

THE USE AND USEFULNESS OF BIG DATA IN FINANCE: EVIDENCE FROM FINANCIAL ANALYSTS

Feng Chi, Byoung-Hyoun Hwang, and Yaping Zheng*

This Draft: March 2024

Alternative data have become a new source of information for investors. This paper examines how this development is reshaping the role of sell-side analysts. Employing textual analysis of analysts' written reports, we show that analysts themselves have begun to adopt alternative data in their analyses. Furthermore, we find that when analysts report having drawn from alternative data, they generate more accurate earnings forecasts and their brokerages subsequently receive higher amounts in trading commissions from investors, suggesting that investors value the adoption of alternative data by analysts.

JEL Classification: G12, G14, G17, G24.

Keywords: Big Data, Alternative Data, Sell-Side Analysts, Investors.

* Chi is affiliated with the Cornell SC Johnson College of Business, Cornell University. Hwang is affiliated with the Nanyang Business School, NTU, Singapore. Zheng is affiliated with the Alberta School of Business, University of Alberta. E-mail: fc354@cornell.edu, bh.hwang@ntu.edu.sg, and yaping2@ualberta.ca. The authors would like to thank the Editor, the Associate Editor, the three reviewers, Andy Call, Clifton Green, Jillian Grennan (discussant), Allen Huang (discussant), Jiekun Huang, Russell Jame, Andrew Karolyi, Ningzhong Li, Ross Lu, Stan Markov (discussant), Joe Pacelli, Marcus Painter (discussant), Hongping Tan, Ye Wang (discussant), Scott Yonker, Chi Zhang (discussant), Frank Zhang, Jingjing Zhang and seminar participants at Korea University, the 2021 EFA Annual Meeting, the 2021 MIT Asia Conference in Accounting, the 2021 China International Conference in Finance, the 2021 CAPANA Conference, the 3rd Future of Financial Information Conference, and the 2021 Conference on Financial Innovation for many valuable suggestions and comments. This research is supported by NTU-Singapore under its start-up grant program.

1. Introduction

Sell-side analysts have long played a pivotal role in financial markets (Bradley, 2023). They gather information on publicly traded companies and share their analyses and opinions with investors. To the degree that analysts pass on unique and value-relevant insights, they can contribute significantly to enhancing market efficiency.

Like many roles in knowledge-driven sectors, the analyst profession is being challenged by technological advances: With the advent of modern information technologies and advances in data analytics, we can increasingly track individuals' and businesses' activities through the digital footprints they leave behind. In the realm of investments, an expanding body of research finds that these digital footprints, commonly referred to as "big data" or "alternative data," are useful in predicting company performance.¹

Alternative data—with its capacity to generate insights in a highly timely, comprehensive, and accurate manner—holds the potential to greatly diminish the usefulness of analysts to the investor community, eventually rendering the analyst profession obsolete.

Yet, there exists a contrasting perspective. Suppose alternative data can provide unique and value-relevant insights and that investors yearn to incorporate these insights into their decision-making. One challenge that may prevent investors from doing so is the substantial cost associated with subscribing to alternative data. Relatedly, extracting meaningful insights from alternative data requires skill and experience. Instead of each individual investor's bearing the cost of subscribing to and learning about alternative data, it is more economical for a few analysts to incur these expenses and subsequently share their insights with investors. Analysts—who have the resources, readiness, and skill to study alternative data—could, therefore, utilize alternative data to maintain (or increase) their usefulness to investors rather than being made obsolete.

Our paper examines this assertion. Our empirical analysis builds on the following idea: Analysts routinely publish written reports detailing their perspectives, along with the data and analyses

¹ Examples include Froot, Kang, Ozik, and Sadka (2017), who consider consumer activity estimated in real-time from mobile phone data; Huang (2018) and Zhu (2019), who consider online consumer activity; and Zhu (2019), Kang, Stice-Lawrence and Wong (2021), Katona, Painter, Patatoukas and Zeng (2023), and Gerken and Painter (2023), who consider satellite images of retail store parking lots.

underpinning these views. We posit that if analysts accessed alternative data and if the consideration of alternative data meaningfully altered their beliefs, they would discuss such influence in the corresponding reports. By parsing analysts' written reports and checking whether they explicitly reference the use of alternative data, we can thus gauge how often they draw from alternative data and in what contexts. We can also study investors' reactions to the issuance of such reports.

To provide more information on our parsing approach, we start with a comprehensive list of in-house data-science teams and external alternative-data vendors. We search for the names of these teams and vendors in analysts' written reports. We then, following prior literature (Hoberg and Moon, 2017; 2019), conduct an iterative keyword search. Among the reports that contain the name of a team or vendor, we extract a list of keywords that analysts use to describe the alternative data. We then use these keywords to search for additional reports that discuss the use of alternative data. At the end of the process, we manually read the relevant passages of each captured report to verify that the report indeed mentions the use of alternative data.

Given the labor-intensive nature of our identification process, we limit our primary analysis to companies listed in the Dow Jones Industrial Average index (DJI) from June 2009 to May 2019. The DJI includes 30 large publicly traded companies. As the composition of the DJI changes over time, our final sample comprises 35 firms. We search for analyst reports on these 35 firms within the Investext database and integrate annual earnings forecast data from IBES. Ultimately, our dataset encompasses 64,018 written reports and their corresponding annual earnings forecasts, issued by 1,002 distinct analysts employed by 55 brokerage firms.²

Our analysis reveals that, by 2009/2010, 11% of the analysts in our sample explicitly reference the use of alternative data in at least one of their reports. By 2018/2019, the corresponding fraction is 28%. Our analysis differentiates eight alternative data categories: app usage, sentiment, employee, geospatial, point of sale, satellite image, web traffic, and others. We find explicit references from all eight categories within the first year of our sample period.

² As we discuss in the main body of the paper, in additional analyses we draw a random sample of 200 companies from the lower half of the size distribution within the CRSP database. We then examine the degree to which our main findings, observed among the largest firms, extend to these smaller firms.

To gauge whether analysts are able to distill unique and value-relevant insights from alternative data, we test whether the annual earnings forecasts from reports that mention the use of alternative data are more accurate than those from reports without alternative data references. We observe within a difference-in-differences specification that annual-earnings-forecast accuracy and alternative-data adoption are positively correlated with each other. Our estimates suggest that the performance improvement accompanying the consideration of alternative data is equivalent to having covered the corresponding firm for 3.6 additional years.

The foundation of the analyst business model rests on soft dollar arrangements, in which research services are funded through trading commissions. As investors perceive greater value in an analyst's research, they increase the number of trades routed through the analyst's brokerage and, as a result, pay higher commissions to the corresponding brokerage (Bradley, 2023). The perceived value of an analyst's research and the magnitude of trading commissions are thus directly linked with each other. To ascertain if analysts' adoption of alternative data enhances their value to investors, we therefore examine whether discussions of alternative data positively correlate with the magnitude of commissions received.

To implement our test, we obtain institutional-investor trade data from ANcerno, a firm that provides transaction-cost analysis to institutional clients. The ANcerno dataset contains specific information for each institutional trade, including the commissions paid to the broker. For each written report, we consider all trades on the corresponding stock within three months after the report was issued and aggregate the total dollar value of the commissions paid to the corresponding broker. We observe within a difference-in-differences specification that institutional investors pay more in trading commissions when an analyst references alternative data. Our finding aligns with our supposition that analysts can enhance their value to investors by adopting alternative data.

Our final test examines the ramifications for investors and financial markets. We hypothesize that barriers to adopting alternative data are more pertinent for smaller or traditional institutional investors, such as mutual funds, than for hedge funds. Consistent with this conjecture, industry reports

document that hedge funds are the primary consumers of alternative data.³ By distilling insights from alternative data and sharing them with non-hedge-fund investors, analysts could play a crucial role in leveling the playing field for hedge funds and non-hedge-fund investors.

To assess this proposition, we again utilize our institutional investors' trade data. Existing research indicates that hedge funds place more profitable trades than other institutional investors (Stulz, 2007). Our paper examines whether this performance gap narrows when analysts incorporate alternative data and share their insights with investors, including non-hedge-fund investors.

To assess the performance of institutional investors, we construct transaction-based calendar time portfolios following the approach of Seasholes and Zhu (2010) and Ben-David, Birru, and Rossi (2019). As we are interested in trades affected by the dissemination of alternative-data insights, we focus on transactions occurring within three months after the issuance of an analyst report that involves the stock mentioned in the report. We employ the methodology outlined by Jame (2018) to distinguish trades by hedge funds from those by other institutional investors.

First, we find that the stocks that hedge funds buy substantially outperform those they sell over the ensuing three months; the outperformance is 10.79% when annualized. The 10.79% figure serves as our benchmark against which we compare the performance of non-hedge-fund institutional investors.

Our initial comparison considers cases where analysts do *not* reference alternative data. In these cases, we find that the stocks that non-hedge funds buy outperform those they sell by 3.85%. The results are noticeably different when analysts reference alternative data. In such cases, the outperformance is 9.85%. The similarity in outperformance for hedge funds and non-hedge funds when analysts share their alternative-data insights is consistent with the notion that analysts' use of alternative data helps level the playing field for hedge funds and non-hedge-fund institutional investors.

Our paper contributes to several research areas. First, our paper relates to the literature on sell-side analysts. Analysts historically served as vital information intermediaries in financial markets and they have been the subject of extensive research. Bradley (2023) observes that the term "analyst" and

³ E.g., <https://www.institutionalinvestor.com/article/2bsxas8ahvzgesgpnq0hs/portfolio/its-time-to-cash-in-on-big-data>

its variations have appeared in the titles and abstracts of nearly 2,500 papers in the top finance and accounting journals over the last two decades.

An emerging body of work shows that the analyst field is suffering from a rise in passive investing and regulatory changes, which have led to a decline in funding for analyst research (Bradley, 2023). The analyst profession is further threatened by technological advancements. There has been a surge in alternative-data vendors (Grennan and Michaely, 2021). These data appear to predict financial performance above and beyond what is forecasted by human analysts (e.g., Froot, Kang, Ozik, and Sadka, 2017; Huang, 2018). Relatedly, Jame, Johnston, Markov, and Wolfe (2016) show that online platforms, which crowdsource investors' earnings forecasts and make the consensus forecasts publicly available, produce predictions that are more accurate than those provided by human analysts. Coleman, Merkley, and Pacelli (2022) study the emergence of “robo-analysts,” which are “*research firms that focus on the use of technology to mass-produce recommendations with limited human involvement*” (page 12). The authors find that the recommendations of Robo-Analysts are superior to those of human analysts. Cao, Jiang, Wang, and Yang (2023) develop an “AI analyst” by leveraging a range of machine-learning toolkits. Like Coleman et al. (2022), the authors find that their AI Analyst outperforms human analysts.

Collectively, the above studies point to a paradigm shift in the way investors can obtain information, with a trend toward advanced data analysis and big data. Our study adds to the analyst literature by providing evidence that some analysts respond to this trend by adopting alternative data themselves. Our results suggest that—by doing so—these analysts have maintained (or enhanced) their relevance to investors.

In that vein, our findings shed light on the broader question of how recent technological advances impact the labor market, particularly in knowledge-intensive industries (e.g., Agarwal, Gans, and Goldfarb, 2019; Frank et al., 2019). Our observations indicate a continued need for human labor, provided that individuals adapt to the evolving technological landscape. One possible reason is that humans can use their intuition to distill insights from machine-generated outputs that, as of now, are lost to machines. As a result, there is value in combining human and machine capabilities, a proposition that aligns with the findings of Cao, Jiang, Wang, and Yang (2023). Cao et al. (2023) find that while

their AI Analyst outperforms human analysts, the most accurate predictions are theoretically achievable by integrating inputs from both AI and human analysts. Another possible reason for the sustained need for human labor is the presence of barriers, which prevent many humans from adopting the technology. As long as these barriers exist, there will be opportunities for certain humans to specialize in these technologies and share their insights with those who have not yet adopted them.

Finally, our study contributes to the body of research that examines the consequences of alternative data for financial markets. The existing literature primarily studies the predictive value of alternative data for companies' future performances.⁴ There is comparatively little research investigating the actual usage of alternative data by financial-market participants. Current studies in this area typically take the initial offering of an alternative dataset as a positive shock to alternative-data availability and compare financial-market outcome variables from before and after the shock (Zhu, 2019). Other studies estimate financial-market participants' use of alternative data for a specific stock by measuring how heavily the stock is covered in an alternative dataset (Dessaint, Foucault, and Fresard, 2023). In contrast, our research considers situations where analysts explicitly describe their use of alternative data. One limitation of our method is that it might overlook situations where analysts consider alternative data but do not mention their use in their reports. The advantage of our approach is that we can document precisely what kinds of alternative data analysts actually use, for which firms, and in what situations. The descriptive evidence we present may provide a valuable reference source for future work on this topic.

2. Analysts' Use of Alternative Data

We begin the main body of our paper with a definition of alternative data. We then describe how we capture analyst reports that reference the use of alternative data.

⁴ Additionally, some studies suggest that econometricians can use alternative data, such as satellite images of retail store parking lots, to create proxies for local store performance (Kang, Stice-Lawrence, and Wong, 2021; Gerken and Painter, 2023). Researchers can then use these proxies to study the degree to which local analysts or local investors rely on, or perhaps overemphasize, local information.

2.1 Alternative Data and Historical Perspective

Alternative data trace the footprints individuals and firms leave behind through their day-to-day activities. These footprints are also known as “exhaust data.” The use of exhaust data is not new to the financial sector. In the past, investment firms dispatched their junior analysts to retail stores to sample the foot traffic; other firms directed their analysts to manufacturing plants to count the number of trucks moving in and out (McMahon and Chu, 2012; Wigglesworth, 2016).

What distinguishes alternative data from the previous exhaust data is that, with the advent of modern information technologies and the rise in computing power, we can now source exhaust data instantaneously, comprehensively, and from a large variety of sources. That is, rather than counting foot traffic for select branches manually over several days, we can now comprehensively track how many consumers visit a merchant’s website.

There are broadly eight alternative data categories: (1) app-usage data, which track the number of active mobile app users and the amount of time they spend on the apps; (2) sentiment data, which include product ratings posted on the Internet and social-media feeds regarding a company’s products and services; (3) employee data, which include online job postings, employee opinions, and manager statements; (4) geospatial data, which contain information about the locations in which a company operates branches; (5) point-of-sale data, which include merchant-level transaction data, product-level purchase data, and pricing data; (6) satellite-image data, which include satellite images of parking lots, manufacturing plants, and construction sites; (7) web-traffic data, which track what terms users search for on the Internet and how frequently and for how long users visit a merchant’s website; and (8) other, which include data that do not fit cleanly into any of the other seven categories (e.g., data on senders and recipients of barrels loaded onto vessels).

2.2 Measuring Reliance on Alternative Data

To gauge analysts’ use of alternative data, we conduct textual analysis of their written reports. We use the Investext database to download analyst reports for constituents of the DJI from June 1, 2009, through May 31, 2019. The Investext database contains active and historical research reports from brokerages, investment banks, and independent research firms around the globe. At any point in time, constituents

of the DJI represent 30 large publicly traded firms. Because the DJI constituent list varies over time, our final sample comprises 35 firms.⁵ For each report, we extract the ticker symbol, company name, report date, analyst names, broker name, report title, and full text. The average report in our sample contains 2,152 words, which is the equivalent of roughly five pages.

We merge our Investext data with annual-earnings-forecast data from the IBES database. We merge these two datasets based on ticker symbols, company names, broker names, analyst names, and dates of forecast issuances.

Of our initial sample of 70,353 Investext reports, we successfully match 65,009 Investext reports ($65,009/70,353 = 92.4\%$). After merging these matches with financial-market data from CRSP and financial-statement data from Compustat, our final sample comprises 64,018 reports and earnings forecasts issued by 1,002 analysts from 55 brokers.

Just as some brokerages are missing from the IBES dataset (e.g., UBS), others are not included in the Investext database (e.g., Goldman Sachs). The non-comprehensive brokerage coverage raises questions regarding the generalizability of our findings. Although we cannot dismiss the potential for non-representative sampling, when we consider analyst reports and earnings forecasts on DJI constituents in the Investext and the IBES databases, respectively, we find no notable differences in underlying brokerage characteristics. The similarities in brokerage characteristics are evident in several key metrics: the average number of firms each analyst covers (10.07 vs. 8.34), the average number of forecasts issued per month (4.84 vs. 4.15), and the average forecast accuracy (-0.54 vs. -0.52). At least based on these observables, there does not seem to be systematic selection bias.

We proceed as follows: We compile a list of in-house data-science teams and external alternative-data vendors from the vendor lists in the J.P. Morgan 2019 Alternative Data Handbook and AlternativeData.org, a platform that connects users to alternative-data providers.⁶ To facilitate research

⁵ The mean and median market capitalization of firms in our sample (as of 2019) are \$250 billion and \$222 billion, respectively. To put these numbers in perspective, the 99th market capitalization percentile among firms in the CRSP/Compustat universe (also as of 2019) is \$144 billion. DJI constituents are thus substantially larger than most firms in the CRSP/Compustat universe. Online Appendix Table A1 compares the industry distribution of the firms in our sample with that of the firms in the CRSP/Compustat universe. Compared with the CRSP/Compustat universe, our sample overweights the Consumer Staples sector and underweights the Health Care sector.

⁶ In-house data-science teams specialize in collecting and analyzing large unstructured data, which analysts can use in their valuation efforts.

on this topic, we report the list of 520 teams and vendors in Online Appendix Figure A1. We use the list of full and abbreviated names as our initial keywords and search for them in analysts' written reports.⁷ For each report identified by these initial keywords, we read the passages surrounding the initial keywords to verify that the report indeed adopts alternative data.

Some analysts explicitly reference the use of alternative data without disclosing their source. To capture these reports, we follow prior literature (Hoberg and Moon, 2017; 2019) and conduct an iterative keyword-search process. In particular, within our first set of analyst reports that reference an in-house data-science team or external alternative-data vendor, we extract a list of keywords that analysts use to describe the alternative data. We then use these new keywords to search for additional reports that incorporate alternative data (but do not reference their sources) and continue expanding our keywords list. Using this iterative process, we arrive at our final set of keywords, which we use to identify reports that reference the use of alternative data. We report our final set of keywords in Appendix 1. In our last step, we (again) read all reports flagged as using alternative data to verify that the analysts indeed discuss the use of alternative data in their analyses.⁸

To illustrate our process by example, one of the alternative data vendors in our sample is "Remote Sensing Metrics," also referred to as "RS Metrics." We first search for reports containing the terms "Remote Sensing Metrics" or "RS Metrics." We find 47 reports that contain these two keywords. The figure below is an excerpt from one such report:⁹

⁷ We convert all names and all text in the reports to lowercase characters.

⁸ To researchers interested in further studying the use of alternative data, we caution that analysts' use of the keywords presented in Appendix 1 is a necessary but not a sufficient condition. We found our final step of carefully re-reading all reports to eliminate false positives crucial for cleanly separating reports that adopt alternative data from those that do not.

⁹ UBS; Neil Currie, Krista Zuber, and David Eads; Walmart Inc; August 12, 2010.

UBS Proprietary National Parking Lot Fill Rate Analysis

We have conducted an analysis with Remote Sensing Metrics, LLC to track parking lot fill rates in order to predict overall US comp-sales performance at Walmart Stores using a sample of between 100 and 150 like-for-like satellite images each month for the past six months. Samples are representative of geographic region, store formats, day of week, and the time of period analysis. All satellite images are usually taken between 10:30am and 1pm to minimize shadows on the images. We believe a traditional grocery trip is less fixed to a certain time of day and thus the time-slot window for imagery results bears less risk than for other more discretionary shopping trips.

Reading the text surrounding the keyword “Remote Sensing Metrics,” we identify two additional keywords related to alternative data: “parking lot fill rates” and “satellite image.” We use these new keywords to search for more reports that adopt alternative data but do not reference “Remote Sensing Metrics” or “RS Metrics.”

The figure below is an excerpt from one such report:¹⁰

- **April Results In +LSD Range:** The McDonald's U.S. sales result for April was +4.0%, which we believe included no meaningful menu price benefit. For April, we had predicted SSS of +3.0%, based widely on our proprietary parking lot fill rate analysis which had suggested lunch trends were positive and in the +3.0% range. We believe that featured core products, breakfast, and beverage platforms marginally contributed to overall U.S. comp trends in line with our expectations. Backing into our quarterly estimate, our +3.3% domestic 2Q11 projection suggests a +3.0% estimated comp for May and June.

In our final step, we read all reports that our procedure flags as discussing the use of alternative data to verify that the analysts indeed incorporate alternative data into their reports. For instance, some firms in our sample provide satellite-related products or employ satellite imagery in their business processes (e.g., oil and gas exploration). We exclude such cases. The figure below is an example of a false positive:¹¹

¹⁰ Piper Jaffary; Nicole Miller Regan and Joshua C. Long; McDonald's Corporation; May 9, 2011.

¹¹ Jefferies Group. John DiFucci, Joseph Gallo, and Howard Ma; Microsoft Corporation; May 11, 2017.

Digital Transformation. Much of the Analyst day was focused on capturing customer relevance which will translate to revenue growth. The company believes that its TAM has increased to \$ 4.5 trillion today. In one example, Microsoft helped Land O'Lakes with a multiple year digital transformation which saw the company use Office 365, Surface, Windows 10, Azure, and Hololens. They were able to take decades of satellite imagery load it into Azure, and then analyze the land, different weather metrics and type of seeds to better layout farms. This process drove yield increases from 130 bushels of corn per acre to now over 500 bushels.

Our primary variable, $I(\text{Alternative Data}_{i,f,t})$, equals one if analyst i 's forecast for the annual earnings of firm f at time t is accompanied by a written report that explicitly references the use of alternative data and zero otherwise.

For $I(\text{Alternative Data})$ to equal one, the following three conditions must be met: (a) Analysts have access to alternative data. (b) Analysts believe the alternative data contain clear signals regarding a company's fundamentals and they incorporate those signals into their forecasts and recommendations. (c) Finally, analysts disclose their reliance in their written reports. In Subsection 4.1, we discuss the implications of measuring this joint effect for the interpretation of our results.

Our variable may also equal one if an analyst pretends to have studied alternative data. While we cannot rule out this possibility, such fabricated use would expose the analyst to serious career and litigation risks.

Online Appendix Figure A2 provides an example of an analyst report that draws from alternative data. As depicted in the figure, discussions of alternative data often include visuals depicting its application, such as satellite images of parking lots and their relation to revenues. After excluding standard boilerplate illustrations, the average number of figures included in analysts' reports is 0.86. When analysts make reference to the use of alternative data, this number rises to 4.16. The high level of detail typically provided by analysts when discussing their use of alternative data suggests further that it is improbable that analysts falsify their use of such data.

2.3 Descriptive Evidence Regarding Analysts' Discussion of Alternative Data

Our first test examines how frequently analysts report using alternative data and, if so, in what manner. Panel A of Table 1 displays summary statistics for $I(\text{Alternative Data})$ across years. Since we have only partial data for 2009 and 2019, we combine the observations in 2009 with those in 2010 and the

observations in 2019 with those in 2018. Panel A reveals that, in 2009/2010, 6% of the analyst forecasts are couched in reports that discuss the use of alternative data. By 2018/2019, the corresponding number is 10%.

The fraction of analysts discussing the use of alternative data for at least one of the firms they cover is naturally greater than the fraction of reports discussing the use of alternative data. In particular, we find that, as of 2009/2010, 11% of the analysts in our sample report using alternative data for at least one of the firms they cover. This fraction increases to 28% by 2018/2019.

In Panel B we report summary statistics for $I(\textit{Alternative Data})$ across industry sectors. Analysts most frequently discuss the use of alternative data for firms operating in the Information Technology sector: the average $I(\textit{Alternative Data})$ is 16%. Alternative data use is also widely discussed for firms operating in Consumer Discretionary (10%), Consumer Staples (10%), Communication Services (9%), Health Care (8%), and Industrials (6%). Analysts infrequently adopt alternative data for firms in the Energy (2%), Financials (1%), and Materials (1%) sectors.

In our study, we manually assign reports mentioning the use of alternative data into the following eight categories: (1) app-usage data, (2) sentiment data, (3) employee data, (4) geospatial data, (5) point-of-sale data, (6) satellite-image data, (7) web-traffic data, and (8) other types of alternative data. Some reports are assigned to more than one category as analysts occasionally reference multiple alternative-data categories. Appendix 1 details how we allocate the reports across the above eight categories. The results reported in Table 2 indicate that, of the 5,639 forecasts from reports that discuss the use of alternative data, 1,944 (34%) are based on web-traffic data. The next most popular categories are other (23%), followed by point of sale (19%), sentiment (19%), employee (10%), and app usage (8%). The least popular categories are geospatial (5%) and satellite image (3%).¹²

Figure 1 displays two timelines. The first timeline indicates when—for our sample—we observe the first analyst report that explicitly references the use of alternative data from a given alternative-data category. The second timeline indicates when we observe the first one hundred such

¹² Online Appendix Figure A3 shows how frequently alternative data from a particular category are discussed across industry sectors. Among others results, the figure reveals that web-traffic data are discussed frequently for firms operating in the Information Technology sector, whereas point-of-sale data are commonly referenced for firms operating in the Consumer Discretionary sector.

reports. The sequence for when we observe the first analyst report is as follows: sentiment (June 11, 2009), web traffic (June 12, 2009), point of sale (August 6, 2009), employee (January 5, 2010), geospatial (January 22, 2010), satellite image (May 3, 2010), other (June 12, 2009) and app usage (July 27, 2010). In other words, within essentially the first year of our sample period, we find that analysts explicitly reference the use of alternative data from all eight categories.

The sequence for when we observe the first one hundred analyst reports is as follows: web traffic (January 19, 2010), point of sale (March 28, 2011), other (August 11, 2011), sentiment (October 10, 2011), satellite image (May 21, 2012), geospatial (June 5, 2012), app usage (September 8, 2014) and employee (August 18, 2015). In other words, by mid-2012, analysts draw extensively from alternative data from six of the eight categories. Only app usage- and employee-level data are not widely adopted until mid-2015.

Overall, our evidence shows that analysts frequently incorporate alternative data into their reports, at least for the largest and economically most meaningful firms.

3. The Usefulness of Alternative Data

In this section, we examine whether institutional investors, the primary clientele of analysts, value analysts' use of alternative data. We also consider the wider implications of analysts' adoption of alternative data on the quality of institutional investors' decision-making and the competitive balance among investors.

3.1 Analysts' Use of Alternative Data and Forecast Accuracy

The hypothesis that institutional investors appreciate analysts' use of alternative data rests on the premise that analysts can derive unique and value-relevant insights from such data. To examine the validity of this assumption, we test whether earnings forecasts from reports that mention the use of alternative data are more precise:

$$Acc_{i,f,t} = \eta_{i,f} + \theta_{f,t} + \beta I(Alternative\ Data_{i,f,t}) + \gamma` Controls + \varepsilon_{i,f,t} \quad (1)$$

The observations are at the analyst/firm/forecast-date level. The construction of $Acc_{i,f,t}$ follows prior literature (e.g., Clement, 1999; Bradley, Gokkaya, and Liu, 2017; Green, Jame, Markov, and Subasi,

2014; Harford, Jiang, Wang, and Xie, 2019). We first compute $AFE_{i,f,t}$ as the absolute value of the difference between analyst i 's annual earnings forecast issued for firm f at time t and the corresponding actual reported annual earnings. We then construct $PMAFE_{i,f,t}$ as the difference between $AFE_{i,f,t}$ and $Avg(AFE)_{f,t}$, scaled by $Avg(AFE)_{f,t}$ to reduce heteroskedasticity. $Avg(AFE)_{f,t}$ is the average absolute forecast error across all analysts covering firm f as of the corresponding forecast period, excluding analyst i and other analysts, who also report drawing from alternative data in their coverage of firm f as of time t . $PMAFE_{i,f,t}$ thus measures analyst i 's forecast accuracy relative to the forecast accuracies of all analysts who cover the same firm at the same time but do not report drawing from alternative data. Negative values, or smaller forecast errors, indicate above-average performance. Positive values, or larger forecast errors, indicate below-average performance. To facilitate interpretation, $Acc_{i,f,t}$ equals $PMAFE_{i,f,t} \times (-1)$.¹³ We provide descriptive statistics regarding Acc and all other variables we use in this paper in Online Appendix Table A3.

We include both analyst-firm, $\eta_{i,f}$, and firm-year, $\theta_{f,t}$, fixed effects. Angrist and Pischke (2008) show that our two-way fixed-effects specification is equivalent to the basic difference-in-differences specification. The estimate of $I(Alternative\ Data)$ thus indicates how much more accurate an analyst becomes in the post-adoption period relative to the pre-adoption period compared with analysts covering the same firm over the same period that do not reference alternative data.

Controls include the following analyst characteristics: *Forecast Age*, *Analyst/Firm Experience*, *Analyst Experience*, *#Firms Covered*, *Forecast Frequency*, and *Broker Size*.¹⁴ We do not control for firm characteristics as our fixed effects subsume them. Since our final sample comprises 64,018 written reports and earnings forecasts, the number of observations on which we estimate regression equation (1) is 64,018. We double-cluster our standard errors at the analyst- and year-month levels.

Difference-in-differences specifications often produce causal evidence. Our setting does *not* lend itself to causal inferences. That is, while the estimate of $I(Alternative\ Data)$ indicates how much more accurate an analyst becomes in the post-adoption period relative to the pre-adoption period

¹³ In robustness checks, we base our analysis directly on $AFE_{i,f,t}$ (Bradley, Gokkaya, and Liu, 2017). As shown in Online Appendix Table A2, the results based on the absolute forecast error are similar to those based on $Acc_{i,f,t}$.

¹⁴ We detail the construction of these variables in Appendix 2.

compared with her non-adopting counterparts, it cannot tell us how much of the abnormal performance improvement is truly caused by alternative data because alternative-data adoption is endogenous. For instance, alternative-data adoption may coincide with an analyst's decision to exert greater effort covering the corresponding firm, which, in turn, leads to improved forecast accuracy ("increased effort channel").

We try to gauge the relevance of the increased-effort channel by constructing various measures of analyst effort used in the literature, including the timeliness of forecasts, the number of forecast revisions, and analyst activity during earnings conference calls (Merkley, Michaely, and Pacelli, 2017; Hwang, Liberti, and Sturgess, 2019; Grennan and Michaely, 2020). We then explore whether adopting alternative data associates with greater effort. As detailed and tabulated in Online Appendix Table A4, alternative-data adoption correlates neither with the timeliness of forecasts nor with the number of forecast revisions. The adoption of alternative data also is not associated with the number of questions asked during earnings conference calls, the number of words spoken, or the types of questions asked; it correlates marginally positively with the number of conference calls attended.

While we generally fail to find empirical support for the increased-effort channel, our tests may lack power. Our point estimate of how much an analyst's forecast accuracy improves after she adopts alternative data should be interpreted with this caveat in mind.

We present our regression results in Table 3. The results reported in column (1) show that the coefficient estimate of $I(\textit{Alternative Data})$ is 0.214 (t -statistic = 6.30). To illustrate the economic significance of this estimate, a 0.214 improvement would move an analyst who is at the median in terms of forecast accuracy to the 62nd percentile.

Another way to gauge the economic significance is to compare the estimate of $I(\textit{Alternative Data})$ with those of our control variables. For instance, the results reported in column (1) show that forecast accuracy increases significantly with the number of years an analyst has been covering a particular firm: The estimate of $\textit{Analyst/Firm Experience}$ is 0.060 (t -statistic = 2.72). Comparing the estimate of $I(\textit{Alternative Data})$ with that of $\textit{Analyst/Firm Experience}$ suggests that the performance improvement accompanying the adoption of alternative data is equivalent to having covered the corresponding firm for 3.6 additional years.

Overall, while we cannot provide strong causal evidence, our results are at least consistent with the premise that analysts can derive distinct and value-relevant insights from alternative data.

3.2 Analysts' Use of Alternative Data and Trading Commissions

The analyst business model is rooted in soft dollar agreements. When institutional investors find greater value in an analyst's research, they allocate a higher volume of trades through the brokerage that employs the analyst and, consequently, pay a greater amount in trading commissions to the corresponding brokerage. A simple test to determine if institutional investors place value on analysts' adoption of alternative data is thus to correlate the reported use of alternative data with the amount in trading commissions received.

We obtain institutional investor trade data from ANcerno, a firm that provides transaction-cost analysis to institutional clients, including investment managers, like Fidelity Investments, and plan sponsors, like CalPERS (Hu, Jo, Wang, and Xie, 2018). While the ANcerno dataset covers a wide set of funds, the coverage is not comprehensive. Hu et al. (2018) estimate that the institutions in the ANcerno dataset account for 15% of all institutional trading volume. As discussed by Hu et al. (2018) and other studies utilizing the ANcerno dataset (Puckett and Yan, 2011; Jame, 2018), the dataset does not seem to suffer from sample selection bias.

For each institution covered by ANcerno, the dataset contains specific information for each of their trades, such as ticker symbol, date, trade direction, number of shares traded, and—crucially—identifiers for the investment manager, the broker executing the trade, and the commissions paid to the broker.

The ANcerno data contain records on all 35 DJI constituents. While the ANcerno data span the period running from 1997 through 2015, as of 2011 ANcerno no longer reports the identifier, which would have allowed us to distinguish trades by individual institutions. Given that our sample of earnings forecasts and analyst reports commences in 2009, our brokerage-commissions analysis thus encompasses the years 2009 and 2010. This period includes a total of 7,634 analyst reports. Of these, 2,877 are issued by brokerages that are not included in the ANcerno database. Our sample for this part of our paper thus comprises 4,757 analyst reports.

To examine how the distribution of commissions is influenced by the dissemination of alternative-data insights, we consider all transactions completed within the first three months following a report's issuance that involve the stock covered in the report. We then total the commissions earned by the brokerage issuing the report. The median commissions earned by a brokerage in the first three months following a report's publication is \$11,578.40. To the degree that ANcerno covers 15% of all institutional trades and the coverage is representative, we estimate that the total median commissions earned by a brokerage in the first three months following a report's publication is \$77,189.33.

We re-estimate regression equation (1) but replace the earnings-forecast-accuracy variable with our new trading-commissions variable. The estimated coefficient β now indicates how much more in trading commissions the analyst's brokerage receives in the post-adoption period relative to the pre-adoption period compared with analysts covering the same firm over the same period but not adopting alternative data. We again double-cluster our standard errors at the broker- and year-month levels.

As reported in Table 4, our estimate suggests that, on average, the adoption of alternative data increases the total commissions paid to the brokerage in the first three months by an additional \$11,858.82 (t -statistic = 1.82), essentially doubling the median commission earned by a brokerage. In additional analyses, we experiment with time frames other than a three-month period. The results of these tests are in line with those reported in this study and are available upon request.

Overall, our findings suggest that institutional investors recognize and value analysts' incorporation of alternative data and compensate them accordingly.

3.3 Analysts' Alternative Data Use and the Playing Field Among Institutional Investors

The adoption of alternative data by sell-side analysts has the potential to reshape not only the relationship between analysts and investors but also the competitive balance among the investors they serve. As alluded to in the introduction, hedge funds have been at the forefront of utilizing alternative data, while other institutional players, such as mutual funds, exhibit slower adoption rates. Additionally, existing research suggests that hedge funds execute more profitable trades than non-hedge funds. This outperformance may arise in part because hedge funds have superior data.

In this part of the study, we investigate whether the dissemination of alternative data insights through analyst reports has leveled the playing field for hedge funds and non-hedge-fund investors. To assess this possibility, we again combine our analyst data with the ANcerno data and explore whether analysts' adoption of alternative data has narrowed the performance gap usually seen between hedge funds and non-hedge funds.

As we are interested in trades affected by the dissemination of alternative-data-driven insights, we again focus on transactions that occur within three months after the issuance of an analyst report and that involve the stock discussed in the report. Analyses based on alternate windows produce results that are similar to those presented here and are available upon request. Applying Jame's (2018) methodology, we separate transactions executed by hedge funds from those made by other institutional investors.¹⁵

We further separate non-hedge-fund trades by whether they are preceded by analyst reports that adopt alternative data or whether they are placed following reports that do not draw from alternative data. To ensure that the investor composition behind the trades associated with alternative data and those without such a linkage is comparable, we restrict our analysis to investors who appear in both subsets of trades.

We assess the performance of investors' trades by creating transaction-based calendar-time portfolios, following the approach of Seasholes and Zhu (2010) and Ben-David, Birru, and Rossi (2019). Stocks purchased (sold) on day t are added to the buy (sell) portfolio at the beginning of day $t+2$. We assume a holding period of three months. We compute DGTW-adjusted returns for each stock/day and construct value-weighted portfolio returns. DGTW-adjusted returns are returns on a stock minus the value-weighted returns on a portfolio of stocks in the same size, book-to-market, and momentum quintiles. We report the time-series average of the buy-minus-sell portfolio returns, whereby each daily portfolio return is weighted by the number of trades contributing to the portfolio on that specific day.

¹⁵ For additional details, we refer the reader to Jame (2018).

We report the results in Table 5. We find that the stocks that hedge funds buy outperform those they sell. The difference in the stocks' ensuing returns is 10.79% when annualized. The statistical significance is modest (t -statistic = 1.54), which largely reflects the high standard errors tied to the comparatively low number of hedge funds and the short sample period.

The performance is noticeably different for non-hedge funds. When analysts do not discuss the use of alternative data, the stocks that non-hedge funds buy outperform those they sell by 3.85% when annualized. The difference in the stocks' ensuing returns of 3.85% for non-hedge funds is similar to the numbers reported by Puckett and Yan (2011) and Busse, Tong, Tong, and Zhang (2019), who also analyze ANcerno data.

Interestingly, when analysts incorporate alternative data, the difference in the stocks' ensuing returns for non-hedge funds improves significantly, reaching 9.85% when annualized. This figure closely matches that achieved by hedge funds. Again, the standard error of the mean is very high because of the comparatively low number of analyst reports discussing the use of alternative data and the short sample period. Still, our results are at least consistent with the proposition that when they employ alternative data analysts can equalize the competitive environment and allow the non-hedge funds that pay attention to alternative-data insights to match the performance typically associated with hedge funds.

Another investor group that may benefit from alternative-data adoption by analysts is retail investors. Although they do not comprise the primary audience for analysts, retail investors can also access analyst reports, depending on their retail broker accounts and the investment-related platforms to which they subscribe.

To study the possible benefits retail investors might derive from analysts' use of alternative data, we utilize Trade and Quote (TAQ) data. We construct a measure of retail order imbalance following the method of Barber, Huang, Jorion, Odean, and Schwarz (2023). Our findings, presented in Online Appendix Table A5, reveal that retail investors are three times more likely to adjust their trades in line with changes in analyst recommendations when these recommendations come from reports that incorporate alternative data. These results imply that at least some retail investors follow analyst recommendations more closely when they are anchored in alternative data. To gauge the future

performance of these trades, we estimate regressions of one-week, one-month, or three-month future returns on retail investor order imbalances. We distinguish between orders executed following the release of an analyst report that discusses alternative data and those following the release of reports that do not include such discussions. As shown in Online Appendix Table A6, retail investor order imbalances more positively predict future returns when the orders are placed in the presence of alternative data—inclusive analyst reports. This finding is consistent with the notion that using alternative data assists retail investors in executing more profitable trades.

In separate analyses, we assess the broader market’s reaction to analysts’ use of alternative data. Specifically, we examine whether the market’s response to changes in analyst research outputs, including earnings forecasts, price-target forecasts, and overall recommendation levels, becomes more pronounced when those outputs are couched in alternative data. We detail our empirical design and our results in Online Appendix Table A7. In short, our results again suggest that broad sections of investors more closely follow analyst outputs when they are anchored in alternative data.

4. Discussion and Additional Analyses

Before concluding, we discuss possible reasons we do not observe more analyst reports mentioning the use of alternative data. We also examine the extent to which our findings regarding the largest U.S. firms apply to smaller firms. Finally, we explore variations in our key independent variable, *I(Alternative Data)* and our main empirical specification.

4.1 Why Do We Not Observe More Analyst Reports Mentioning the Use of Alternative Data?

The seeming benefits accompanying the adoption of alternative data, such as higher trading commissions, raise the question why analysts do not incorporate alternative data into their reports more frequently.

Our key independent variable, *I(Alternative Data)*, can equal zero for at least three reasons. First, analysts may not have the resources or willingness to study alternative data, causing *I(Alternative Data)* to be zero (“resource limitations”). Second, even if analysts access and study alternative data,

they may not always uncover clear, unique, or value-relevant insights. Analysts are unlikely to discuss such “failed” uses of alternative data in the limited space that is available to them in their reports (“intermittent usefulness”). Finally, analysts may identify relevant signals but choose strategically not to detail these insights in their reports (“strategic considerations”).

4.1.1 *The Relevance of Resource Limitations and Intermittent Usefulness*

To gauge the relevance of the resource-limitations and intermittent-usefulness perspectives for explaining why $I(\text{Alternative Data})$ does not equal one more frequently, we estimate the following probit regression:

$$I(\text{Alternative Data}_{i,f,t}) = \alpha + \beta' X_{i,f,t} + \delta' \text{Controls}_{i,f,t} + \varepsilon_{i,f,t} \quad (2)$$

The observations are at the analyst/firm/forecast-date level. *Controls* include the following analyst- and firm-level characteristics: *Analyst/Firm Experience*, *Analyst Experience*, *#Firms Covered*, *Forecast Frequency*, *Broker Size*, *Size*, *M/B*, and *Momentum*. We detail the construction of these variables in Appendix 2. We double-cluster our standard errors at the analyst- and year-month levels.

4.1.1.1 *Variables Tied to Resource Limitations*

Our first key independent variable is an indicator of whether an analyst’s brokerage has an in-house data-science team. Our second key variable is the number of an analyst’s colleagues working in the same city who have already discussed the use of alternative data in at least one of their reports as of the analyst’s forecast date. The expenses necessary to access and analyze alternative data are presumably reduced when analysts are supported by an in-house data-science team. Moreover, suppose that individuals in knowledge-based industries owe much of their success to their colleagues (Hwang, Liberti, and Sturgess, 2019). In that case, analysts surrounded by peers who already utilize alternative data should find the learning curve less steep. If resource limitations are an important reason that $I(\text{Alternative Data})$ does not equal one more frequently, we should observe more frequent alternative-data adoptions when there are fewer resource limitations; we should thus observe positive coefficient estimates for the first two key independent variables.

4.1.1.2 Variables Tied to Intermittent Usefulness

Similarly, if intermittent usefulness is an important reason that $I(\text{Alternative Data})$ does not equal one more frequently, we should observe a positive correlation between the frequency of alternative-data adoptions and measures of its usefulness.

Alternative data offer three key advantages. First, alternative data provide immediate indications of a company's performance, allowing analysts to draw from current rather than stale information. Second, the origin of alternative data from independent third-party vendors mitigates the risk of bias or distortion often associated with data produced by company managers. Third, the detailed nature of alternative data—for instance data based on breaking down information into specific products or branches—provides analysts with a more detailed and nuanced understanding of a company's performance.

Building on these considerations, we propose that alternative data are more useful and, as a result, become more frequently discussed in analyst reports when (1) receiving instantaneous signals regarding a company's performance is critical, (2) concerns of misrepresentation are urgent and traditional data are ambiguous, and (3) analysts lack access to granular data.

Receiving instantaneous signals regarding a company's performance becomes critical when there are relatively few company announcements and when the uncertainty regarding a company's performance is high. We thus conjecture that the usefulness of alternative data and the likelihood of alternative-data adoption are higher when a firm files relatively few Form 8-Ks, when stock return volatility is high, and when the absolute value of earnings surprises is high. For a description of the precise construction of these and the subsequent key independent variables outlined below, we direct the reader's attention to Appendix 2.

That managers and firms are removed from the "alternative data-generating process" becomes particularly relevant when misrepresentation concerns are urgent and traditional data are ambiguous. As measures for concerns of misrepresentation and the ambiguity of traditional data, we consider whether a company has had to restate its earnings (Wilson, 2008) and the absolute value of discretionary accruals (Bhattacharya, Desai, and Venkataraman, 2013).

Finally, private meetings with management are one of the key channels through which analysts can obtain a more nuanced perspective on a company's performance (Green, Jame, Markov, and Subasi, 2014; Soltes, 2014; Solomon and Soltes, 2015; Brown, Call, Clement, and Sharp, 2015; Bengtzen, 2017). Not all analysts are granted private meetings with management, however, putting them at a significant disadvantage.

To measure whether analyst i has preferential access to the management of firm f , we consider whether analyst i works for a broker that hosts an investor conference in which firm f participates. Green, Jame, Markov, and Subasi (2014) argue that broker-hosted investor conferences, which provide selected investors opportunities to interact with senior corporate managers, offer insights into whether a particular analyst has preferential access to the firms participating in the conference.

4.1.1.3 Results Regarding The Relevance of Resource Limitations and Intermittent Usefulness

We report the results of the analyses of resource limitations and intermittent usefulness in Table 6. We find that analysts incorporate alternative data into their reports less frequently when their brokerages do not deploy in-house data-science teams and when analysts do not have colleagues that already draw from alternative data. These results are consistent with the resource-limitations perspective.

Consistent with the intermittent-usefulness idea, we find that analysts adopt alternative data more frequently when a company provides relatively few announcements, when stock-return volatility and the absolute value of earnings surprises are high, when the firm has had to restate its earnings and the absolute value of discretionary accruals are high, and when analysts lack preferential access to management through private meetings.¹⁶

Overall, our results suggest that resource limitations and intermittent usefulness are important determinants of $I(Alternative\ Data)$ and help to explain why analysts do not discuss alternative data in their reports more frequently.

¹⁶ In Online Appendix Table A8, we describe analyses showing that the incremental improvements in earnings-forecast accuracy associated with discussions of alternative data also strengthen when the absolute value of earnings surprises are high and when the firm has had to restate its earnings.

4.1.2 *The Relevance of Strategic Considerations*

Analysts may choose not to disclose their reliance on alternative data for two strategic reasons. First, some analysts may strive to continuously provide fresh perspectives. An analyst reporting the use of alternative data in one year may thus avoid mentioning her continued reliance in the following years to avoid repetition.

To investigate this possibility, we consider cases where analysts report using alternative data for a particular firm in a particular year. We find that these analysts discuss their reliance on the same alternative-data category in the subsequent year 55% of the time.

It appears likely that an analyst who incorporated alternative data into her report in one year could easily obtain an updated version of the data in the ensuing year. The fact that in 45% of the cases analysts do not reference the same alternative data category in the subsequent year suggests either that the perceived value of alternative data decreases rapidly (“intermittent usefulness”) or that analysts avoid repeating their continued reliance (“strategic considerations”).

To shed light on the relevance of these possibilities, we modify our earnings-forecast-accuracy regression. Our adjusted independent variable equals one if an analyst does not discuss the use of alternative data in year t but did so in year $t-1$. Suppose analysts no longer mention the use of alternative data because they find the data no longer useful. Suppose further that analysts’ assessments of the diminished usefulness are accurate. In that case, our revised independent variable should no longer be significantly associated with earnings-forecast accuracy.

Our analysis yields a coefficient estimate of 0.070 (t -statistic = 2.57). This robust estimate is consistent with the idea that analysts continue to benefit from alternative data insights but to avoid repetition opt not to reiterate their reliance. At the same time, the estimate of 0.070 is markedly lower than that of our main specification (coefficient estimate = 0.214, t -statistic = 6.30), suggesting that, in many cases, the utility of alternative data does diminish over time. This result again points to intermittent usefulness as an important reason that $I(\text{Alternative Data})$ does not equal one more frequently.

Another strategic consideration that may dissuade analysts from disclosing their use of alternative data is that they fear that disclosing their methodologies and data sources would make it

easier for competing analysts to imitate their outputs. This could result in the erosion of the competitive advantage the original analyst initially held.

We find limited evidence for the relevance of this particular consideration. In our sample, we detect only 92 cases where the initial adoption of alternative data is replicated by another analyst in the concurrent year for the identical company. In these situations, we find that the original analyst receives, on average, \$20,059.55 in commissions (for trades in the corresponding stock occurring within the first three months of the report publication). Subsequent to the imitation, the total three-month average commission of the original analyst stands at \$25,537.25, suggesting that the financial repercussions from imitation are limited.¹⁷

4.2 The Use and Usefulness of Alternative Data Among Small Firms

Our tests so far have involved constituents of the DJI. To explore the prevalence and impact of alternative data within the small-cap sector, we analyze a random sample of 200 companies drawn from the lower half of the size distribution in the CRSP database. The sample period is 2009–2019. Of these companies, 143 have earnings forecasts and analyst reports data in IBES and Investext. The median market capitalization of our small firms is \$1.013 billion. In comparison, the median market capitalization of DJI constituents is \$11.854 billion. We retrieve a total of 13,123 analyst reports for our 143 small firms over the 2009–2019 period, compared with 64,018 reports for the 35 firms in the DJI. Our small firms thus receive substantially less extensive coverage, which is no surprise given that analysts' key clientele, institutional investors, invest primarily in larger, more liquid stocks.

Following the same procedure described in Subsection 2.2, we identify the fraction of reports that discuss the use of alternative data. We find that for our small firms, analysts discuss alternative data in only 2% of their reports, compared with 9% for the DJI constituents. It thus appears that analysts use alternative data significantly less frequently among small firms. One factor that might explain this result is that analysts are reluctant to bear the financial and learning costs associated with alternative-data

¹⁷ Concerns over imitation might still affect analysts' reporting behavior if they are unaware that imitations are rare and do not affect an original analyst's trading commissions.

adoption for firms that attract limited attention from institutional investors. It is also possible that the quality of alternative data is lower for smaller firms.

We find that references to alternative data are accompanied by more accurate earnings forecasts even among our small firms. As reported in Online Appendix Table A9, the coefficient estimate is 0.197, suggesting that an analyst at the median accuracy level improves to the 60th percentile. The magnitude of the association is similar to that for DJI constituents.

We also re-run our analyses of the trading commissions received by analysts' brokerages and the performance differential between hedge funds and non-hedge-fund institutional investors. Since institutions trade significantly less extensively in smaller firms, the sample sizes for these additional analyses are substantially smaller. As shown in Online Appendix Table A10, we find evidence that institutional investors continue to reward analysts who discuss alternative data by directing more trades through their brokerages. The coefficient estimate is 2,818.23 (t -statistic = 3.38). Compared with the estimate of 11,858.82 for DJI constituents, this estimate is substantially smaller. The weaker economic significance suggests that there is a weaker incentive for analysts to adopt alternative data for smaller firms, which could explain why the fraction of reports that reference alternative data is much lower for smaller firms than for DJI constituents.

When we repeat our transaction-based calendar-time portfolio analysis for the smaller firms, we find that analysts' discussions of alternative data coincide with a narrower performance gap between institutional investors and hedge funds. However, the standard errors are so large that none of the average returns are statistically different from zero.

4.3 Differences in Usefulness by Alternative Data Types

4.3.1 Variation Across Alternative Data Categories

While our results suggest that alternative data contain unique and value-relevant insights, the usefulness of the data may vary by type. To explore this possibility, we first replace $I(\text{Alternative Data})$ with eight indicator variables, each denoting whether an analyst uses alternative data from a particular alternative-data category. We then re-estimate our earnings-forecast-accuracy regression (1). The results reported in Table 7 show that the adoption of alternative data from six of the eight categories is associated with

statistically significant performance improvements. Within those six, the ranking in descending order based on the magnitude of the coefficient estimates is as follows: (1) app usage, (2) sentiment, (3) employee, (4) other, (5) point of sale, and (6) web traffic.

Unlike the adoption of alternative data from the above six categories, our results show that (7) geospatial data and (8) satellite image data are not associated with more accurate earnings forecasts. As shown in Table 2, these are also the two categories that analysts report using least frequently. In 2017, Ernst & Young Global Limited surveyed hedge funds and asked which datasets, in their experience, have been the *least* accurate and *least* insightful.¹⁸ The two datasets that are by far the most frequently mentioned are “geolocation” and “satellite.” While the survey conducted by Ernst & Young Global Limited represents a one-time snapshot of investors’ opinions, we nevertheless find the overlap between the survey results and our regression results revealing.¹⁹

In another test, we replace $I(\text{Alternative Data})$ with the number of distinct alternative data categories discussed in an analyst’s report. We then re-estimate our regressions within the subset of cases where an analyst reports the use of alternative data from at least one category. In short, we detect a strong positive correlation between the number of alternative data categories and the accuracy of earnings forecasts. This finding implies that, when analysts’ views are supported by a wider range of data types, their predictions are particularly accurate (Table 7).

4.3.2 Variation Across More or Less Proprietary Alternative Data

Among our eight alternative-data categories, four categories—(1) sentiment data, (2) employee data, (3) geospatial data, and (4) web-traffic data—are comparatively more accessible to the public and less dependent on external data vendors (“accessible data”). For example, investors can acquire variants of sentiment data through social media platforms, catch a glimpse of employee-related information through public websites (e.g., Glassdoor.com), gather geospatial data through a combination of Google

¹⁸ The survey results are viewable at <https://alternativedata.org/stats/>.

¹⁹ In separate tests, we also gauge the usefulness of alternative data categories across industries. Specifically, we re-estimate our regressions separately for each Global Industry Classification Standard Sector. We present the results in Online Appendix Figure A4. Among other findings, our results suggest that app-usage data are particularly beneficial for forecasting the performance of firms that operate in the Information Technology sector, while point-of-sale data is particularly advantageous for predicting the performance of firms that operate in the Consumer Staples sector.

Maps and publicly available demographic data, and retrieve web-traffic data using tools such as Google Trends. In contrast, data from the remaining four categories—(5) app-usage data, (6) point-of-sale data, (7) satellite image data, and (8) other unspecified types of data—are generally not accessible to the general public (“proprietary data”).²⁰

In separate tests, we gauge whether accessible data are as useful to analysts for predicting a company’s performance as proprietary data are. We find that, among analyst reports that discuss alternative data, 64.2% base their analysis on accessible data, while 49.3% rely on proprietary data.²¹ Proprietary data are adopted more frequently when analysts work for larger brokerages or brokerages with internal data-science teams; they are also adopted more frequently when more colleagues in the same locale lean on alternative data in their reports.

Regarding the implications for forecast accuracy, results reported in Table 7 show that references to proprietary data (coefficient estimate = 0.243, t -statistic = 6.82) and accessible data (coefficient estimate = 0.188, t -statistic = 4.07) come with similar increases in the precision of earnings forecasts; the two coefficient estimates are not statistically different from each other (F -statistic = 0.99, p -value = 0.32).

4.4 Alternative Data Use and Earnings-Forecast Accuracy: Instrumental Variable- and Matching Analyses

To study the relationship between alternative-data use and earnings-forecast accuracy, our primary analysis adopts a standard difference-in-differences method and assesses how much more accurate an analyst becomes in the post-adoption period relative to the pre-adoption period compared with analysts covering the same firm over the same period that do not reference alternative data. As alluded to in Subsection 3.1, an analyst’s decision to adopt alternative data is not exogenous and, instead, could reflect the decision to exert greater effort. To address this concern, we experiment with an instrumental-variable approach. The basic premise behind our instrumental-variable approach is that an individual

²⁰ We emphasize that this is a relative statement. Even data types we deem more accessible to the public, such as web-traffic data, encompass specific details such as the duration of user visits to a merchant’s website, which are typically not accessible to the general public.

²¹ The sum of the two percentages exceeds 100%, as some analyst reports base their analysis on both accessible and proprietary data.

analyst's ability to use alternative data depends on her firm's infrastructure and technology investment ("relevance"). At the same time, certain brokerage-level decisions, which allow analysts to access alternative data, are made independently of a single analyst's preferences or actions ("exclusion restriction").

We consider two instruments. The first instrument, *First Time Use*, is an indicator variable that equals one when an analyst's colleague, working in the same city, adopts alternative data for the first time. The second, *Software Budget*, refers to the allocated budget for software purchases at the broker-year level, sourced from Aberdeen's Computer Intelligence Technology Database. We assess the validity of these instruments through the Hansen J -test for overidentifying restrictions and confirm their efficacy with underidentification tests and F -tests in the first-stage regression. Our regression results pass all diagnostic tests, implying that our instrumenting strategy is valid.

We report the results from the instrumental-variable approach in Online Appendix Table A11. In column (1) we report the results from the first-stage test. The coefficient estimates of *First Time Use* and *Software Budget* are positive and significant at the 1% level, suggesting that our instruments are highly correlated with the endogenous variable. The results reported in column (2) show that the second-stage coefficient estimates of the instrumented $I(\textit{Alternative Data})$ are positive and significant. The estimates are similar in magnitude to those obtained from the original regression specification.

We also experiment with a matching-sample analysis. We pair each alternative-data report within our sample with a non-alternative-data report written by analysts covering the same firm during the same forecast period. We ensure that the analysts with the matched reports work for brokerages that fall into the same size quintile (based on the number of analysts employed) and the same firm-specific experience quintile (based on the number of years of having covered the respective firm). If multiple reports satisfy the above criteria, we select the report with the forecast horizon that is most similar to that of the alternative data report. We then repeat our forecast accuracy analysis within this matching sample.

The results obtained with this matching-sample approach are presented in Online Appendix Table A12. The coefficient estimate for $I(\textit{Alternative Data})$ is 0.204, with a t -statistic of 3.83. Again, this estimate is very similar to that obtained in our primary analysis.

4.5 Other Key Performance Indicators

Our final test considers revenue-forecast accuracy as another crucial analyst output. Revenue forecasts are sourced directly from the IBES database and the accuracy of these forecasts is measured in the same manner as earnings-forecast accuracy. Not all analysts have revenue forecasts in the IBES database, though, so the sample size drops to 27,661.

In addition to revenue-forecast accuracy, one might also consider the accuracy of cost forecasts. Unfortunately, the IBES database does not contain cost forecasts. Instead, we compute “residual forecasts” by taking the difference between revenue-per-share forecasts and earnings-per-share forecasts and constructing a measure of accuracy based on these residual forecasts. The resulting measure may be seen as capturing any improvement in earnings-forecast accuracy that cannot be tied to more accurate revenue forecasts.

We report our results in Online Appendix Table A13. In column (1) and column (2) we present results based on revenue-forecast accuracy and residual-forecast accuracy, respectively. The results reported in column (1) reveal that discussing alternative data is associated with significantly more accurate revenue forecasts. Conversely, the findings reported in column (2) suggest that the use of alternative data does not lead to more accurate residual forecasts, given that the estimate of $I(\textit{Alternative Data})$ is close to zero. Put differently, our results indicate that, once any improvement tied to more accurate revenue forecasts is removed from the analysis, the adoption of alternative data no longer leads to more accurate earnings forecasts. Such a suggestion seems reasonable, considering that most alternative data pertain to a company’s sales, not profits.

5. Conclusion

Our study documents the dynamic role of sell-side analysts. Rapid advancements in information technology and the emergence of alternative-data sources are creating a wealth of information which, from an investor’s point of view, might dominate those provided by sell-side analysts, eventually rendering the analyst profession obsolete.

Our paper points to a more nuanced picture. We propose that the complexity and costs associated with accessing and interpreting alternative data represent a significant challenge for many

investors. This challenge opens a window of opportunity for analysts to serve as essential conduits of information. Instead of investors each bearing the financial and learning costs associated with studying alternative data, it is more efficient for a few analysts to absorb these expenses and become proficient in parsing alternative data. Analysts can then share their insights with investors, who, in turn, can integrate these signals with other information to determine a stock's value.

Our research findings are in line with this perspective. We find that analysts have begun to adopt alternative data and that investors value analysts' alternative-data-driven insights and reward them accordingly. Our study thus highlights the enduring importance of human expertise in financial analysis, even in an era of increasing reliance on big data and technology. Simultaneously, our findings indicate that maintaining this relevance requires a transformation in the function and approach of human analysts.

References

- Agrawal A, Gans JS, Goldfarb A (2019) Artificial intelligence: the ambiguous labor market impact of automating prediction. *J. Econom. Perspectives.* 33(2): 31-50.
- Angrist JD, Pischke JS (2008) *Mostly harmless econometrics* (Princeton University Press, Princeton, NJ).
- Barber BM, Huang X, Jorion P, Odean T, Schwarz C (2023) A (sub) penny for your thoughts: Tracking retail investor activity in TAQ. *J. Finance.* Forthcoming.
- Ben-David I, Birru J, Rossi A (2019) Industry familiarity and trading: Evidence from the personal portfolios of industry insiders. *J. Financial Econom.* 132(1): 49-75.
- Bengtzen M (2017) Private investor meetings in public firms: the case for increasing transparency. *Fordham J. Corporate & Financial Law* 22 (1): 33-132.
- Bhattacharya N, Desai H, Venkataraman K (2013) Does earnings quality affect information asymmetry? Evidence from trading costs. *Contemp. Accounting Res.* 30(2): 482-516.
- Bradley D (2023) Financial Analysts. *Handbook of Financial Decision Making.* (Edward Elgar Publishing, Northampton, MA)
- Bradley D, Gokkaya S, Liu X (2017) Before an analyst becomes an analyst: does industry experience matter? *J. Finance* 72(2): 751-792.
- Brown LD, Call AC, Clement MB, Sharp NY (2015) Inside the “black box” of sell-side financial analysts. *J. Accounting Res.* 53(1): 1-47.
- Busse JA, Tong L, Tong Q, Zhang Z (2019) Trading regularity and fund performance. *Rev. Financial Stud.* 32(1): 374-422.
- Cao S, Jiang W, Wang JL, Yang B (2023) From man vs. machine to man + machine: The art and AI of stock analyses. Working paper, National Bureau of Economic Research, Cambridge, MA.
- Clement MB (1999) Analyst forecast accuracy: Do ability, resources, and portfolio complexity matter? *J. Accounting Econom.* 27(3): 285-303.
- Coleman B, Merkley K, Pacelli J (2022) Human versus machine: A comparison of robo-analyst and traditional research analyst investment recommendations. *Accounting Rev.* 97(5): 221-244.
- Dessaint O, Foucault T, Frésard L (2023) Does Alternative Data Improve Financial Forecasting? The horizon effect. *J. Finance.* Forthcoming.
- Frank MR, Autor D, Bessen JE, Brynjolfsson E, Cebrian M, Deming DJ, Feldman M, Groh M, Lobo J, Moro E, Wang D (2019) Toward understanding the impact of artificial intelligence on labor. *Proceedings of the National Academy of Sciences.* 116(14): 6531-9.
- Froot K, Kang N, Ozik G, Sadka R (2017) What do measures of real-time corporate sales say about earnings surprises and post-announcement returns? *J. Financial Econom.* 125(1): 143-162.
- Gerken WC, Painter MO (2023) The value of differing points of view: Evidence from financial analysts’ geographic diversity. *Rev. Financial Stud.* 36(2): 409-449.
- Green TC, Jame R, Markov S, Subasi M (2014) Access to management and the informativeness of analyst research. *J. Financial Econom.* 114(2): 239-255.
- Grennan J, Michaely R (2020) Artificial intelligence and high-skilled work: Evidence from analysts. Working paper, Duke University, Durham, NC.
- Harford J, Jiang F, Wang R, Xie F (2019) Analyst career concerns, effort allocation, and firms’ information environment. *Rev. Financial Stud.* 32(6): 2179-2224.

- Hoberg G, Moon SK (2017) Offshore activities and financial vs. operational hedging. *J. Financial Econom.* 125(2): 217-244.
- Hoberg G, Moon SK (2019) The offshoring return premium. *Management Sci.* 65(6): 2876-2899.
- Hu, G., Jo, K.M., Wang, Y.A. and Xie, J., (2018). Institutional trading and Abel Noser data. *Journal of Corporate Finance*, 52, pp.143-167.
- Huang J (2018) The customer knows best: the investment value of consumer opinions. *J. Financial Econom.* 128(1): 164-182.
- Hwang BH, Liberti JM, Sturgess J (2019) Information sharing and spillovers: evidence from financial analysts. *Management Sci.* 65(8): 3624-3636.
- Jame R (2018) Liquidity provision and the cross section of hedge fund returns. *Management Sci.* 64(7): 3288-3312.
- Jame R, Johnston R, Markov S, Wolfe MC (2016) The value of crowdsourced earnings forecasts. *J. Accounting Res.* 54(4): 1077-1110.
- Kang JK, Stice-Lawrence L, Wong YTF (2021) The firm next door: Using satellite images to study local information advantage. *J. Accounting Res.* 59(2): 713-750.
- Katona Z, Painter M, Patatoukas PN, Zeng J (2023) On the capital market consequences of alternative data: Evidence from outer space. *J. Financial Quant. Anal.* Forthcoming.
- Kothari SP, Leone AJ, Wasley CE (2005) Performance matched discretionary accrual measures. *J. Accounting Res.* 39(1): 163-197.
- Livnat J, Mendenhall RR (2006) Comparing the post-earnings announcement drift for surprises calculated from analyst and time series forecasts. *J. Accounting Res.* 44(1): 77-205.
- McMahon D, Chu K (2012) Clampdown in China on corporate sleuthing. *Wall Street Journal*. Retrieved from <https://www.wsj.com>
- Merkley K, Michaely R, Pacelli J (2017) Does the scope of the sell-side analyst industry matter? An examination of bias, accuracy, and information content of analyst reports. *J. Finance.* 72(3): 1285-1334.
- Puckett A, Yan X (2011) The interim trading skills of institutional investors. *J. Finance.* 66(2): 601-633.
- Seasholes MS, Zhu N (2010). Individual investors and local bias. *J. Finance.* 65(5): 1987-2010.
- Solomon D, Soltes E (2015) What are we meeting for? The consequences of private meetings with investors. *J. Law & Econom.* 58(2): 325-355.
- Soltes E (2014) Private interaction between firm management and sell-side analysts. *J. Accounting Res.* 52(1): 245-272.
- Stulz RM (2007) Hedge funds: Past, present, and future. *J. Econom Perspectives.* 21(2): 175-194.
- Wigglesworth R (2016) Investors mine big data for cutting-edge strategies. *Financial Times*. Retrieved from <https://www.ft.com>
- Wilson WM (2008) An empirical analysis of the decline in the information content of earnings following restatements. *Accounting Rev.* 83(2): 519-548.
- Zhu C (2019) Big data as a governance mechanism. *Rev. Financial Stud.* 32(5): 2021-2061.

Appendix 1. Continued.

Category (1)	Definition (2)		Keywords (3)	Example from Analyst Report (4)
Point of Sale	These data include merchant-level transaction data, product-level purchase data, and pricing data.	1010Data* airbnb + listing compared online prices discount tracker financial rate monitor First Data SpendTrend* footlocker.com footwear scrapes hotel tracker hotel tracking listing monitor MasterCard Advisors* Nielson* online price survey online pricing study our proprietary datasets Point-Of-Sale*	price comparisons price intelligence price monitoring price observations pricing monitor pricing study pricing tracker property listing SG2* spend tracker Standard Media Index* SuperData* vehicle listing web analytics web mining web scraping zillow	<p><i>“CS Proprietary Home Pricing Tracker (median home price trends on an individual store basis across entire store base) shows similar trends in HD/LOW markets.”</i></p> <p>[Issued by CREDIT SUISSE on 02/19/16 for HOME DEPOT]</p>
Satellite Image	Satellite images can be used to track consumer traffic and gauge inventory levels as well as production activities at mines, construction sites, and plants, among other facilities.	Orbital Insight* parking lot fill rate parking lot traffic proprietary satellite data remote sensing	Remote Sensing Metrics* RS Metrics* satellite analysis satellite image traffic analysis	<p><i>“The satellite analysis points to a y/y Q1 parking lot fill rate change of +0.4%, however, the y/y change in fill rate became progressively worse over the quarter. Based on this data and the headwinds faced in Q1 from cash for appliances and weather, we feel comfortable with our Q1E comp of +1%.”</i></p> <p>[Issued by PIPER SANDLER on 05/12/11 for HOME DEPOT]</p>

Appendix 1. Continued.

Category (1)	Definition (2)		Keywords (3)	Example from Analyst Report (4)
Sentiment	These data include social-media feeds and news flow that help gauge consumer sentiment regarding products and services.	brand sentiment CMS Data* consumer sentiment customer rating customer review customer satisfaction rating customer satisfaction trend facebook analysis facebook data facebook like facebook likes facebook post facebook track facebook user guest sentiment instagram data instagram engagement instagram follower Internet World Stats* Investing Analytics* Medicare Plan Finder* Merchant Centric* net sentiment NetBase*	online customer review online review Prosper Insights* ratings on tripadvisor review analytics scoring released by cms sentiment analysis sentiment data social media analysis social media engagement social media follower star(s) rating tracking on twitter tripadvisor ratings twitter analysis twitter data twitter purchase intent twitter sentiment web analytics web mining web scraping yelp	<p><i>“In this report, we introduce our proprietary consumer sentiment analysis, using information from Merchant Centric, a company that works with multi-location brands across consumer and service industries to “help them manage and learn from guests’ online feedback.” For our purposes, Merchant Centric tracks location-specific, user-generated reviews across multiple social media platforms. The reviews are user-generated and tied to a specific location, and then sourced from the following social media sites: Facebook, Google, Yelp, Trip Advisor, Superpages, and CitySearch. Our specific Jefferies data set utilizes reviews on a representative sample of some 500 McDonald’s locations across the country, as well as reviews for just over 2,400 Bojangles, Burger King, Del Taco, Dunkin’ Donuts, Jack in the Box, Sonic, Taco Bell, and Wendy’s units located within the same zip code. We have chosen to exclude independent/local operators, and focus on a sample of national and regional competitors.”</i></p> <p>[Issued by Jefferies on 12/05/17 for MCDONALD’S CORP]</p>
Web Traffic	These data track what users search for on the Internet and how frequently/for how long users visit given websites.	baidu analysis baidu data baidu search data baidu search index baidu search volume ComScore* daily traffic google search analysis google search trend google trend google-searched iphone monitor	iphone tracker Scrapehero* search interest search trend search volume smartphone tracker Thinknum* traffic analysis traffic monitor web hit activity web search	<p><i>“The AlphaWise Smartphone Tracker has been developed by Morgan Stanley’s AlphaWise using multi-country web search analysis using Google Trends. The approach accounts for different search criteria in multiple countries, as well as the differential between search and sales data seasonality, where appropriate. The in-sample period consists of 2008-2011 for Apple and 2010-2012 for Samsung Galaxy.”</i></p> <p>[Issued by MORGAN STANLEY: on 09/18/13 for APPLE]</p>

Appendix 1. Continued.

Category (1)	Definition (2)	Keywords (3)	Example from Analyst Report (4)	
Other	These include alternative data which do not fit cleanly into any of the above categories (e.g., data on senders and recipients of barrels loaded onto vessels).	BuildFax* climatology ClipperData* Collateral Verifications* Dodge* Drillinginfo* Dun & Bradstreet* Edmunds* Entgroup* EPFR* evidence lab macro Flightglobal* formulary coverage home improvement tracker IFI Claims* Innovata*	lower end spending m2m macro-to-micro network traffic lab nowcast One Click Retail* Ookla* OpenSignal* Root Metrics* Rystad* STR Data* wait time monitor Wards Automotive* weather monitor	<p><i>“The most effective way to measure the integrated advantage is our proprietary use of ClipperData, which can track the shipper ID of barrels loading onto vessels in the Gulf of Mexico. For this analysis, we do not isolate the loadings to export barrels, but look at Jones Act activity as well, as the ability to move its crude production anywhere is the proper reflection of the business model’s advantage, in our view.”</i></p> <p>[Issued by WOLFE RESEARCH on 09/28/18 for EXXON MOBIL CORP]</p>

Appendix 2
Variable Description

Variables	Definition
<i>Acc</i>	We first calculate the proportional mean absolute forecast error, <i>PMAFE</i> , as the difference between the absolute forecast error of an analyst's annual earnings forecast and the average absolute forecast error across all analysts, scaled by the average absolute forecast error. When we compute the average absolute forecast error, we exclude all analysts who (also) report drawing from alternative data in their coverage of the corresponding firm in the respective forecast period. Since negative (positive) values of <i>PMAFE</i> indicate above (below) average performance, <i>Acc</i> is defined as $PMAFE \times (-1)$.
<i>I(Alternative Data)</i>	An indicator variable that equals one if the corresponding analyst issues an earnings forecast explicitly supported by alternative data and zero otherwise.
<i>Forecast Age</i>	The logarithm of one plus the number of calendar days between the forecast date and the corresponding earnings report date.
<i>Analyst/Firm Experience</i>	The number of years since the corresponding analyst first issued a forecast for the corresponding firm.
<i>Analyst Experience</i>	The number of years since the corresponding analyst first issued a forecast for any firm in the IBES database.
<i>#Firms Covered</i>	The logarithm of one plus the number of firms the corresponding analyst covers in the corresponding year.
<i>Forecast Frequency</i>	The logarithm of one plus the number of forecasts made by the corresponding analyst in the corresponding year.
<i>Broker Size</i>	The number of analysts working at the corresponding analyst's broker in the year of the forecast.
<i>Trading Commissions</i>	For each written report, we consider all trades involving the stock discussed in the report within three months after the report was issued and aggregate the total dollar value of the corresponding commissions paid to the corresponding broker.
<i>I(In-House Data Science Team)</i>	An indicator variable that equals one if the corresponding analyst works for a broker that has an in-house data-science team.
\sum <i>Colleagues</i> <small><i>Alternative Data</i></small>	The number of analysts that draw from alternative data and work for the same broker in the same city as the corresponding analyst during the previous year.
<i>Number of 8-Ks</i>	The total number of Form 8-Ks filed during the previous annual forecast period.
<i>Return Volatility</i>	The standard deviation of daily stock returns during the previous annual forecast period.
<i>Earnings Surprise</i>	The most recent earnings surprise in the previous forecast period based on diluted earnings per share, excluding extraordinary items, and applying a seasonal random walk (Livnat and Mendenhall, 2006).
<i>I(Earnings Restatement)</i>	An indicator variable that equals one if a given firm has issued restatements in the past.

Appendix 2. Continued.

<i>Discretionary Accruals</i>	The most recent discretionary accruals based on the Modified Jones model matched to another from the same industry and year with the closest ROA (Kothari, Leone, and Wasley, 2005).
<i>I(Lack of Preferential Access to Management)</i>	An indicator variable that equals one if the corresponding firm did not participate in a conference hosted by the corresponding analyst's broker over the previous year.
<i>Size</i>	The market capitalization of the corresponding firm at the end of the previous fiscal year in billions.
<i>M/B</i>	The market value of equity divided by the book value of equity at the end of the previous fiscal year.
<i>Momentum</i>	Buy-and-hold return on the corresponding stock over the previous six months.
<i>I(Category = X)</i>	An indicator variable that equals one if a given analyst explicitly references the use of alternative data from alternative-data category X.
\sum <i>Categories</i>	The number of distinct alternative-data categories an analyst references in her report.
<i>I(Source = Proprietary Data)</i>	An indicator variable that equals one if a given analyst explicitly references the use of proprietary alternative data.
<i>I(Source = Accessible Data)</i>	An indicator variable that equals one if a given analyst explicitly references the use of accessible alternative data.

Figure 1

This figure displays two timelines, indicating when—for our sample of firms in the Dow Jones Industrial Average Index—we observe the first analyst report, or, the first one hundred analyst reports, explicitly referencing the use of alternative data from a particular alternative-data category. We describe our alternative data categories in Subsection 2.3.

First Time an Analyst Report In Our Sample Mentions the Use of Alternative Data from a Particular Category

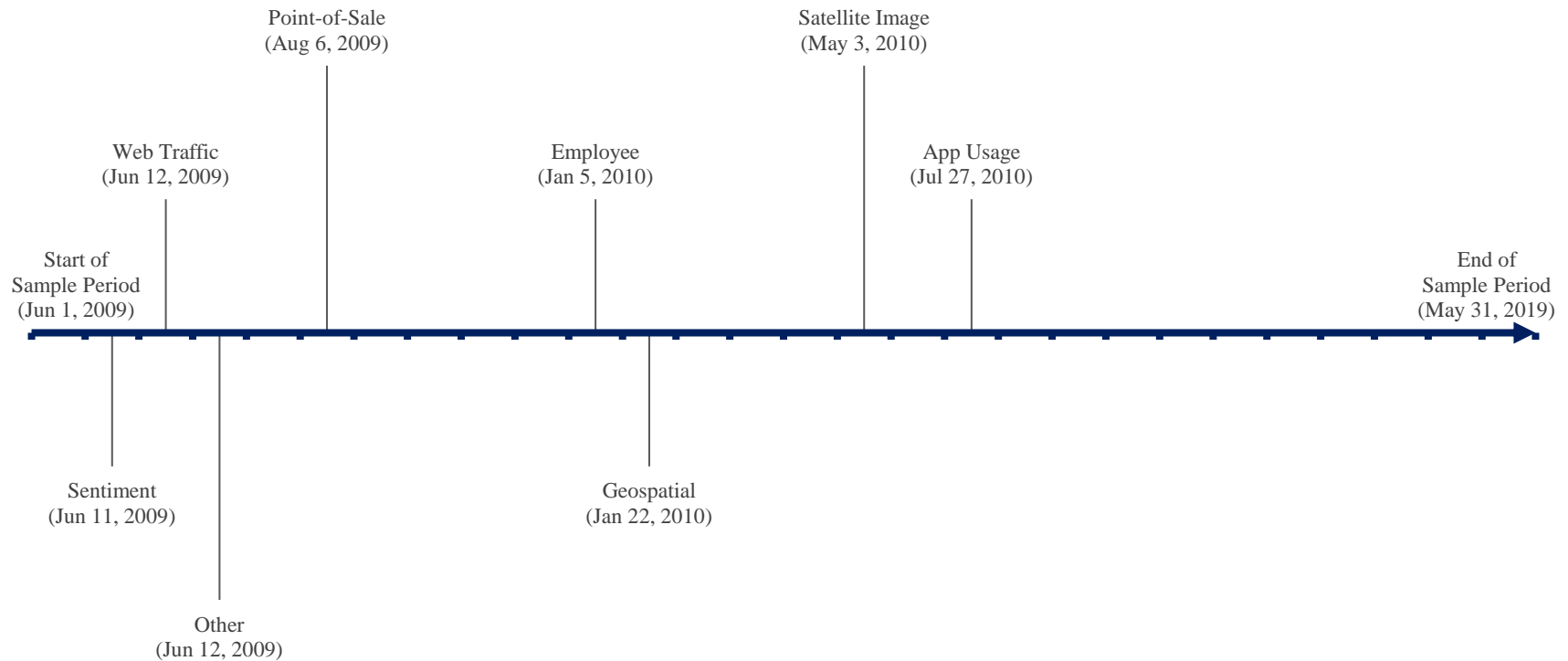


Figure 1. Continued.

First Time More Than 100 Analyst Reports In Our Sample Mention the Use of Alternative Data from a Particular Category

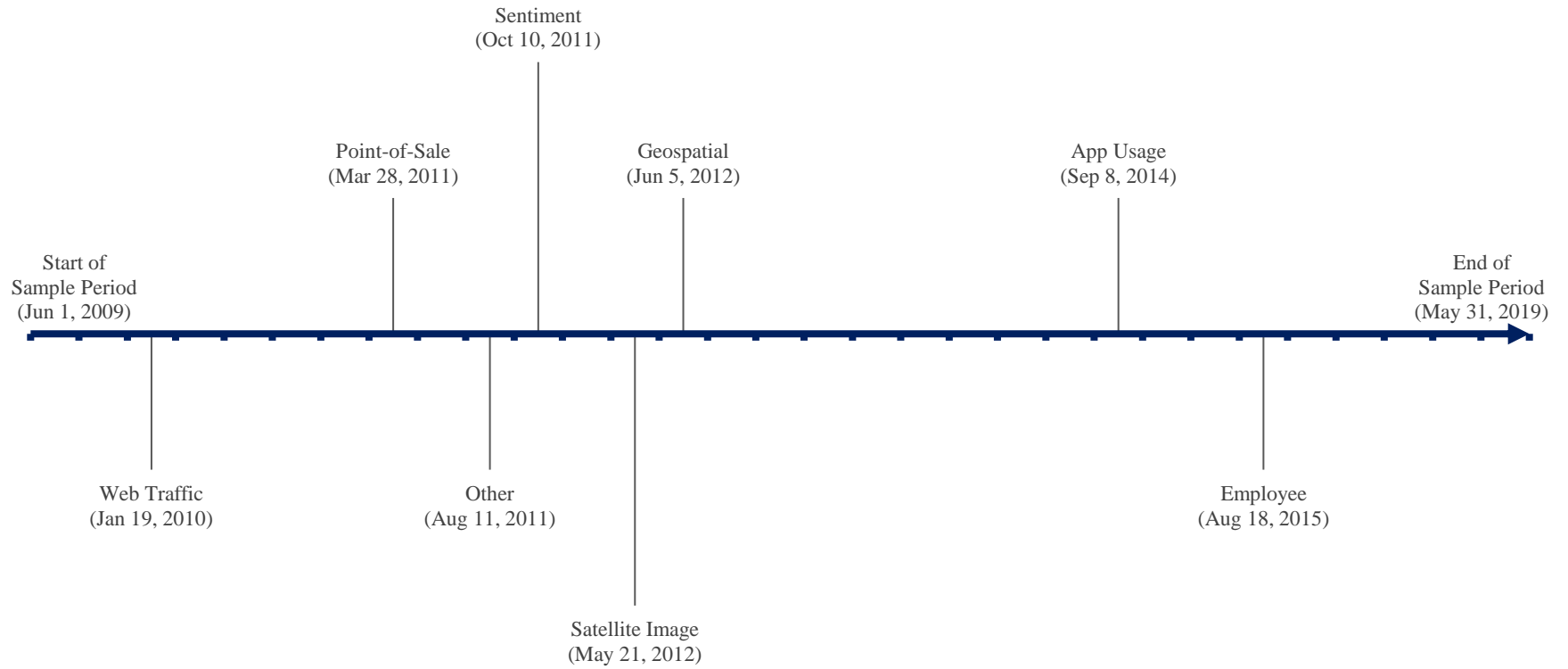


Table 1
Numbers and Fractions of Analyst Forecasts Explicitly Supported by Alternative Data

In this table, we present the numbers and fractions of analyst forecasts explicitly supported (not explicitly supported) by alternative data. Our sample contains all Dow Jones Industrial Average Index firms from June 1, 2009, through May 31, 2019. We combine the years 2009 and 2010 and the years 2018 and 2019 as we have only partial data for the years 2009 and 2019.

	Number of Forecasts ...		Fraction of Forecasts...
	... Explicitly Supported by Alternative Data	... Not Explicitly Supported by Alternative Data	... Explicitly Supported by Alternative Data
<i>Panel A: By Year</i>			
2009/2010	515	7,616	6%
2011	615	6,239	9%
2012	488	6,769	7%
2013	490	6,348	7%
2014	497	6,058	8%
2015	694	5,998	10%
2016	729	5,691	11%
2017	659	5,444	11%
2018/2019	952	8,216	10%
2009 - 2019	5,639	58,379	9%
<i>Panel B: By Industry Sector</i>			
Energy	40	2,512	2%
Materials	29	2,596	1%
Industrials	580	8,398	6%
Consumer Discretionary	443	4,194	10%
Consumer Staples	841	7,199	10%
Health Care	661	7,487	8%
Financials	103	8,082	1%
Information Technology	2,513	13,597	16%
Communication Services	429	4,314	9%

Table 2
 Numbers of Analyst Forecasts Explicitly Supported by Data from a Particular Category

In this table, we present the numbers and fractions of analyst forecasts explicitly supported by data from a particular alternative-data category. We describe our alternative data categories in Subsection 2.3. Since a given analyst report may draw from multiple alternative data categories, the sum of the numbers of forecasts reported in Table 2 exceeds the total number of forecasts explicitly supported by alternative data reported in Table 1; the fractions do not add up to 100% for the same reason.

Alternative Data Category	Number [Fractions] of Forecasts Explicitly Supported by Alternative Data
App Usage	476 [8%]
Employee	543 [10%]
Geospatial	257 [5%]
Point of Sale	1,080 [19%]
Satellite Image	171 [3%]
Sentiment	1,062 [19%]
Web Traffic	1,944 [34%]
Other	1,322 [23%]

Table 3
Alternative Data and Forecast Accuracy

In this table we report coefficient estimates derived from regressions of forecast accuracy on whether an analyst explicitly references the use of alternative data in her corresponding written report. The observations are at the analyst/firm/forecast-date level. To construct the dependent variable, *Acc*, we first compute *PMAFE* as the difference between the absolute forecast error of an analyst and the average absolute forecast error across all analysts not referencing the use of alternative data, scaled by the average absolute forecast error. Since negative (positive) values of *PMAFE* indicate above (below) average performance, we define *Acc* as $PMAFE \times (-1)$. $I(\text{Alternative Data})$ equals one if the corresponding analyst's earnings forecast is explicitly supported by alternative data as described in Subsection 2.2 and zero otherwise. We define all remaining variables in Appendix 2. We report *t*-statistics in parentheses. We double-cluster our standard errors at the analyst- and year-month levels. *, **, and *** denote statistical significance at the 10%, 5%, and 1% levels, respectively.

	(1)
<i>I(Alternative Data)</i>	0.214*** (6.30)
<i>Forecast Age</i>	-0.251*** (-12.34)
<i>Analyst/Firm Experience</i>	0.060*** (2.72)
<i>Analyst Experience</i>	0.056 (1.00)
<i>#Firms Covered</i>	0.029 (0.52)
<i>Forecast Frequency</i>	0.033 (1.11)
<i>Broker Size</i>	-0.000 (-0.85)
Analyst-Firm Fixed Effects	Yes
Firm-Year Fixed Effects	Yes
<i>N</i>	64,018
Adjusted R^2	0.238

Table 4
Alternative Data and Trading Commissions

In this table we report coefficient estimates derived from regressions of trading commissions paid to a brokerage on whether an analyst explicitly references the use of alternative data in her corresponding written report. The observations are at the brokerage/firm/forecast-date level. We consider trades and the associated trading commissions in the ANcerno database. As we are interested in how the allocation of trades and commissions are affected by the dissemination of alternative-data insights, we focus on transactions that occur within three months after the issuance of an analyst report and which involve the stock discussed in the report. For each analyst report, we compute the total commissions earned by the brokerage issuing the report and use this as our dependent variable. $I(\text{Alternative Data})$ equals one if the corresponding reports explicitly describe the use of alternative data as described in Subsection 2.2 and zero otherwise. We define all remaining variables in Appendix 2. We report t -statistics in parentheses. We double-cluster our standard errors at the broker- and year-month levels. *, **, and *** denote statistical significance at the 10%, 5%, and 1% levels, respectively.

	(1)
<i>I(Alternative Data)</i>	11,858.82* (1.82)
<i>Forecast Age</i>	-3,393.30 (-1.17)
<i>Analyst/Firm Experience</i>	1,416.47 (0.63)
<i>Analyst Experience</i>	1,233.77* (1.88)
<i>#Firms Covered</i>	5,091.09 (0.22)
<i>Forecast Frequency</i>	-2,909.48 (-0.33)
<i>Broker Size</i>	-720.16*** (-4.23)
Broker-Firm Fixed Effects	Yes
Firm-Year Fixed Effects	Yes
<i>N</i>	4,757
Adjusted R^2	0.589

Table 5
Alternative Data and Portfolio Returns

In this table we report the performance of transaction-based buy-and-sell calendar-time portfolios. We consider trades in the ANcerno database. As we are interested in trades affected by the dissemination of alternative-data insights, we focus on transactions that occur within three months after the issuance of an analyst report and that involve the stock discussed in the report. Applying Jame’s (2018) methodology, we separate trades into those made by hedge funds and those made by non-hedge funds. We further separate non-hedge-fund trades by whether they are preceded by an analyst report that adopts alternative data or whether they are placed following reports that do not draw from alternative data. We assess the performance of trades by creating transaction-based calendar-time portfolios, following the approach of Seasholes and Zhu (2010) and Ben-David, Birru and Rossi (2019). Stocks purchased (sold) on day t are added to the buy (sell) portfolio at the beginning of day $t + 2$. We assume a holding period of three months. We compute DGTW-adjusted returns for each stock/day (= return on a stock minus the value-weighted return on a portfolio of stocks in the same size, book-to-market, and momentum quintiles) and construct value-weighted portfolio returns. We report the time-series average of the buy-minus-sell portfolio returns, whereby each daily portfolio return is weighted by the number of trades contributing to the portfolio on that specific day. We also report the annualized difference. *, **, and *** denote statistical significance at the 10%, 5%, and 1% levels, respectively.

	Buy – Sell [Daily] (1)	Buy – Sell [Annual] (2)
Hedge funds	0.04% (1.54)	10.79%
Non-Hedge Funds without Alternative Data	0.02% ** (2.19)	3.85%
Non-Hedge Funds with Alternative Data	0.04% (1.27)	9.85%

Table 6
Variation in the Use of Alternative Data

In this table we report coefficient estimates derived from regressions of alternative-data adoption on various analyst- and firm-level characteristics. The observations are at the analyst/firm/forecast-date level. In column (1) we report the results derived from a probit model, while in column (2) we report the results derived from a linear probability model to allow for consistent parameter estimation while including fixed effects. The dependent variable equals one if the analyst's earnings forecast is explicitly supported by alternative data and zero otherwise. $I(\text{In-House Data Science Team})$ equals one if the analyst's brokerage has an in-house data-science team. $\sum \text{Colleagues}_{\text{Alternative Data}}$ is the number of colleagues working in the same city as the corresponding analyst and drawing from alternative data for one of the firms they cover during the previous year. *Number of 8-Ks* is the total number of Form 8-Ks filed during the previous annual forecast period. *Return Volatility* is the standard deviation of daily stock returns during the previous annual forecast period. *Earnings Surprise* is measured as in Livnat and Mendenhall (2006). $I(\text{Earnings Restatement})$ equals one if the corresponding firm has had to restate its financial accounts. We compute *Discretionary Accruals* as in Kothari, Leone, and Wasley (2005). $I(\text{Lack of Preferential Access to Management})$ equals one if the corresponding firm did not participate in a conference hosted by the corresponding analyst's broker over the previous year. To facilitate a comparison of the coefficient estimates, we convert $\sum \text{Colleagues}_{\text{Alternative Data}}$, *Number of 8-Ks*, *Return Volatility*, *Earnings Surprise*, and *Discretionary Accruals* into quintile-rank variables, ranging from one if the corresponding realization is in the bottom quintile of its distribution to five if the corresponding realization is in the top quintile. We define all remaining variables in Appendix 2. We report z-statistics in parentheses. We double-cluster our standard errors at the analyst- and the year-month levels. *, **, and *** denote statistical significance at the 10%, 5%, and 1% levels, respectively.

	(1)	(2)
<i>I(In-House Data Science Team)</i>	0.411*** (3.14)	0.088*** (2.77)
<i>Rank</i> ($\sum \text{Colleagues}_{\text{Alternative Data}}$)	0.123*** (3.54)	0.015** (2.58)
<i>Rank</i> (<i>Number of 8-Ks</i>)	-0.105*** (-3.18)	-0.014** (-2.24)
<i>Rank</i> (<i>Return Volatility</i>)	0.134*** (4.18)	0.013*** (2.63)
<i>Rank</i> (<i>Earnings Surprise</i>)	0.060*** (2.60)	0.010*** (3.07)
<i>I(Earnings Restatement)</i>	0.281*** (2.71)	0.070*** (3.16)
<i>Rank</i> (<i>Discretionary Accruals</i>)	0.053** (2.05)	0.011*** (2.95)
<i>I(Lack of Preferential Access to Management)</i>	0.188* (1.76)	0.032** (2.03)
<i>Analyst/Firm Experience</i>	-0.004 (-0.53)	-0.001 (-0.55)
<i>Analyst Experience</i>	0.004 (0.63)	0.001 (1.02)
<i>#Firms Covered</i>	0.134 (0.87)	-0.007 (-0.32)
<i>Forecast Frequency</i>	-0.025 (-0.33)	0.006 (0.62)
<i>Broker Size</i>	0.001 (0.52)	0.000 (1.13)

Table 6. Continued.

	(1)	(2)
<i>Size</i>	0.191*** (2.96)	0.044*** (3.12)
<i>M/B</i>	0.011 (1.03)	0.001 (0.53)
<i>Momentum</i>	0.284* (1.70)	0.040 (1.46)
Analyst-Firm Fixed Effects	No	No
Firm-Year Fixed Effects	No	No
Industry Fixed Effects	No	Yes
<i>N</i>	55,464	55,464
Pseudo R^2	0.130	0.133

Table 7
Differences in the Usefulness by Alternative Data Types

In this table we report results derived from analyses similar to those associated with Table 3. For column (1), we replace $I(\text{Alternative Data})$ with $I(\text{Category} = X)$, which equals one if the analyst explicitly references the use of alternative data from alternative-data category X. For column (2), we replace $I(\text{Alternative Data})$ with $\sum \text{Categories}$, which is the number of distinct alternative data categories an analyst references in her report. For column (3), we replace $I(\text{Alternative Data})$ with $I(\text{Source} = \text{Proprietary Data})$ and $I(\text{Source} = \text{Accessible Data})$, which equals one if the analyst explicitly references the use of proprietary and accessible alternative data as described in Subsection 4.3.2, respectively. We report t -statistics in parentheses. We double-cluster our standard errors at the analyst- and year-month levels. *, **, and *** denote statistical significance at the 10%, 5%, and 1% levels, respectively.

	(1)	(2)	(3)	(4)
$I(\text{Category} = \text{App Usage})$	0.310*** (4.04)			
$I(\text{Category} = \text{Sentiment})$	0.230*** (2.88)			
$I(\text{Category} = \text{Employee})$	0.212*** (3.09)			
$I(\text{Category} = \text{Geospatial})$	-0.033 (-0.31)			
$I(\text{Category} = \text{Point of Sale})$	0.183*** (3.44)			
$I(\text{Category} = \text{Satellite Image})$	0.053 (0.53)			
$I(\text{Category} = \text{Web Traffic})$	0.137** (2.09)			
$I(\text{Category} = \text{Others})$	0.187*** (3.88)			
$\sum \text{Categories}$		0.175* (1.75)		
$I(\text{Source} = \text{Proprietary Data})$			0.243*** (6.82)	
$I(\text{Source} = \text{Accessible Data})$				0.188*** (4.07)
<i>Forecast Age</i>	-0.250*** (-12.36)	-0.181*** (-3.67)	-0.251*** (-12.19)	-0.252*** (-12.31)
<i>Analyst/Firm Experience</i>	0.060*** (2.68)	-0.017 (-1.24)	0.059** (2.60)	0.061*** (2.79)
<i>Analyst Experience</i>	0.055 (0.98)	0.378 (1.50)	0.059 (1.03)	0.057 (1.02)
<i>#Firms Covered</i>	0.030 (0.54)	-0.295 (-1.10)	0.029 (0.53)	0.034 (0.60)
<i>Forecast Frequency</i>	0.033 (1.09)	-0.101 (-0.96)	0.033 (1.08)	0.029 (0.96)
<i>Broker Size</i>	-0.000 (-0.91)	0.004 (1.46)	-0.000 (-0.84)	-0.000 (-0.86)

Table 7. Continued.

	(1)	(2)	(3)	(4)
Analyst-Firm Fixed Effects	Yes	Yes	Yes	Yes
Firm-Year Fixed Effects	Yes	Yes	Yes	Yes
<i>N</i>	64,018	5,639	64,018	64,018
Adjusted R^2	0.239	0.406	0.237	0.236