative process where the analyst iteratively selects the data to be visualised, modifies any parameters on the data or visualisation and then starts again with a new dataset (Purchase et al., 2008).

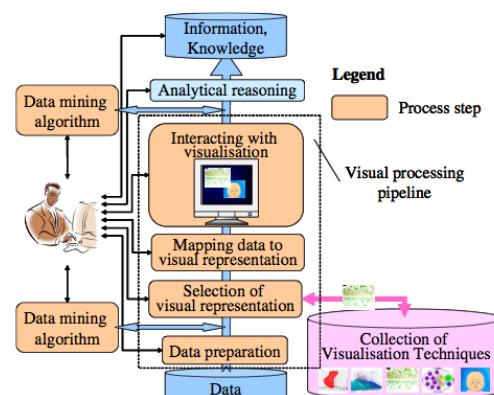


Figure 1: Diagram showing the VDM process

### **Overview first, Filter and Zoom, Details on demand**

A refined approach to this iterative Visual Data Mining exploration process that puts visualizations on the centre is the mantra by Shneiderman's (Keim et al., 2008): 'Overview first, Filter and zoom, Details on demand.'

First the user gets an overview of the data getting a wide idea of the details of the entire dataset, identifying interesting patterns to focus on. Then, to analyse these patterns he might be able to zoom to a specific area of its interest to obtain the details on the data. For the interaction to be successful it is important to keep an overview (in some sort of visualization) of the entire dataset when the user is focusing on a smaller set. This mantra is based on an iterative behaviour and the analyst is likely to go back and retrieve the initial data (Ye & Rey, 2011) (Keim, 2002) (Keim & Zhang, 2001).

As the reader might recall the current research question is to test whether such visual data mining approach can improve the stock-picking process. The prototype that will be built thus will need to facilitate this feedback loop where the analyst can explore the data freely.

## **2.3. OTHER EQUITY SCREENERS**

When designing a new product or application it is always compulsory to evaluate whether other similar tools exist on the market and how they solve the task. What needs to be evaluated is whether other equity screeners exist and if they have any visual data mining capabilities.

There are several different tools used for screening equities available on the market. The researcher analysed a wide range of different screeners by trying to answer the following questions: does the equity screener have a visualization solution? What visualization idioms are used?

Most of the current equity screeners that were analysed used a simple table to display the result of the screening process (The Telegraph, n.d.) (Financial Times, n.d.) (Thomson Reuters, n.d.) (Marketwatch, n.d.) (EquityMaster.com, n.d.) (Fool.com, n.d.). Even premium equity screeners even where customers had to pay a monthly quote to access the service did not integrate any charting capabilities (MorningStar, n.d.). Such is the case of Thomson Reuters competitor, Bloomberg, which offers an equity screener without any visualization.

On the other hand, some equity screeners were more focused on the technical side of screening (4Traders, n.d.) (TMX Money, n.d.) thus offering mainly technical charts (charts focused mainly on price movement) such as the depicted below.



Figure 2: Technical Analysis chart

Another set of different screeners make use of visualizations to improve the process. Such is the case of Panopticon's stock screener example (PanopticonSoftwareAB, n.d.). This software is different from the previous screeners as the filtering process is performed using the visualization as the main source of information (rather than simple tables). Although the Finviz.com screener is a table based screener (Finviz.com, n.d.), the website authors have created a series of interesting visualizations that show the overall state of the stock market with different Treemap-based visualizations (Finviz.com, n.d.). However these visualizations do not relate to the stock-picking process directly.

From this wide range of equity screeners that have been evaluated none take a direct advantage on visualization to help the stock-picking process. Therefore, no of the evaluated equity screeners have visual data mining capabilities.

The current research project thus, devises a new and innovative approach to screening by bringing visualization to the centre of the equity screening process.

### 3. METHODS

To backbone the current research project the researcher decided to follow the Design Study Methodology (DSM). This methodology is used in problem-driven visualization research. DSM unifies several different approaches to visualization research providing a holistic view of the process and guidance on each step. This approach seemed a valid and good fit for the project in hand being a problem-driven research.

One of the main reasons for such selection was due to the exposure from the researcher to such method. Tamara Munzner came to give a master class on such method at City University by presenting several success cases where this methodology was applied. Although these success stories were similar to the current research the main difference was the time frame; normally DSM projects lasted much longer. Yet, the agile nature of the methodology seemed that allowed a degree of freedom in terms of the project length.

#### 3.1. THE DESIGN STUDY METHODOLOGY

DSM allows researchers to tackle visualisation research in a systematic way. The definition of DSM can be drawn from Munzner, which explains the methodology and the potential pitfalls the researcher might be faced with when using this approach:

*"A design study is a project in which visualization researchers analyse a specific real-world problem faced by domain experts, design a visualization system that supports solving this problem, validate the design, and reflect about lessons learned in order to refine visualization design guidelines."* (Sedlmair et al., 2012)

Although DSM seemed a suitable and appropriate approach, not all real-world problems are suitable to this framework.

##### 3.1.1. SUITABILITY

DSM focuses on two different aspects of the problem being faced to design a visualization that fulfills the requirements: the data and the tasks; what data will the intended user use and what task will he perform with this data. The degree of suitability depends mainly on these two factors. First, how clear the tasks the visualization needs to support are defined and second, where the information to solve these tasks is located (Sedlmair et al., 2012).

Regarding the information location, we can have two extremes. In one hand, it might happen that all the information required to complete the task is actually on the domain expert's head. That is, information is implicit and relies solely on the expert's knowledge. In this case no matter if the task is well-defined or not, a visualization solution will prove useless in helping the expert to complete his task. On the other hand, it could happen that all the knowledge required to complete the task might be located on the computer itself. In this case, if the task to be sorted out is sufficiently well defined we could potentially write an algorithm that automates the task. In this case visualization is not needed either.

## Suitability of the current research project

In the case of the current research, the tasks in hand are not crisply defined. There is some room for further elaboration on what the exploration of equities is, and how is to be performed. Additionally not all the information for the task to be completed is actually on the expert's head. The ability to discern a good equity from another equity is partly on the expert's head, and partly on the data itself. Thus, the expert still needs to interact with the tool to extract and discover the interesting equities.

Consequently, the current problem we are trying to solve is effectively a DSM-suitable question as it can be allocated within the DSM space.

## 3.2. THE 9-STAGE FRAMEWORK

The Design Study Methodology (DSM) defines a 9-stage framework that vertebrates the entire methodology. This framework does not define a really strict approach but a starting point, a scaffold for the process of design (Sedlmair et al., 2012).

The methodology defines 9-stages. Cycling between the different stages is encouraged at any point of the flow. The designer can jump back to a previous stage to refine the design or to enhance and re-evaluate the knowledge on the problem. For example, we could go back from the implementation stage to the discover stage to learn about a new visualization that might prove better suited for a new use case.

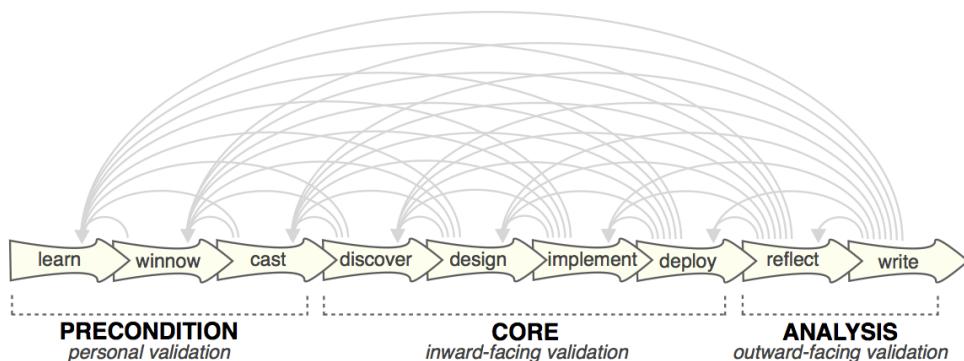


Figure 3: The 9-stages framework depicted

## 3.3. PRECONDITION PHASE

The precondition phase should set the foundations for a good problem-driven DSM research project. It has to allow the researcher to detect if the project is really suitable for a DSM-based approach. Not only that, it also allows the researcher to balance out if the research project is sufficiently interesting and that it can benefit the researcher. Additionally it should be a phase to allow the researcher to equip him with enough knowledge to be successful on the following phases of the framework. There are three main activities that take place during this phase.

### 3.3.1. LEARN

DSM emphasizes the need for the researcher to have a good knowledge of the visualization literature and its 'visual encoding and interaction techniques, design guidelines and evaluation methods' (Sedlmair et al., 2012). DSM does not require the researcher to be an expert on visualization, yet it strongly recommends a good basic of knowledge on the visualization field as it will be useful all throughout the project.

Before being able to get in contact with a domain expert to kick-start the project there was some time where the researcher could not start the project. There was room to effectively focus on the learn phase. The researcher was able to spend some time re-learning the concepts of visualization and researching the most novel visualization techniques. This phase lasted one and a half weeks approximately from around the August 25<sup>th</sup> to September the 3<sup>rd</sup> where the first meeting with the domain expert took place.

### 3.3.2. WINNOW

The winnowing stage tries to identify the suitability for a visualization project and discard the less useful projects. DSM poses a set of questions to help define if a given project is both feasible and attractive to the researcher (Sedlmair et al., 2012).

The winnowing stage was useful in trying to detect any problems that could jeopardize the project from the start such as the lack of access to the data or the inability to contact domain experts. The researcher considered the answers to these questions and was able to act upon any problems detected at this stage.

**Does real data exist? Is it enough, can I have it?**

Yes. Thomson Reuters provides the researcher with access to real stock data. Two different ways to access the data where possible: via exporting an Excel file or dynamically via a Web Service querying the stock data.

**How much time do they have for the project? How much time can I spend on their environment?**

The time allocated for the project was from the end of August till December, enough time to be able to gather requirements and develop a visualization prototype.

The great benefit of the internship is the fact that the researcher can work on site with close contact with domain experts. This allows an easier discussion and faster feedback loops informing the design.

**Is there an interesting visualisation research question?**

Solving a data exploration problem by devising novel interactive visualization is an interesting approach for the researcher.

**Is there a real need or they can use the existing tools? Are there any existing approaches good enough?**

Currently, only a small set of indicators can be evaluated visually. A new way to visualise the data is required to enable analysts to have a better overview of the data.

***Is this a real task? How many people are going to use it?***

The equity screener is in fact one of the most used applications on Thomson Reuters Eikon and a really important piece of software on the financial analyst's workflow. Finding investment opportunities by exploring the equity space is obviously a real task.

### 3.3.3. CAST

Two main roles can be found on the DSM methodology: the front-line analyst and the gatekeeper. The front-line analyst is the user doing the actual data analysis, the final user. The gatekeeper is the person that can approve or block the project as well as to authorize people to spend time on the project.

In the case of the current internship at Thomson Reuters, the roles defined by the DSM are slightly different.

Although the gatekeeper still exists it can be understood as two different roles, the work supervisor and the academic supervisor. The work supervisor can authorize people to spend time on the project as well as facilitate the resources required by the researcher. The academic supervisor can approve or block the project in terms of its academic quality. But the situation is different from a normal DSM project on a company where the funding and access to resources takes a more important role than in the current project.

A major requirement for DSM is the need to establish contact with at least a front-line analyst and obtain approval from the gatekeeper. In fact, one of the first things facilitated by the work supervisor was to facilitate the contact with a domain expert. The domain expert the work supervisor facilitated was a Product Manager for the Assets Management and Investment group at Thomson Reuters.

Thomson Reuters Eikon is organized in an app-based approach. Eikon offers its functionality on several different applications. For example the Equity Screener is in fact an application that resides within Eikon. The Reuters News platform that lists the latest news offered by the agency is an application too. Each application has one or multiple Product Managers that control what are the features that will be included on the next release: they are in contact with customers to know what features could help customers when using the tool, control the feature requests as well as envision new and innovative functionalities that can be incorporated on the app to make it more attractive. To carry on with such task Product Managers need to have a wide knowledge on their field, not only theoretically but also in practice (most of them have been traders, investors, portfolio managers in the past so they do know the tools available on the market). Being able to have a front-line analyst of such nature in the project is of great benefit to the outcomes of the project. The only problem with this connection was the fact that there was a considerable time difference (8 hours), as this front-line expert was based in San Francisco, US and all communications needed to be done by phone.

Front-line analysts in the case of the current research are not real-world users. At Thomson Reuters it is complex getting hold of real-world users due to company policy. When accessing real-world users, the new features that are being tested and that might be included in Eikon are exposed and could potentially be copied. Additionally, if customers have access to a piece of software that fulfills their needs better than the current solution they might want that feature implemented immediately on Eikon. This is virtually impossible and might cause the customer to be unhappy with the product. To avoid such problems

testing is normally performed internally with Thomson Reuters employees and product managers act as ‘proxies’ for real world users, deciding what features are needed.

## 3.4. CORE PHASE

Once the initial precondition phase takes place and the major issues amended, the core phase starts. Although it might seem we are following a cascade-like fashion, the nature of DSM allows revisiting previous steps to improve the overall design, so the precondition phase could potentially be revisited later on the project.

This phase is where the major design work takes place. It starts by learning about the domain, identifying the needs the visualization needs to fulfill, designing a visualization solution that supports the tasks, implement it, deploy it and test it with users. Again, the inherently iterative nature of DSM allows revisiting the previous phases making the design an iterative process.

### 3.4.1. DISCOVER

In this phase the general problem we want to solve with an appropriate visualization needs to be abstracted into a set of both tasks and data. The process of characterizing and abstracting the problem is cyclic. The expert explains the domain problem; the researcher listens to this explanation, abstracts the problem presents back what he has gathered and obtains feedback from the expert on the abstraction (Sedlmair et al., 2012).

During this abstraction process the researcher learns about the domain, the practices, needs and problems directly from domain experts. In doing so he can discover if and how visualization can enable insight and discovery (Sedlmair et al., 2012).

Abstracting both tasks and data means translating broadly defined tasks specific to a domain into a more abstract and low-level set of tasks. These task abstractions can then be addressed through visualizations.

Tasks can be derived from the requirements elicited. Requirements analysis takes place during this phase (Sedlmair et al., 2012). The researcher not only has to learn about the domain but also has to translate the user needs to requirements that the visualization should fulfill. These in turn are further analyzed and the tasks that need to be fulfilled are abstracted.

By creating a correct characterization of the problem, the researcher is achieving a shared understanding with domain experts, something crucial when designing a complex visualization (Sedlmair et al., 2012).

To abstract the problem, learn about the domain and elicit the requirements for the visualization a series of meetings were scheduled with the Product Manager. As mentioned before no access to real-world users was possible so Product Managers (PM) acted as proxies for real world users.

#### Kick-start Interview

Interviews are the most common technique for requirements gathering, mainly due to the fact that allow collecting a great amount of information about the domain quickly (Zowghi & Coulin, 2005). The requirements elicitation started with an initial unstructured interview.

Although this format is informal it helps the researcher to better understand the problem in hand by asking the expert about specific areas of interest leaving the domain expert expand on these areas as he wishes. However, the main drawback of this sort of interviews is that they don't follow an agenda and might end up focusing too much on some areas while neglecting others (Zowghi & Coulin, 2005).

The initial meeting was held on the September 3<sup>rd</sup> following such format. It was facilitated by the work supervisor and had the objective of gathering general information about the domain. The PM explained the common workflow of a financial analyst using the tool; the environment where the tool is used and the most common tasks performed using it.

### 3.4.2. USE CASES

From the initial meeting the researcher could distil a rich set of use cases for functionalities the visualization should fulfill. These use cases were somewhat detailed descriptions transcribed from the conversation with the domain expert on how the visualization was ought to be used. These were not complete use cases but a lighter version destined to be a starting point for eliciting the requirements and abstracting the tasks.

### 3.4.3. REQUIREMENTS

In many software projects requirements can come from multiple sources (Zowghi & Coulin, 2005). The domain expert was indeed the only source of requirements for the visualization solution as no other expert was available to the researcher.

From the use cases elicited, the researcher extracted a set of requirements the visualization prototype should fulfill. Not only functional but also non-functional requirements were elicited based on the information extracted.

### 3.4.4. DESIGN

The DSM suggests designing visualization solutions taking into account three main questions for each of the requirements the visualization should fulfil:

- **Why** the user intends to use the visualization
- **What** data the user will see
- **How** the visualization conveys this information concerning the interaction and the encoding choices.

A combination of answers for these questions is known as an instance. An instance is thus a triad of a *task-data-idiom*. A really simple visualization solution can consist of only one instance whereas much more complex visualizations can have multiple combinations of these instances.

To effectively design a visualization that can fulfil the requirements we will use this framework to analyse the requirements and abstract the tasks the visualization should fulfil.

## Task Analysis: What?

The first part of the triad **what**-why-how. What data are we actually dealing with? We need to know what data we have available and can show to the user. There are several different types of data that can be visualized.

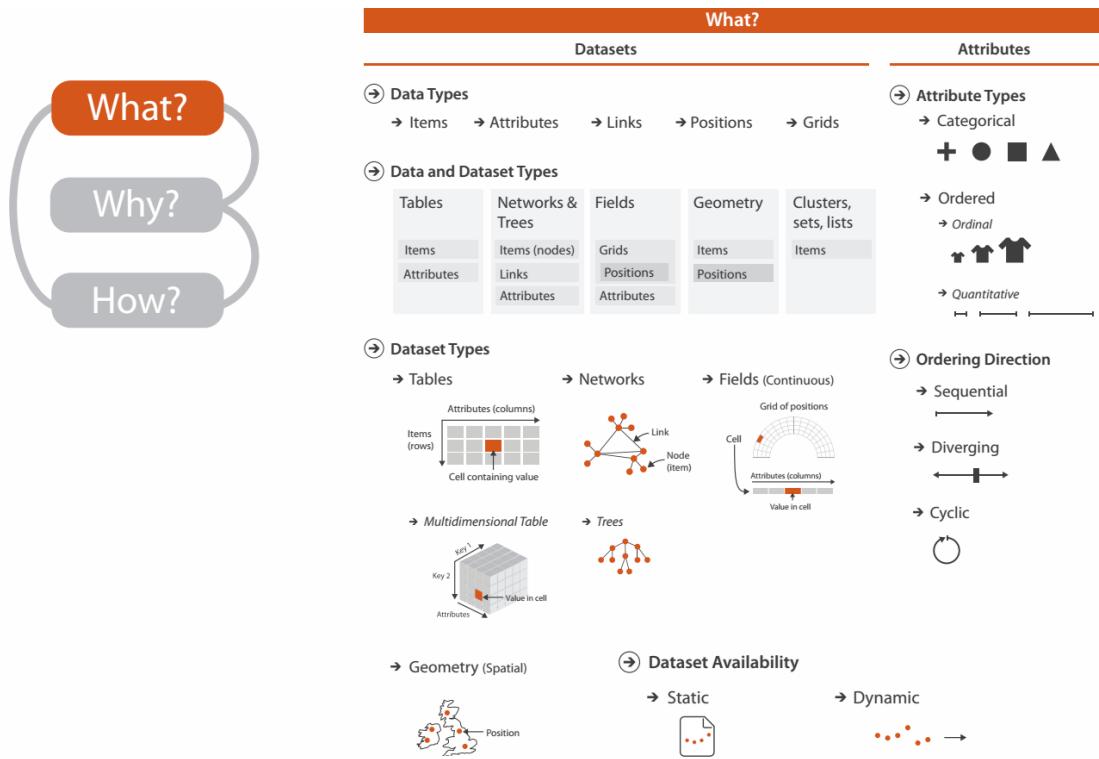


Figure 4: The different aspects of answering the *What?* question

Data is normally organized in datasets. The dataset that is available in the Equity Screener is a table-based dataset. All data is organized in rows and columns, every row corresponding to a unique stock and with a set of different attributes that are the columns of the table.

## Task Analysis: Why

We now know about the data most task idioms will be based on, now we can analyse what is that we want to do with it. That is, what will be this data used for. For example, data can be consumed to discover interesting insights, or used to produce new information.

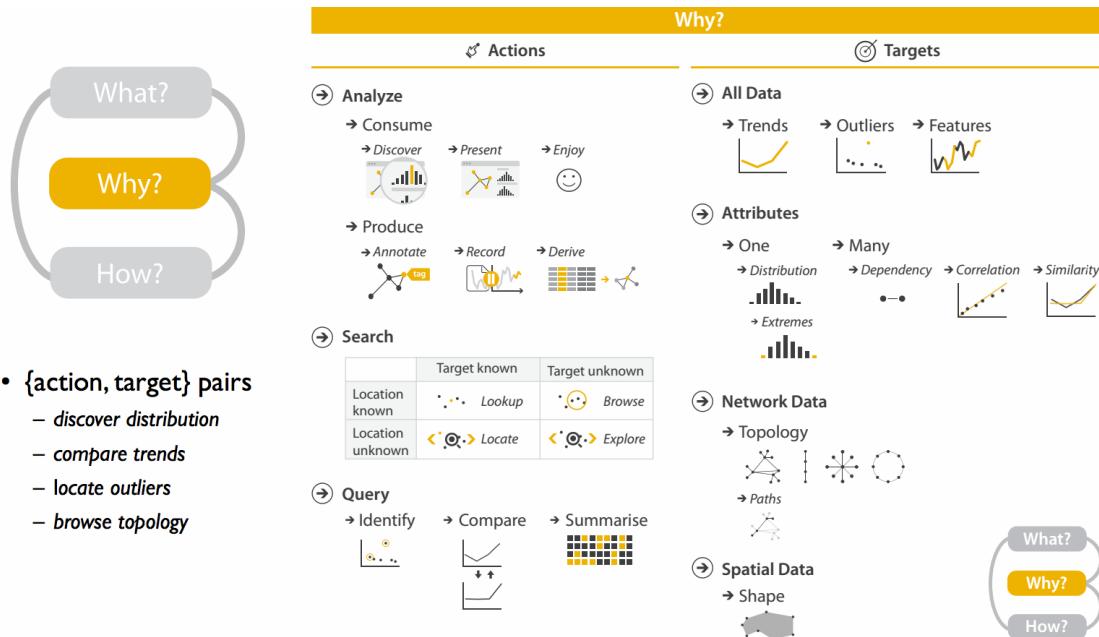


Figure 5: The Why? Part of the triad

From the initial use cases, requirements were elicited, and from these tasks could be abstracted. To do so a simple process of translation from the tasks was followed. This process required evaluating the tasks against two different aspects: the actions and targets of the task.

## Actions

Actions represent the user goals for the task. That is, what the user wants to achieve by doing the task. Three levels can be distinguished in terms of the actions: Use, search and query.

- Use: Why the visualization is being used for e.g. consume information to discover interesting patterns.
- Search: Whether the target is known or not.
- Query: Whether the objective is to identify, compare or summarize the data.

## Targets

Targets represent what the user wants to understand from the data such as locating trends and outliers on the data, finding the extreme values, dependencies between attributes, correlations etc. The picture below depicts the different targets the user might want to discover from the data:

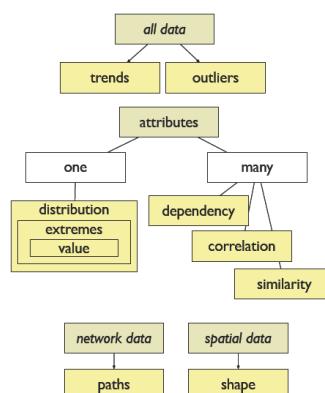


Figure 6: The different targets of the analysis

## Abstracting tasks in Practice

To abstract the tasks, each single requirement was analyzed and both the data and the task derived into a domain-independent version. This activity mainly took place between the initial kick-off meeting where the requirements were elicited from the domain expert and the tasks abstracted. Additionally the data that could be used was also being analyzed at that time.

The requirements gathering process is a critical stage of any software project. It not only occurs once but continuously (Zowghi & Coulin, 2005). That is, requirements elicitation happens all throughout the project life. This idea fits within the DSM approach as DSM embraces cycling through the process making the task / data abstraction a continuous process.

## Task Analysis: How

Once tasks are abstracted, the focus is placed on designing the visualization idioms that can be used to support these tasks with the data we have available. Doing so means selecting the best visualization idiom to fulfill each task from a wide range of different design choices.

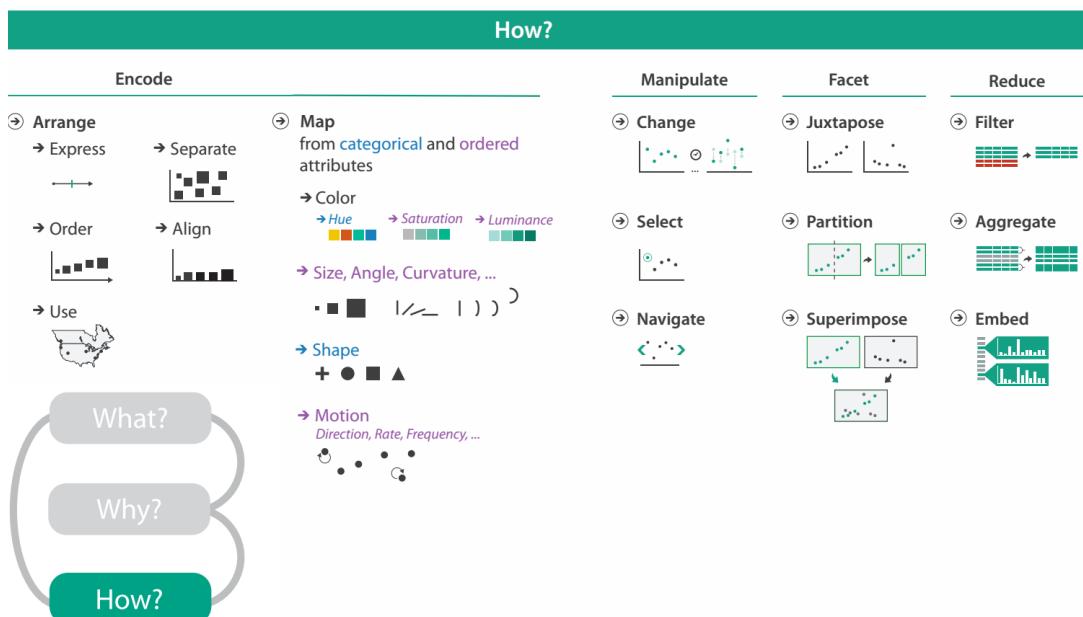


Figure 7: The how? Part of design

Two main design choices need to be addressed for each task, first how the data is encoded visually and secondly how this can be manipulated. Selecting both the encoding and interaction idiom for each task means solving the last part of the what-why-how triad.

Committing too early visualization idiom without investigating the wider range of solutions is a key pitfall in DSM (Sedlmair et al., 2012) and is best avoided by starting with a wide set of different solutions and iteratively narrowing down the proposal space. This is best achieved using a form of low-level prototyping, which helps presenting the different solutions to domain experts without spending too much time in their elaboration.

Does a visual data mining approach improve stock picking for equity analysts?

The process of refining had to take into account the timespan of the project. Although some solutions seemed valid at first, due to the short implementation time some were too complex to be implemented and discarded.

To avoid this main design issue a parallel-prototyping (DOW et al., 2010) approach was taken. Developing several different prototypes allows considering a wider range of solutions. These prototypes where then validated with domain experts and the less optimal solutions discarded.

The investigation of the different visualizations took place all throughout the design phase. Between the initial kick-off meeting (that took place on the 3<sup>rd</sup> September) and a series of different meetings until around the 29<sup>th</sup> of September most project time was dedicated to investigating the visualization design space in search for visualizations that could fulfill the tasks in hand.

## Keynote prototyping

Initially, the validation of the tasks/data and visualization idiom could be easily conveyed with simple non-interactive prototypes. These initial prototypes consisted in mock-ups of the visualization elaborated using Keynote (the Mac OS equivalent for Microsoft Office). These mock-ups lacked interaction but enabled the researcher to easily explain and present the features of the prototype of visualization. When the interaction was more difficult the researcher moved into implementation to be able to show the interactions.

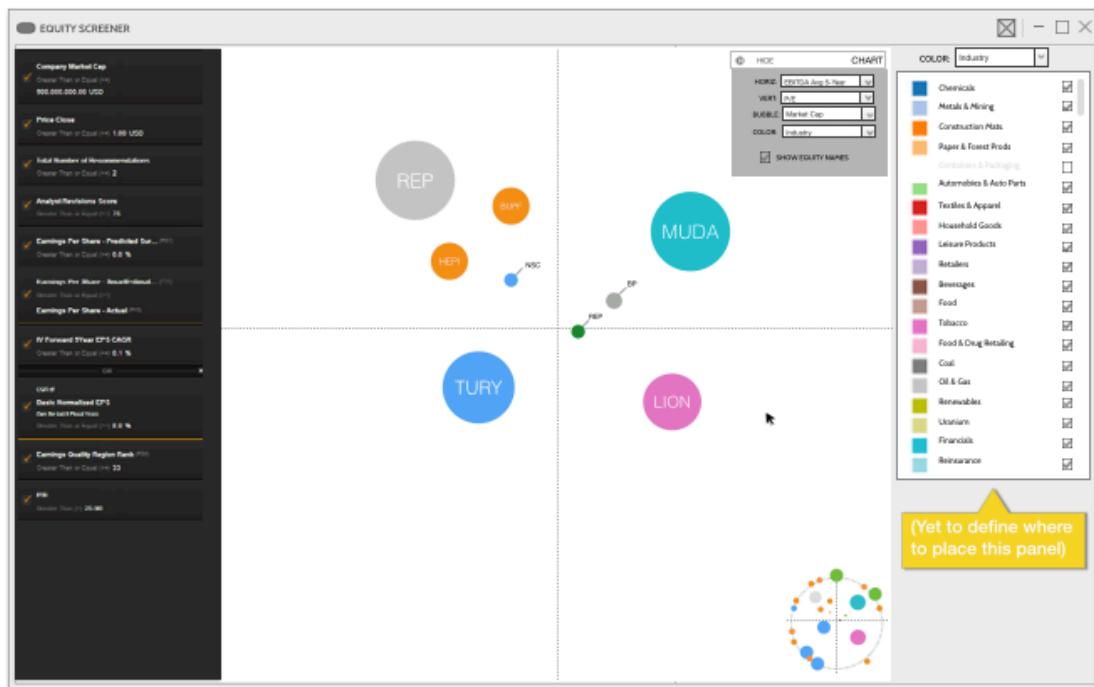


Figure 8: A slide of the mock-ups presented to the domain expert

## Visualization Prototypes with real data

Visualizations are best validated with real users, real tasks, and large, complex datasets (Carpendale, 2008). In order to validate the visual encoding idiom the researcher implemented the encoding of key parts of the visualization with real, complex data

generated from the Equity Screener. The objective of this during the design phase was to verify if the given visualization would work when real data is charted.

Sometimes visualization researchers make the mistake of working with fictional data while designing a visualization solution. When real data is used as the input to the visualization it is possible it might not work at all due to problems on the encoding.

### 3.5. THE 4 LEVELS OF DESIGN

DSM defines four major levels of visualization design. These levels range from the highest level; the problem abstraction, to the lowest level; the computer algorithm used to compute the visualization. The levels are nested in a cascade-like fashion. This means the output of a higher level is used as an input on the level below.

- Domain situation: Problem is understood and requirements gathered
- Data/task abstraction: Both tasks and data are abstracted into a domain-independent version.
- Encoding/interaction idiom: The visualization is designed and its interaction defined.
- Algorithm: The computation algorithm for both the data and the visualization is implemented.

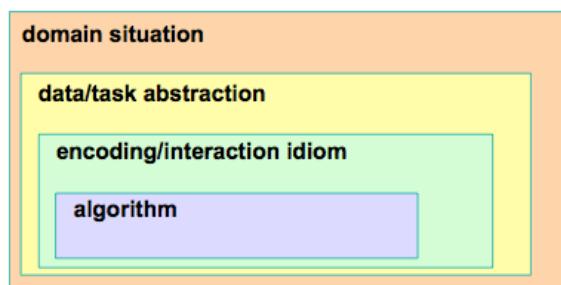


Figure 9: The 4 levels of design

Inevitably making a wrong decision on a higher level cascades to lower levels (Munzner, 2014). The positive aspect of this approach is that you can analyse whether each level has been addressed correctly and amend it if needed to refine the visualization design.

#### Level Validation

Each level has different threats that the researcher has to overcome to design a successful visualization (Munzner, 2009) (Munzner, 2014). Threats can be understood as wrong decisions taken on each design stage and that end up producing an inappropriate visualization for the problem we are trying to solve. Specific threats need special attention from the researcher at each level and DSM offers different methods for validating each stage to prevent these from happening.

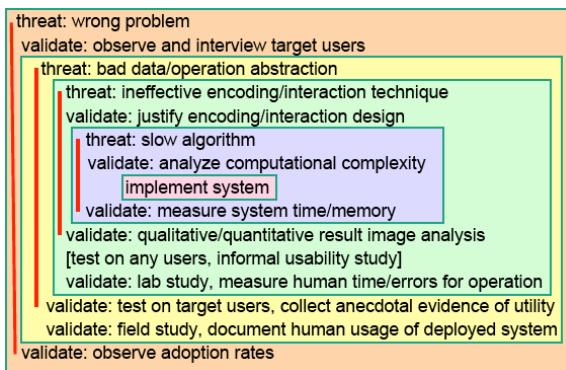


Figure 10: Validation diagram of the 4 nested levels of design

Although this nested model suggests a sequential validation in a *downstream* fashion the DSM allows going back to a higher level and re-evaluating the design. Moving *upstream* (back to a higher level) is possible and sometimes required when validation fails at a lower level and changes need to be done to the visual encoding, abstractions or understanding of the problem.

### 3.5.1. DOMAIN SITUATION VALIDATION

The correctness of the problem abstraction is validated at this level. This is tightly related to the Discovery phase where the researcher learns about the domain and extracts the particular features of the problem. The desired outcome of the validation is a clear understanding of the needs of the user and the ability to translate these needs into questions about what the user actually needs.

### 3.5.2. DATA/TASK ABSTRACTION VALIDATION

Considering the problem is understood correctly and the researcher has a clear idea of the needs, the next level of validation is the why-what abstraction level. This can be seen too as the task-data abstraction level (note here the idiom is missing as on the task-data-idiom as this belongs to the next validation level).

The abstraction of the needs into tasks and data needs to be validated with domain experts to make sure the correct needs are being fulfilled. Abstracting from a domain-focused language such as the one used in the higher level allows the researcher to make sure this translation is correct when verified with domain experts.

Validating this level means making sure the tasks that have been elicited are correct and that the data that is to be shown to users is the right data too (Munzner, 2014).

### 3.5.3. VISUAL ENCODING AND INTERACTION IDIOM VALIDATION

This level is intended to verify the design choices are correct in fulfilling the tasks. DSM encourages considering not only the visualization per se (what is known as the visual encoding idiom) but also how users control what they see (the interaction idiom). Information visualization in its deepest couples visualization and interaction to enable users to perform their tasks. The best interaction is that where the user forgets about the system and only

focuses on the tasks in hand (Fekete et al., 2008). This makes both idioms to best be evaluated and validated together rather than separately.

This validation level is important as the design space for visualization is big and when we include interaction to the mix it becomes enormous. DSM encourages following the nested model strictly by identifying and validating accurately the task and data abstractions on the previous level. This allows the designer to narrow down the range of different visualization and interaction idioms to a smaller set of feasible solutions at this level.

### 3.5.4. ALGORITHM VALIDATION

DSM differentiates the design of the visualization and interaction idioms from the computational issues related to its implementation (Munzner, 2014). In some way though, these two levels are related, as some interactions might not be possible if the computation of the result takes minutes or hours rather than seconds.

All throughout the design and implementation of the visualization the author has taken into account computational costs of the design choices, always discarding the options that required too much computational cost.

DSM encourages measuring the computational cost in a practical way (e.g. measure the processor and memory time of the application). The author discarded this validation, as the software being developed was not a final version but an experimental prototype therefore performance measures had little to no importance.

## 3.6. THE 4 LEVELS OF DESIGN IN PRACTICE

During all the timespan of the project the researcher spent time acquiring more information and learning about the domain. This not only happened on the initial stages of the project but all along. The fact that stock-picking is a very complex process made the validating task abstractions, domain assumptions and design decisions a constant during the duration of the project.

Initially the researcher had access to only one domain expert and all validation levels were effectively validated with the help of this only unique expert. Yet, during the duration of the project more domain experts were made available to the researcher which could help on refining this validation by providing feedback on the different levels.

## Process Overview

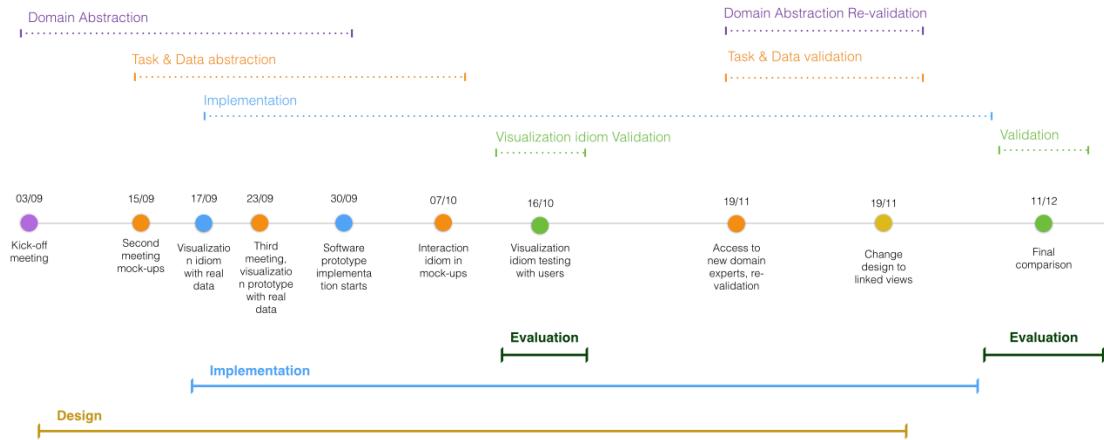


Figure 11: Overview of the entire project

### Initial domain situation abstraction

The initial kick-off meeting followed with another meeting that took place on September the 15<sup>th</sup>. The objective of this meeting was to validate that the problem abstraction was clearly understood by the researcher. The researcher presented a set of mock-ups with sketches of a set of different visualizations that could be used to support the tasks abstracted from the first meeting. On this presentation the researcher elaborated on what tasks could the prototype support and what was the main objective the tool should be used for.

Creating a mock-up of the different visualizations enabled the researcher to better communicate the task and data abstractions as well as the possible visualization idioms designed. Instead of focusing only on tasks and data, having the complete set of triads developed helped validating the different levels of design. During this meeting, feedback was gathered about the problems detected by the domain expert on the domain abstraction elaborated by the researcher.

The validation of the tasks and data took place over the two initial meetings until the third meeting where the visual encoding took a more important role.

### Visual Encoding Validation

Once the initial task and data abstractions have been validated with the domain experts, the focus was placed on validating the visual encoding and interaction.

At this stage, using Keynote as the main prototyping tool was not proving to be enough. The complexity of including interaction to simulate the interaction idiom was too high. This, together with the fact that task/data abstractions were stable enough, suggested the researcher could move into implementing the software prototype. In fact, in design studies the implementation process is tightly interleaved with the design process, starting the implementation at this point was the next logical step defined on DSM (Sedlmair et al., 2012).

### 3.7. IMPLEMENTATION

Although Thomson Reuters Eikon is a complex financial platform conglomerating many different technologies there is a movement towards web-based technologies. This is due to the need of traders and investment professionals needing a constant connection to Eikon from their laptops, tablets and smartphones (Thomson Reuters, 2013). This fact, in addition to the wide experience of the author in web-based technologies (HTML, CSS and JavaScript) made a browser-based visualization prototype the preferable solution.

There are a great number of JavaScript-powered visualization frameworks available (DATAVISUALIZATION.CH, 2014). The author researched what frameworks could be a better fit for the solution. Initial evaluation considered the following open source frameworks: D3.js, Processing.js, Degrafa and Raphaël. This evaluation showed that D3.js was the best candidate in terms of functionality and usability. One of the key factors was the higher number of plugins and sample code (Bostock, 2015) that could help the author implement a working prototype within the short project time frame.

### 3.8. USER TESTING FOR VISUAL ENCODING VALIDATION

Due to the fact that some visualization idioms chosen to support the abstracted tasks were quite novel, especially to the financial domain, the researcher chose to perform a testing with users. The objective of this testing was to validate the visual encoding idiom to make sure the visualization and interaction could be understood and the tasks achieved.

Performing an evaluation during the design process is a common practice (Lam et al., 2012) and is also encouraged in DSM as it can be used to inform the design before its implementation and deployment. A common way to test and validate the visualization idiom does work is performing a lab test.

#### DECIDE Framework

To perform the user evaluation the DECIDE framework (Preece et al., 2002) was used as a guide to create a correct evaluation study. This framework helps novel user evaluators on elaborating user tests in a structured manner by focusing on the main points to guide an evaluation correctly. The framework helps by placing a focus on the following points:

1. Determine the overall goals the evaluation addresses
2. Explore the questions that need to be answered
3. Choose the techniques that are needed to answer the questions
4. Identify practical issues that need to be addressed (lab availability, participant selection...)
5. Decide how to act towards ethical issues
6. Evaluate, interpret and present the data

## **1) Determine the overall evaluation goals**

The objective of the evaluation will be to test the efficiency of the visualization idiom. To do so both qualitative and quantitative measures will be collected during the test.

This evaluation is needed to validate the visualization idiom and confirm that it can be used on the visualization prototype to support the user needs.

Additionally, the goal of the evaluation is to highlight the major usability problems on the prototype that might cause the tasks not to be fulfilled. These usability errors might be causing the interaction idiom work incorrectly.

## **2) Explore the Questions that need to be answered**

The question that needs to be answered is: Does the envisioned visualization idiom work effectively to enable users perform the abstracted tasks? Answering this question will validate the devised visualization or invalidate it, thus requiring a re-design.

## **3) Choose the techniques and paradigm**

To validate the visual encoding, a series of different tasks involving the use of the visualization idiom where designed and various metrics captured. These tasks would be performed with a functional software prototype with enough capabilities to enable users on fulfilling their tasks.

## **4) Practical Issues**

### **Equipment & facilities**

Thomson Reuters has a really well equipped usability laboratory that is specifically used to perform usability tests. The Research and Strategy team manages such facilities and was able to give the researcher access to them.

### **Users**

Access to final users during this project is not feasible. Domain experts act all throughout the project as proxies for final users. Although this might seem an impediment in validating the visualization idiom, it might not be such case. Visual encoding and interaction idiom validation concerns the ability of users in performing tasks with the visualization. This concerns testing the user's abilities in terms of the visual perception and memory by measuring error rates and performance, something that is independent from the knowledge of the field. Thus there is no need to have final users test the visualization solution with users that are familiar with the process of stock-picking.

The Research and Strategy team organizes a series of usability testing sessions where volunteers from all over the organization can join and help the designers test their creations. The researcher was given the possibility to use one of these testing sessions with the visualization prototype with users.

## **5) Ethical Issues**

The users being tested are volunteers on this process. In any case their names are asked nor are considered important for the testing evaluation. Before the testing takes place volunteers are informed of the test and a participation sheet given explaining how the test works and how will they be informed. No other ethical issues are important in this testing procedure.

## 6) Data

This test should measure the impact of specific design choices by measuring quantitative indicators such as task performance, error rates and/or qualitative indicator such as user opinions on a series of tasks selected by the researcher (Munzner, 2014).

The tasks being tested were abstracted from the requirements. By doing so we were testing how users could use the visualization idiom to fulfill the requirements.

The test was more focused on measuring qualitative measures rather than quantitative ones. Still, the error rates were an important measure used to validate the fact that the visualization could effectively be used to fulfill the tasks. Qualitative measures such as user feedback and the user's perceived difficulty on fulfilling the task would be used to validate or reject the visualization idiom.

The qualitative data would be evaluated and used to improve the design of the prototype.

The quantitative data should be used to validate whether the visualization could be used to fulfill the tasks (correct idiom) or it failed on enabling users in fulfilling their tasks (incorrect idiom).

### Performing the test

The test took place the 16<sup>th</sup> of October and 3 participants were invited to take the test, feedback from users was gathered and used to validate the visualization idiom.

## 3.9. RE-VALIDATION OF FUNCTIONALITIES

The researcher was given the opportunity by the work supervisor to promote his work within the company intranet, in a group dedicated to Data Visualization. The researcher created an internal blog entry that was published on October the 27<sup>th</sup>. This blog entry had a really good welcome raising interest on the current piece of research.

Additionally, the researcher was invited by the work supervisor to participate on a visualization-related conference call that takes place every two months and given the opportunity to share the work of this research too.

These promotion efforts that were facilitated by the help of the work supervisor enabled access to a wider number of domain experts. Now the researcher could contact these domain experts to validate the prototype being built.

With the objective of validating the functionalities implemented on the prototype the researcher interviewed seven different domain experts with the objective of gathering feedback this validation took place on the week of November the 19<sup>th</sup>. During these interviews the domain expert was taken through each of the features of the prototype. Feedback and suggestions for improvement for each individual functionality was recorded.

## 3.10.FINAL COMPARISON TEST

To objectively answer the research question we need to measure if the prototype does help analysts on developing insight and new knowledge relevant to the stock picking process. To do so we will use the current Equity Screener that does not have visual data mining capabilities to be used as a baseline to compare to.

Evaluating visualization tools is complex because such tools normally are used in complex analytical processes. Evaluating how a visualization supports analysis and reasoning is even more difficult due to the fact that visual reasoning is an ill-defined process (Lam et al., 2012). However Lam et al., 2012 suggest a specific evaluation approach that can be used to objectively measure the degree in which visualization can support this process. This evaluation approach (or scenario if we are using their suggested name) was coined as Evaluation of the Visual Data Analysis and Reasoning (VDAR).

By evaluating the VDAR of both the old and the new tool the researcher should be able to infer how both visualization tools support the generation of knowledge and insights on the data (Lam et al., 2012). Because evaluating the VDAR involves testing the entire workflow and how well the tool supports discovering new insights on the data, it requires stable and reliable software to be able to evaluate it properly. Usability problems will block the user in executing the workflow properly and should be solved before being able to evaluate the VDAR.

Usually to evaluate VDAR we will define tasks that measure how well the tool supports the following processes (Lam et al., 2012):

- *Data Exploration:* How well the tool supports seeking information, searching, filtering, reading and extracting information.
- *Knowledge Discovery:* How well it supports schematizing information and re-analysing theories.
- *Hypothesis Generation:* How it supports the ability to generate and test hypothesis interactively.
- *Decision Making:* How it enables the communication and application of the analysis results.

The best approach to evaluate VDAR is to do it on its intended environment with real users and real tasks (Lam et al., 2012). Unfortunately it is not possible to evaluate the prototype on its real environment because the access to customers is too complex for testing an early prototype. What can be done instead is test the VDAR capabilities of the prototype against the current Equity Screener in a controlled environment.

For the evaluation of VDAR we will be testing both the old Equity Screener against the new prototype. We will thus be testing how well each screener handles VDAR-related tasks.

### 3.10.1. DECIDE FRAMEWORK

#### 1) Determine the overall evaluation goals

Does a visual data mining approach improve stock picking for equity analysts?

The objective of the evaluation will be to measure the Equity Screener compared with the prototype in terms of VDAR capabilities.

This evaluation is needed to test the research question.

## **2) Explore the Questions that need to be answered**

We need to translate the goals for the evaluation into questions that can be answered through the evaluation (Preece et al., 2002). The main question that needs to be answered is: What tool provides a higher degree of VDAR capabilities?

## **3) Choose the techniques and paradigm**

To test the hypotheses both qualitative and quantitative data will be collected during the different lab sessions. By collecting both types of data the researcher should be able to validate the results in a more precise manner than compared with only measuring a type of indicators. The task design for this evaluation test is complex and will be explained on section Designing the tasks.

## **4) Practical Issues**

### **Equipment & facilities**

As with the previous user evaluation, the Research and Strategy team at Thomson Reuters provided access to the user testing facilities available on the company.

### **Users**

Domain experts will be evaluated instead of final users. This is due to the fact that access to final users is not possible in the current project.

The researcher managed (with the help of the work supervisor) to have three domain experts available for this comparison testing. These users were specialists with a high degree of expertise in investment banking and portfolio management.

### **Testing conditions**

An evaluation looks to answer a given hypothesis we are testing against; these can be seen as relationship between two events, also known as variables. Hypotheses are tested by manipulating these variables, and evaluating whether users perform better in one condition or another (Preece et al., 2002).

In the current evaluation two independent variables will be manipulated. One is whether the evaluation tasks will be performed using the Prototype or the current Equity Screener. The other independent variable will be the dataset being used to fulfil the task. Two different yet similar datasets will be used to test the task.

When multiple independent variables are used (Preece et al., 2002) participants need to be distributed in such a way that all conditions are tested. If this is not possible the distribution must be done in such a way that most of the conditions are tested[Appendix: Task

Distribution]. In such cases counter-balancing is important: make sure that the order of tasks does not affect the result.

## 5) Ethical Issues

Before proceeding to the testing participants were informed about the nature of what they were being tested against and how by providing a participation information sheet. Additionally an informed consent document was also handed out to each participant too.

## 6) Data

During this test we will be collecting both quantitative as well as qualitative information. Both types will be used to inform the comparison.

The quantitative data we will be measuring is both the error rates (based on the questions that need to be answered) and the time performing each task.

### 3.10.2. DESIGNING THE TASKS

Evaluating visualizations is a complex undertaking, not only we are testing the tool itself but the complex processes that the tool supports (Lam et al., 2012). The tasks need to be designed in such a way that enable us to compare how the two tools support common activities performed by the potential users. Obviously we should be testing tasks that can be performed in the two tools; otherwise no comparison would be possible.

These activities in fact have been elicited on the initial part of the project in the form of use cases. From these use cases, requirements were extracted. And these are, in fact, the features that are required by equity analysts to better support the stock-picking process. The objective of the task design is thus designing tasks that can test how well both tools support the domain expert's requirements.

#### Training tasks

When comparing a well-established piece of software against a new solution, it is likely that users will be much more familiar with the current software rather than with the new solution. This might skew the final results (Carpendale, 2008) (North et al., 2011). In this case it is recommended to provide some training on the visualization such that users become more familiar with the new tool and the evaluation is more balanced.

As the prototype is a new tool for the users being evaluated, a set of training tasks will be devised in order to familiarise users with the prototype functionalities. This should be done before performing the set of measured tasks.

#### Benchmark Tasks

Benchmark tasks are those designed to compare how efficient is the user when utilizing the visualization to achieve the results he wants. Benchmark tasks are the most usual ways of evaluating visualizations (North, 2006).

What we are looking for is measuring simple tasks and compare the metrics recorded when using both visualizations. Simple benchmark tasks are those with the following objectives: to compare, contrast, associate, rank, cluster, correlate or categorize data (Carpendale, 2008)

Although benchmark tasks are the most usual method to evaluate visualizations (P. et al., 2005) they have some fundamental drawbacks (North, 2006):

- They are predefined, so there is no degree of freedom for the user to discover new features of the data.
- They don't provide enough time for the user to experiment with the visualization.
- No qualitative insights on the data are evaluated, just quantitative.
- They require simple answers
- Might focus on tasks that are not on the typical user workflow
- They make users

Whilst benchmark tasks can help comparing how easy is to actually perform a specific task it will not prove whether the prototype of visualization helps to take more informed decisions on the stock picking process. Moreover, when evaluating using benchmark tasks we are assuming that more complex tasks build on these simple tasks. The difference between these simple tasks and more complex workflow tasks can be quite large, so a different method for measuring VDAR capabilities is needed.

## Defining insights

There is no formal definition for what is an insight, probably because it might fall too short or be too general to be useful (North, 2006). Although insights are not formally defined, insights are:

- *Complex*: There is no limit on what amount of data is required (maybe all the data or maybe a small set is enough).
- *Deep. They build over time*.
- *Evolving*: Can evolve into further hypotheses which once answered can generate more insight.
- *Qualitative*. Not exact, can be uncertain and subjective.
- *Unexpected*. Unpredictable, serendipitous, and creative.
- *Relevant*. They connect the data with the domain. It is more than simple data analysis as it gives truly relevant meaning to the data.

To better evaluate insights that are the key on the VDAR process, benchmark tasks are not enough; we need a more sophisticated way to measure them.

## Medium-insight tasks

While benchmark tasks evaluate essential analyst procedures, they do not 'provide a firm basis for supporting the kinds of knowledge-making activities that people seek to perform every day' (Amar & Stasko, 2005). An evolution of such tasks is what we will know as

*medium-insight* tasks. These types of tasks add a degree of uncertainty to benchmark tasks. By adding some uncertainty, we are allowing a higher degree of freedom for the user that will allow us to evaluate how he understands the data (North, 2006).

To do so, we'll force users to interpret the visualization by designing more complex cognitive tasks. Amar & Stasko define a set of complex prototypical activities analysts perform when exploring the data and that need to be evaluated to measure the degree of insight: developing an understanding of data trends, uncertainties, and causal relationships, predicting the future, or learning a domain (Amar & Stasko, 2005).

By designing a set of tasks that test these analytical activities we will evaluate the level of understanding of the users' mental model of the visualization, as users need to understand it clearly to perform such tasks correctly (North, 2006).

Measuring such insights is a more difficult method compared to simple benchmark tasks. It requires creating more complex tasks, longer tasks and greater variability in task times and correctness. An additional step in evaluating VDAR is eliminating benchmark tasks completely.

### **Completely insight-based**

Medium-insight tasks still have some of the issues of benchmark tasks but to a lesser degree (North, 2006). Yet, for a deeper evaluation of the usefulness of the visualization a more suitable solution is to use a completely open protocol and observe how users gain insights on their own. With this approach, researchers do not participate on the tasks they rather observe what insights the users get from the data. In this kind of evaluations, users are instructed to explore the data (normally with a small set of initial questions) and report their insights until they feel that they have learned all that they can from the data.

### **Mixture of different evaluations**

For the purpose of evaluating the visualization a mixture of the three different approaches will be used. This will enrich our evaluation by bringing all benefits of the three approaches. Benchmark tasks will enable us to highlight differences between the two products at a low level; the other evaluations will provide us with a higher knowledge on how users understand and interpret the data and how they facilitate insight generation (North et al., 2005).

When performing a combination of the three approaches it is recommended to have the more open-ended tasks first and finish with the less open-ended portion. This way the user will not become constrained by the way of proceeding on the benchmark tasks when moving to the open ended part (North, 2006).

### **Final evaluation**

The final comparison test took place in the week of the December the 11<sup>th</sup>. Three domain experts were called to perform the comparison testing using this mixture of evaluations, testing the different tools.

## 4. RESULTS

The results of the different phases of the DSM study together with the results of the evaluations that took place during the project are included in this section. The results are presented from the initial requirements to the final evaluations of the visualization solution.

### 4.1. REQUIREMENTS

#### 4.1.1. USE CASES

Listed below are the use cases derived from the domain experts. These use cases are the result of multiple validations of the data abstraction. These final use cases define activities portfolio managers or equity analysts perform and that need to be supported by the visualization. These use cases have been refined all throughout the project with the input from the different domain experts.

##### **UC01: Similar Stocks & Stock Comparison**

As an investor want to find a similar stock to the one I have in my portfolio and has gathered good results. For example, I might have AAPL stock and want shares of a company with similar results. This also means I should somehow be able to compare a small number of different stocks.

##### **UC02: Monitor Portfolio**

I want to **monitor** the status of the stocks that are in my portfolio and its tendency. I want to see if their price is dropping, if the company is doing great or poor (I want to control the stocks I have and add any stocks I want to the visualisation to see which are likely to drift).

##### **UC03: Value investing highlights**

I'm looking for companies that are undervalued so I set a range of filters to remove the companies that are overvalued. I want to see the companies that best match the filters I've set highlighted in the visualisation. This requires the idea of strength defined for the company, how well does it match the filters.

##### **UC04: Equities Grouping & Differentiation**

I want to see how the different sectors are performing, there is no grouping now whatsoever so I want to know what industry are the stock related to. I might also want to see the country of the stock for example to differentiate the stock and make sure I can actually consider it as an option for investment.

##### **UC05: Equity Exploration**

I want to be able to explore the equity space by applying filters, selecting stocks, querying for their characteristics to find a equity which is worth investing. The tool should allow me to discover new insights of the dataset, enable hypothesis generation and testing and improve the decision making on what equities are worth investing.

##### **UC06: Industry and overall Mean**

When evaluating equities I want to see a company or a set of companies behave in terms of the industry. For that as an equity analyst I normally compare the values against the median

(or the average) for the overall or the industry. Maybe I can even use a index such as the S&P or a given stock as the baseline for comparing.

#### **UC07: Adding an equity**

I'd like to be able to add an arbitrary equity to discover how it compares to the set I've filtered down. For example, if I've got a set of different oil companies on screen and I've currently have an oil company in my portfolio, I'd like to know how it compares to the selection, maybe because the company I own is better to the selection I have or maybe it is indeed worse.

### **4.1.2. FINAL REQUIREMENTS**

Listed below the reader will find the full list of the requirements that have been extracted from the use cases. A link is provided to what specific use case was used as the base as well as a rationale derived from the use case.

#### **ID: FR01**

**Description:** The visualisation should enable **comparing multiple stocks** in a range of **different metrics**.

**Rationale:** The user wants to know how different companies compare one to each other, and not only in a specific metric but a set of metrics.

**Use case:** UC01 (Similar Stocks and Stock Comparison)

#### **ID: FR02**

**Description:** The visualisation should allow displaying **trends** that show the movement of the equities.

**Rationale:** As an investor I want to see how my portfolio is doing, or how a given equity is getting on through time.

**Use case:** UC02 (Monitor portfolio).

#### **ID: FR03**

**Description:** The visualisation should display the company / companies that **ace** the specified filters. That is, the companies with the highest strength should be clearly visible.

**Rationale:** The equity space can be really high, so a fast scan should enable users to easily discern what equities better match their filters than others that do not match these so appropriately. The calculation would be done via quants (this is a method that assigns a rating value based on a weight for a set of attributes chosen for the equity, it is effectively the "Quick Rank" of the old screener).

**Use case:** UC03 (Value investing highlights)

#### **ID: FR04**

**Description:** The visualisation should allow analysts to spot **outliers**, that is, companies that are behaving differently, in one or more of the metrics defined.

**Rationale:** The outliers on the visualisation effectively represent investment opportunities for the equity analyst as they are companies that are behaving differently.

**Use case:** UC03 (Value investing highlights)

#### **ID: FR05**

**Description:** The visualisation should allow **grouping** companies by **industry, country** or some other **categorical** attribute.

**Rationale:** Group stocks by industry to be able to see what interesting patterns / groupings can be found. This highlighting allows to better understand the data in terms of how equities are behaving.

**Use case:** UC04 (Equities grouping and Differentiation)

**ID: FR06**

**Description:** The system shall allow the user to **add** an arbitrary stock to the visualisation and display why it does not match the filters (highlight specific filters that are failing).

**Rationale:** The user wants to know why a filter does not actually match the filters defined on the first place when he expects it to appear. This requirement should allow the user to answer the question: Why my stock (e.g. IBM, AAPL...) does not appear on my filtering?

**Use case:** UC07 (Add an equity).

**ID: FR07**

**Description:** The visualisation should allow **selecting** what equity **groups** are to be displayed.

**Rationale:** I want to be able to focus on specific industries that seem interesting, so I need an easy way to exclude groups that might not seem not interesting.

**Use case:** UC04 (equity grouping)

**ID: FR08**

**Description:** The visualisation should allow me to **locate** a specific share.

**Rationale:** I want to see where a company is actually positioned compared to other companies, because I might want to track it or use it as a reference.

**Use case:** UC05 (Data exploration)

**ID: FR09**

**Description:** The visualisation should allow me to compare how a given company compares to the **median** of the industry.

**Rationale:** As an equity analyst I consider the relation of a given equity with the industry or the overall set of equities available.

**Use case:** UC06 (Industry and overall mean)

**ID: FR10**

**Description:** The visualisation should enable **filtering down** progressively to reduce the universe of stocks into a more manageable set.

**Rationale:** As an equity analyst I start considering a wide universe of stocks and by playing with the filters and discarding equities I'm not interested in. I reduce the number of equities I want to invest in.

**Use case:** UC05 (Data Exploration)

## Non-functional requirements

Although only one non-functional requirement was elicited, it is of upmost importance. Data exploration in the current Equity Screener is made difficult due to the fact that the filtering process lags back the process by performing such tasks really slowly.

**NFR01:** Any filtering or categorising actions should be almost immediate

**Rationale:** The system shall allow exploring the data in a seamless way. The experience of

exploring the data shall be immediate; any modifications on filtering the output data should be immediate to ensure this.

**Use Case:** UC05

**Fit criteria:** The most common interactions should not take more than 5 seconds. By common interactions we understand operations such as retrieving company data, the company trend.

## Summary requirements analysis

The process of gathering requirements from the domain expert was done initially via unstructured interviews where the researcher was able to acquire a greater understanding of the domain problem. Thanks to these initial meetings the researcher produced a solid number of use cases from which to extract the requirements.

Additionally, the process of re-validation of the prototype features that occurred halfway through the project was useful to trigger new use cases and thus new requirements.

## 4.2. DATA

The **what-why-how** triad framework recommends focusing first on what data we can use in our visualization. The Equity Screener works with data that is organized in a table-like fashion. DSM defines how data is organized as datasets. The dataset we will be working basically takes the shape of a **table** datatype. In this table, we'll have stocks as items and a series of economic indicators as its attributes.

### Data Abstraction

The objective of the data abstraction is determining if the dataset we are using needs to be transformed visualized. We are dealing with stock data, each stock item we are going to visualize has many thousands of attributes related to the company it belongs to. We will have an overview of the shape these attributes can take.

- **RIC (Reuters Identifying Code):**

The identifier of the stock, it can uniquely identify the stock on the table.

- **Ratios and financial indicators:**

The majority of ratios that we can filter on the screener are generally quantitative i.e. numbers that can be compared but do not have a specific ordering (in contrast with ordinal data (Munzner, 2014)). There are some cases that are interesting such as the case of the Quick Rankings. Given a set of economic indicators we can order the companies on the table by how they do in terms of these indicators. The better they do, the higher the ranking.

Other attributes can be time-based (but instead of being continuous they are sampled by the Equity Screener). This specific type will be evaluated on the next section.

- **Categorical data:**

The stock might have some categorical attributes associated to it. For example there is the "Economic Sector" attribute used to classify the stock within a specific sector.

Some categorical data indicators might be hierarchical. Thomson Reuters has a classified industries by their economic sector, industry and many other categories. Such classifications normally are organised in a hierarchical way. For example, Shell Inc., the oil company in its 'Industry' category follows this hierarchy: 'Energy' Economic sector > 'Fossil Fuels' business sector >'Oil and Gas' industry group > 'Oil & Gas Exploration and Production'.

Quantitative data that depends on time (time-based) is important for some of the tasks elicited so we'll be focusing more on it.

### Time-series datasets

Some attribute values depend on time. When they are processed by Thomson Reuters and stored for the use of the Equity Screener they are sampled and stored (the smaller unit of time the data is stored for is a single Day). When the Equity Screener wants to show this value it requires the user to sample it to a specific point in time. For example, 'Market Capitalization' is an indicator that varies every moment because it depends on the price the company is being traded at, still to be able to display it on the table we need to specify what point in time we want the value for.

Two main types of data can be found, the data that is specific to a period such as a Fiscal Year, and the data that does not depend on a period (does not refer to a Fiscal year/quarter...). To use this data it still needs to be sampled and specify a Day, Week or Month to obtain its value.

	Period-dependent	Period-independent
Description	This data always references a specific period and can't be considered without a period.	This data can be considered individually, without a specific period.
Time units	Fiscal Year, Fiscal Semi-annual, Fiscal Quarter, Fiscal Interim...  For estimates: Next Twelve Months (NTM), Next Fiscal Quarter...	Last trading Day, Last Week Day, Week, Month...
Example	Income of the company in a Quarter 'Revenue Smart Estimate NTM': The estimated revenue for the next twelve months.	'Company Market Cap' which depends on the volume of shares being traded and the price of the stock.

### Period-dependent

Some data is only available at the end of a given company's financial period, such as 'Yearly Income'. Some other data is available every quarter 'Quarterly Income'.

Yet, some other data is projected on the future because it refers to estimates of the indicator. Thomson Reuters calculates forecasted values for the most important indicators of a company. These are values that evolve through time but are not continuous. Their values are relative to similar units to normal period dependent metrics.

### Summary data analysis

The data the Equity Screener works with is really vast and can be quite complex. Following the DSM methodology we can perform a top down analysis of the data by evaluating how data is organized, what types of data we are working with and judge whether we need to modify this data model in order to fulfill the tasks.

#### 4.2.1. TASK ABSTRACTIONS

The DSM approach encourages abstracting requirements tasks into high/mid-level tasks into lower level task definitions, part of the what-why-how triad. These tasks are pairs of: {action, target}. The different tasks extracted from the requirements can be found on the Appendix [Appendix: Abstracted Tasks].

### Summary of the task analysis

The translation from use cases to requirements and then to data and task abstractions can be complex at times. Even if requirements are defined granularly it can be difficult to abstract it into a task into such low level.

Although this translation was difficult it helped a deeper analysis of each requirement. This is due to the need to defining the requirement so well that it can be translated to a low-level task.

Additionally not all requirements can be abstracted into tasks and data as sometimes they only refer to manipulations of the data.

### 4.3. DESIGN AND IMPLEMENTATION

#### 4.3.1. CURRENT VISUALIZATION

The current Equity Screener only includes a scatterplot to display the equities resulting from the initial filtering. The other option available on the screener is displaying the data in a table format; users normally switch from one format to the other to evaluate the data.

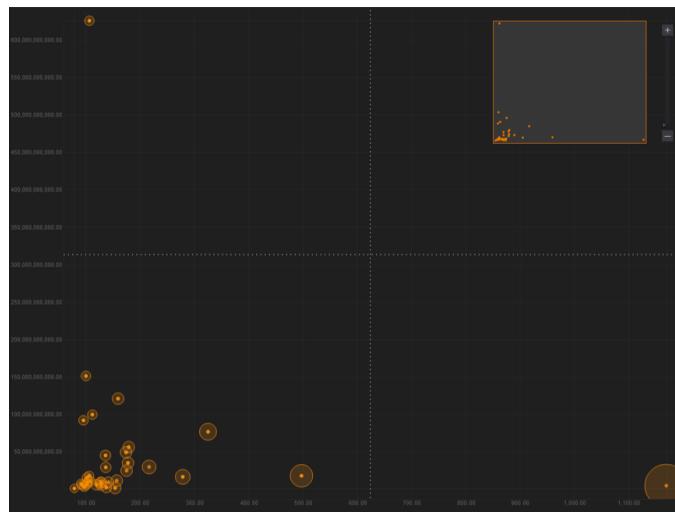


Figure 12: The scatterplot on the Equity Screener

Thanks to this visualization users can see three dimensions for each of the equities: the user can choose what is charted on the X-axis, Y-axis and what data is used for the size of the bubble. Up to 5000 equities can be drawn using this approach, although this seems a technology related limit, when a great number of equities are charted the visualization becomes too crowded and difficult to understand.

#### 4.3.2. TASK ENCODING

Following the DSM methodology each of the tasks abstracted should be analyzed to solve the “how” part of each task triad. For the researcher it was in fact easier to focus not only on individual tasks trying to find an encoding and interaction for each one independently, but to focus on the overall set of tasks to think on what would be the best visualizations to that could be used to fulfill the tasks.

In fact, in this section DSM was not followed strictly, it was rather used as a guideline to knowing what could be the different options to encode the data as well as what interaction idiom could be used.

#### 4.3.3. INTERACTION CODING

At a given point in the project, once the visualization idioms where validated by the domain experts, the interaction idioms was much more complex to validate. Using Keynote some basic interactions could be validated but the more complex interactions where really difficult to simulate. This was the point where implementation started, not because everything on the design was definitive but because the interaction idiom needed to be tested with the prototype itself.

#### 4.3.4. DESIGN FEATURES

##### Representing multiple dimensions

The normal workflow of an equity analyst consists in evaluating a considerable number of attributes for a given company. This number is variable but is normally much higher than only

the two or three dimensions we can actually draw with a scatterplot. A new method to display multiple dimensions in a comprehensible way was needed to better support this evaluation (this is related to FR01 & FR04).

An option to chart multiple datasets is to use dimensionality reduction so that data can be charted using 2D scatterplots, interactive 3D scatterplots, or Scatterplot Matrices (also known as SPLOMs) (Sedlmair et al., 2013). Dimensionality reduction can be used to create a lower-dimensional version of the data which preserves much of the information available yet preserving the most important features of it. Although this at first seemed an interesting approach the high dimensionality of the data (easily more than 10 indicators can be considered for each company) made this unsuitable.

A different approach was required to chart highly multidimensional data without having to visualize a version of the data with fewer dimensions. At this point various works were studied to see how multidimensional data can be charted. Long's thesis pointed out that the most common way to chart multidimensional data is using Scatterplot Matrices or by using Parallel Coordinates (Long, 2009, pp.4-6). Another interesting report evaluated 11 different techniques to chart multidimensional data focusing especially on financial data (Marghesku, 2007). For each of these visualizations they highlighted the capabilities to fulfill data mining tasks such as detecting outliers, finding correlations or describing the data. Finally, the article that helped choose what was the most suitable visualization was Keim's and Kriegel (Keim & Kriegel, 1996) that highlighted Parallel Coordinates as a powerful visualization to chart multidimensional datasets and at the same time supporting a great number of visual data mining tasks.

Parallel Coordinates (Inselberg, 1985) allow displaying multidimensional data using lines. Each dimension is displayed as a vertical axis, known as coordinates and each data point is a line crossing each vertical axis on the chart. Parallel coordinates allow comparing between different companies as well as the detection of relationships between the dimensions (Marghesku, 2007). Normally, vertical dimensions can be reordered making it easier to detect such relationships.

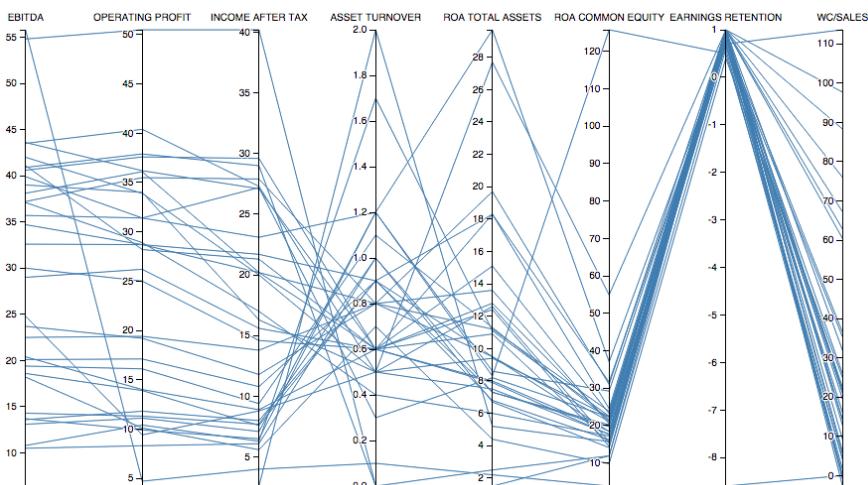


Figure 13: An example of parallel coordinates displaying 8 different dimensions

Parallel coordinates can implement a powerful interaction known as 'brushing'. Brushing is the most basic interaction yet powerful interaction and it consists in highlighting a set of lines

on the parallels to focus on this specific data (Siirtola & Raiha, 2006). Normally this selection is done in a single dimension at a time, but these highlighting can be overlaid so that the parallels will only highlight the data that complies with all selections. The power of this interaction relies on the fact that by simply brushing a dimension we are filtering down to focus only on the data lines that comply with the filters.

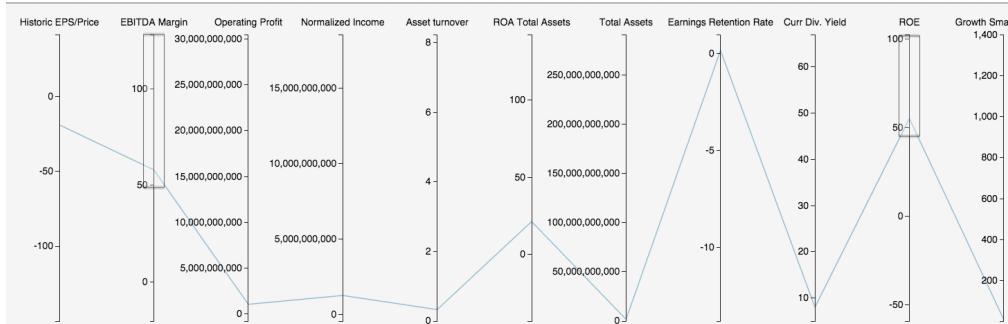


Figure 14: Brushing can be represented as drawing boxes on specific dimensions  
in this case highlighting data lines with 'EBIDTA Margin' >50 and 'ROE' >45

Parallel coordinates have been used in several different data exploration problems in and in different domains as a complementary or as a main visualization to explore and gain insight on the data (Kumasakaa & Shibatab, 2008) (Adrienko & Adrienko, 2001) (Xiang et al., 2012). However, they are not lacking of problems. One of the biggest problems is the difficulty to find relationships between non-adjacent variables (Huh & Park, 2008). The solution to this problem is reasonably easy: instead of using straight data lines, using smooth curves allows doing this task much easier.

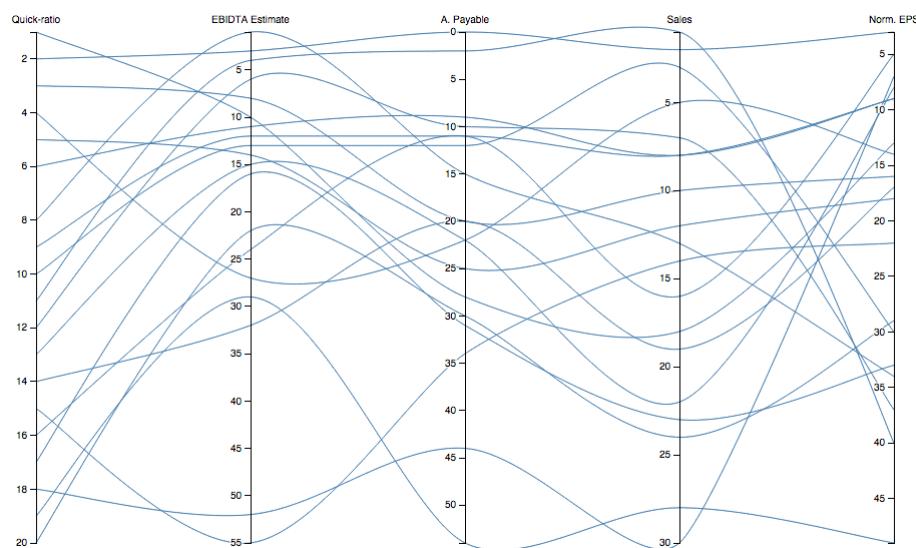


Figure 15: An example of parallel coordinates with smooth curves,  
note it is much easier to follow a curve rather than a straight line.

Parallel coordinates, with these modifications can fulfill a considerable number of the elicited requirements (FR01, FR04, FR10).

## Scatterplot enhancements

The scatterplot is a visualization equity analysts are used to (it is incorporated in the product and they are familiar with it). So the intention of the researcher was to enhance this visualization in order to incorporate the newly devised features. Additionally, scatterplots have some benefits such as they are able to show three dimensions in one chart and can be used to spot outliers and tendencies from a data set.

One of the requirements gathered was that the visualization should allow grouping equities by a specific criteria, for instance the industry or economic sector the company belongs to (FR05). If we evaluate the different ways suggested by Munzner on how categorical attributes such as these can be encoded we could use color hue or shape. Applying this reasoning to the scatterplot we could use hue or shape to show the different industries. Because the number of industries can be quite large, encoding different shapes did not seem an option so the hue was used as a way to encode this information. In fact, Scatterplots are often augmented by using color-coding to show additional attributes (Munzner, 2014, p.139). There is a limitation regarding encoding information using color, human perception can distinguish a small number of different colors (around seven), so the number of categories that can be encoded is considerably small.

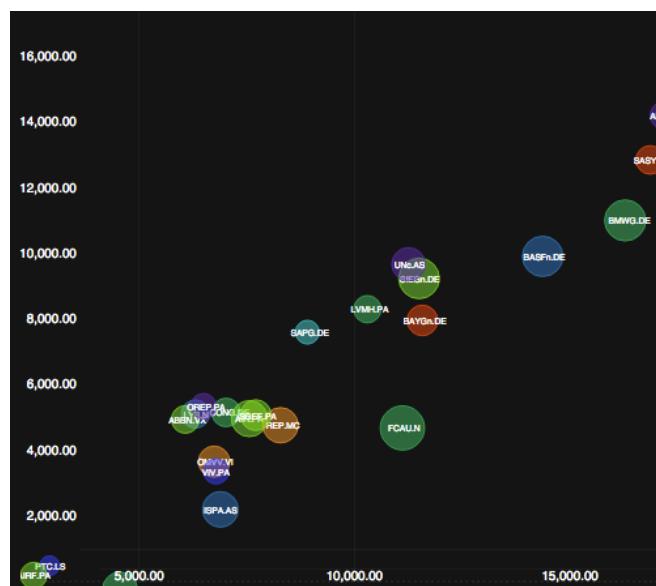


Figure 16: The scatterplot bubbles with different hue encoding the information of their corresponding category.

Generally, equity analysts do not consider equities individually, they compare their performance against their industry or other equities. To fulfill such feature the scatterplot the median could be charted in the scatterplot visualization in order to be able to compare how equities compare against the median (FR09). Additionally, the user can choose what industries to display on the scatterplot by simply checking/unchecking them, thus helping to focus only on a given industry removing the noise of the other equities (FR07).

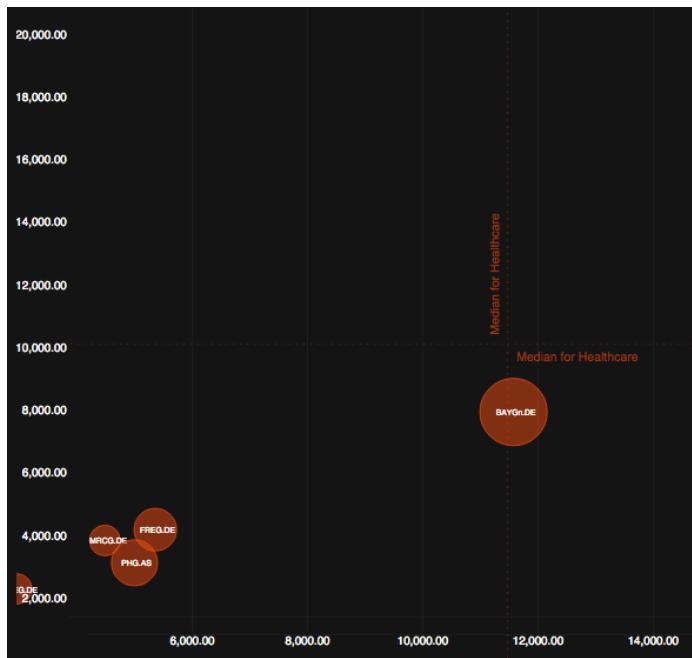


Figure 17: Only companies of the 'Healthcare' sector are displayed,  
note the median line too.

The process of enhancing the scatterplot does not finish here: more requirements can effectively be fulfilled by making modifications to the visualization. We could also find a way to show the trend of the equities using the scatterplot. However it now seems difficult to encode more information on this visualization as we already have four dimensions (x-axis, y-axis, size of the bubble, color). By evaluating the remaining encoding options we find this trend could be encoded by animating the bubbles. Animation is an extremely salient way of encoding information and must not be abused (Munzner, 2014, p.214), it should be used to draw the user's attention rather than encode information such as the trend of the equity.

Robertson et al. suggested a different approach to showing trends by using static depictions of the trends rather than using an animation-based approach (Robertson et al., 2008). They tested that this approach helps when comparing trends of different bubbles with higher accuracy and thus conveys trend information more appropriately than the animated counterpart. In practice they overlay the trends on the scatterplot. To do so, they suggest using a series of bubbles linked with a trend line, having the oldest bubble on the sequence with the lowest opacity and the newest with the highest opacity (because there is no animation is difficult to see what is the oldest). This approach has been slightly modified on the final design by labeling the bubbles with the year the data refers to.

Does a visual data mining approach improve stock picking for equity analysts?

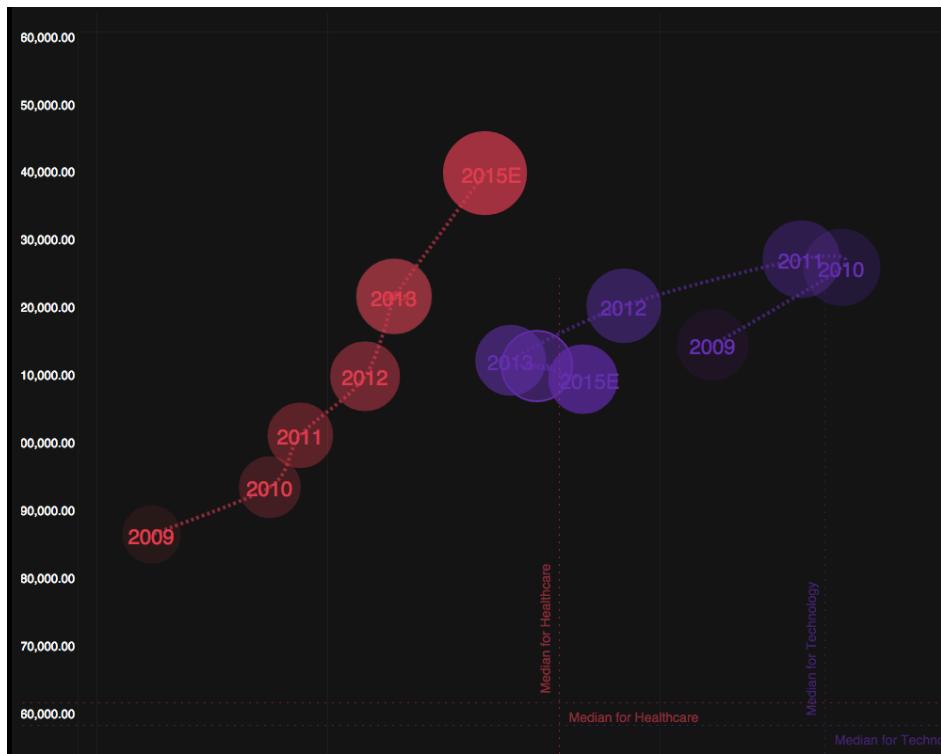


Figure 18: Depiction of the trend of two different companies, the suffixed 'E' indicates the value is estimated.

When scatterplots show a large number of items they appear too overcrowded and they don't convey information as efficiently. To avoid this from happening the researcher used a clustering algorithm to hide the more dense areas of equities into artificial groups. This was designed to help reduce the noise and make it easier to spot outliers.



Figure 19: An example of a cluster that holds 12 equities inside



Figure 20: The equities inside the cluster are expanded in a circular shape around the cluster

This clustering approach together with the categorisation of posed an interesting challenge, if equities are hidden behind a cluster, how would the analyst know how many equities of a given category are inside the cluster? To solve such issue an arc of the colour of the category and proportional to the number of equities inside the cluster was devised.

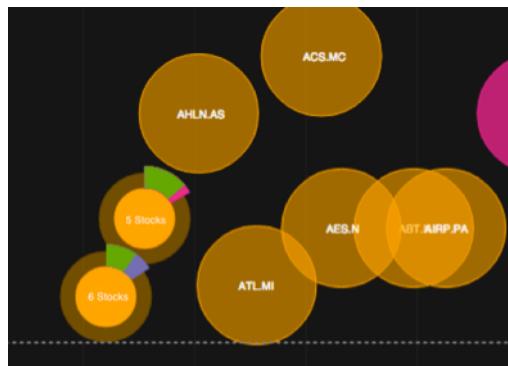


Figure 21: The arcs are proportional to the number of stocks of the category

Finally, to show companies that ace the filters defined by the user (FR03), that is, companies that have a great ranking the researcher decided to encode the information by the luminance level of the equities on the same scatterplot. Because we are already using colour for the categories, to show the ranking (QuickRank) information I decided to add a mode button that allowed changing from the QuickRank mode (depicted below) to the normal mode.

Does a visual data mining approach improve stock picking for equity analysts?

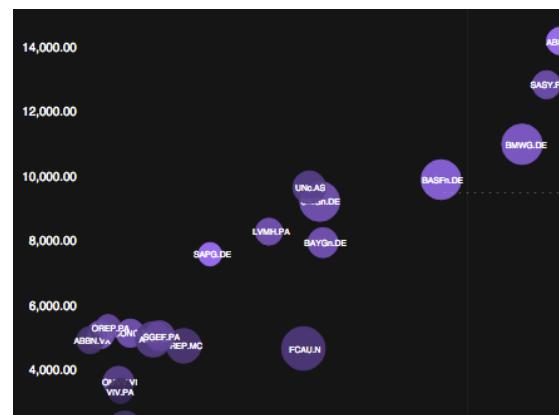


Figure 22: The QuickRank mode

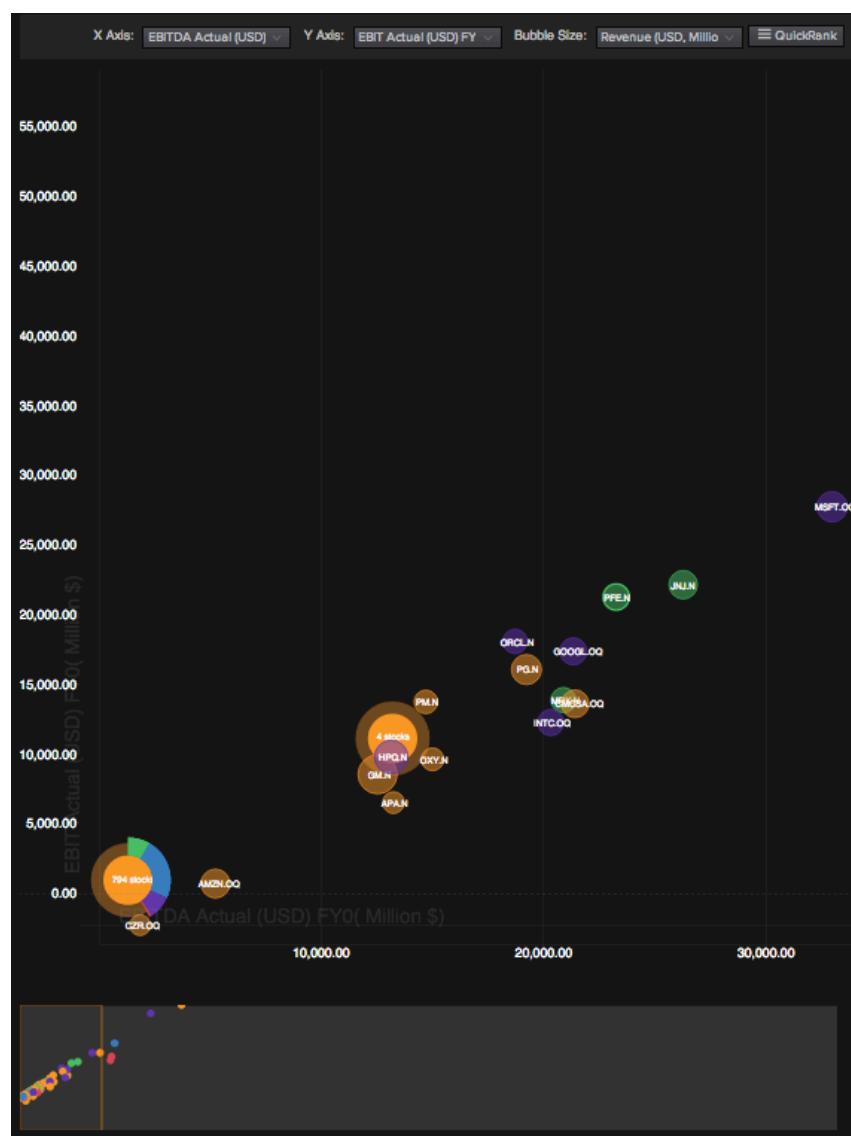


Figure 23: The final scatterplot design

## 4.4. INITIAL DESIGN TESTING

As pointed out before, the implementation of the prototype happened during the design stage in order to validate the selected interaction idioms. But implementing the prototype the decisions taken during design also allowed validating the entire visualization.

Parallel coordinates are sometimes seen as ‘being difficult to understand, expert-only representations’ (Siirtola & Raiha, 2006). This was a concern for the author who to test it out designed a test with real users to evaluate whether users would be able to understand and use parallel coordinates or a different approach should be used to chart multidimensional data. Under the methods section the reader will find how this testing was performed, here the results of the test will be summarized.

### Results

The testing was performed in order to find usability issues and to evaluate if design decisions were correct.

Overall the majority of participants were successful on using both the parallel coordinates and the scatterplot. They completed the tasks that were designed to test both visualizations. Regarding the parallel coordinates, it took more time to figure out what information the visualization was actually showing, but once this initial barrier was overcome, they could interact with it fairly easily.

The third subject being tested was the most familiar with the Equity Screener. In fact his job was training customers on using Eikon. He was really enthusiastic about the parallel coordinates seeing how useful the brushing interaction would be when highlighting and filtering down by multiple dimensions. During the testing the participant even commented: “I can see how customers would be using it, it needs a bit of work but it is a really good visualization”.

The test helped uncover a number of usability problems on the prototype that needed to be sorted in order to improve the interaction, problems such as difficulties when trying to highlight a specific category on the scatterplot. But it also helped validate the parallel coordinates as a visualization that can be understood by users who have not been exposed to this type of visualization before.

## 4.5. LINKED VIEWS

The validated design proved that both visualizations could be used independently to perform basic tasks performed by analysts. However, many visual data mining tools use an approach called linked views (Poulet, 2002) (Schulz et al., 2006) (Ladstadter et al., 2010) (Adrienko & Adrienko, 2001) (Broeksema et al., 2013). Linked views consist on several visualizations that are interconnected thanks to a high degree of interaction. This means that when we are interacting with a visualization (selecting items, moving the mouse over its items, modifying its parameters etc.) the other visualizations on the screen react to this interaction. Linked views allow having alternative viewpoints on the data, to have a clearer view of the underlying structure of the data and the ability to compare items more precisely (Roberts, 2005).

In order to implement this approach, the author worked on how the enhanced scatterplot could be coupled with the parallel coordinates in order to achieve a flawless interaction between the two visualizations. This resulted into a set of powerful interactions that enabled having different views on the data and more powerful visual data mining capabilities. This was considered to be the final design and fully was fully implemented.

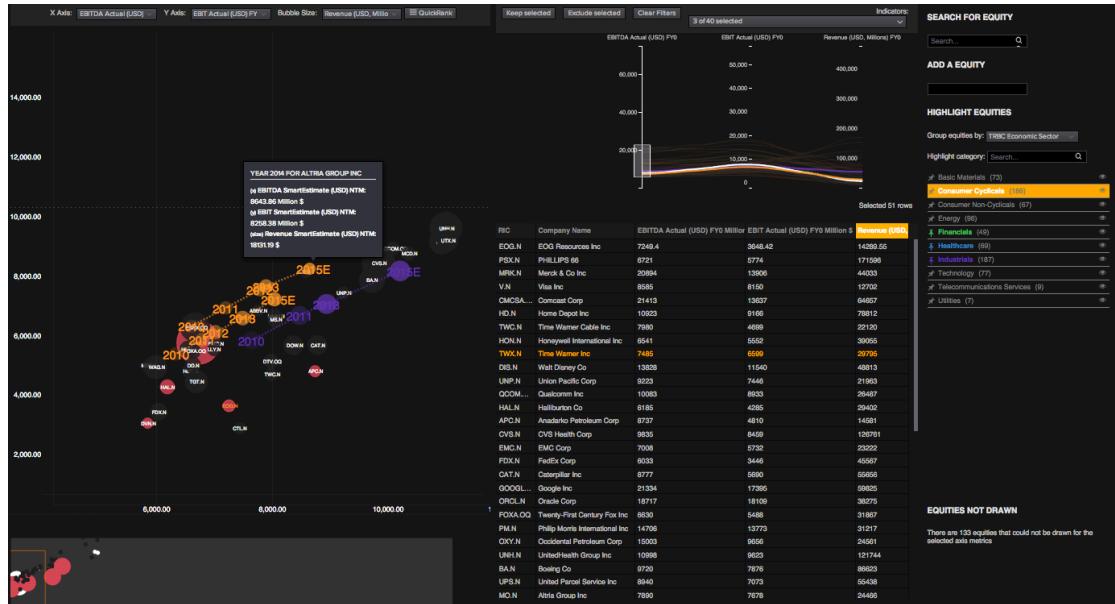


Figure 24: Screen capture of the final design of the prototype implementing the linked views approach

## 4.6. RE-VALIDATION OF THE FUNCTIONALITIES

Towards the end of the project the researcher gained access to more domain experts thanks to the internal promotion of the design work and the help of the work supervisor.

On the interviews with the seven different domain experts the prototype gathered really positive feedback. Although some issues were raised during the interviews, the great majority of the functionalities were successfully re-validated by the domain experts. Furthermore, most of the feedback gathered were in fact new functionalities to expand the functionality of the screener [Appendix: Functionality revalidation Feedback Summary].

The problems raised did not affect the domain validity nor the task/data abstractions, in most cases these were usability problems.

The most important problem was related to the Quick-Rank functionality, where the encoding of the ranking as luminance was considered to be wrong.

## 4.7. FINAL COMPARISON TEST

The table below depicts the match between tasks defined for the final comparison test and the requirements. (Green means the task is testing the requirement)

	FR01	FR02	FR03	FR04	FR05	FR06	FR07	FR08	FR09	FR10	
TASK	Compare multiple metrics			Ace the filters	Find outliers	Group by attribute	Add equity	Filter by groups	Locate a share	Compare medians	Filter progressively
BT1											
BT2											
BT3											
BT4											
BT5											
BT6											
MI1											
MI2											
MI3											

There is no task for the requirement FR06. Two main reasons exist for this fact. First, the functionality itself does not exist on the current Equity Screener, thus we would be testing a complete new feature without being able to compare with the old screener. The second and more important is that the feature was not stable enough to be tested.

## 4.8. FINAL TESTING RESULTS

### 4.8.1. BENCHMARK TASKS

The summarized results for the comparison test are listed below. To preserve the anonymity of each participant their names have been replaced by Subject A, B and C. Three domain experts on investment banking and portfolio management participated voluntarily on the test. Two of these tests took place at the Thomson Reuters offices. The remaining test was conducted via an online web conferencing tool that allowed the researcher to communicate with the participant as well as see his screen and mouse interaction.

A full list of the final scores and additional information on the tests can be found on [Appendix: Final Comparison Results Summary]

Initially, the test was much more complex with up to 9 tasks instead of 6. When performing the tests the author realized too much time was spent on some tasks and amended the tasks by eliminating duplicates. The results of the initial task duplicates were discarded.

A think-aloud protocol was used to gather feedback from users whilst performing the tasks. Some gave interesting insights on the task execution. Below the reader can find a table with the summary of the results, and for further analysis of the testing results each task is analyzed separately.

Task	Subject A		Subject B		Subject C		Prototype	Screener
	Time	Tool	Time	Tool	Time	Tool		
B1	03:41	Prototype	02:45	Prototype	04:11	Screener	03:13	04:11
B2	03:00	Prototype	04:51	Prototype	04:04	Screener	03:55	04:04
B3	00:58	Prototype	01:45	Prototype	02:19	Screener	01:21	02:19
B4	02:11	Screener	09:00	Screener	03:34	Prototype	03:34	05:35
B5	04:23	Screener	01:50	Screener	05:41	Prototype	05:41	03:06
B6	04:03	Screener	03:25	Screener	03:48	Prototype	03:48	05:35

### **Task B1: Find three companies with the highest values on 3 metrics**

Some participants struggled here to find how to use the visualization to actually find the 3 companies with such high values with both the prototype and the Screener.

With the prototype the given participant was trying to find the three companies with the table, trying to sort for each 3 fields (a functionality which is not supported on the screener nor on the prototype, so was unsuccessful at first. Later on he realized he could find the companies with such values by charting these 3 on the scatterplot. Another participant focused on the parallels and wanted to chart the 3 metrics yet because the screen that he was using was too small, all the indicators were not visible. It took some time to realize the screen could be resized to accommodate all indicators.

The subject using the Screener performed the task by filtering down from a big number of companies to a smaller set by playing with the values of the three metrics. This was performed a bit randomly (play with the filter values, run the screen, wait for the results, modify the filters).

Although it was a seemingly simple task, task times were considerably high due to the problems experienced by the participants.

### **Task B2: Trend evaluation**

The main problem for Subject B using the prototype was the problems when navigating the scatterplot chart. This participant was using a integrated mouse on a laptop and the degree of precision required to navigate the scatterplot and see the values of the trends is quite high. This meant a considerable amount of time was spent on navigating rather than evaluating the trends themselves.

Subject A also used the prototype and was trying to make sense of the parallel coordinates to try use them to chart the trends. These suggested the visualization purpose was not being understood at first. The participant finally realized the fact that the trend could be more easily seen using the scatterplot than with the parallel coordinates.

The participant using the screener (Subject C) had to evaluate the metrics one by one, so it took time to gather enough conclusions to answer the research question.

### **Task B3: Evaluate specific metric of a company**

Two participants where tested on this task using the prototype. The first participant only took 58 seconds on fulfilling it. To do so he simply added the new metric on the visualization and focused on the specific company the test was asking for, he got the result straight away. However, the second participant struggled a bit more, it took him time to understand the linking between the two visualizations, by a trial-and-error fashion ended up charting the metric and evaluating each company independently. In this case the behavior of the first participant was the expected whereas the second was not ideal.

The participant using the Screener spent a considerable amount of time trying to find the different companies we where asking for on the task, searching for the metric on the table and mentally comparing the values for the three different companies. During the process no visualization was involved, the participant only used the data grid to perform the task.

### **Task B4: Company over median for an economic sector**

For this task two participants where tested using the Screener. The great difference on the time it took (Subject A & B) performing the tasks relies in the user's proficiency on performing the tasks. The first participant was more used to performing such task on the screener and it did not take too much time to; categorize the stocks depending on their economic sector, set the calculation of the median and finally evaluate what companies where over that median.

Subject B was not so familiar with this functionality and was struggling to see how it could be performed on the screener. He started by grouping equities but then did not know how to calculate the median. Finally, he created a complex filter where the output would be only companies over the median. This produced the correct results, showing how the filtering process the Screener is equipped with is really powerful yet sometimes complex to use.

Subject C was using the prototype to perform the task. Although It did not take time until the participant understood how to filter by sector, it did take time until the user navigated the entire equity space on the scatterplot visualization trying to look for a companies that would be over the calculated median of the two economic sectors.

### **Task B5: Companies with a high Quick-Rank value**

With the screener this task was a fairly simple task. The ranking criterion was already set-up so that users could easily sort the data grid to find the equities between quick-rank values. Subject A was not familiar with the concept of Quick-rank so the researcher needed time to explain such concept, thus the great amount of time.

On the other hand, Subject C using the prototype to perform the task struggled to give a correct answer to the question. The problem was that, as the participant pointed out "it is

really difficult to differentiate on the ranking of the companies because all them seem quite similar". This problem was already highlighted by domain experts and reappeared here.

### **Task B6: Complete filtering process**

Subject A completed this task with the prototype by adding the different metrics on the parallel coordinates and started by creating filters. The participant was engaged on the task trying to give the best possible answer at a given time the participant pointed out "I can see there are quite a lot of efficient companies that comply with this strategy but I'm refining to get the best ones."

Subject B used the prototype too. To perform the filtering the participant also used the parallel coordinates and filtered down to a set of interesting equities quite readily. He then used both the data grid and the scatterplot to further evaluate the equities. The filtering process seemed easy for the participant as he commented: "this is much easier on the prototype than on the Screener, to make it that easy to change the filters you should use something like scrollbars on the Screener".

## **4.8.2. MEDIUM INSIGHT TASKS**

Medium insight tasks are designed to test whether users understand the structure of the data properly. Listed below there is a list of the different medium insight tasks and the time it took to answer the question.

Task	Subject A			Subject B			Subject C		
	Time	Dataset	Tool	Time	Dataset	Tool	Time	Dataset	Tool
MI1	01:27	Europe	Prototype	-	USA	Screener	03:05	USA	Prototype
MI2	03:13	Europe	Prototype	-	USA	Screener	04:17	USA	Prototype
MI3	01:49	Europe	Prototype	-	USA	Screener	03:05	USA	Prototype

Unfortunately, due to the amount of time that was spent setting up the remote environment for Subject B (the reader might recall the testing was done remotely) the researcher could not perform the testing with this participant.

### **MI1: Describe relations between two metrics**

Subject A could clearly see the tendency on the scatterplot. So much so, he commented "every analyst knows that these metrics are related but you can see it clearly on the scatterplot". For the second part of the question where two metrics not directly related are evaluated the participant stated: "mainly depends on the company". To reach this conclusion he started playing with the scatterplot visualization and evaluating the metrics for each company.

The other participant also performed this task with the prototype. Both participants used a combination of the scatterplot and the parallels to reach their final conclusion and answer the question correctly.

### **MI2: General difference between industries**

The participant that performed this task with the prototype was Subject A and stated “more or less the same, just the number of companies changes”. He mistakenly misjudged the scaling on the data and provided an incorrect answer. This might well be due an error on the scaling of the prototype solution.

### **MI3: Performance difference between industries**

Subject A was now used to evaluating the differences between industries and performed the task pretty fast, pointing out a some deeper insights such as that in one of the industries there where less outliers than the other, which suggested these would be fairly similar.

#### **4.8.3. FULL-INSIGHT TASKS**

Two questions where asked to the participants, not to be answered and recorded but as a stating point for the analysis, the focus of the results will be on the ‘open exploration’ where the researcher observed analysts exploring the data freely.

Subject A started with the parallel coordinates by adding a set of different metrics that measure how efficient a given company is such as ROE, Net Profit Margin and metrics that indicate if a company is undervalued (its price under its intrinsic value) EV/EBIT. By drawing on the parallel coordinates in just a couple of minutes ended up with an equity he found to be interesting. It was Take Two Interactive. This particular equity was actually a good pick because the company was efficient and was undervalued by the market. To further verify it was a company worth investing he added the metric ‘Analysts Revision Score’ this indicates how financial analysts value the company; the higher the better the company is considered. He evaluated how Take Two interactive was valued, in 98 a number close to 100 the maximum. This was a stock worth investing in. The entire analysis process only took around 4 minutes.

Subject B started using the parallels too but with a different approach, he added ‘EPS’ the ‘Beta’ which is a measure for how volatile is a company and ‘Price to Sales’. He then started evaluating what industries had a higher variation on the beta (more volatile). By doing so he concluded Industrials where a volatile industry at the moment because he could see the data was distributed more dispersed in this industry, he commented “Interesting... Industrials seem quite volatile at the moment, if you fancy gambling you’d invest in this industry now”. Then he included the ranking of the equities on the parallels and filtered down to the top equities, discovering Cedar Fair Ltd. a company that seemed interesting to him. This process took about 10 minutes.

Subject C was entirely focused on comparing the different industries. Instead of only using the parallel coordinates he started playing with the different categories and comparing equities against the median. The user did not end up with an equity worth buying but with a good overview of the industry. This process took about 10 minutes.

## 5. DISCUSSION

The main activities of the project are discussed in this section compared against the research objectives, evaluating whether these were met or not.

### Requirements Gathering & Validation

The requirements elicited during the initial phase of the project were validated with domain experts in several occasions during the project. The DSM enforces different validations during the 4 levels of design (Munzner, 2009) making sure each step is correct. The main advantage resulting from following such approach is that requirements are refined continuously. This leaves less room for abrupt changes and has less probability of implementing features that are not really needed. When requirements validation is performed with the help of parallel prototyping, it allows exploring a wider range of possible requirements and their solutions to these requirements. By following this process it is possible to detect and discard features that are not needed or that maybe need to be redefined to be correct.

The set of requirements elicited from the domain expert has derived a set of really useful features that could help equity analysts on their tasks. For example a participant on the final research commented during the testing "I can see the data now, for example when focusing on an industry I know how the stocks are behaving straight away, this is really cool stuff for an analyst" (regarding requirements FR05, FR07). Another example regarding the ability of comparing trends between stocks (FR01, FR02) "it's so much easier comparing stocks in the prototype than with the screener". Another participant commented on the ability to perform the filtering process on the prototype: "visually doing the filtering is so much better than manually" (FR10).

All this feedback proves RO1 has been achieved as the set of requirements elicited can in fact help equity analysts on performing their tasks.

### Design

A great majority of advanced stock data solutions that were evaluated on the initial review of the literature (Chapter 2: Context) were based in complex visualization and interaction idioms (Lei & Zhang, 2010) (Schreck et al., 2007) (Ziegler et al., 2010). The prototype, product of the current research, is aligned with the efforts of the current reviewed literature. A powerful interactive visualization with multiple linked views that allow exploring the data from various perspectives is not an uncommon approach. Yet, the current research is indeed a more comprehensive approach than the previous literature as the current prototype has been thoroughly validated, with domain expert feedback and then by evaluating the degree of insights achieved with the solution. The DSM approach has been used to design and validate the design resulting in a prototype incorporating visual data mining features fulfilling RO2.

One of the most surprisingly positive aspects of the project was the degree of good feedback gathered on the re-validation of the domain as well as during the final comparison test. All domain experts that were faced with the tool at first understood all the functionalities really easily and could easily envision how they could help on the stock-picking process. Furthermore, not only they understood and validated the functionalities, but

also even started to think about new features that could be included on top of the ones on the prototype. From the feedback gathered we could start defining even more requirements and new features. For example, regarding the trend feature (that draws a line showing how the stock has progressed through time) a domain expert commented, “we could plot many other variables, such as the size of the fleet for an airline, this way you can compare both financial and non-financial trends”. Or considering the median functionality another user pointed out “would be useful to be able to chart an index such as the Dow Jones or S&P500 as the median is in order to compare visually”. This feedback highlights an interesting fact, rather than giving feedback to correct or solve the prototype functionalities they were extending the features on top of the prototype.

The main problem on the design of the visualization was the Quick-Rank functionality. The author took an incorrect approach when encoding the rank information of the equities. The reader might recall this functionality where a set of different stocks is ranked (1<sup>st</sup> better, last worse) depending on how well they match a set of characteristics such as having low debt and high revenue. This information was encoded on the luminescence of the stock representation. This was highlighted not only on the first round of feedback with domain experts but up to the final comparison test. The users highlighted it was really difficult to know the rank just by looking at the color “It is difficult to distinguish between rank position; the best ranking (rank 1) has a similar color to rank 15”. This feature was incorrectly designed and should be redesigned if it were to be useful.

## Implementation

The main product of this research has been a prototype implementing the elicited requirements, which was designed and validated using the DSM approach, incorporating visual analytics features. The proof that the functionalities of the prototype were stable enough is the fact that its functionality was tested against the current Equity Screener, which is a stable piece of software used by financial analysts everyday. Being able to carry out this test means the RO3 was indeed fulfilled.

## Evaluation

An evaluation test was designed to compare between the Screener and the prototype. Both tools were evaluated by using a set of different tasks in order to measure how they supported visual data reasoning. Task times and error rates were recorded in order to compare both software applications.

Designing the tasks for evaluating insights was indeed a complex task, as literature suggests (North, 2006). Furthermore, collecting and analyzing the data was a complex undertaking too. Collecting data was difficult because following the analytical reasoning of an expert on the domain can be quite complex, and much more when you are not familiar enough to follow the reasoning processes. Evaluating the resulting data is a time consuming task. Both quantitative and qualitative data was collected and analyzed. Analyzing qualitative data is always more time consuming than quantitative data.

From the analysis of the quantitative data we can see that the mean task times are smaller when performed using the prototype in comparison to the equity screener. These quantitative measures indicate an improvement on common processes and tasks carried out by analysts. It proves that a powerful interactive visualization can in fact be useful on the

stock-screening process. But it does not offer any proof that the equity screening process can be improved by a visual data mining approach.

An interesting result when evaluating both solutions is the low degree of error. This might happen due to two different matters. First, might be the complexity of the tasks being performed was indeed lower than common analysis tasks. This seems improbable otherwise completion times would be much lower than expected. Second, it might be that the prototype and the Equity Screener are indeed both good enough performing the required tasks. However, task times are lower when using the prototype, it offers a better performance.

Finally, the time to get interesting insights was really small when performing the insight evaluation. Although the main problem of this approach is that some analysts are keener on exploring the data than others, so measuring the number of insights is not a good metric. However, when a participant was keen to explore the data, sometimes it was rather difficult recording all required insights that were generated when using the prototype.

From the evidence collected, the feedback, the benchmark-focused testing, and the insight-based testing there is enough evidence to offer a proof that visual-data mining can in fact enhance the equity screening process, thus fulfilling the research objective RO4.

## 5.1. VALIDITY AND GENERALISATION OF THE RESULTS

### 5.1.1. VALIDITY

Common visualization evaluation is focused on testing users performing simple tasks with the visualization (North, 2006). This is acceptable when we are testing how visualization can in fact support a set of well-defined tasks. However, for the current project a more complex approach was needed because what was being tested is whether a tool can actually improve a complex analytical process, as it is the stock picking process. The benchmark tests assume simple tasks build up to more complex tasks (North, 2006). The medium and full-insight approaches are focused on testing how the visualization facilitates the understanding the data and the process itself. They help getting an insight on how it helps the analysis process by allowing the analyst to reflect on the process whilst doing it.

The highest degree of validation for the results would be achieved not by performing a comparison test between the two tools but rather measuring the usage of the tool in the real-world environment with real users (Munzner, 2014, p.71). Measuring the adoption rate would be possible if both the Screener and the Prototype were launched as separate applications within Eikon and its usage rate could be measured. This would answer if we have been successful on developing a useful visualization: greater usage rates are a sign of successful design (Munzner, 2014, pp.71-73).

### 5.1.2. GENERALIZATION

The final comparison test was performed with only three participants. Although this number of users seems really small, their value was that they were real experts on the field. In any company workers are normally under heavy time pressure to meet deadlines, this poses a

major problem when recruiting participants: every hour spent with the evaluator means they are not directly working on their tasks (Sedlmair et al., 2010). Additionally, Thomson Reuters being such a large company has a well-structured hierarchy (compared to smaller companies where hierarchy is normally not so well defined), this makes somewhat more difficult to gain access to domain experts. An interesting fact that has been useful on recruiting users is the one-hour limit (Sedlmair et al., 2010). Limiting the evaluations, meetings and feedback sessions to only one hour does have a positive impact on the recruiting process

Having a small population does not make results less valid or generalizable. A given result is generalizable to the extent to which it can apply to other people or other situations (Carpendale, 2008). Equity analysis is a really specialized job, thus the number of experts on the matter is considerably small. Taking this into consideration, even having access to three domain experts on the subject, it gives enough room to be able to generalize. Additionally, if results are consistent, as such is the case of the current research, we could consider generalizing these to other Equity Analysts.

The tool that has been built would be useful for any equity analyst when performing the stock picking process. However, when it comes to generalizing into any data exploration problem generalizing is not so easy. The initial problem was a data exploration problem, being able to gather insight from a complex multidimensional dataset. Generalizing the results obtained into other data exploration problems is not that easy because it is in fact a really niche problem with quite a specific dataset.

## 6. EVALUATION, REFLECTIONS AND CONCLUSIONS

The main aspects of research will be evaluated in this chapter reflecting on what has gone right and wrong through each project phase.

### 6.1. LITERATURE REVIEW

The initial review of how visualization is being used in finance was helpful as it gave the researcher an overall idea of what is being done. Yet, the initial review was not enough as it was restricted to academic papers and the academic environment, and was considered to be too restrictive. The author then decided to widen the review and evaluate how financial software and especially commercially available equity screeners make use of visualizations. A review of the different stock screening tools was produced collecting information on how other equity screener tools make use of visualization.

Regarding the methodology that has been followed, DSM, the literature available was not enough to fully understand it (Munzner, 2009) (Sedlmair et al., 2012). However, the author was able to access a draft version of Tamara's Munzner's book "Visualization Analysis and Design" which was publicly available on its website before publication. Having access to the draft has been key to understanding the DSM methodology as the book clearly explains the methodology and how to face each step, as well as a wide collection of visualization idioms that can be used to solve the tasks.

All throughout the project, academic literature has been used to support design decisions, investigate different visualization idioms, learn how to implement visual analytics features, inform how to evaluate visualizations etc. This provides a high degree of validity on the most important decisions. Compared to taking decisions to just relying on personal knowledge and guesswork, relying on well-established academic knowledge helps base decisions on tested and accepted methods.

### 6.2. REFLECTION ON TOPIC AND OBJECTIVES

Choosing the research topic was not a simple task. As this research project took place during an internship, the outputs of it should be an actual research question and a set of different objectives that should benefit not only the researcher but Thomson Reuters too. In order to define an interesting research topic for both parts a series of meetings where held between the researcher, the academic supervisor and the work supervisor. From these meetings a series of topic suggestions where produced, yet the final research topic was not defined until the author knew more about the data exploration problem and the domain. The initial research topic suggested on the project proposal (Can a design study approach be used to improve data exploration?) was found to be too difficult to prove and thus rejected. The main reason is that in order to judge the efficiency of a methodology you need several different cases where the methodology was applied and consistent measures of its success. Thus this was found to be not a suitable research topic.

During the initial stages of the project the selection of objectives defined in the project proposal where also refined. These initial research objectives where solely based on the

Does a visual data mining approach improve stock picking for equity analysts?

description of the internship and based on the initial research topic. They were redefined together with the new research topic.

### 6.3. REFLECTION ON CHOSEN METHODS

The Design Study Methodology has helped guide the project execution. The cyclic nature of the methodology does indeed support real-world work. In design work it is not possible getting everything completely correct at first, so having a structured design process that allows being able to go back on the process to find out where the mistake was made, amend it and cascade the changes is a great approach to support the real-world visualization design work. Additionally, approaching the visualization design in four different levels helps separate the work in stages and allows validating each level in isolation (note this does not mean independently, there is a correlation between the levels).

The commitment to use a Visual Analytics, and specially Visual Data Mining supports the tasks was driven by the fact that these technologies help users to gain a higher degree of insight on the data (Lei & Zhang, 2010). The results of the insight test do in fact prove that visual data mining can help get a better insight on the data, showing what is difficult to discover without the use of a visualization approach.

### 6.4. REFLECTION ON PLANNING

The great majority of software projects have several stakeholders participate in different levels of the process. It is also the case of the current project work. Obtaining the description of the problem needed the participation of a domain expert, gathering feedback on the abstractions created by the designer too... On a working environment it is not possible to have the domain expert always available for the researcher. This can cause delays and work needs to be re-scheduled until the domain expert is available. This sometimes takes more than desired and time has to be managed wisely to be able to follow the project.

The initial time schedule suggested in the project proposal was not precise at all. For instance the time allocated for the Domain Problem Characterization was only 6 days, nothing further than reality, in fact this phase took almost a month to complete and was not completed independently it happened together with the design of the visualization and the analysis of the data to be displayed.

### 6.5. WHAT HAS BEEN ACHIEVED

Thomson Reuters is one of the largest information companies in the world, providing data to the great majority of the financial markets. At the time of the writing, Thomson Reuters Eikon is used by 185.000 financial professionals around the world and is used by some of the biggest companies in the financial sector such as HSBC or UBS. Thomson Reuters are data and information providers with multiple sources of information ranging from financial markets to their all important Reuters Agency. They produce data and information that is delivered to professionals around the world. Yet, the use of visualizations to convey this complex information is not generalized. Many applications in Eikon could benefit not only how this information is conveyed but of a new workflow approach by using visualization and interaction. The present work has presented the concept of visual analytics to a part of the organization. The ability to explore data thanks to a set of well-thought and interconnected visualizations has been a welcomed concept, gaining a considerable support from both

product managers and even management. In fact, the blog entry written by the author gathered more than 200 views, a number well above the average of 50 views on the group where it was posted.

Product managers are key people at Thomson Reuters. We could say Thomson Reuters is a fairly Product-manager centric business. They are who devise the features that are included in the application they own in the Eikon platform. In the different occasions where the author has been able to gather feedback from product managers about the prototype, the experience has been really positive.

Furthermore, towards the end of the project the author presented the prototype to the design team as well as the majority of members from the investment team. To present the prototype appropriately the author created a slideshow summarizing the design process and showing the major features of the prototype. The feedback arising from this presentation was really positive, the head of investment and management even commented: "some of the features designed could go into the product", and "this prototype will be taken into account for the set of features being designed for the upcoming year". This clearly signs that the features devised for the screener and the execution excelled.

## 6.6. WHAT I'VE LEARNED

Visualization design is a complex process. It involves not only fulfilling a set of requirements but also researching about many more aspects of the solution such as human computer interaction, human perception and cognitive reasoning. Additionally there is a part of creativity involved on the process. To create new visualizations there is a need to experiment with different visualizations. DSM encourages starting with a broad set of different solutions and narrow down these into a manageable set of solutions that can be implemented. In the case of the current project, although being an internship, working with Thomson Reuters, the main focus should be on the research question and on envisioning the features that could go on the product straight away. It is not easy to find the balance between the features that can go on the product and the features that are needed to support the visual analysis.

Another complex aspect is the fact that the filtering down from a broad set of visualization solutions to a smaller set is not a simple process. Choosing what features for the visualization need to be implemented has to take into account both the time and resources, quite limited in this project as well as in reality.

The series of validations encouraged by the DMS during the design phases helps make sure we are taking the appropriate decisions on each of the steps of the process of creating the visualization. We validate whether we understand the problem, if the tasks that need to be completed and with which data, the visualization and the interaction idiom and finally if the implementation works as intended to fulfill all these design levels.

Promoting the visualization work internally is an important part of the success of the project. By promoting work internally the author was able to gather new connections to domain experts within Thomson Reuters. These connections were facilitated with the help of the work supervisor yet one of the triggers to gain a good degree of involvement was thanks to this internal promotion. An important aspect of this internal promotion is that aesthetics and usability should not be forgotten, they are tools that can be used to gain acceptance of what you are building.

Validating early design decisions improves the design process. Making sure the design of the visualization is correct before moving into implementation improves the visualization design process, as it is much easier to modify a mock-up than a working piece of software code.

## 6.7. CONCLUSIONS

The project as a whole can be considered a successful collaboration between Thomson Reuters, City University and the researcher.

The researcher has been able to produce an interesting visualization prototype that not only shows innovative features that could potentially be included on the current Equity Screener, but also brings into light a new workflow for equity screening. This workflow is based on the use of visual data mining functionalities to enhance the equity screening process. If included on the platform, such approach could become a competitive advantage for Thomson Reuters against other Equity Screeners on the market. Additionally, the visual data mining approach suggested by this research could be useful in other applications of the platform. Or maybe even the disclosing of specific visualizations such as the parallel coordinates might be an interesting addition to some of the applications on the platform.

The Design Study Methodology has proved to be a valid methodology on the execution of this short span project. This shows that the methodology is not only suitable for lengthy projects that span over years but can be used too with projects with a short timespan.

Designing visualizations is an iterative endeavor, validating design choices early and often leads to better visualizations. Being able to do this requires forging connections with different domain experts and the promotion of visualization work (even in somewhat early stages of the design) can help in building these connections.

Evaluating software products is difficult because designing tasks that simulate complex workflows is not an easy. Yet, designing tasks to evaluate visualizations is even more complex as normally visualizations support deep analytical thinking processes which are really difficult to measure.

## 6.8. FURTHER WORK

The current prototype could be developed into a fully functional application included in Eikon providing a new way to explore and analyze the stock market data.

Additionally some of the new features envisioned by product managers when asked to give feedback about the prototype could be implemented on the prototype to test if they enhance the analyst's workflow.

Although the final evaluation comparing both tools suggests that the visual analytics approach can improve the equity analysis process, a more thorough evaluation is possible. This evaluation should involve real-world users performing real equity analysis.

Does a visual data mining approach improve stock picking for equity analysts?

Finally, the visual analytics approach could be beneficial in many other applications of Thomson Reuters Eikon: its benefits are not obviously restricted to Equity Screening. Many other applications within the platform could be modified to include such approach.

## 7. WORKS CITED

- Ye, X. & Rey, S., 2011. A framework for exploratory space-time analysis of economic data. *Springer-Verlag* .
- VIAU, C., 2012. HYBRID VISUALIZATIONS FOR DATA EXPLORATION. In *HYBRID VISUALIZATIONS FOR DATA EXPLORATION*. Montreal.
- Wattenberg, M., 1999. Visualizing the Stock Market. *Dow Jones & Co. (SmartMoney Magazine)* , 99, pp.188-89.
- Xiang, Y. et al., 2012. Visualizing Clusters in Parallel Coordinates for Visual Knowledge Discovery. *PAKDD*, pp.505-16.
- Ziegler, H., Jenny, M., Gruse, T. & Keim, D.A., 2010. Visual Market Sector Analysis for Financial Time Series Data. In *IEEE Symposium on Visual Analytics Science and Technology*. Utah, 2010.
- Zowghi, D. & Coulin, C., 2005. Requirements Elicitation: A Survey of Techniques, Approaches, and Tools. *Engineering and Managing Software Requirements*, pp.19-46.
- 4Traders, n.d. *4Traders Equity Screener*. [Online] Available at: <http://www.4-traders.com/top-records/ratings>.
- ACM, 1987. Definition of Visualization. *SIGGRAPH Comput. Graph.*, 21(6), p.3.
- Adrienko, G. & Adrienko, N., 2001. Exploring spatial data with dominan attribute map and parallel coordinates. *Computers, Environment and Urban Systems*, 25, pp.5-15.
- Ankerst, M., Berchtold, S. & Keim, D.A., 1998. Similarity Clustering of Dimensions for an Enhanced Visualization of Multidimensional Data. In *Proceedings IEEE Symposium on Information Visualization*., 1998.
- Amar, R.A. & Stasko, J.T., 2005. Knowledge Precepts for Design and Evaluation of Information Visualizations. *IEEE TRANSACTIONS ON VISUALIZATION AND COMPUTER GRAPHICS*, 11(4), pp.432-41.
- Bostock, M., 2015. *Github D3 Example Gallery*. [Online] Available at: <https://github.com/mbostock/d3/wiki/Gallery>.
- Broeksema, B., Telea, A.C. & Baudel, T., 2013. Visual Analysis of Multi-Dimensional Categorical Data Sets. *COMPUTER GRAPHICS forum*, pp.158-69.
- Carpendale, S., 2008. Evaluating Information Visualizations. In *Information Visualization: Human-Centered Issues and Perspectives*. Berlin: Springer-Verlag.

Does a visual data mining approach improve stock picking for equity analysts?

Chang, R. et al., 2007. WireVis: Visualization of Categorical, Time-Varying Data From Financial Transactions. In *IEEE Symposium on Visual Analytics Science and Technology*. Chicago, 2007.

CNBC, M.-, 2014. *CNBC Markets*. [Online] Available at: <http://www.cnbc.com/id/101563401> [Accessed 2 November 2014].

EquityMaster.com, n.d. *EquityMaster Stock Screener*. [Online] Available at: <https://www.equitymaster.com/research-it/company-info/stock-screener-india.asp>.

Data Driven Documents (D3), n.d. *Data Driven Documents (D3) Home*. [Online] Available at: <http://d3js.org/>.

DATAVISUALIZATION.CH, 2014. *DataVisualization Tool Selection*. [Online] Available at: <http://selection.datavisualization.ch/> [Accessed 2014].

Dillow, C., 2009. Stock-picking works! *Investors Chronicle*, 1 September.

DOW, S.P. et al., 2010. Parallel Prototyping Leads to Better Design Results, More Divergence, and Increased Self-Efficacy. *ACM-TRANSACTION*, 17(4).

Fekete, J.-D., Wijk, J.J.v., Stasko, J.T. & North, C., 2008. The Value of Information Visualization. In Stasko, J.T., Fekete, J.-D. & North, C. *Information Visualization: Human-Centered Issues and Perspectives*. Berlin: Springer-Verlag.

Finviz.com, n.d. *FinViz.com Screener*. [Online] Available at: <http://finviz.com/screener.ashx>.

Finviz.com, n.d. *S&P Heatmap*. [Online] Available at: <http://finviz.com/map.ashx>.

Financial Times, n.d. *Financial Times Screener*. [Online] Available at: [markets.ft.com/screener](http://markets.ft.com/screener).

Fool.com, n.d. *Stock Screener*. [Online] Available at: <http://caps.fool.com/screener.aspx>.

Graham, B., 1945. *The Intelligent Investor*. New York: HarperCollins.

Investopedia, 2012. *Investopedia - Stock Picking*. [Online] Available at: <http://www.investopedia.com/terms/s/stockpick.asp> [Accessed 22 October 2014].

Investopedia, 2013. *Fundamental Analysis: Qualitative Factors - The Company*. [Online] Available at: <http://www.investopedia.com/university/fundamentalanalysis/fundanalysis2.asp> [Accessed August 2014].

Infogengineering, 2013. *Understanding Information Overload*. [Online] Available at: <http://www.infogengineering.net/understanding-information-overload.htm>.

Inselberg, A., 1985. The Plane with Parallel Coordinates. *Special Issue on Computational Geometry*, pp.69–91.

Does a visual data mining approach improve stock picking for equity analysts?

Islam, A., Zaman, H. & Ahmed, R., 2009. Automated Fundamental Analysis for Stock Ranking and Growth Prediction. In *Computer and Information Technology (ICCIT 2009)*. Dhaka, 2009.

Huh, M.-H. & Park, D.Y., 2008. Enhancing parallel coordinate plots. *Journal of the Korean Statistical Society*, 37, pp.129-33.

Hagstrom, R.G., 2005. *The Warren Buffet Way*. 2nd ed. Hoboken, New Jersey, USA: John Wiley and Sons.

Kumasakaa, N. & Shibatab, R., 2008. High-dimensional data visualisation: The textile plot. *Computational Statistics and Data Analysis*, 52.

Keim, D. & Zhang, L., 2001. Solving Problems with Visual Analytics: Challenges and Applications. In York, A.N., ed. *roceedings of the 11th International Conference on Knowledge Management and Knowledge Technologies*. New York, 2001.

Keim, D.A., 2002. Information Visualization and Visual Data Mining. *IEEE TRANSACTIONS ON VISUALIZATION AND COMPUTER GRAPHICS*, 8(1), pp.100-07.

Keim, D. et al., 2008. Visual Analytics: Definition, Process, and Challenges. In *Information Visualization: Human-Centered Issues and Perspectives*. Springer.

Keim, D.A. & Kriegel, H.-P., 1996. Visualization Techniques for Mining Large Databases: A Comparison. *IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING*, 8(9).

Keim, D.A. et al., 2008. Visual Analytics: Scope and Challenges. In S.-V.B. Heidelberg, ed. *Visual Data Mining: Theory, Techniques and Tools for Visual Analytics*. Berlin. pp.76-90.

Kent, P., 1985. An efficient new way to represent multi-dimensional data. *The Computer Journal*, 28(2), pp.184-92.

Ladstadter, F. et al., 2010. Exploration of Climate Data Using Interactive Visualization. *American Meteorological Society*, April.

Lam, H. et al., 2012. Empirical Studies in Information Visualization: Seven Scenarios. *IEEE Transactions on Visualization and Computer Graphics* , 18(9), pp.1520-36.

Larman, C. & Basili, V.R., 2003. Iterative and Incremental Development: A Brief History. *IEEE Computer Society*.

Lei, S.T. & Zhang, K., 2010. A Visual Analytics System for Financial Time-Series Data.

Long, T.V., 2009. *Visualizing High-density Clusters in Multidimensional Data Phd*. PhD Thesis. Jacobs University.

North, C., 2006. Toward Measuring Visualization Insight. pp.6-10.

North, C., Saraiya, P. & Duca, K., 2005. Knowledge Precepts for Design and Evaluation of Information Visualizations. *IEEE TRANSACTIONS ON VISUALIZATION AND COMPUTER GRAPHICS*, July/August. pp.432-42.

Does a visual data mining approach improve stock picking for equity analysts?

North, C., Saraiya, P. & Duca, K., 2011. A comparison of benchmark task and insight evaluation methods for information visualization. *Information Visualization*, 10(3), pp.162–81.

Munzner, T., 2009. A Nested Model for Visualization Design and Validation. *IEEE TRANSACTIONS ON VISUALIZATION AND COMPUTER GRAPHICS*, 15.

Munzner, T., 2014. In D.o.C.S.U.o.B. Columbia, ed. *Visualization Analysis and Design*.

Marghesku, D., 2007. *Multidimensional Data Visualization Techniques for Financial Performance Data: A Review*. Technical Report. Turku, Finland: Turku Centre for Computer Science TUCS University.

Marketwatch, n.d. *MarketWatch Stock Screener*. [Online] Available at: <http://www.marketwatch.com/tools/stockresearch/screener/>.

Mnkandla, E., 2005. A Thinking Framework for the Adaptation of Iterative Incremental Development Methodologies. pp.315-16.

MorningStar, n.d. *MorningStar Equity Screener*. [Online] Available at: <http://screen.morningstar.com/AdvStocks/Selector.html>.

P., S., C., N. & K., D., 2005. An insight-based methodology for evaluating bioinformatics visualizations. *IEEE Trans Vis Comput Graphics*, 11(4), pp.443-56.

Purchase, H.C., Andrienko, N., Jankun-Kelly, T.J. & Ward, M., 2008. Theoretical Foundations of Information Visualization. In S.-V.B. Heidelberg, ed. *Information Visualization*. Springer-Verlag Berlin Heidelberg. pp.44-64.

PanopticonSoftwareAB, n.d. *Screening Equities with Treemap and Scatter Plot Data Visualizations - Youtube*. [Online] Available at: <https://www.youtube.com/watch?v=z8RgaoWONV0>.

Philips, M., 2013. *How the Robots Lost: High-Frequency Trading's Rise and Fall*. [Online] Available at: <http://www.businessweek.com/articles/2013-06-06/how-the-robots-lost-high-frequency-tradings-rise-and-fall> [Accessed 28 July 2014].

Poulet, F., 2002. FULL-VIEW: A VISUAL DATA-MINING ENVIRONMENT. *International Journal of Image and Graphics*, 2(1), pp.127-43.

Preece, Rogers & Sharp, 2002. In I..6.T.A.N.Y. John Wiley & Sons, ed. *Interaction Design: Beyond human-computer interaction*. John Wiley & Sons, Inc. pp.340-60.

Roberts, J.C., 2005. *Exploration through Multiple Linked Views (MLV) - University of Kent*. [Online] Available at: [http://www.roe.ac.uk/~rgm/sc4devo/sc4devo4/sdmiv2\\_17\\_roberts.pdf](http://www.roe.ac.uk/~rgm/sc4devo/sc4devo4/sdmiv2_17_roberts.pdf).

Robertson, G. et al., 2008. Effectiveness of Animation in Trend Visualization. *IEEE TRANSACTIONS ON VISUALIZATION AND COMPUTER GRAPHICS*, 14(6), pp.1326-30.

Does a visual data mining approach improve stock picking for equity analysts?

Symanzik, J., 2001. *Visual Data Mining - Techniques and Examples*. [Online] Utah State University, Logan, UT Available at: [http://www.math.usu.edu/~symanzik/talks/2001\\_byu.pdf](http://www.math.usu.edu/~symanzik/talks/2001_byu.pdf) [Accessed 1 October 2014].

Schulz, H.-J., Nocke, T. & Schumann, H., 2006. A Framework for Visual Data Mining of Structures. In Society, A.C., ed. *Twenty-Ninth Australasian Computer Science Conference*. Hobart, 2006. Australian Computer Society.

Schulz, H.-J., Nocke, T. & Schumann, H., 2006. A Framework for Visual Data Mining of Structures. In Technology, C.i.R.a.P.i.l., ed. *Twenty-Ninth Australasian Computer Science Conference*. Hobart, 2006.

Schreck, T., Tekusov, T., Kohlhammer, J. & Fellner, D., 2007. Trajectory-Based Visual Analysis of Large Financial Time Series Data. *SIGKDD Explorations*, 9(2), pp.30-37.

Sedlmair, M., Isenberg, P., Baur, D. & Butz, A., 2010. Evaluating Information Visualization in Large Companies: Challenges, Experiences and Recommendations. In *Proceedings of the 3rd BELIV'10 Workshop*., 2010. ACM.

Sedlmair, M., Munzner, T. & Tory, M., 2013. Empirical Guidance on Scatterplot and Dimension Reduction Technique Choices. *IEEE TRANSACTIONS ON VISUALIZATION AND COMPUTER GRAPHICS*, 19(12), pp.2634-44.

Sedlmair, M., Meyer, M. & Munzner, T., 2012. Design Study Methodology: Reflections from the Trenches and the Stacks. *IEEE TRANSACTIONS ON VISUALIZATION AND COMPUTER GRAPHICS*, 18(12), pp.2431-40.

Sedlmair, M., Meyer, M. & Munzner, T., 2012. Design Study Methodology: Reflections from the Trenches and the Stacks. *IEEE TRANSACTIONS ON VISUALIZATION AND COMPUTER GRAPHICS*, 18(12), pp.2431-41.

Siirtola, H. & Raiha, K.-J., 2006. Interacting with parallel coordinates. *Interacting with computers*, 18, pp.1278–1309.

Simoff, S.J., Böhnen, M.H. & Mazeika, A., 2008. Visual Data Mining: An Introduction and Overview. In *Visual Data Mining: Theory, Techniques and Tools for Visual Analytics*. Bozen-Bolzano, Italy : Springer-Verlag Berlin. pp.11-21.

The Economist, 2010. *One big, bad trade*. [Online] Available at: [http://www.economist.com/blogs/newsbook/2010/10/what\\_caused\\_flash\\_crash](http://www.economist.com/blogs/newsbook/2010/10/what_caused_flash_crash) [Accessed 25 July 2014].

The Telegraph, n.d. *The Telegraph Equity Screener*. [Online] Available at: [shares.telegraph.co.uk/stockscreener](http://shares.telegraph.co.uk/stockscreener).

Thomson Reuters, 2012. *Thomson Reuters Business Classification*. [Online] Available at: [http://thomsonreuters.com/products/financial-risk/01\\_281/trbc-quick-guide.pdf](http://thomsonreuters.com/products/financial-risk/01_281/trbc-quick-guide.pdf) [Accessed September 2014].

Does a visual data mining approach improve stock picking for equity analysts?

Thomson Reuters, 2013. *Thomson Reuters Annual Review*. [Online] Available at: <http://ar.thomsonreuters.com/#chapter-3>.

Thomson Reuters, n.d. *Reuters.com Screener*. [Online] Available at: <stockscreener.us.reuters.com>.

Thomson, S., 2013. Stock-picking marvels. *Investors Chronicle*, 16 January.

TMX Money, n.d. *TMX Money Screener*. [Online] Available at: <http://web.tmxmoney.com/>.

## 8. APPENDIX

### 8.1. ORIGINAL PROJECT PROPOSAL

Can a design study approach be used to improve data exploration?

#### Introduction

The current project proposal is the starting point of a project developed as part of an internship programme offered by Thomson Reuters. Thomson Reuters (TR from now on) is leading multinational media and information firm aiming to provide professionals with quality, trustful information. TR proud themselves in delivering high quality information to professionals in a wide variety of fields ranging from finance, legal, pharmacy & life sciences to intellectual property. Since its inception in 1851 TR has been providing stock market quotations to brokers and, whilst it has widened the range of information products, it has kept clear ties with the stock market.

Eikon is one of TR's flagship products. It provides financial professionals with all the information they need to discover investment opportunities, track assets, commodities and stocks amongst other functionalities. It is an extremely powerful tool that combines multiple information sources to offer a deep insight on market information. One of the utilities that Eikon offers is a tool to analyse stock values. It allows screening and analysing an impressive amount of stocks coming from the most important stock exchanges in the world. This tool offers a powerful insight on the stock market allowing financial professionals to make the right investment decisions. Put in other words it allows financial analysts to evaluate stocks and pick which are worth investing in.

Investors need a vast amount of information to take the write decisions. There are multiple sources this information can come from and that influence not only investment decisions but complete stock market movements. Each communication channel investors use adds more complexity to the overall investing activity. Classic sources of data like the financial reports, growth trends etc. are now just a starting point whilst investors are bombarded with news, opinions from financial writers, analysts, market strategists... No wonder why investment is becoming increasingly complex making much more difficult for investors to make any profit. Some are hard pressed even to continue as stock prices rocket or plummet with no apparent reason (Hagstrom, 2005). The current project proposal will focus on improving the stock picking process using TR's stock analysis utility making this process easier to financial professionals.

There are two main approaches for evaluating whether a public company is worth investing in. In the one hand a more long-term approach that evaluates the value of a company and its prospects for growth. This approach is called Fundamental Analysis (FA). Its main objective is to estimate a company's *intrinsic* value by evaluating a series of indicators of the company to deduce its price. This value is then compared to the price it's trading at to decide if it would be a good idea to buy shares (if the price is under the intrinsic value, undervalued) or if walk away (the intrinsic value is lower than the price is trading at, overvalued). The other main approach targets solely on the price movements of the market and is known as technical analysis. Technical analysis focuses on the short-term income by trying to predict if the stock

## Does a visual data mining approach improve stock picking for equity analysts?

price will increase or decrease and purchase shares accordingly. This type of analysis can be automated using complex algorithms that execute trades according to a set of rules, this is known as high frequency trading (HFT). Although HFT is still an important part of the stock orders executed all over the world it seems there has been a decline in trading volumes (Philips, 2013) probably due to its unintended consequences as in several occasions the stock market to plummeted without any apparent reason (The Economist, 2010). This has played against technical analysis and in favour of tools that take advantage of FA such as Eikon's stock market analysis tool.

Although there have been some attempts to automate FA (Islam et al., 2009), it is much more difficult to automate as it involves evaluating both qualitative and quantitative factors that affect the company's value (Investopedia, 2013). These factors can be either macroeconomic (status of the economy and industry conditions) or company-specific (financial reports and management). But FA is more of a subjective task that involves the evaluation of several dozens of indicators of a company to try grasping the intrinsic value of a company. Evaluating the indicators for a single company is complex, but evaluating these in relation to the industry makes this task extremely challenging and time expensive (Investopedia, 2013). Additionally sometimes the greatest rewards come from investing in companies that are not behaving like others in their industry. Stock-picking then becomes an exploration task where the objective is to find the stocks for companies that are performing differently than the competitors.

Information visualization (InfoVis) can be defined as "the use of computer-supported, interactive visual representations of data to amplify cognition" (VIAU, 2012). The last three words of their definition communicate the ultimate purpose of visualization, to amplify cognition: to help understand. InfoVis is about developing insights from data (Fekete et al., 2008). It does so by using human perception as the basis for inducing new insights from this data. Vision acts as a very fast filter, if it perceives a pattern it probably is because there is actually a pattern in the data structure (Fekete et al., 2008). Thus, InfoVis systems are best applied for exploratory tasks that involve large datasets such as stock market data. Sometimes, the user using the InfoVis system may not have a specific goal in mind. Instead, he might simply be examining the data to learn more about it, to make new discoveries, or to gain greater insight (Fekete et al., 2008).

This project intends to gain advantage of the exploration capabilities of information visualization. As we've highlighted InfoVis can present complex financial data in an easy way for financial analysts to gain a greater insight on this data. TR required their stock analysis tool to offer exploration capabilities, in such a way that patterns, outliers and other behaviours could be detected.

## Aims and objectives

The objective of the current project is to develop a tool that can take advantage of data visualization techniques to improve the stock picking process. This visualization should empower users with new ways to explore the complex universe of financial data involved in stock picking. To allow users to explore the financial data in novel ways two new dimensions will be introduced in the current charting solution.

In one hand the new chart will allow charting time-series data, that is, time is introduced as a new dimension. By introducing this new variable analysts will be able to follow the stocks

Does a visual data mining approach improve stock picking for equity analysts?

trajectory throughout time. This will allow evaluating the possible trend of a given stock, getting insight of what stocks are moving in the same direction through time etc.

On the other hand, a colour scale will be computed for similar companies. The idea is to group companies under a similarity scale according to a given field. For example, if we are given Philips and Intel they might seem quite disparate companies but they have great similarities when the investment on R&D is taken into account as both companies invest a considerable percentage of their income on the department, thus they will be coloured similarly within a colour scale

Several objectives can be derived from this:

- Learn how data is interpreted in the process of stock picking. There are several different ways in which stocks can be evaluated but this research will focus on a specific type of fundamental analysis: Growth At a Reasonable Price (GARP) which looks for companies that will grow and that are under their intrinsic value.
- Understand how to find correlations between complex data indicators. Learn how visualization techniques can be used to highlight correlations between data (such as color-coding depending on the degree of correlation).
- Build knowledge around how the different company and industry indicators can be used to forecast stock movements thanks to TR's financial experts.
- Understand how a design study can help structure the design process.
- Learn how to interact with complex web-based visualization APIs to plot complex data.
- Acquire the knowledge related to Visual Analytics required to create a useful visualization platform for exploiting complex data
- Validate the utility of the prototype by evaluating it with financial experts.

### **Products of this research**

The product of this research will be a web-based prototype that will exploit the vast amount of stocks TR's data feed contains by taking into consideration two new dimensions that might affect the stock-picking process: time and stock similarity.

The visualization implemented on this tool will allow financial analysts to consider time as a variable that can influence the stock picking process. Surprisingly this dimension is now completely ignored in TR's stock picking utility. Additionally a grouping algorithm will be implemented that will compare and classify a set of companies depending on one or multiple criteria.

### **Beneficiaries**

Researchers that are working on information visualization as well as financial analysts working on new ways to understand any company's financial data.

The new visualization approach will be based on adding new dimensions to the visualization of already complex data. The outcome of this approach (whether positive or negative) might prove useful in several other fields. Other situations where the characteristics of data might be similar to the complex financial dataset we are working with might also benefit from the discoveries uncovered by the present work. Ultimately the present project is not only

Does a visual data mining approach improve stock picking for equity analysts?

intended to develop a prototype for the sole benefit of TR's, but also to add to the body of knowledge on visualization in general and on financial data visualization particularly.

## **Scope and definition**

Although it might seem simplistic at first, stock picking is an extremely complex and risky process. Thus the current project will only focus on how visualization can be used to improve and facilitate this only process. There is no intention to widen the scope of the project to evaluate how visualization can help in other types of financial-related activates.

The prototype being developed will only focus on stock picking and will provide a reduced number of visualization approaches as that is proven to be the most effective.

## **Critical Context**

Probably one of the most well-known visualization techniques for financial data is the Map of the Markets (Wattenberg, 1999). It is a Tree Map representation of the whole stock market that in a simple graph gives a feeling of the overall status of the market. This visualization cleverly combines stock tendencies (color-coded squares), market capitalization (depicted by size) and the grouping by industry (by wider borders) all in a single graph. Although this initial contribution in financial data visualization seemed promising, financial analysts have stuck to the traditional charts for their analysis. Most studies highlight that the traditional one-dimensional trend line chart is the most commonly used visualization for stock picking (Ziegler et al., 2010) (Schreck et al., 2007). Furthermore, very few solutions have been developed which can handle stock market data information effectively in order to gain enough insight to understand how the market works (Ziegler et al., 2010).

Information visualization has been applied successfully in many other scenarios related to financial data. In fact, the field of financial data analysis can clearly benefit from information visualization and its outgrowth field of visual analytics (Schreck et al., 2007). Visual analytics focuses not only in data visualization but also in providing a great level of interaction with the visualization to enhance the data analysis. An interesting example on how visualization can dramatically improve the analysis of vast amounts of financial data comes from the collaboration of the University of North Carolina together with The Bank of America. This collaboration produced a system that uses a highly interactive visualization to allow the detection of fraudulent financial transactions from millions of transactions (Chang et al., 2007). This research offers a clear example of the benefits that applied visualization can bring.

The great majority of research focuses on representing the stock market in general to give a sense of how the market is trading. Lei and Zhang have conducted an interesting research on this direction. They suggested several different visualizations that should help financial analysts. These visualizations such as stock variation plot (which highlights which sector is attracting the most market orders), a ring chart (similar to the TreeMap, but circular) that gives an overview of the market or even a pattern matching visual tool (Lei & Zhang, 2010). Another interesting work by Schrek et al. uses visualization to find if there are any tendencies of the market data. Their research focuses on trying to increase the ability of financial analysts to discover patterns on data (Schreck et al., 2007). Ziegler et al. use visualization not only to find patterns on data but also to give an overview of the market and the industry. Their tool allows comparing the features of companies, assets and even entire countries. This

Does a visual data mining approach improve stock picking for equity analysts?

allows analysts to explore and discover interesting facts such as patterns that indicate that markets or countries go into turbulence (Ziegler et al., 2010).

This research sets the ground in which build the current project to improve Eikon's stock picking utility. A new and improved visualization that allows exploring the data interactively seems the best approach to fulfil TR's requirements.

## Methodology

Due to the particular characteristics of the current research question being faced it seems a problem-driven approach research methodology would be the best fit. Design studies are a form of such problem-driven research. A problem-driven research is such where the goal is to work with real users to solve their real-world problems (Sedlmair et al., 2012).

A design study can be defined as a project in which researchers analyse real-world problems, design a visualization that helps to solve the problem and reflect about the lessons learned to refine the field of visualization design (Sedlmair et al., 2012).

A design study approach defines a process with several phases that dictate how to conduct the research. Each part of the study can be quite deep and involve a complex analysis but for the current proposal an overview will be provided for the major steps involved (Munzner, 2009).

**Analysis:** Is the process of translating domain-specific data and tasks into abstractions (elements belonging to information systems) that can be translated into a problem that can be solved thanks to visualization. The analysis will consist in capturing the main indicators used in stock picking to focus on those that really add value to the problem being solved. Not only data needs to be abstracted, tasks also play an important role on a design study. Tasks have to be elicited to understand at what specific point the visualization can be used to improve their execution.

Within fundamental analysis there are several different techniques that can be applied to stock picking such as Value Investing, Growth Investing, GARP Investing etc. These techniques involve evaluating a set of different indicators for a given company (or set of companies) and probably comparing these to the industry, to the overall market etc. These are tasks that have to be analysed to try abstract workflows where visualization can help improve the process.

**Real-world problem:** A proper design study requires real users and a real set of data. In the current piece of work we will be working with real data thanks to TR's data feed on stock information. This guarantees that domain experts (which need to be involved in the design study compulsory) can understand and make sense of the data visualized and judge whether the approach taken improves the process or not. TR has domain-experts

**Design:** This is the most creative part of the design study process. The main objective is to explore the vast amount of different data visualization techniques available to iteratively reduce this number to the most suitable ones (Munzner, 2009). At the same time the researcher will further define the tasks and data needs that should be clearer as the process offers a deeper insight on the problem. Knowledge on the vast number of possibilities information visualization offers is a requirement for this step: only knowing the vast number

of techniques that can be used one can judge which are more suitable than others. The time for this project is quite limited so there is not much room for a fully iterative approach. An interesting technique that will be used in this study is Parallel Prototyping (DOW et al., 2010), a technique that is based on exploring several different propositions at once instead of trying one per iteration leading to better design results. The produced prototypes will then be tested with real-world domain experts in an iterative process with the objective of refining the most valuable prototypes whilst discarding the less useful.

**Validation:** This step is directly correlated with both the Design and the Analysis. Both steps of the process will need to be validated. For example, once activities are abstracted these need to be validated with the domain experts to make sure these abstractions are correct are part of the typical workflow.

**Reflection:** This main objective of a design study is to reflect the lessons learnt on solving the real-world problem. One of the objectives of this project is to add to the body of knowledge on financial data visualization, thus this part of the project is important.

As can be observed the flow within the different parts of the design study is not sequential, it looks more like an iteration-based approach. This is because the degree of complexity of the problem that is being faced. The researcher has little to no knowledge on how financial markets work compared to domain experts. But instead of working separately the idea is to combine the expert knowledge of both and work together towards a solution.

## **Software Development Methodology**

Any software development project requires a certain methodology to be followed. This project has the objective to develop a functional prototype that uses data visualization to improve stock picking. The software methodology that will be used is Incremental Development (Larman & Basili, 2003). The development of the final prototype will be divided in small iterations. For each iteration a new set of prototypes will be developed, validated with the domain experts and only the best will be refined. Using multiple prototypes ensures the researcher is exploring multiple alternatives, although it requires more work it leads to better results (DOW et al., 2010). Usually this approach in developing InfoVis systems requires a process of shared learning between the researcher and the domain experts. Almost in every iteration new requirements are created, discovered or elicited and implemented to be validated on the next iteration with the creation of a new prototype this approach is known as iterative Prototyping (Mnkandla, 2005).

## **Project Feasibility**

To properly execute the project, knowledge in the following fields is required: financial analysis, statistics, visualization, big data and software development. As a computer scientist the researcher is familiar with statistics and software development. Additionally having spent 4 years programming at a professional level I have developed the skills required to be able to face the programming challenges that might arise.

Regarding financial knowledge, the great benefit for the researcher is that TR has extremely talented financial experts in-house that can be inquired about any financial information required. Having these experts to participate in the design process is crucial for the success

Does a visual data mining approach improve stock picking for equity analysts?

of the project and as it has been assured by TR the researcher will have access to the experts anytime.

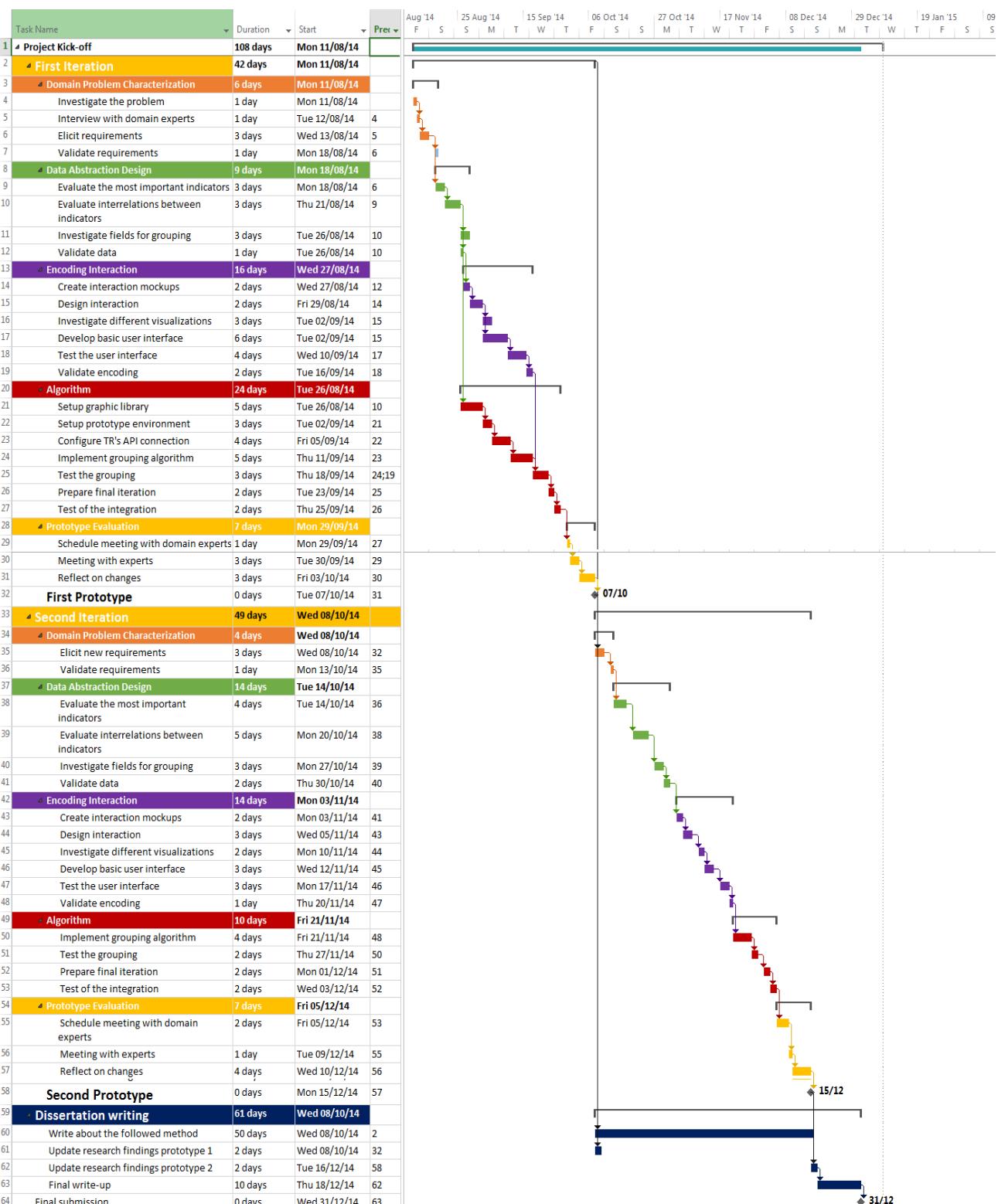
Does a visual data mining approach improve stock picking for equity analysts?

## Risk Register

Nature of the risk	Description	Impact	Likelihood	Triggers	Control	Mitigation	Contingency Plan
<b>Data availability</b>	Not being able to access the financial data	Very High	Very Low	Not enough data / No data to build the visualization on	Halfway through the project	Request access to the financial data API as soon as possible	Talk to the employer to gain access to the financial data they have
<b>No access to employer's experts</b>	Being unable to contact the employer's experts on the field that can contribute to the project	Medium	High	The only source of information contributed for the project is bibliography		Try to arrange informal meetings with team members	Bring the employer's supervisor support to help schedule the meetings
<b>Excessive pressure on project result</b>	There might be tension between the academic objective of the project and the employer's interests	High	Low	Not being allowing time to write for the dissertation as the employer is pushing to meet the deadline	Track the progress of the document elaboration against the timeline	Make sure the employer understands the importance of the project document, not only of working software	Discuss with the employer the main objective of the project, bring the supervisor in if required
<b>Time Constraints</b>	Not able to deliver the project on time	High	Medium	Way behind the projected schedule	Weekly self-evaluations of the milestones of the project	Be realistic on the project scope as time is quite limited	Try reduce the project scope
<b>Technical abilities</b>	Not enough technical skills to implement a solution	High	Low	The progress made with the project is much more smaller than expected	Review the functionality of the project prototype against the feature list	Review the project scope and make sure there is enough time allocated for technically complex tasks	Check whether the feature is completely required for the project to progress.
<b>Loss of work</b>	Loss of code or work might occur due to a failing computer for instance	Very high	Medium	Work/code lost	Make sure the backup is up to date (daily)	Set up a backup/repository for both the work done and the code	Restore the latest backup from the server

Does a visual data mining approach improve stock picking for equity analysts?

## Schedule



## 8.2. FUNCTIONALITY REVALIDATION FEEDBACK SUMMARY

Domain expert names have been made anonymous.

- Scatterplot
  - **R & G:** I'm missing a 'select' equities functionality to select some equities for further analysis.
- Clusters
  - **E & R & Ed.:** Clusters are useful for reducing the noise.
  - **A & R & G:** The clustering works but when you expand, the positioning of elements is arbitrary.
  - **E & R:** Mini-map is useful and helps navigate and see all the equities at once.
  - **AB:** I think the listing of equities is better than the 'declustering' of equities because I can't really see where the equities were actually located, can't differentiate from the equities on the background.
- Quick Rank
  - **E:** Quick rank is OK, works fine.
  - **R & G & Ed. & AB:** It is difficult to distinguish between rank position. Rank 1 colour is similar to rank 3.
  - **Ed.:** How about show me the top 15, or the top 10 and then highlight those on a different colour? Now is not clear at all how can we display that.
  - **AB:** It is a really useful feature. Would be great being able to rank by any metric, like when using multi-factor ranks.
- Trends & Medians
  - **A & E & G & Ed. & AB:** Trends are useful, as well as medians for comparing companies
  - **E:** Trend is great, we could plot many other variables and to compare between stocks. For example the size of the fleet for airlines, or oil barrels per day. You could also draw the trend for the entire industry, or the economy.
  - **R:** You might want to know where ALL the companies were in the past (not only the one you select). Also you might need a play button to see how they evolve. (Gonzalo not sure about this, it might make the interaction too complex for little benefit).
  - **G & Er. & AB:** We could have a bubble that shows the aggregate for the industry (or another benchmark). Or a index such as the S&P500 or FTSE100. We could allow users to select the benchmark. We actually have the data for the additive values per industry.
  - **G:** The trend is yearly but you have metrics that are end-of-fiscal year for example that won't match. We should then interpolate.
  - **Er.:** Might be useful to draw a bubble representing the directly competing peers and their trend.

- **Ed.:** I think it would be useful to chart a **regression line** for the entire industry. Also, couldn't we add the **average** instead of only the median of the industry?
  - **Ed. & AB:** Would be useful to see how the metrics move through time, like for example how the company was doing last year in terms of 'Net Income' and how it has progressed through time.
- Parallels
  - **A & Er.:** Parallels look as time-series at first you should give a hint they are not.
  - **E & R & G & Er. & Ed.:** They seem really useful and match the workflow. You could filter down pretty easily from a wide range to a lower range. By setting some metrics up you can easily filter down to the best performers.
  - **R:** You might normally look for composite metrics such as EV/EBITDA for example, normally not direct EBITDA measures.
  - **R:** Have these been tested? What users think about them? I think once they understand how they work they are powerful.
  - **Ed.:** I'd like to have the zoom-in feature so that to focus on a set of. Also when zooming in it should scale.
  - **Ed.:** Would be really helpful to see the progression for the several indicators through time.
  - **AB:** When I click on a table row I expect the parallel to activate on a 'comparison mode'. Like what actually happens when I click the scatterplot, but when clicking on the row. (usability issue mainly)
  - **Ed. & AB:** I expect to be able to Keep/Exclude the data that I'm evaluating.
- Categorisation
  - **E & R & G & Er.:** Categorising data is useful, both for filtering as well as knowing the data distribution of your selection.
  - **A & E:** I'm missing the complete hierarchy of categories.
  - **A & E & R & G & Er.:** The comparison is good, you sometimes compare different companies between them or against the industry.
  - **A & Ed.:** Might be good to have stats of the industry (trends or financial ratios such as P/E) when you move over the type of category (Something like Avg EBITDA for the industry).
  - **E & A:** Hierarchical categorisation might be more useful in some cases but the simplification (by sector) is still quite useful.
  - **Ed.:** It would be even more useful being able to categorise by a statistical feature of the data such as 'Best P/E' or 'Top Net Margin'. This is a sort of statistic analysis of the input data to create the categories. For example a use case would be: 'Which high revenue companies are not considered to be good by the analysts?'.
- Workflow
  - **E:** Workflow is better with the parallels as you can refine really easily the filters by just drawing boxes, and then updating the initial filter.
  - **E:** Analysts sometimes use a given equity as a baseline to compare with all other equities.
  - **A:** The approach you've taken consolidates different investment approaches in one. It also enhances the interaction (breaking the: refine the filter –wait for server data - evaluate – refine loop)

Does a visual data mining approach improve stock picking for equity analysts?

- **R:** Key point here is the exit points, and how this visualization can be fitted within the workflow (linking with the information of company research for example).
- Others
  - **Er.:** I see many analysts looking to add an arbitrary equity to compare it with the current equities you have on the scatterplot.
  - **A & R & G:** Might not look good on a small screen.
  - **G:** Everything is connected, linked and gives place some really powerful interactions.
  - **Er.:** Complete new interaction level (bubbles as industries and you can drill down to lower levels).
  - **G:** Should be really useful to have a WebEx as a walkthrough the features of the prototype.

### 8.3. ABSTRACTED TASKS

**TASK01** derived from **FR01:**

{Compare, Similarity}

**TASK02** derived from **FR02:**

{Present, Trends}

**TASK03** derived from **FR03:**

{Identify, Extremes}

**TASK04** derived from **FR04:**

{Explore, Outliers}

**TASK05** derived from **FR05:**

{Discover, Features}

**TASK06** derived from **FR07:**

{Identify, Features}

**TASK07** derived from **UR08:**

{Lookup, Distribution}

**TASK08** derived from **UR09:**

{Present, Correlation}

### 8.4. QUESTIONS FOR THE COMPARISON TEST

#### 8.4.1. BENCHMARK TEST

1. What 3 are the companies with the highest values on 'EBITDA Actual', 'Revenue' & 'Market Capitalisation'?

USA	EU
-----	----

Does a visual data mining approach improve stock picking for equity analysts?

a) MSFT.OQ, JNJ.N, WMT.N b) WMT.N, XO.MN, AAPL.OQ c) AAPL.OQ, T.N, GE.N	a) ROG.VX, NESN.VX, DIAGn.DE b) ENI.MI, ROG.VX, VOWG_p.DE c) SIEGn.DE, UNc.AS, BASFn.DE
<b>Correct:</b> b	<b>Correct:</b> b

2. From the following 3 companies listed below, what company do you think has grown more from 2009-2014 in terms of both 'EBITDA' and 'Revenue'?

*The company that has grown more is:*

For USA: (WMT.N, AAPL.OQ, XOM.N)

For EU: (TEF.MC, VOWG\_p.DE, EN.MI)

USA	EU
a) WMT.N b) AAPL.OQ c) XOM.N	a) TEF.MC b) VOWG_p.DE c) ENI.MI
<b>Correct:</b> b	<b>Correct:</b> b

3. Now, tell me from these 3 companies what is the one with the highest 'Dividend Yield'?

*The company with the highest 'Dividend Yield' is:*

For USA: (WMT.N, AAPL.OQ, XOM.N)

For EU: (TEF.MC, VOWG\_p.DE, EN.MI)

USA	EU
a) WMT.N b) AAPL.OQ c) XOM.N	a) TEF.MC b) VOWG_p.DE c) ENI.MI
<b>Correct:</b> c	<b>Correct:</b> c

4. Select the company that is over the median for 'Revenue' for both 'Consumer Cycicals' and 'Technology'?

USA	EU
a) PSX.N b) ORCL.N c) GM.N	a) ENI.MI b) ASC.MC c) VOWG_p.DE

Does a visual data mining approach improve stock picking for equity analysts?

d) AMZN.OQ	d) REP.MC
Correct: b	Correct: c

5. Can you please tell me 2 companies that have a high Quick-Rank?
- 

USA – Possible Answers	EU – Possible Answers
GOOGL.OQ, PCL.OQ, DAL.N, GS.N, ORCL.N, DAL.N ...	SAPG.DE, GRMN.OQ, EUP.PA, CONG.DE, DAE.S, SCGFI...
Correct: c	<b>Correct: a</b>

6. Now we assume a Price-to-sales strategy. We will be looking for companies that match the following features.

*NOTE: Consider the sectors we had before: 'Consumer Cyclicals' and 'Technology'.*

THESE CHANGE DEPENDING ON THE DATASET	
USA	EU
'Price to Sales' < 5.0 '3 Year % Change of Net Profit Margin' > 0.2% '3-year CAGR (Compound Average Growth Rate) Basic Normalized EPS' > 0.5%.	'Price to Sales' < 2.0 '3 Year % Change of Net Profit Margin' > 0% '3-year CAGR (Compound Average Growth Rate) Basic Normalized EPS' > 0.5%
'Free Cash Flow, Unlevered' > 0	'Free Cash Flow, Unlevered' > 0

Use the mechanisms provided to filter with these parameters and select 2 companies that appear among the highest '3-year CAGR Basic Normalized EPS' in these terms.

---

USA	EU
MAN.N, DAN.N, TSO.N, CNW.N, RX.N, WY.N...	VOWG_p.DE, HOLN.VX, AFLG.DE, REP.MC, SFLG.MI, KUNN.SE...

## 8.4.2. MEDIUM INSIGHT TEST

- Can you describe me how is the data distributed when charting with 'EBIT', 'EBIDTA' & 'Market Cap'?
  - These two metrics seem positively related (the higher one, the higher the other).
  - They seem negatively related (the lower one, the higher the other).
  - There seems not to be any relation

Does a visual data mining approach improve stock picking for equity analysts?

And how about 'EBITDA' and 'Revenue'?

- a) These two metrics seem positively related.
- b) These seem to be negatively related.
- c) There seems not to be any relation

Correct Answers: a & a

2. Now we will consider the two following industries: 'Consumer Cyclicals' and 'Consumer Non-Cyclicals'. How is data distributed in terms of 'Net Profit Margin' and 'ROA'.
  - a) In 'Consumer Cyclicals' there is higher variation, 'Non-Cyclicals' have the same behaviour.
  - b) In 'Consumer Cyclicals' they have the same behaviour, 'Non-Cyclicals' there is higher variation.
  - c) Both industries have similar behaviour

Correct: b

3. How does 'Healthcare' compared to 'Technology' in relation to the 'EBITDA' and 'Revenue'.

USA	EU
<ul style="list-style-type: none"><li>a) Healthcare sector does better</li><li>b) Healthcare sector does worse</li><li>c) They are really similar</li></ul>	<ul style="list-style-type: none"><li>a) Healthcare sector does better</li><li>b) Healthcare sector does worse</li><li>c) They are really</li></ul>
Correct: a (c maybe too)	Correct: b

#### 8.4.3. INSIGHT EXPLORATION TASK

1. Focus on 'Industrials' & 'Consumer Cyclicals'. Tell me what thoughts you have and what can actually discover from these two different industries.
2. Filter down to a set of equities that seem interesting in any sector. Feel free to explore the data until you know it thoroughly.

### 8.5. FINAL COMPARISON TEST RESULTS

#### 8.5.1. TASK DISTRIBUTION

	EU dataset, Prototype	EU dataset, Screener	USA dataset, Prototype	USA dataset, Screener
B1	Subject A	Subject B	Subject C	
B2	Subject A	Subject B	Subject C	
B3	Subject A	Subject B	Subject C	
B4		Subject A	Subject B	Subject C
B5		Subject A	Subject B	Subject C
B6		Subject A	Subject B	Subject C

Does a visual data mining approach improve stock picking for equity analysts?

MI1	Subject A	Subject B	Subject C	
MI2	Subject A	Subject B	Subject C	
I1	Subject A		Subject B	Subject C
I2	Subject A		Subject B	Subject C

### 8.5.2. TEST A RESULTS

RESULTS TEST A				
Task	Answer	Correct	Time	Comments / Feedback
B1	B	Yes	03:41	It took time to figure out this was better performed with the scatterplot.
B2	B	Yes	03:00	Was trying to figure out how the parallels worked before being able to actually see the trends
B3	C	Yes	00:58	Adding the indicator and clicking on the equity he got the result straight away.
B4	A	Yes	04:43	I can tell by looking at the numbers, or I could create a percentage to know how much it has grown, it is much more difficult.
B5	A	Yes	01:49	I don't know how to reduce to 30, I can tell you the 30 first in terms of a given metric but don't know how to reduce it.
B6	C	No	02:11	With the old screener you can do this numerically fairly easily, but the time to calculate it, reload the table etc. makes the task lengthy
B7	LLTC.OQ, TTWO.OQ, PCLN.OQ	Yes	02:09	It is difficult to say, the colouring is not ideal. It is tough to differentiate between the different rankings.
B8	OTHER (MOODYS,PL C.OQ)	Yes	04:23	He struggled to see companies with the highest rank and started panning on the scatterplot trying to figure out which one was the best.
B9	CONN.NN, MODN.N	Yes	04:03	"I can see there are quite a lot of efficient companies that comply with this strategy but I'm refining to get the best ones."  He was playing with the boxes to try find the best ones, spent some time refining the results to get the best answer.
OVERALL TEST TIME				
MI1	A	Yes	01:27	a) You can see the line, the tendency.

Does a visual data mining approach improve stock picking for equity analysts?

				Every analyst knows that they are related but you also can see that on the scatterplot.  b) "Between EBIT and Revenue, it depends mainly on the company."  He started playing with the trends and figuring out what was the relation between the measures. "I can't see a clear relation between those two measures." It mainly depends on the specific company
MI2	C	No	03:13	"More or less the same, it just changes on the number of companies available."
MI3	A		1:49	"By grouping by sector on the old screener it can pretty easily differentiate for these two metrics on the two sectors."  Did not take long to take the result.
MI4	-	-	-	-
I1	03:15	Revenue and Net Sales is obviously related, it's a linear relationship. Free Cashflow is a reverse relation, you can clearly see that on the parallels.		
I2		I'd like to see the 'Net Income' of these industries, it would be really useful to have that.		
I3		I don't know how to filter down into 10 equities. What I'll do is explore the data and find the 10 I think are the best.		
I4	04:08 to finding an equity worth buying	"I'm adding ROE, Net Profit Margin, EV/EBIT"  Now I'm looking for undervalued companies. Companies that are efficient but are a multiplier of their value 'undervalued'.  He started drilling down with the parallels (did not work with the scatterplot at all).  "I've found Take2 Interactive is undervalued, let's see what analysts say about it"  Added analyst revision score and:  "Ah, so its now undervalued (for a multiplier of its EV/EBIT) and analysts think is actually is really good. So it's a company worth		

Does a visual data mining approach improve stock picking for equity analysts?

		purchasing"
--	--	-------------

### 8.5.3. TEST B RESULTS

RESULTS TEST B				
Task	Answer	Correct	Time	Comments / Feedback
B1	B	Yes	2:45	Really difficult to know not seeing the figure.  <b>NOTE &gt;</b> It was a small screen so you could not actually see the column that he added.
B2	B	Yes	4:51	Using the mouse pad is difficult to navigate the scatterplot
B3	C	Yes	1:45	It is hard to tell, using the combination of scatterplot and parallels.
B4	B / D / C	No, no clear answer	9:00	First started evaluating the different medians for the metrics but did not find out how to get the ones that matched both metrics. Then he created a filter on the equity screener to figure out what companies were over the median.
B5	TW.N, X.N	Yes	1:50	Order by the 'Quick Rank' easily.
B6	TWC.N, RL.N	Yes	3:25	This is really easy on this prototype. To translate this on the current screener, you could actually use sliders. These sliders might help you emulate the same as you have here.
OVERALL TEST TIME				
MI1				
MI2	Could not actually do this because we spent too much time setting up the screener remotely.			
I1	If I'm doing industry / regional I'd use a different app, because now the screener is not really prepared for that. I can add the colouring to the categories and do the filtering on the prototype now.			
I2	To do the same analysis, you'll have to modify the screen parameters, run the screen and go back again. On the prototype you just compare information and between industries visually.  I only see that in some extreme situations where the metrics are really complex that this filtering process might not really work. For example for extremely complex equations such as complex formulas on the screener this might need rethinking.  Industrials seem to have a much higher Beta, more volatility. I can see there is a			

Does a visual data mining approach improve stock picking for equity analysts?

	higher EPS growth straight away on Industrials.
	There is an interesting company CEDAR Fair, it seems as an outlier to me.

### 8.5.4. TEST C RESULTS

RESULTS TEST C				
Task	Answer	Correct	Time	Comments / Feedback
B1	B	Yes	4:11	Using the filtering process and filtering typing more or less randomly filter values
B2	B	Yes	4:04	
B3	B	Yes	2:19	Finding the required companies and comparing by hand each value on the table
B4	C	Yes	3:34	Struggling to find the companies that where over the median, exploring all the space until found them.
B5	SAP.G, CFR.VX	Yes	5:41	Struggling to see the difference on the rankings. "Can't see the value of the ranking, they are really similar"
B6	DUEG.DE, OTEPr.AT	No	3:48	Started adding the metrics to the screener and
OVERALL TEST TIME				
MI1	3:05			
MI2	4:17	<p>It is fairly even between industries.</p> <p>If you actually want to invest in a sector with a considerable volatility you might invest on the Industry sector.</p> <p>If you want a more stable approach you can use 'Consumer Cyclical'. So if you want to play it self do invest so.</p> <p>Industrials seem much more volatile.</p>		
MI3	03:05			
I1	10:50	<p>If I was doing this analysis on the old screener I could not actually do it, I could only play with 3 dimensions. I can't add the categories and compare them. I can only add them separately and compare them.</p> <p>The screener is pretty good at what it does but if you want to do more powerful things such as comparing against the median or the average is much more difficult / can't be done. You can only see 3 indicators against the many you have here.</p>		
I2	07:48			

Does a visual data mining approach improve stock picking for equity analysts?

## 8.6. FINAL COMPARISON RESULTS SUMMARY

### Task performance

Task	Subject A				Subject B				Subject C						
	Time	Dataset Used	Tool	Time	Dataset Used	Tool	Time	Dataset Used	Tool	Average Prototype	Max	Min	Average Screener	Max	Min
B1	03:41	Europe	Prototype	02:45	USA	Prototype	04:11	Europe	Screener	03:13	04:11	02:45	04:11	04:11	04:11
B2	03:00	Europe	Prototype	04:51	USA	Prototype	04:04	Europe	Screener	03:55	04:04	03:00	04:04	04:04	04:04
B3	00:58	Europe	Prototype	01:45	USA	Prototype	02:19	Europe	Screener	01:21	02:19	02:43	02:19	02:19	02:19
B4	02:11	Europe	Screener	09:00	USA	Screener	03:34	Europe	Prototype	03:34	03:34	03:34	05:35	09:00	02:11
B5	04:23	Europe	Screener	01:50	USA	Screener	05:41	Europe	Prototype	05:41	05:41	05:41	03:06	04:23	01:50
B6	04:03	Europe	Screener	03:25	USA	Screener	03:48	Europe	Prototype	03:48	03:48	03:48	05:35	04:03	03:25

### Error Rates

Task	Subject A			Subject B			Subject C					
	Correct	Dataset Used	Tool Used	Correct	Dataset Used	Tool Used	Correct	Dataset Used	Tool Used	Errors Prototype	Errors Screener	
B1	Yes	Europe	Prototype	Yes	USA	Prototype	Yes	Europe	Screener	0	0	
B2	Yes	Europe	Prototype	Yes	USA	Prototype	Yes	Europe	Screener	0	0	
B3	Yes	Europe	Prototype	Yes	USA	Prototype	Yes	Europe	Screener	0	0	
B4	Yes	Europe	Screener	No	USA	Screener	Yes	Europe	Prototype	1	0	
B5	Yes	Europe	Screener	Yes	USA	Screener	Yes	Europe	Prototype	0	0	
B6	Yes	Europe	Screener	Yes	USA	Screener	No	Europe	Prototype	0	1	
										TOTAL: 1	TOTAL: 1	



Does a visual data mining approach improve stock picking for equity analysts?