



## Exploring the predictive value of Google search data in finance

Erasmus University Rotterdam  
Erasmus School of Economics  
Bachelor thesis, Finance Group

Student: William Beckwith (433908wb)

Supervisor: Xia Shuo

### Abstract:

The publication of search frequency data by Google Trends has opened many opportunities for analysing and predicting human behaviour. One of these opportunities is the use of Google search data to test the hypothesis that an increase in investor attention causes a price increase. This paper provides evidence that the hypothesis is consistent in The Netherlands from 2013 to 2018. Furthermore, Google search data is used to create a hypothetical investment strategy based on economic search terms that generated a 127% return from 2004 until 2008. Although the search data of economic terms is unlikely to have a direct impact on stock prices, it may be able to give indications of certain macroeconomic activities.

Keywords: Google, Trends, search, predict, returns

## Table of content

1. Introduction	3
2. Data	5
2.1. Google Trends	5
2.2. Other methods of attention	5
2.3. Returns	6
3. Methodology and results	7
3.1. Comparing measures of attention	7
3.2. Impact of search frequency on abnormal returns	9
3.3. Using search terms to predict investor sentiment	12
4. Conclusion	15
5. References	16
6. Appendices	17

# 1. Introduction

The unique set of search frequency data that was first released by Google in May 2006 has the potential of being a beneficial predictive power. Vosen & Schmidt (2011) found that Google Trends data outperformed survey-based indicators when predicting private consumption. Askitas & Zimmermann (2009) used Google Trends data to forecast unemployment. A lot of other successful research has been done to anticipate for example video game and film sales, the rank of songs on top charts, tourism, and even disease prevalence. It is of no doubt that time series data on search frequency can be used to predict human behaviour.

The attention hypothesis described by Barber & Odean (2007) suggests that individual investors are net buyers of stocks with higher than usual attention. The idea is that individual investors do not have the time to investigate every stock in the world and therefore base their portfolio on stocks that have recently caught their attention. As individual investors are often bounded by short sell constraints or high transaction costs for short selling, negative attention has less impact on a stock return than positive attention. This means that in general, the average abnormal return of a stock that has had a recent increase in attention is significantly positive.

Unfortunately, attention is something that cannot be measured exactly and therefore proxies are used. Barber & Odean (2007) used variables such as the number of times a stock was mentioned in the news, abnormal trading volume and extreme daily returns. These proxies each have their own disadvantages. Investors get their news from many different sources. Whether it be television, radio, word of mouth or social media, a lot of investors get their news from other sources than regular news articles. The other two proxies obviously capture attention. When a stock has an extremely positive or extremely negative return people want to find out why. Positive abnormal trading volume indicates many people were buying and selling and thus paying attention to that stock. However, this does not go to say that a stock experiencing an attention increase always has a positive abnormal trading volume or a previous day extreme return.

Recent research argues that search frequency data is a more direct proxy than those mentioned previously. "Return or turnover can be driven by factors unrelated to investor attention and a news article in the Wall Street Journal does not guarantee attention unless investors actually read it" (Engelberg & Gao, 2011). Using empirical data from firms in the Dow Jones Index, Engelberg & Gao (2011) find that Google Trends captures attention in a timelier fashion than other attention proxies and likely measures the attention of individual investors. This leads to the first

hypothesis: 'Is Google Trends search data a viable proxy for attention when testing the attention hypothesis?'

Other research by Preis, Moat & Stanley (2013) takes a different approach to predict stock prices by using Google Trends data to analyse human behaviour. Specifically, they found that the Dow Jones Industrial Average (DJIA) would tend to go down when people in the U.S. had searched on 'debt' more than usual in the week before. A hypothetical portfolio they created, returned 326% from 2004 to 2011 by buying when 'debt' had a low search frequency in the previous week and selling otherwise. It is important to note that they made the unrealistic assumption of no transaction costs, however, it is still a remarkable result as a simple 'buy and hold' strategy would have only secured a 16% return over the same period. This leads to the second hypothesis: 'Can investment strategies based solely on Google Trends data generate a significant excess return?'

## 2. Data

### 2.1 Google Trends

Google Trends provides search frequency data on a relative basis for a given search term, time range, region, and category. The output will show numbers between 0 and 100 for each day, week or month for the given query. For example, if the output shows the search frequency was 100 for a specific day, this means that the maximum search frequency was achieved on this day. If the output was 50 on a specific day, this means that the search frequency was half that of the maximum search frequency within the set constraints. Whether the output is shown per day, week or month depends on the set time range. In this paper, different time ranges and consequently different scales will be used for different purposes. The set region will be worldwide for comparing measures of investor attention and exploring the impact of search frequency on abnormal returns. For attempting to predict investor sentiment, the set region will be narrowed down to The Netherlands. The set category is left blank as this is not of any relevance for the research. Search frequency as described in this paragraph will from now on be referred to as *SF*.

In this paper, the focus will be stocks in the Amsterdam Exchange index (AEX) as of the 1<sup>st</sup> of March 2018 to keep the data collection and cleaning task manageable. Companies that have stock data of less than 5 years will be left out which eliminates ABN Amro, Altice, Galapagos and NN Group. From the remaining 21 stocks search data will be downloaded from Google Trends on a monthly scale from March 2013 until March 2018. To ensure the majority of searches captured are actually searches for financial reasons, the full company name is used. For example, 'Heineken International' is used to filter out people who are searching for 'Heineken' beer. Similarly, 'Royal Dutch Shell' is used to exclude the data from people searching for a simple shell. A list of search terms used for all companies will be presented in Appendix A. Of course, one cannot exclude that people have searched the company name for non-financial reasons such as looking for a job. As the main interest is investor attention, potential noise should be taken into consideration.

### 2.2 Other measures of attention

The news data is obtained from Factiva for the same period and scale as the Google Trends data. The value of the variable *News* for a certain company is the number of times that company was mentioned in an article from the Dow Jones Newswire in a specific month. The filter for articles from the Dow Jones Newswire was used to ensure all news mentions were finance related. The

most recent two years were available on a monthly basis, however for the three years prior to that the number of news mentions had to be looked up manually for each month. The variable ‘Analysts’ describes the number of times an analyst made a recommendation about a stock according to I/B/E/S. Again, the data is distributed monthly from March 2013 until March 2018. It is important to note that all three variables (search frequency, news, analysts) could describe positive and negative attention. Although negative attention might not have a positive effect on the stock price, it should not be excluded from the data as the hypothesis is about attention in general.

### 2.3 Returns

Abnormal returns are retrieved from Datastream for the 21 AEX companies that will be analysed from March 2013 until March 2018. The calculation of abnormal returns uses the market model with an estimation period of 550 trading days. The chosen benchmark is the S&P Europe as it contains a wide variety of industries and covers 350 stocks throughout Europe. The  $\alpha$  and  $\beta$  of each stock are calculated by taking the intercept and the slope respectively, between the stock return and the market return in the estimation period. The abnormal returns during the evaluation period for each stock are then calculated as shown below, where  $R_t$  is the stock return in period  $t$  and  $RM_t$  is the market return for period  $t$ .

$$AR_t = R_t - \alpha - \beta * RM_t$$

For the second hypothesis, returns from AEX will be used. The prices of the AEX from January 2004 until June 2018 are downloaded from Yahoo Finance. As the AEX is representative of 25 different companies and has characteristics of a market index, the returns will not be abnormalized.

### 3. Methodology and results

#### 3.1 Comparing measures of attention

First, the correlations between the different measures of investor attention are calculated. The logarithms of these measures are taken to reduce skewness. Because the variables *Analysts* and *News* have observations that are 0, each observation will get +1 before taking the logarithm to prevent undefined values. The correlations will be taken for each stock individually and then averaged across stocks, the results are presented in Table 1. The highest correlation is between  $\log(SF)$  and  $\log(News + 1)$ , a likely explanation for this is that investors search companies that appear in the news to find more information. The lowest correlation is between  $\log(SF)$  and  $\log(Analysts + 1)$ , which could imply that analysts do not necessarily base their recommendations on stocks that have a lot of investor attention.

**Table 1: Correlation between measures of attention**

*The logarithm is taken of the search frequency, the number of analyst recommendations and the number of news mentions. Of these three logarithmic variables, the correlation is calculated and presented below.*

	Log(SF)	Log(Analysts +1)
Log(Analysts +1)	0.0662	
Log(News +1)	0.2690	0.1356

Next, a vector autoregression (VAR) is used to examine the weekly lead-lag relation among measures of attention. A time trend is included as an exogenic variable to account for a possible increase in general investor attention over time. The VAR is run for each individual stock and the coefficients are averaged over all stocks and reported in Table 2 together with corresponding p-values. The p-values are computed using a bootstrapping method to account for cross-sectional correlation in the error terms. For each variable, the panel of coefficients is bootstrapped to construct 1000 replication panels which are all t-tested under the null hypothesis that the average coefficient is 0. As can be seen in Table 2, the significant coefficients belong to the lagged independent variables in the regressions where the current dependent variable is the same. In other words, the previous week search frequency has predictive value for the following week search frequency. The same goes for the number of analysts and the number of news mentions, even when a time trend is included.

**Table 2: Vector autoregression (VAR) on measures of attention**

The VAR is run for each individual stock and the coefficients are averaged over all stocks and a time trend is included as an exogenic variable. The p-values are computed using a bootstrapping method to account for cross-sectional correlation in the error terms. For each variable, the panel of coefficients is bootstrapped to construct 1000 replication panels which are all t-tested under the null hypothesis that the average coefficient is 0. P-values are presented in parentheses under the coefficient and \*, \*\*, and \*\*\* represent significance at the 10%, 5%, and 1%, respectively.

	log(SF)	log(Analysts + 1)	log(News + 1)	R <sup>2</sup>
log(SF)	0.3515 (0.000)***	-0.0401 (0.669)	-0.0626 (0.263)	43.09%
log(Analysts + 1)	0.0421 (0.604)	0.7371 (0.000)***	0.0377 (0.412)	82.23%
log(News + 1)	0.3407 (0.270)	0.0975 (0.849)	0.1174 (0.021)***	28.49%

Finally, the relationship between  $SF$  and the other measures of attention is examined in a set of regressions. The dependent variable will be abnormal  $SF$ , which is calculated by subtracting the logarithm of the median of the previous  $\Delta t$  months from the logarithm of the current month:

$$ASF_t = \log(SF_t) - \log\{Med(SF_{t-1}, \dots, SF_{t-\Delta t})\}$$

The median of the previous  $\Delta t$  months is taken to capture the ‘normal’ attention a stock gets. Some stocks may get more attention in certain seasons, so taking a time window of, for example, two months, makes the  $ASF$  robust to these seasonality’s. A high  $SF$  does not necessarily mean a high increase in attention, whereas a high  $ASF$  does. The independent variables are  $\log(Analysts + 1)$  and  $\log(News + 1)$ , they are regressed on  $ASF$  for a  $\Delta t$  of 2, 4 and 8 months in separate regressions. Because the dataset consists of panel data where longitudinal observations exist for the same measurements of attention, a linear model with fixed effects is used to represent the measurement-specific means. Standard errors are robust and clustered by firm. The results are reported in Table 3.



**Table 3: Regression on ASF**

The dependent variable in each regression is *ASF*, the independent variables are  $\log(\text{Analysts} + 1)$  and  $\log(\text{News} + 1)$ . *ASF* is calculated with a  $\Delta t$  of 2, 4 and 8 months separately. Each regression is a linear model containing fixed effects and the robust standard errors are clustered by firm. P-values are presented in parentheses under the coefficient and \*, \*\*, and \*\*\* represent significance at the 10%, 5%, and 1%, respectively.

	$\Delta t = 2$ months	$\Delta t = 4$ months	$\Delta t = 8$ months
Constant	-0.0530 (0.000)***	-0.0495 (0.000)***	-0.0425 (0.001)***
$\log(\text{Analysts} + 1)$	-0.0086 (0.040)***	-0.0125 (0.017)***	-0.0207 (0.005)***
$\log(\text{News} + 1)$	0.0386 (0.000)***	0.0388 (0.000)***	0.0402 (0.000)***
Observations	1281	1281	1281
Clusters (firms)	21	21	21
R <sup>2</sup>	0.0166	0.0194	0.0205

As can be seen, *ASF* is negatively related to  $\log(\text{Analysts} + 1)$  and positively related to  $\log(\text{News} + 1)$ . The p-values of the coefficients of both variables are significant at the 1% level for all three regressions. The negative relation between abnormal search frequency and the number of analyst recommendations implies that investors do not suddenly start Googling a stock that was recommended by an analyst, which could be because the analyst already provides the information needed for the investors to take a decision. It also implies that investors do not necessarily base their recommendations on stocks that investors have Googled a lot recently, but more likely on ‘hidden gems’ they found. The positive relation between abnormal search frequency and number of the news mentions is more obvious. An increase in news mentions could result in investors Googling the stock to find out more, or both could be related to a certain event which caused the stock to get more attention. Furthermore, the R<sup>2</sup> of these regressions is between 1.66% and 2.05%, suggesting that the number of analyst recommendations and the number of news mentions only explain a small amount of the variation in the abnormal search frequency. It is possible that some variation could be driven by other proxies for attention or by noise.

### 3.2 Impact of search frequency on abnormal returns

According to the attention hypothesis, individual investors are net buyers of stocks that have gained recent attention. This would cause a price increase investor attention increases followed by a price decrease when investor attention eventually drops to normal. To test this hypothesis with Google Trends search data, *ASF* will be regressed on abnormal returns. Both

$\log(\text{Analysts} + 1)$  and  $\log(\text{News} + 1)$  will be used as control variables to ensure the coefficient of  $ASF$  solely describes the effect of an increase in search frequency. Because  $\log(\text{Analysts} + 1)$ ,  $\log(\text{News} + 1)$  and  $ASF$  are distributed monthly, the abnormal returns will be cumulated to get monthly cumulated abnormal returns ( $CAR$ ). Firstly, the measures of attention will be regressed on the following month  $CAR$  to find out whether an increase in investor attention causes a price pressure. Secondly, the same regression will be done on the month after the following month  $CAR$  to see if the price pressure is consistent or there is any kind of reversal. Finally, the regression will be done on the  $CAR$  of the ten months after the following two months to test if the price pressure due to investor attention causes an eventual price decrease. With similar reasoning as for the regressions on  $ASF$ , a linear model with fixed effects is used to represent the measurement-specific means. Standard errors are robust and clustered by firm. The results are reported below in Table 4.

**Table 4: Regressions on  $CAR$**

The dependent variable in each regression is  $CAR$ , the independent variables are  $ASF$ ,  $\log(\text{Analysts} + 1)$  and  $\log(\text{News} + 1)$ . The  $CAR$  is calculated for the following month, following 2<sup>nd</sup> month and the following 3<sup>rd</sup> to 12<sup>th</sup> month. Each regression is a linear model containing fixed effects and the robust standard errors are clustered by firm.  $P$ -values are presented in parentheses under the coefficient and \*, \*\*, and \*\*\* represent significance at the 10%, 5%, and 1%, respectively.

	Following month $CAR$	Following 2 <sup>nd</sup> month $CAR$	Following 3 <sup>rd</sup> to 12 <sup>th</sup> month $CAR$
Constant	-0.0031 (0.900)	0.0175 (0.345)	-0.1123 (0.733)
$ASF$ ( $\Delta t = 2$ )	0.0810 (0.014)**	-0.0957 (0.117)	-0.3920 (0.343)
$\log(\text{Analysts} + 1)$	-0.0163 (0.323)	-0.0212 (0.148)	-0.2483 (0.374)
$\log(\text{News} + 1)$	0.0169 (0.211)	0.0078 (0.274)	0.3043 (0.338)
Observations	1260	1239	1029
Clusters (firms)	21	21	21
$R^2$	0.0009	0.0000	0.0068

As can be seen in Table 4, the only significant variable in the set of regressions is the  $ASF$  on the following month  $CAR$ . This is consistent with the attention hypothesis that an increase in investor attention is often followed by a price increase. Also, the attention hypothesis states that especially an increase in individual investor attention causes price pressure. Engelberg & Gao (2011) found empirical proof that Google Trends search data captures mainly individual investor attention. The fact that Google is used by almost all regular people whereas institutional investors often use

more complex databases, makes their findings quite logical. Furthermore, price pressure due to increased investor attention is known to be larger for smaller companies. Therefore, to test the robustness of the result in Table 4, the same regression will be done for the smallest 11 firms, the middle 11 firms and the largest 10 firms of the dataset based on market capitalization. The results are presented in Table 5.

**Table 5: Robustness check for regressions on CAR**

The dependent variable in each regression is CAR of the following month, the independent variables are ASF,  $\log(\text{Analysts} + 1)$  and  $\log(\text{News} + 1)$ . The regression is done for the smallest 11 firms, middle 11 firms and largest 10 firms based on market capitalization. Note that there is an overlap between small and middle and middle and large, as there are only 21 firms in the dataset. Each regression is a linear model containing fixed effects and the robust standard errors are clustered by firm. P-values are presented in parentheses under the coefficient and \*, \*\*, and \*\*\* represent significance at the 10%, 5%, and 1%, respectively.

	Smallest 11 firms	Middle 11 firms	Largest 10 firms
Constant	0.1145 (0.519)	0.0498 (0.005)***	0.001 (0.983)
ASF ( $\Delta t = 2$ )	0.1071 (0.024)**	0.0689 (0.207)	0.0276 (0.327)
$\log(\text{Analysts} + 1)$	-0.1164 (0.482)	-0.0545 (0.059)*	-0.0107 (0.449)
$\log(\text{News} + 1)$	0.0265 (0.276)	0.0282 (0.259)	0.0059 (0.452)
Observations	660	660	600
Clusters (firms)	11	11	10
R <sup>2</sup>	0.0017	0.0041	0.0061

The coefficient of ASF turns out to be significant at the 5% level for the regression on the CAR of the following month for the smallest 11 firms and not significant for the middle 11 firms or largest 10 firms. This result is not surprising; the smaller firms are less well known and when they experience an increase in attention, more people will know about them and consequently more people will consider buying them. This is less the case with larger firms because most investors already know about them and therefore might already have considered buying them. The set of regressions from Table 5 was also done for the following 2<sup>nd</sup> month and the following 3<sup>rd</sup> to 12<sup>th</sup> month CAR as described in Table 4. As none of these coefficients were significant at the 5% level, they will not be reported. No evidence has been found that the price pressure from an increase in investor attention is followed by an eventual price decrease. This could be because the dataset is too small, or the firms analysed are too well known to experience price pressure from investor attention. Nonetheless, there is evidence that individual investor attention has a

positive effect on the following month cumulative abnormal returns, especially for firms with a relatively small market capitalization.

### 3.3 Using search terms to predict investor sentiment

As mentioned in the introduction, Preis et al. (2013) found certain search terms to have predictive value for stock prices. Especially the investment strategies they implemented for the search terms ‘debt’, ‘inflation’, ‘share’, ‘economy’ and ‘unemployment’ generated significant hypothetical returns from January 2004 until February 2011 for the DJIA. Their results were quite remarkable and therefore this paper will test the robustness of their results for different time periods in the Netherlands. The time periods used are January 2004 until December 2008, January 2009 until December 2014 and January 2015 until June 2018, the index that will be used to calculate returns is the AEX.

The investment strategy used to generate the hypothetical returns involves short selling when the search terms were Googled more than normal and buying when the search terms were Googled less than normal. Specifically, when the preceding week  $ASF$  is positive, the return for the current week is calculated for a short position and when the preceding week  $ASF$  is negative, the return is calculated for a long position. This strategy is used for the three previously mentioned time periods, the five previously mentioned search terms and for  $\Delta t$  equal to 1 until 12. An overview of the exact calculations is given below.

$$ASF_t = \log(SF_t) - \log\{Med(SF_{t-1}, \dots, SF_{t-\Delta t})\} > 0 \rightarrow \frac{p_t - p_{t+1}}{p_{t+1}}$$

$$ASF_t = \log(SF_t) - \log\{Med(SF_{t-1}, \dots, SF_{t-\Delta t})\} \leq 0 \rightarrow \frac{p_{t+1} - p_t}{p_{t+1}}$$

Furthermore, for each period the total return of a ‘buy and hold’ position in the AEX is taken as a comparison. Also, a random strategy is created by where for each week there is a 50% chance of going short and a 50% chance of going long. For each period, the returns are calculated for this random strategy 10,000 times. The returns are then averaged out and the standard deviation is taken. With this data, a two-sided t-statistic is calculated for the 180 hypothetical returns generated earlier: 3 different time periods, 5 different search terms and  $\Delta t$  from 1 until 12. The results for the search term ‘debt’ will be presented in Table 6 and the results for the rest of the search terms will be presented in Appendix B.

**Table 6: Investment strategy based on ‘debt’ search frequency**

Cumulated returns from investment the investment strategy based on the search term ‘debt’ are presented below as the amount the starting portfolio is multiplied (a value of 1 means the portfolio has not increased nor decreased during the period). The strategy is used for the periods January 2004 until December 2008, January 2009 until December 2014 and January 2015 until June 2018 and the index used is AEX. The return is calculated for a short position when the previous week ASF is positive and for a long position when the previous week ASF is negative. ASF is calculated for 12 different values of  $\Delta t$ : 1 to 12. The weekly returns are cumulated over the whole period. A ‘buy and hold’ position in the AEX is presented as a comparison. Also, a random strategy is created by where for each week there is a 50% chance of going short and a 50% chance of going long. For each period, the returns are calculated for this random strategy 10,000 times. Next, the returns are averaged out and the standard deviation is taken. With this data, a two-sided t-statistic is calculated for the 36 hypothetical returns and presented in parentheses. \*, \*\* and \*\*\* represent significance at the 10%, 5% and 1%, respectively.

	2004-2008	2009-2013	2014-2018
$\Delta t = 1$	2.3358 (2.570)***	1.2365 (0.499)	0.4350 (-1.667)*
$\Delta t = 2$	1.8569 (1.653)**	0.854 (-0.300)	0.8607 (-0.414)
$\Delta t = 3$	2.1406 (2.196)**	0.9372 (-0.127)	0.8266 (-0.514)
$\Delta t = 4$	2.3668 (2.630)***	1.0108 (0.027)	0.9316 (-0.206)
$\Delta t = 5$	2.1238 (2.164)**	0.3518* (-1.350)	0.9057 (-0.282)
$\Delta t = 6$	1.9143 (1.763)**	0.4536 (-1.137)	0.7035 (-0.877)
$\Delta t = 7$	1.6295 (1.217)	0.4090 (-1.230)	0.7241 (-0.816)
$\Delta t = 8$	1.5144 (0.997)	0.5318 (-0.974)	0.7641 (-0.698)
$\Delta t = 9$	1.5173 (1.002)	0.5797 (-0.873)	0.7719 (-0.676)
$\Delta t = 10$	1.5384 (1.043)	0.6833 (-0.657)	0.8061 (-0.575)
$\Delta t = 11$	1.5975 (1.156)	0.4920 (-1.057)	0.7485 (-0.744)
$\Delta t = 12$	1.7976 (1.539)*	0.5439 (-0.948)	0.7777 (-0.658)
Buy and hold	0.7517	1.5101	1.3657
Random	0.9938	0.9978	1.0015
Standard deviation	0.5221	0.4786	0.3399

The investment strategy based on search frequency of the term ‘debt’ produces significant returns for the period 2004-2008. This implies that an increase in Google searches on ‘debt’ is often followed by a decrease in the price of the AEX and a decrease in Google searches on ‘debt’ is often followed by an increase in the price of the AEX, which is consistent with results found by Preis et al. (2013). However, the significance is less apparent when observing the periods 2009-

2013 and 2014-2018. From the 24 hypothetical strategies for these periods, only 2 are significant at the 10% level which is probably a coincidence. Nonetheless, the question remains how a hypothetical investment strategy based simply on Google search frequency of the word 'debt', could multiply the starting portfolio by 2.37 (using  $\Delta t = 4$ ) in a period where both the 'buy and hold' and the random strategy would have generated a loss.

The most likely explanation is that there is a correlation between the search frequency of 'debt' and the economic crisis of 2008. As the crisis became more apparent, people got more interested in it and started to search for information about it. The crisis also had the effect that almost all stock prices and index prices such as the AEX experienced a loss. These two facts must have led to the hypothetical investment strategy going short during the period when the AEX price was decreasing and long when the price was recovering, which is obviously an optimal way of investing. For the strategies based on the search terms 'inflation', 'share', 'economy' and 'unemployment', only 10 out of 144 were significant and thus they provided little predictive value. It seems that search terms have little predictive value for short terms stock prices, however certain search terms including 'debt', and other potentially undiscovered terms could have predictive value for macroeconomic events. Further research could be done to provide more evidence for this statement.

## 4. Conclusion

Google Trends search data is a viable proxy for attention when used correctly. It is important to use search terms that capture investor attention as exclusively as possible. Search terms that could be used to search for the product of the company or something completely different should be avoided. A negative relation was found between the abnormal search frequency and the number of analyst recommendations which could imply these proxies measure different types of attention. The relation found between the abnormal search frequency and the number of news mentions was positive and significant as expected.

An increase in attention measured by Google search frequency has a positive effect on the following month cumulated abnormal return. This is especially the case for smaller firms, which is consistent with the attention hypothesis. The attention hypothesis also states that a price reversal should occur when investor attention returns back to normal, however, no evidence was found for this. More research could be done to investigate this in different countries. Google Trends search data is a good measurement to use for the attention hypothesis due to its timely fashion and high representativeness of individual investors.

The 326% hypothetical return generated by Preis et al. (2013) were not reoccurring for later time periods in The Netherlands. The search term 'debt' must have had a correlation with the economic crisis in 2008 which led to the investment strategy based on that search term to go short at the right moments. Although it is unlikely that Google search frequency data of economic terms has predictive value for next week stock prices, there may be other search terms that give an indication of macroeconomic events such as the term 'debt' did for the 2008 crisis. Further research could be done to explore this.

## 5. References

- Askatas, N., & Zimmermann, K. F. (2009). Google econometrics and unemployment forecasting. *Applied Economics Quarterly*, *55*(2), 107-120.
- Barber, B. M., & Odean, T. (2007). All that glitters: The effect of attention and news on the buying behavior of individual and institutional investors. *The Review of Financial Studies*, *21*(2), 785-818.
- ENGELBERG, J., & Gao, P. (2011). In search of attention. *The Journal of Finance*, *66*(5), 1461-1499.
- Preis, T., Moat, H. S., & Stanley, H. E. (2013). Quantifying trading behavior in financial markets using Google Trends. *Scientific reports*, *3*, 1684.
- Vosen, S., & Schmidt, T. (2011). Forecasting private consumption: survey-based indicators vs. Google trends. *Journal of Forecasting*, *30*(6), 565-578.



## 6. Appendices

### Appendix A: Search terms used for company names

'Aalberts Industries'

'Aegon Nederland NV'

'Ahold Delhaize'

'AkzoNobel'

'ArcelorMittal'

'ASML Holding'

'Boskalis'

'DSM'

'Gemalto'

'Heineken International'

'ING Group'

'KPN'

'Philips'

'Randstad Holding'

'RELX Group'

'Royal Dutch Shell'

'SBM Offshore'

'Unibail-Rodamco-Westfield'

'Unilever'

'Vopak'

'Wolters Kluwer'

## Appendix B: Investment strategies based on different search terms

**Table B1: Investment strategy based on 'inflation' search frequency**

Cumulated returns from investment the investment strategy based on the search term 'inflation' are presented below as the amount the starting portfolio is multiplied (a value of 1 means the portfolio has not increased nor decreased during the period). The strategy is used for the periods January 2004 until December 2008, January 2009 until December 2014 and January 2015 until June 2018 and the index used is AEX. The return is calculated for a short position when the previous week ASF is positive and for a long position when the previous week ASF is negative. ASF is calculated for 12 different values of  $\Delta t$ : 1 to 12. The weekly returns are cumulated over the whole period. A 'buy and hold' position in the AEX is presented as a comparison. Also, a random strategy is created by where for each week there is a 50% chance of going short and a 50% chance of going long. For each period, the returns are calculated for this random strategy 10,000 times. Next, the returns are averaged out and the standard deviation is taken. With this data, a two-sided *t*-statistic is calculated for the 36 hypothetical returns and presented in parentheses. \*, \*\* and \*\*\* represent significance at the 10%, 5% and 1%, respectively.

	2004-2008	2009-2013	2014-2018
$\Delta t = 1$	0.4581 (-1.026)	0.5651 (-0.904)	1.007 (0.017)
$\Delta t = 2$	0.3294 (-1.272)	1.1369 (0.291)	1.1469 (0.428)
$\Delta t = 3$	0.4124 (-1.113)	0.7743 (-0.467)	1.6058 (1.778)**
$\Delta t = 4$	0.4782 (-0.987)	0.9430 (-0.114)	1.5644 (1.656)**
$\Delta t = 5$	0.4790 (-0.986)	1.2523 (0.532)	1.0695 (0.200)
$\Delta t = 6$	0.4698 (-1.004)	1.1739 (0.368)	0.9919 (-0.028)
$\Delta t = 7$	0.4208 (-1.097)	1.3828 (0.804)	1.0965 (0.280)
$\Delta t = 8$	0.4652 (-1.013)	0.9485 (-0.103)	1.182 (0.531)
$\Delta t = 9$	0.4572 (-1.028)	0.9753 (-0.047)	0.9204 (-0.239)
$\Delta t = 10$	0.4086 (-1.121)	0.9996 (0.004)	0.9723 (-0.086)
$\Delta t = 11$	0.4835 (-0.977)	0.9674 (-0.063)	1.0716 (0.206)
$\Delta t = 12$	0.4872 (-0.97)	0.8499 (-0.309)	1.0630 (0.181)
Buy and hold	0.7517	1.5101	1.3657
Random	0.9938	0.9978	1.0015
Standard deviation	0.5221	0.4786	0.3399

**Table B2: Investment strategy based on 'share' search frequency**

Cumulated returns from investment the investment strategy based on the search term 'share' are presented below as the amount the starting portfolio is multiplied (a value of 1 means the portfolio has not increased nor decreased during the period). The strategy is used for the periods January 2004 until December 2008, January 2009 until December 2014 and January 2015 until June 2018 and the index used is AEX. The return is calculated for a short position when the previous week ASF is positive and for a long position when the previous week ASF is negative. ASF is calculated for 12 different values of  $\Delta t$ : 1 to 12. The weekly returns are cumulated over the whole period. A 'buy and hold' position in the AEX is presented as a comparison. Also, a random strategy is created by where for each week there is a 50% chance of going short and a 50% chance of going long. For each period, the returns are calculated for this random strategy 10,000 times. Next, the returns are averaged out and the standard deviation is taken. With this data, a two-sided t-statistic is calculated for the 36 hypothetical returns and presented in parentheses. \*, \*\* and \*\*\* represent significance at the 10%, 5% and 1%, respectively.

	2004-2008	2009-2013	2014-2018
$\Delta t = 1$	1.2152 (0.424)	0.9003 (-0.204)	0.6214 (-1.118)
$\Delta t = 2$	1.3605 (0.702)	0.5385 (-0.960)	0.7114 (-0.854)
$\Delta t = 3$	1.3762 (0.732)	1.0549 (0.119)	0.8946 (-0.314)
$\Delta t = 4$	1.0878 (0.18)	0.8641 (-0.279)	0.7956 (-0.606)
$\Delta t = 5$	1.1220 (0.245)	1.1259 (0.268)	0.6687 (-0.979)
$\Delta t = 6$	1.1795 (0.356)	1.5105 (1.071)	0.6564 (-1.015)
$\Delta t = 7$	1.0261 (0.062)	1.6568 (1.377)*	0.6556 (-1.018)
$\Delta t = 8$	0.9877 (-0.012)	1.7923 (1.660)**	0.6676 (-0.982)
$\Delta t = 9$	1.0946 (0.193)	1.7508 (1.573)*	0.7592 (-0.713)
$\Delta t = 10$	1.2186 (0.431)	1.9366 (1.962)**	0.5736 (-1.259)
$\Delta t = 11$	1.2637 (0.517)	1.5911 (1.24)	0.6608 (-1.002)
$\Delta t = 12$	1.2231 (0.439)*	1.6150 (1.290)*	0.7042 (-0.875)
Buy and hold	0.7517	1.5101	1.3657
Random	0.9938	0.9978	1.0015
Stantard deviation	0.5221	0.4786	0.3399

**Table B3: Investment strategy based on ‘economy’ search frequency**

Cumulated returns from investment the investment strategy based on the search term ‘economy’ are presented below as the amount the starting portfolio is multiplied (a value of 1 means the portfolio has not increased nor decreased during the period). The strategy is used for the periods January 2004 until December 2008, January 2009 until December 2014 and January 2015 until June 2018 and the index used is AEX. The return is calculated for a short position when the previous week ASF is positive and for a long position when the previous week ASF is negative. ASF is calculated for 12 different values of  $\Delta t$ : 1 to 12. The weekly returns are cumulated over the whole period. A ‘buy and hold’ position in the AEX is presented as a comparison. Also, a random strategy is created by where for each week there is a 50% chance of going short and a 50% chance of going long. For each period, the returns are calculated for this random strategy 10,000 times. Next, the returns are averaged out and the standard deviation is taken. With this data, a two-sided t-statistic is calculated for the 36 hypothetical returns and presented in parentheses. \*, \*\* and \*\*\* represent significance at the 10%, 5% and 1%, respectively.

	2004-2008	2009-2013	2014-2018
$\Delta t = 1$	0.6386 (-0.680)	1.3227 (0.679)	0.9381 (-0.187)
$\Delta t = 2$	1.2415 (0.474)	0.9742 (-0.049)	0.8976 (-0.306)
$\Delta t = 3$	0.6374 (-0.683)	1.3811 (0.801)	1.0339 (0.095)
$\Delta t = 4$	0.8589 (-0.258)	0.9597 (-0.080)	1.1889 (0.551)
$\Delta t = 5$	0.9892 (-0.009)	0.9024 (-0.199)	0.8395 (-0.477)
$\Delta t = 6$	1.1488 (0.297)	0.7904 (-0.433)	0.7569 (-0.720)
$\Delta t = 7$	0.9305 (-0.121)	0.7762 (-0.463)	0.7654 (-0.695)
$\Delta t = 8$	1.3745 (0.729)	0.7810 (-0.453)	0.8669 (-0.396)
$\Delta t = 9$	1.2518 (0.494)	0.7989 (-0.416)	0.8352 (-0.489)
$\Delta t = 10$	1.2192 (0.432)	0.8137 (-0.385)	0.6910 (-0.913)
$\Delta t = 11$	1.8482 (1.636)*	0.7236 (-0.573)	0.928 (-0.216)
$\Delta t = 12$	1.8464 (1.633)*	0.9040 (-0.196)	0.9235 (-0.23)
Buy and hold	0.7517	1.5101	1.3657
Random	0.9938	0.9978	1.0015
Stantard deviation	0.5221	0.4786	0.3399

**Table B4: Investment strategy based on ‘unemployment’ search frequency**

Cumulated returns from investment the investment strategy based on the search term ‘unemployment’ are presented below as the amount the starting portfolio is multiplied (a value of 1 means the portfolio has not increased nor decreased during the period). The strategy is used for the periods January 2004 until December 2008, January 2009 until December 2014 and January 2015 until June 2018 and the index used is AEX. The return is calculated for a short position when the previous week ASF is positive and for a long position when the previous week ASF is negative. ASF is calculated for 12 different values of  $\Delta t$ : 1 to 12. The weekly returns are cumulated over the whole period. A ‘buy and hold’ position in the AEX is presented as a comparison. Also, a random strategy is created by where for each week there is a 50% chance of going short and a 50% chance of going long. For each period, the returns are calculated for this random strategy 10,000 times. Next, the returns are averaged out and the standard deviation is taken. With this data, a two-sided t-statistic is calculated for the 36 hypothetical returns and presented in parentheses. \*, \*\* and \*\*\* represent significance at the 10%, 5% and 1%, respectively.

	2004-2008	2009-2013	2014-2018
$\Delta t = 1$	0.6026 (-0.749)	0.7251 (-0.57)	0.9716 (-0.088)
$\Delta t = 2$	0.4342 (-1.072)	1.3891 (0.818)	1.0345 (0.097)
$\Delta t = 3$	0.4476 (-1.046)	1.4778 (1.003)	1.3799 (1.113)
$\Delta t = 4$	0.5028 (-0.940)	1.2114 (0.446)	1.0979 (0.284)
$\Delta t = 5$	0.7405 (-0.485)	0.9920 (-0.012)	0.8666 (-0.397)
$\Delta t = 6$	1.1340 (0.269)	0.6325 (-0.763)	0.6955 (-0.9)
$\Delta t = 7$	1.1311 (0.263)	0.5820 (-0.869)	0.6772 (-0.954)
$\Delta t = 8$	1.0376 (0.084)	0.7708 (-0.474)	0.7259 (-0.811)
$\Delta t = 9$	1.0030 (0.018)	0.8410 (-0.328)	0.6897 (-0.917)
$\Delta t = 10$	0.8174 (-0.338)	0.9021 (-0.200)	0.7908 (-0.620)
$\Delta t = 11$	0.8383 (-0.298)	0.8244 (-0.362)	0.7001 (-0.887)
$\Delta t = 12$	0.8674 (-0.242)	1.0461 (0.101)	0.7674 (-0.689)
Buy and hold	0.7517	1.5101	1.3657
Random	0.9938	0.9978	1.0015
Stantard deviation	0.5221	0.4786	0.3399