

Introduction

Investor behaviour has long played a role in asset price formation, with recent research focusing on the measurable effects of investor attention. As digital search activity becomes a critical component of information gathering, tools like Google Trends provide an opportunity to quantify public interest in real time.

This analysis will answer the question “Does Google Trends Search Interest Predict Stock Returns in the Short Term Among Fortune 25 Companies?”. Using a fixed-effects panel regression framework, I will analyse the relationship between google search interest and quarterly returns over the period between 2010 and 2025. Through a combination of search and financial data, I aim to test the value of search interest as an indicator of price movement, even among the most mature firms in the economy.

Literature Review

Investor attention has emerged as a non-traditional factor in asset-pricing, particularly the use of online data. Da, Engelberg, and Gao (2011) pioneered this approach by introducing Google Search Volume Index (SVI) as a proxy for investor attention, observing an association between search interest and short-term return predictability. Similarly, Bank, Larch, and Peter (2011) showed that abnormal search volume can forecast excess returns in European markets. These findings provide a crucial background; however, results are often sensitive to sector, time horizon, and sentiment. This analysis will build on past studies through sectoral-level analysis, using fixed-effects regression to isolate firm-specific dynamics.

Methodology

Data Selection & Collection

To analyse the correlation between google trends data and stock prices, I decided to choose the top 25 companies currently on the fortune 500. (Fortune, 2024). As these companies represent the largest proportion of market-share and encompass several sectors.

Google trends data was downloaded directly from their website, encompassing the period from 01/01/2010 to 01/01/2025. (Google Trends, 2025). Google Trends data consists of a 0–100-point scale measuring ‘search interest’.

Historical stock market data was sourced from Yahoo Finance, and was taken from the same period. (Yahoo Finance, 2025). It is important to note that not every company on the list has been trading publicly for the entire period.

I then stitched together all the data for each company into a single file using python. (See appendix)

Data Analysis

Analysis was entirely conducted using python libraries including Pandas, NumPy and StatsModels. The stock price and search interest data was loaded and reshaped into long format to facilitate time-series analysis.

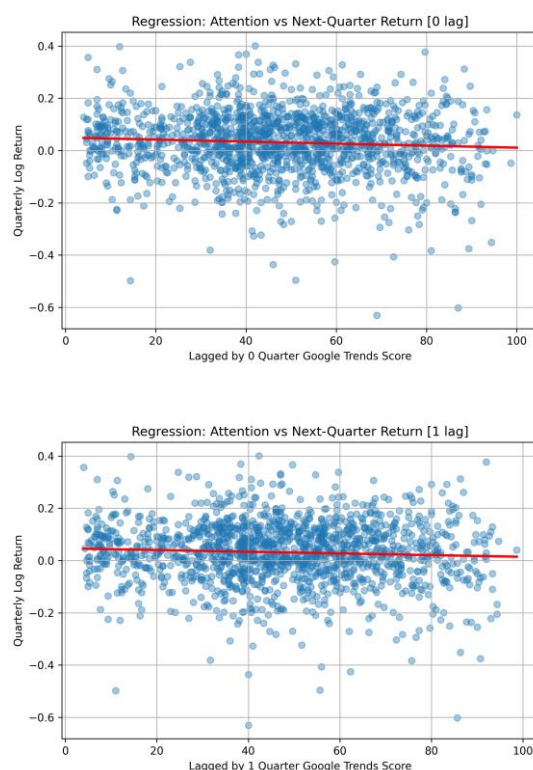
To stabilise variance, reduce skew, and compare performance across stocks with various price scales, I then calculated quarterly log returns by taking the difference in the log of the stock price quarter to quarter within groups defined by their stock ticker.

To identify the delayed impacts of search interest on stock returns, analysis was conducted using a lagged attention variable. Essentially comparing log return to search interest with various levels of offset. This was carried out for lag ranging from 0-4 quarters.

I then used the Ordinary Least Squares (OLS) regression model to incorporate stock-specific and time-specific fixed effects, with the aim of controlling for unobserved heterogeneity and temporal trends. Additionally,

clustered standard errors were implemented to correct for autocorrelation within individual stock tickers, thus enhancing the robustness of the statistical inference. (All code can be found in the appendix)

General Analysis with Varied Lag



Above are the plots for quarterly log return against lagged google trends scores for the two most statistically significant lag levels. Refer to the ‘Graphs’ section of the Appendix for all plots.

lag	coef	pval	r2
0	-0.000598	0.000064	0.2624
1	-0.000399	0.003224	0.2603
2	-0.000358	0.005834	0.2586
3	-0.000199	0.177993	0.2576
4	-0.000359	0.013713	0.2587

Above is a table showing the results of the regression test at each level of lag, there is the strongest correlation at zero lag. Across all levels of lag, it is observable that there is some negative impact of search interest on stock prices.

Observations

From the data, it is clear that there is a general correlation between search interest and stock price changes. The most significant windows are immediate and within 1 quarter. Both of which have a p value below 0.01, indicating

statistical significance at a 1% level. Their R^2 values are similar at roughly 0.26, indicating a moderate impact, as roughly 26% of variance in stock price is explainable by search interest.

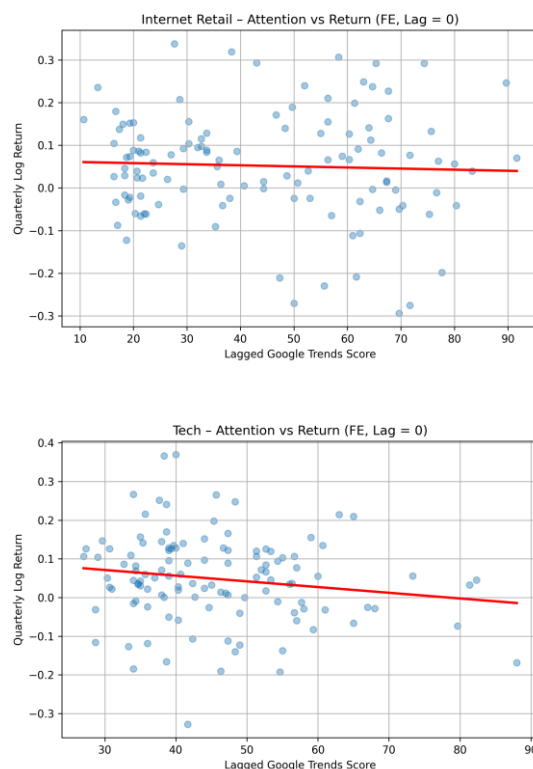
It is interesting that their coefficients are both negative, indicating a negative correlation between search interest and stock price. This may be explained by the fact that search interest is increased by negative coverage of companies more than it is by positive coverage. This is a valuable insight, it shows that consumer/investor interest in a company is not necessarily positive interest.

Methodology Continued

After exploratory analysis identified that a lag of 0 provided the strongest and most statistically significant predictive power, I selected this lag for a detailed sector-by-sector evaluation. The goal of a sector specific evaluation was to enhance the specificity and applicability of insights. Stocks were mapped to industry sectors and sectors with fewer than 30 observations or fewer than three unique tickers (companies) were excluded to ensure sufficient statistical validity.

Similar OLS regression modelling to earlier was used, maintaining consistent controls and clustering strategies. The results of sector-by-sector analysis are displayed in the following section. (Code found in appendix)

Sector Analysis with Zero Lag

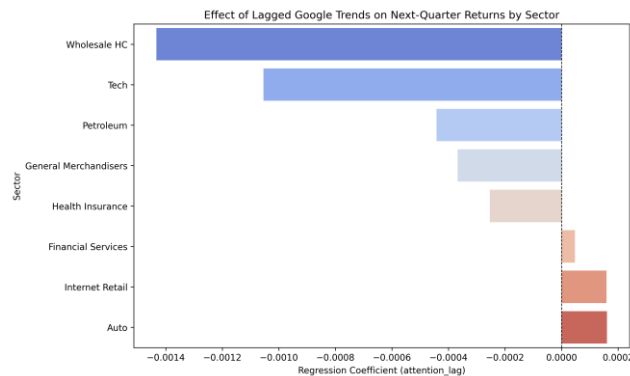


Above are the plots for quarterly log return against zero-lag google trends scores for the two most statistically significant sectors, Technology, and Internet Retail. Plots for other sectors may be found in the appendix.

sector	coef	pval	r2
Tech	-0.00105	~0.000000	0.7567538980017445
Internet Retail	0.00016	~0.000000	0.7749495267751321
Petroleum	-0.000440	0.001609	0.7525893660281782
Health Insurance	-0.000250	0.390977	0.6646704565440116

Financial Services	0.000047	0.521184	0.7523698148483239
General Merchandisers	-0.000370	0.728089	0.3980707428062471

Above is a table with the results of the regression calculations for each sector. Technology and Internet Retail are the sectors with the most significant correlation between trends activity and stock price. While there are other sectors with larger coefficients, Internet Retail has a low p-value and according to its R^2 value, the largest proportion of change in stock price explainable by trends data.



Above is a bar chart showing the coefficients for several sectors, visualising the large difference in search interest impact across sectors. Note that wholesale HC and Auto are included on the bar chart, but were not included in the table of summary statistics as they did not produce a p-value.

Observations

From the data presented, there is visibly a large disparity in impact across sectors. Technology and Internet Retail seem to be the most impacted by search interest in a zero-lag timeframe. Internet Retail is the only sector with a statistically significant positive correlation between search interest and stock price, this may be due to the fact that search interest in an online retailer has a higher chance of being genuine interest in consumption (indicated by its R^2 of 0.775, the highest); search interest in other industries is more likely attributed to news or other factors.

Petroleum is also worth discussing, the sector showed a strong R^2 value and a sub-5% significance level, so while its coefficient was not as high as Technology or Wholesale HC, its variance can be strongly attributed to search interest. Its negative regression coefficient strongly suggests that search interest is a negative for petroleum companies, likely due to climate scrutiny and other scandals.

General Merchandisers and Financial Services each have very high p-values of 0.73 and 0.52 respectively, this indicates that the variable the regression was performed on (search interest) was unlikely to be the cause behind any variance. Observations were likely due to chance or confounding factors not considered in this analysis.

General vs Sector Specific

When compared in the same time-frame of 0 lag, there are clear differences between our general analysis and sector-specific analysis.

Primarily, the disparity in R^2 values. General analysis returned a value of ~0.26 in this time frame, a moderate proportion of variance attributable to search interest. But each sector analysed returned values much higher. This can be attributed to the homogenous nature of each sector, patterns can be more consistently explained by search

interest; In general analysis, it is less easy to explain variation through search interest due to the diverse nature of the dataset.

Additionally, there is a disparity in p values. Apart from Technology and Internet Retail, every other sector analysed has a much higher p value than that of the general analysis. The disparity between the other sectors and the general analysis can be explained by the fact that the dataset in the general analysis is much larger. Technology and Internet Retail have a very high correlation between search interest and stock price changes, it is large, stable, and significant. This means they have extremely low p values, and are likely the main drivers of the low general p value.

There is no clear, large difference between the general regression coefficient and the sector coefficients as a whole; they likely average out to form the general regression coefficient. However, it is notable that Internet Retail is distinct in its positive coefficient (apart from Financial Services, however that sector has a high p value and is thus not statistically significant). This indicates that while in general search interest predicts a decrease in stock price, Internet Retail is unique as a sector. Likely due to Internet Retail being retail on the internet, any consumption is associated with a google search. This naturally brings us onto the limitations of our model.

Conclusion

In this analysis, I have explored the predictive relationship between search interest (through Google Trends data), and short-term stock returns among the Fortune 25. Using a fixed-effects panel regression framework, it was found that search interest is a statistically significant predictor of returns in general, with stronger effects in specific sectors such as Technology and Internet Retail. It is notable that in general analysis, I found that the effect is strongest within one quarter, and that that effect was a negative return.

However, the direction and strength of the correlation varied by sector, and most sectors did not show significant correlation in isolation. This may be due to the limited number of companies in the dataset. But it could also be due to confounding variables, an endogenous relationship between stock price and search interest, or a lack of sentiment analysis. This approach was limited, but still yielded statistically significant insights into the link between search interest and stock returns.

APPENDIX: Code

CSV Stitcher for Trends, Similar used for Stock	Code for General Regression & Plotting
---	--

```

1 # === Google Trends Data Aggregation Script ===
2 # This script loads monthly Google Trends CSVs, extracts data from the correct header row,
3 # cleans and aligns them, averages values by quarter, and combines them into a single DataFrame.
4
5 import pandas as pd
6 import glob
7 import os
8
9 csv_dir = r"C:\Users\#####\Documents\Q1 Project stuff\trends CSVs"
10 csv_files = glob.glob(os.path.join(csv_dir, "*.csv"))
11
12 trend_data = {}
13
14 # Loop through each file in the folder
15 for file in csv_files:
16     # Use filename (minus .csv) as the column name
17     trend_name = os.path.splitext(os.path.basename(file))[0]
18
19     # Identify the line that contains the actual header
20     header_line_index = None
21     with open(file, 'r', encoding='utf-8') as f:
22         for i, line in enumerate(f):
23             if "Month" in line:
24                 header_line_index = i
25             break
26     if header_line_index is None:
27         continue # skip file if no "Month" found
28
29     # Read the CSV, starting from the correct header line
30     df = pd.read_csv(file, sep=";", skiprows=header_line_index, engine="python")
31     df.columns = df.columns.str.strip() # clean column names
32
33     if "Month" not in df.columns:
34         continue # skip if structure is unexpected
35
36     # Rename for consistency
37     df.rename(columns={"Month": "Date", inplace=True)
38     df["Date"] = pd.to_datetime(df["Date"], format="%Y-%m", errors="coerce")
39
40     # Identify the search interest column (assumes only one besides Date)
41     trend_col = [col for col in df.columns if col != "Date"]
42     if not trend_col:
43         continue
44
45     trend_col = trend_col[0]
46     df[trend_col] = pd.to_numeric(df[trend_col], errors="coerce")
47     df.set_index("Date", inplace=True)
48
49     # Resample monthly data into quarterly averages
50     quarterly_avg = df.resample("Q").mean()
51     quarterly_avg.rename(columns={trend_col: trend_name, inplace=True)
52
53     # Store cleaned and averaged data
54     trend_data[trend_name] = quarterly_avg[trend_name]
55
56 # Combine all trends into a single DataFrame, indexed by quarter-end
57 combined_df = pd.DataFrame(trend_data)
58 combined_df.index.name = "Quarter_End"
59
60 # Export to Excel for use in regression analysis
61 combined_df.to_excel("combined_trends_quarterly_avg.xlsx")

```

Sector-By-Sector Regression and Plotting 1

```

1 # === Sector-By-Sector Analysis of Attention and Returns ===
2 # This script runs a fixed-effects regression for each sector to estimate whether
3 # lagged Google Trends attention predicts next-quarter stock returns.
4 # It also outputs sector-level plots and a summary bar chart.
5
6 import pandas as pd
7 import numpy as np
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10 import statsmodels.formula.api as smf
11
12 # =====
13 # CONFIGURATION
14 # =====
15 lag = 0 # Number of quarters to lag attention scores
16 min_points = 30 # Minimum data points required to include a sector
17 dpi = 300 # Resolution for saved plots
18 output_dir = "sector_plots" # Folder to save individual sector plots
19 # =====
20
21 # 1. Load pre-processed quarterly attention and price data
22 trends = pd.read_excel("combined_trends.xlsx")
23 prices = pd.read_excel("combined_stocks.xlsx")
24
25 # 2. Reshape to long format for merging
26 trends_long = trends.melt(id_vars="Quarter_End", var_name="ticker", value_name="attention")
27 prices_long = prices.melt(id_vars="Quarter_End", var_name="ticker", value_name="price")
28
29 # 3. Merge and sort chronologically
30 df = pd.merge(trends_long, prices_long, on=["Quarter_End", "ticker"], how="inner")
31 df["Quarter_End"] = pd.to_datetime(df["Quarter_End"])
32 df = df.sort_values(["ticker", "Quarter_End"])
33
34 # 4. Compute log returns and lagged attention values
35 df["log_return"] = df.groupby("ticker")["price"].transform(lambda x: np.log(x) - np.log(x.shift(1)))
36 df["attention_lag"] = df.groupby("ticker")["attention"].shift(lag)
37
38 # 5. Drop rows with missing values
39 df_clean = df.dropna(subset=["log_return", "attention_lag"]).copy()
40
41 # 6. Map company tickers to defined sectors
42 sector_map = {
43     "WMT": "General Merchandisers",
44     "COST": "General Merchandisers",
45     "HD": "General Merchandisers",
46     "KRC": "General Merchandisers",
47     "WDC": "Internet Retail",
48     "GOODO": "Internet Retail",
49     "AMPL": "Tech",
50     "AAPL": "Tech",
51     "MSFT": "Tech",
52     "WU": "Health Insurance",
53     "CVS": "Health Insurance",
54     "CNC": "Health Insurance",
55     "TLP": "Health Insurance",
56     "BAC": "Financial Services",
57     "WFC": "Financial Services",
58     "JP": "Financial Services",
59     "AXP": "Financial Services",
60     "C": "Petroleum",
61     "XOM": "Petroleum",
62     "CVX": "Petroleum",
63     "PFE": "Pharmaceuticals",
64     "MRK": "Pharmaceuticals",
65     "NVS": "Pharmaceuticals",
66     "T": "Auto",
67     "GM": "Auto",
68 }
69
70 df_clean["sector"] = df_clean["ticker"].map(sector_map).fillna("Other")
71 df_clean["quarter_start"] = df_clean["Quarter_End"].dt.to_period("Q").start_time

```

```

1 # === General Analysis: Google Trends vs Stock Returns ===
2 # This script runs a fixed-effects regression to test whether lagged Google Trends
3 # scores predict next-quarter log returns across all Fortune 25 companies.
4
5 import pandas as pd
6 import numpy as np
7 import statsmodels.formula.api as smf
8 import seaborn as sns
9 import matplotlib.pyplot as plt
10
11 # Set lag length to use in the regression
12 lag = 4 # Number of quarters to lag attention data
13
14 # Load combined quarterly data
15 trends = pd.read_excel("combined_trends.xlsx")
16 prices = pd.read_excel("combined_stocks.xlsx")
17
18 # Convert from wide to long format for merging
19 trends_long = trends.melt(id_vars="Quarter_End", var_name="ticker", value_name="attention")
20 prices_long = prices.melt(id_vars="Quarter_End", var_name="ticker", value_name="price")
21
22 # Merge Trends and Price data
23 df = pd.merge(trends_long, prices_long, on=["Quarter_End", "ticker"], how="inner")
24 df["Quarter_End"] = pd.to_datetime(df["Quarter_End"])
25 df = df.sort_values(["ticker", "Quarter_End"])
26
27 # Compute log returns and lagged attention scores
28 df["log_return"] = df.groupby("ticker")["price"].transform(lambda x: np.log(x) - np.log(x.shift(1)))
29 df["attention_lag"] = df.groupby("ticker")["attention"].shift(lag)
30
31 # Drop rows with missing values (usually early quarters)
32 df_clean = df.dropna(subset=["log_return", "attention_lag"])
33
34 # Create quarter identifiers for fixed effects
35 df_clean["quarter_str"] = df_clean["Quarter_End"].dt.to_period("Q").astype(str)
36
37 # Run fixed-effects regression with firm and time dummies, clustered by ticker
38 model = smf.ols("log_return ~ attention_lag + C(ticker) + C(quarter_str)", data=df_clean)
39 results = model.fit(cov_type="cluster", cov_kwds={"groups": df_clean["ticker"]})
40
41 # Extract key statistics
42 coef = results.params.get("attention_lag", np.nan)
43 stderr = results.bse.get("attention_lag", np.nan)
44 pval = results.pvalues.get("attention_lag", np.nan)
45 r2 = results.rsquared
46
47 # Format and save regression summary
48 summary_text = (
49     f"===== Regression Summary (Lag = {lag}) =====\n"
50     f"Coefficient for attention_lag: {coef:.6f}\n"
51     f"Standard Error: {stderr:.6f}\n"
52     f"P-value: {pval:.6f}\n"
53     f"R-squared: {r2:.4f}\n"
54     "===== "
55 )
56 with open(f"regression_summary_lag_{lag}.txt", "w") as f:
57     f.write(summary_text)
58
59 # Plot raw scatterplot with regression line (note: this does not reflect fixed effects)
60 plt.figure(figsize=(8, 5))
61 sns.regplot(
62     x="attention_lag",
63     y="log_return",
64     data=df_clean,
65     scatter_kws={"alpha": 0.4},
66     line_kws={"color": "red"},
67     ci=None
68 )
69 plt.xlabel(f"Lagged by {lag} Quarter Google Trends Score")
70 plt.ylabel("Quarterly Log Return")
71 plt.title(f"Raw Scatter: Attention vs Return (Lag = {lag})")
72 plt.grid(True)
73 plt.savefig(f"regression_plot_lag_{lag}.png", dpi=300)
74 plt.show()

```

Sector-By-Sector Regression and Plotting 2

```

1 # 7. Iterate output directory
2 os.makedirs(output_dir, exist_ok=True)
3
4 # 8. Run regressions and generate sector plots
5 sector_results = []
6
7 for sector in df_clean["sector"].unique():
8     sector_data = df_clean[df_clean["sector"] == sector]
9
10     if len(sector_data) < min_points:
11         print(f"Skipping sector (only {len(sector_data)} rows)")
12         continue
13
14     unique_tickers = sector_data["ticker"].unique()
15     if len(unique_tickers) < 2:
16         print(f"Skipping sector (only {len(unique_tickers)} unique tickers)")
17         continue
18
19     # Fit fixed-effects regression with firm and time dummies
20     model = smf.ols("log_return ~ attention_lag + C(ticker) + C(quarter_str)", data=sector_data)
21     results = model.fit(cov_type="cluster", cov_kwds={"groups": sector_data["ticker"]})
22
23     coef = results.params.get("attention_lag", float("nan"))
24     pval = results.pvalues.get("attention_lag", float("nan"))
25     r2 = results.rsquared
26
27     print(f"Sector {sector} | Coef: {coef:.6f} | SE: {stderr:.6f} | P: {pval:.6f} | R2: {r2:.4f}")
28     sector_results.append({"sector": sector, "coef": coef, "pval": pval, "r2": r2})
29
30 # Plot scatter with OLS trendline (not fixed effects)
31 plt.figure(figsize=(8, 5))
32 sns.regplot(
33     x="attention_lag",
34     y="log_return",
35     data=sector_data,
36     scatter_kws={"alpha": 0.4},
37     line_kws={"color": "red"},
38     ci=None
39 )
40
41 plt.title(f"Sector {sector} | Raw Scatter: Attention vs Return (Lag = {lag})")
42 plt.xlabel(f"Lagged Google Trends Score")
43 plt.ylabel("Quarterly Log Return")
44 plt.grid(True)
45
46 save_name = sector.replace(" ", "_").replace(".", "")
47 plt.savefig(f"{output_dir}/{save_name}_regression.png", dpi=dpi)
48 plt.close()
49
50 # 9. Save summary of sector-level results
51 summary_df = pd.DataFrame(sector_results)
52 summary_df = summary_df.sort_values("pval")
53 summary_df.to_csv("sector_regression_summary.csv", index=False)
54 print(f"Saved regression summary to 'sector_regression_summary.csv'")
55
56 # 10. Generate summary bar chart of coefficients by sector
57 summary_df = summary_df.sort_values("coef")
58
59 plt.figure(figsize=(10, 6))
60 sns.barplot(
61     data=summary_df,
62     x="sector",
63     y="coef",
64     palette="coolwarm",
65     order=s
66 )
67
68 plt.grid(True)
69 plt.title("Effect of Lagged Google Trends on Next-Quarter Returns by Sector")
70 plt.xlabel("Regression Coefficient (attention_lag)")
71 plt.ylabel("Sector")
72 plt.tight_layout()
73 plt.savefig("sector_regression_bar_chart.png", dpi=300)
74 plt.show()

```

Bibliography

- Bank, M., Larch, M. and Peter, G. (2011) ‘Google Search Volume and Its Influence on Liquidity and Returns of German Stocks’, *Financial Markets and Portfolio Management*, 25(3), pp. 239–264.
- Da, Z., Engelberg, J. and Gao, P. (2011) ‘In Search of Attention’, *The Journal of Finance*, 66(5), pp. 1461–1499.
- Fortune (2024) *Fortune 500*. Available at: <https://fortune.com/ranking/fortune500/> (Accessed: 1 May 2025).
- Google Trends (2025) *Google Trends*. Available at: <https://trends.google.com/> (Accessed: 3 May 2025).
- Yahoo Finance (2025) *Historical stock data*. Available at: <https://finance.yahoo.com/> (Accessed: 29 April 2025).