

# November 1, 2021: Status of *The Economist* Leader-Sentiment-Earnings Project

Joshua Y. Levy \*

November 1, 2021

---

\*Joshua Y. Levy (joshua.levy@chicagobooth.edu), Stigler Center, Booth School of Business, University of Chicago

# 1 Introduction to this Document

Since the beginning of August, I have been working on a project that tries to identify the role business-facing publications in the “revolving door” phenomenon between government and the private sector. In an effort to summarize the results of this project to date, I have produced this document which does the following. First it outlines the motivating argument underlying the project, and identifies how the project has evolved over the past three months. Subsequently, it identifies some key results that have driven the methodological evolution of the project and the challenges associated with those various methods. Finally, in this document, I make the case that this project — as it is currently conceived — likely cannot be completed accurately with any degree of assurance that the underlying results would be 1) accurate; 2) replicable; or 3) externally valid.

## 2 Hypothesis and Argument

Modern democracies exhibit strong “revolving door” phenomena whereby government officials often leave their posts for more lucrative opportunities in the private sector. Taken to its natural extension this phenomena raises the question: what do world leaders do after they leave office?

Anecdotally, there is plenty of evidence of former presidents and prime ministers establishing NGOs, joining boards of private corporations, becoming consultants, writing books, and racking up large fees on the speaking circuit. A subsequent question arises: why can some leaders make more than others?

In the revolving door setting, we suppose that business-facing publications (which cater to the business elite) serve as an important intermediary. Publications like *The Economist*, *Financial Times*, and *Wall Street Journal* provide informative signals to business leaders as to the pro-business *bona fides* of world leaders. In turn, businesses are willing to pay more in consulting fees and board compensation to former world leaders they perceive as being more business-friendly or more amenable to the objectives of the corporation.

We consider the former leaders of only substantively democratic nations for two reasons. First, substantive democracies tend to be more transparent — suggesting that the press has a bigger role in those nations relative to autocratic ones. Indeed, in many settings, freedom of the press is considered one of the necessary constituent parts of substantive democracy. Additionally, when considering former leaders of autocratic or authoritarian nations, businesses may be inclined to access or retain former leaders not because of the supposed competence of the leader in question as conveyed by the press, but because of the perception that doing so is either necessary or advantageous for practicing some kind of (illicit) corruption. For the question at hand, it would be difficult to isolate the effect of this “illicit corruption” channel from the one of interest — the role of business-facing news publications.

**HYPOTHESIS:** We suppose that the sentiment expressed by *The Economist* is a predictor of the post-term (out-of-office) earnings of former democratic world leaders. That is, the more positively *The Economist* writes about a leader, the more that leader is able to command (or the more businesses are willing to pay) in compensation for out-of-office services.

### 3 Sentiment Analysis Tools

After data collection, to date, work on this project has primarily focussed on conducting sentiment analysis on *The Economist*. This work can be divided into three periods generally defined by the sentiment analysis tools used: SentimentR, VADER, and BART.

#### 3.1 SentimentR

The first tool that I used to evaluate sentiment is called SentimentR. SentimentR relies on a valence-shifter-aware dictionary-based approach to sentiment analysis. This is one degree more sophisticated than a naive dictionary approach which takes a (weighted) average of pre-assigned “polarity scores”

for the words in a sentence.<sup>1</sup> Valence shifters are important because they can have significant effects on the expressed sentiment. For instance, the word **not** is a “negator” that “negates” the sign on the sentiment score assigned to a sentence like: **I do not like ice cream**. Other classes of valence shifters include amplifiers, de-amplifiers, and adversative conjunctions.

Importantly, SentimentR scores text at the sentence-level. Even an article that exclusively profiles a single leader is unlikely to identify that leader in every sentence. Instead, for the purposes of style, copy writers might use:

- Pronouns and titles: For example **She** instead of **Angela Merkel**;
- Metonymy: For example **the Blue House** instead of **Moon Jae-In**;<sup>2</sup> or
- Synecdoche: For example **Washington** instead of **the administration of Donald Trump**.<sup>3</sup>

Important, sentiment-bearing information may be included in such sentences. In order to address this problem I use a Natural Language Processing (NLP) technique called coreference resolution (more details below). This process is intended to maximize the number of sentences that could be tagged as a “leader-sentence” (i.e. a sentence that implicitly or explicitly mentions a leader of interest.).

The process of segregating leader-sentences from “ambient-sentences” is important for extracting sentiment expressed by *The Economist* about the leader in particular as opposed to the sentiment expressed by *The Economist* about the article in general. This becomes particularly acute in the cases of an “event- ” or “disaster- ” bias (more details in Section 5). In addition to computing sentence-level sentiment, I aggregate to the “entity-level” to have two numbers statistic that represent the sentiment expressed by an article about the leader of interest, and the “ambient” sentiment expressed by the article. SentimentR returns scores ranging from  $-1$ , most negative to  $1$ , most positive.

---

<sup>1</sup>These dictionaries often include common idioms that have scores assigned at the idiom-level rather than the word-level.

<sup>2</sup>Metonymy is the linguistic device of invoking one thing to refer to something related.

<sup>3</sup>Synecdoche is subset of metonymy, often using a part a part of a thing as a stand-in for the broader concept.

### 3.1.1 AllenNLP Coreference Resolution

Coreference resolution is an NLP technique for identifying repeated expressions or references of a single entity in a large (i.e. multi-sentence) body of text. This is often an important step in higher-level NLP tasks such as text summarization, question answer, and other reading comprehension problems.

In our case, we use it narrowly for the purposes of cluster and span replacement. The particular model that I use, SpanBERT, identifies “clusters,” or entities, within a body of text and associates “spans,” or text-segments, with those clusters. If an identified cluster matches the leader of interest, I can replace all of the identified spans with the leader’s name to increase the number of leader-sentences available for sentiment analysis.

## 3.2 VADER

Valence Aware Dictionary and sEntiment Reasoner (VADER) is similar in approach to SentimentR. This was additionally used as a sort of robustness check to ensure that sentiment analysis results produced by SentimentR were consistent.

Because VADER also evaluates text at the sentence-level, the same coreference resolution process as described above is used prior to evaluation. Following sentence-level evaluation, the same entity-level aggregation process is used to determine the average leader sentiment and average ambient sentiment. Additionally, VADER also returns scores ranging from  $-1$  to  $1$ .

## 3.3 BART

Finally, I used a pre-trained BART model for sentiment analysis. Unlike SentimentR and VADER, BART represents a machine-learning approach. It is of critical importance to note that BART was not trained for sentiment analysis, but instead trained for text-classification.

That is, in being trained on a hand-labelled corpus of text, BART develops a (very) high dimensional representation of the word embeddings that constitutes those labels. The high-dimensionality of this space permits for implicit “topics” or “clusters” to develop even for concepts that were not in the labelled training set. In prediction, BART attempts to identify the cluster that is most closely associated with the text of interest.

For our purposes, I attempt to “classify” text segments as being associated with the candidate labels/clusters: “Positive,” “Neutral,” and “Negative.” Notably these topics do not appear in the set of labels that are used for training purposes. BART returns a three-tuple of scores on the unit interval that represents estimated “closeness” to a cluster.

Notably, BART is not limited to evaluating text segments at the sentence level. Consequently I perform evaluation at the article-, paragraph- and sentence-level. Results are discussed below.

## 4 Various Intermediary Results Detailed

Below follows some results from each stage of analysis that demonstrate each approach’s strengths and shortcomings, as well as the logic for how the use of each tool developed.

### 4.1 SentimentR

Below, Figure 1, is an example of SentimentR output based on articles that mention prime ministers of Australia. The imposed lines are lowess-smoothed lines of entity-level sentiment at each date of observation (issue publication date). Each point in the scatter plot is a sentence-level observation. Figure 1 illustrates two challenges. First, note that on average, SentimentR evaluates *The Economist* as have almost perfectly neutral sentiment about most all leaders over time. The absence of heterogeneity is concerning for two reasons — first that it is unlikely the case that *The Economist* is truly neutral about these leaders given that it has a well-publicized *ex ante* ideological position, and second that the absence of heterogeneity will make any regression analysis particularly sensitive to

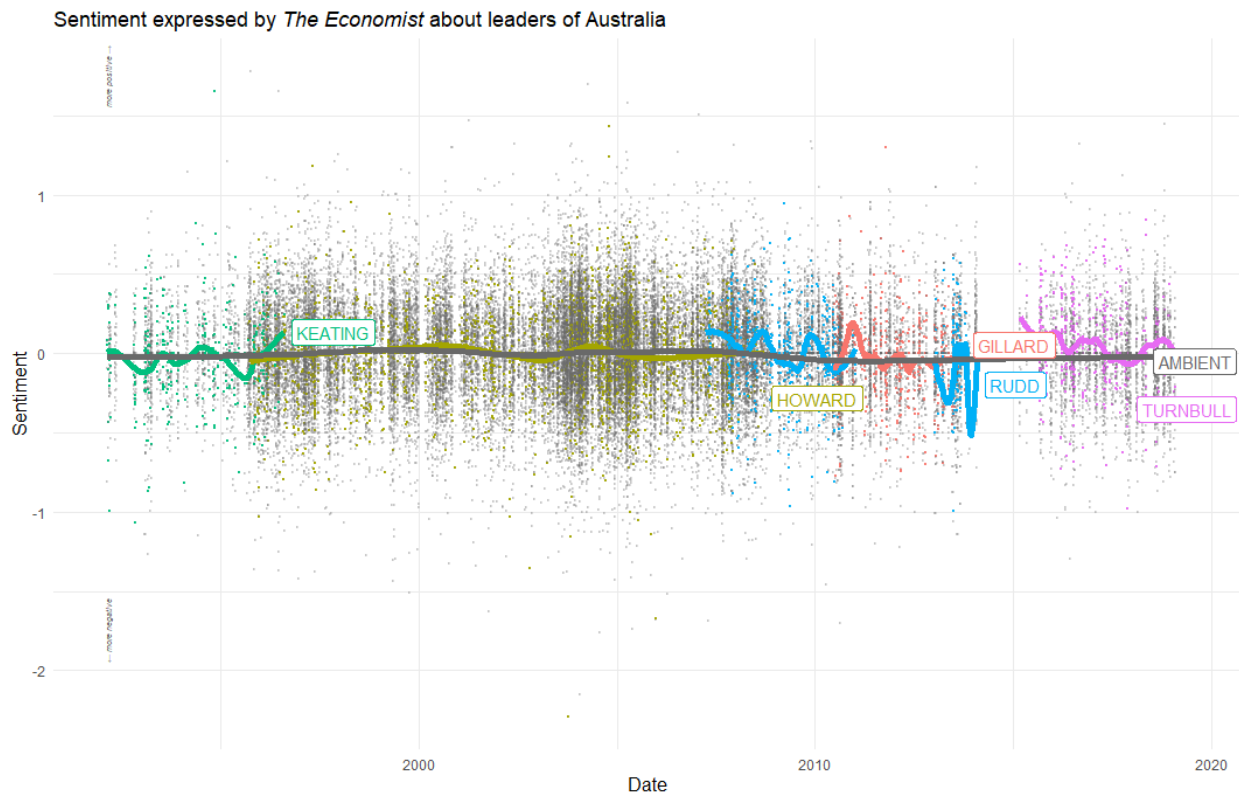


Figure 1: Sentiment expressed by *The Economist* about Australian prime ministers

outliers.<sup>4</sup>

## 4.2 VADER

In light of the poor SentimentR results we turned to VADER to see if the valence-aware dictionary-based approach was worthwhile. In effect, I was testing to see if VADER and SentimentR produced different scores for the same sentences.

Below, Figure 2 plots the correlation between VADER and SentimentR scores for selected Australian leaders. One of the features of VADER over SentimentR is that it preserves extreme values more often than does SentimentR. However, this property is symmetric at the sentence-level and entity-level observations are similarly clustered around 0 (perfectly neutral.)

A closer inspection of SentimentR and VADER results suggests that these tools are inadequate for

---

<sup>4</sup>Reading *The Economist* is evidence enough that these very neutral scores are unrealistic and/or inaccurate.

the task at hand. Consider, for example, the following sentence from the July 7th, 2020, edition of the economist:

Brazil’s president Jair Bolsonaro, who has downplayed the threat of COVID-19, flouted social-distancing guidelines and said that his “athletic history” would protect him, tested positive for the coronavirus.

SentimentR scores this sentence positively (0.426), largely due to the presence of the word “positive,” when this sentence — to a human reader — clearly does not express a positive sentiment.<sup>5</sup> SentimentR and VADER lack knowledge of domain-specific idioms that can heavily weight the sentiment expressed in a sentence to the extremes. Because *The Economist* writes about a wide variety of domains with the expectation that the reader is “an educated layman,” SentimentR and VADER are inappropriate.

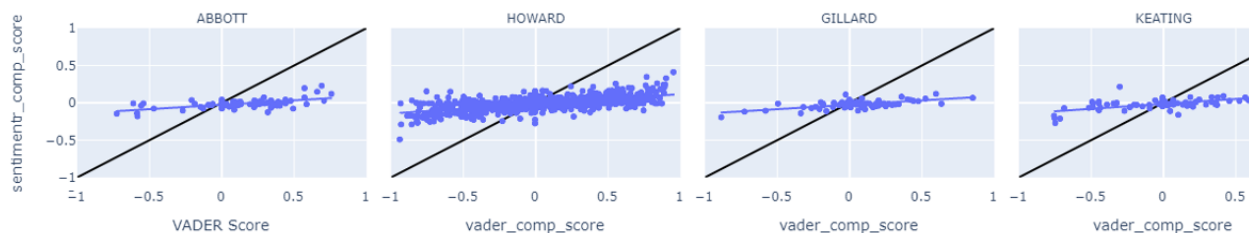
### 4.3 BART

Consequently we turned to BART, which is considerably more computationally intensive than are SentimentR and VADER, in the hope that more advanced techniques would produce results similar to an expert human reader.

As noted above, at the article-, paragraph-, and sentence-level, BART returns a three-tuple of scores on the unit interval indicating the probability that the text segment of interest is closest to either “Positive,” “Neutral,” or “Negative.” To test whether BART was returning accurate sentiment evaluations, I assigned members of the research group 20 articles to score by hand so that results

---

<sup>5</sup>Even without context for what COVID-19 is, a human reader can identify the word “flout” as having a negative connotation. Additionally, in the medical context, testing “positive” for any kind of disease carries a negative connotation.





could be compared. Readers were asked to mark each article as being “Very negative” “Negative,” “Neutral,” “Positive,” or “Very Positive.”<sup>6</sup> For purposes of comparison, I coerced BART scores into similar categories.

Below, Table 1 contains pairwise Cohen’s Kappa statistics that measure the degree of agreement between readers and/or BART. Additionally, because of the difficulty of interpreting the Kappa statistic, table (INSERT TABLE HERE) outlines generally agreed-upon thresholds for what each statistic means.

	BART FIVE	LUIGI	TANO	FABIO	FEDRICO	RACHEL	UTSAV	KHWAJA	JOSHUA	RP	AVG. (FIVE)
BART FIVE	1.0000				0.1028	0.2105	0.1638	0.1419	0.0841		0.0315
LUIGI		1.0000									
TANO			1.0000								
FABIO				1.0000							
FEDRICO					1.0000	0.5238	0.1696	0.3506	0.2809		0.6601
RACHEL						1.0000	0.2105	1.0000	1.0000		0.7368
UTSAV							1.0000	0.3728	0.2715		0.4946
KHWAJA								1.0000	0.6689		0.6564
JOSHUA									1.0000		0.5235
RP AVG. (FIVE)											1.0000

Table 1: Cohen’s Kappa Statistic ( $\kappa$ ) – Inter-rater agreement

Kappa value ( $\kappa$ )	Interpretation
<0	Less than chance agreement
0-0.2	Slight agreement
0.2-0.4	Fair agreement
0.4-0.6	Moderate agreement
0.6-0.8	Substantial agreement
0.8-1	Almost perfect agreement

Table 2: Cohen’s Kappa Statistic ( $\kappa$ ) – Inter-rater agreement Coefficient Interpretation

These results are discussed at greater length below.

## 5 Challenges

As evidenced by the inter-rater agreement statistics presented in Table 1, evaluating the sentiment of article published by *The Economist* is difficult, even for human readers. That is, because the sentiment expressed by *The Economist* is indicated by subtle cues, it is hard for humans to regularly

<sup>6</sup>In some instances sentence- and paragraph-level text segments were also scored in the same fashion.

and consistently agree upon the sentiment expressed about a leader.

In particular, the purpose of this project is to identify something that is very abstract and very diffuse: the sentiment of *The Economist's* evaluation of a leader. These layers of abstraction require an extraordinary amount of context about the issue domain that is the subject matter of a given article, *The Economist's* regular editorial position, any deviation that the *The Economist* might take from its regular position, domain-specific idioms, and a familiarity with the authorial voice and style used by *The Economist*.

Individually, these problems might be tractable. Taken together however, they pose a challenge that I have identified in five parts. The remainder of this section will focus on interpreting BART results (even if the challenges identified also apply to the SentimentR/VADER approach).

## 5.1 Higher-resolution segment (un)certainty

<sup>7</sup> A *prima facie* inspection of BART results might be encouraging. At the article-level, the highest scores (used to identify the “categorical” score for a text segment) tend to reflect the general tone of the article. Additionally, these scores are often returned with a considerable degree of certainty. That is, BART evaluates an article as being considerably more likely to belong to a single cluster than to the other two.<sup>8</sup>

However, when evaluating text segments at the paragraph- or sentence-level, BART scores moderate across the three categories (i.e. BART is much less sure of which cluster a paragraph or sentence might belong to) or tend to the extremes without obvious (to humans) reason. With respect to the first problem, consider a sentence in an article about Emmanuel Macron published on September 30, 2017:

Having had no seats, Alternative for Germany, a disruptive and polarising force, is now

---

<sup>7</sup>For ease of reading, this section might best be read with `BART_SCORING_SMOOTH_COLORS.html` open.

<sup>8</sup>A notable exception to this is the article about Spanish Prime Minister Jose Luis Rodriguez Zapatero, about which BART thinks the article is either quite positive or quite negative be definitely not neutral.

the Bundestag’s third largest party.

Because BART lacks domain specific context and knowledge about *The Economist’s* position on the AfD, BART is very uncertain about which category this sentence belongs to. By contrast, a human reader would likely mark this sentence as “negative” “very negative.”

On the other hand, in an article published in 2002 about European political shifts, the sentence “What is true of Germany is true of the rest of Europe” is scored very confidently positively.

A desirable property of a good sentiment analysis tool should be consistency, something that BART lacks.

## **5.2 Affect bias**

## **5.3 Disaster/event bias**

## **5.4 Leader reference ambiguity**

## **5.5 Leader-policy-colleague metonymic ambiguity**

# **6 The Case to End this (version of the) Project**