# Sinclair Broadcasting and Racial Resentment: Evidence from Google Trends

April 9, 2021

# Aknowledgements

# Abstract

# Contents

*This paper uses Google search terms containing offensive language as a proxy measure for racial animus. Following the practice of similar works (Stephens-Davidowitz, 2014; Chae, Clouston, Hatzenbuehler, et al., 2015; Chae, Clouston, Martz, et al., 2018; Isoya and Yamada, 2021), I use coded language to refer to these terms. The words themselves can be found in Table **??***

# 1  Introduction

Can media coverage influence racial resentment? Empirical studies suggest that when white Americans understand welfare policies to threaten their privileged status in the U.S. Social hierarchy, their resentment of minorities increases and their support for welfare decreases (Willer, Feinberg, and Wetts, 2016; Wetts and Willer, 2018). This tendency seems to be weaponed by conservative media institutions and politicians, who seem to encourage racial animus to erode support for social programs, as in the infamous example of the "Welfare Queen" narrative, a racial stereotype employed to undercut support for the Aid to Families with Dependent Children (AFDC). However, the link between traditional media coverage and racial resentment has not been extensively studied.

I propose a study exploiting the expansion of the Sinclair Broadcasting Group from 2004-2021 to understand how conservative media messaging impacts racial resentment in a media market. The expansion of the Sinclair Broadcasting Group during this period provides the basis for a difference-in-differences analysis estimating the effect of Sinclair Media purchasing a station on racial animus in an area. I use the concentration of racially charged Google searches a proxy for the intensity of racial animus in an area. This measure has previously been used to measure the contributions of racial resentment towards African-American mortality (Chae, Clouston, Martz, et al., 2018; Chae, Clouston, Hatzenbuehler, et al., 2015), election outcomes (Stephens-Davidowitz, 2014), and economic inequality (Connor et al., 2019), and does not suffer from the same social censoring issues that confound traditional measures.

# 2  Theoretical Framework

# 3  Background

## 3.1  Sinclair Broadcast Group

During the period of 2004-2021, the Sinclair Broadcasting Group bought or sold stations in 67 media markets. I use this pattern of expansion as the basis for a difference-in-differences analysis.

I extract a record of which markets Sinclair has a presence in from the company's yearly filings with the Securities and Exchange Commission, which report each station owned by the company at the end of the time of filing a map of which regions Sinclair moved into / out of can be seen in figure 4a. This strategy has previously been used by Miho (2018) to track Sinclair Expanison over the period of 1995-2017. SEC filings have an advantage over the typical sets data used to track network expansion maintained by the Nielsen Corporation as they are publicly available, so the findings can be easily reproduced.

## 3.2  Does Sinclair Have a Racial Bias?

This analysis assumes that Sinclair stations differ from non-Sinclair stations on coverage of racial issues: if Sinclair stations deploy the same coverage on racial issues as other stations, then we would expect to see no difference in levels of racial resentment between Sinclair and non-Sinclair media markets. So, it must be asked: does Sinclair coverage have a racially-conservative bent? Admittedly, there have been no large-scale analyses of the effects of Sinclair ownership on

coverage of racial issues. Nonetheless, I submit that there is evidence to show Sinclair Stations cover racial issues in a more conservative light than they would absent Sinclair ownership.

Previous research has demonstrated that Sinclair acquisition of a network is associated with a sharp rightwards shift in its coverage (Martin and McCrain, 2019). I argue this effect extends to its framing of racial issues.

This year, the Sinclair corporation drew ire for a series of "must-run" segments on police violence following the murder of George Floyd pushing the "black-on-black violence" canard and advocating for a military response to the protests (Pleat and Savillo, 2020; Pleat, 2020).

In analyzing the effects of news coverage on racial animus, it might be natural to examine the expansion of Fox News, which occupies a position in the public consciousness as among the most conservative stations on racial issues, and is the most trusted media outlet among Republican and Republican-leaning respondents in many polls (Mark Jurkowitz and Walker, 2020). However, I choose to use Sinclair over Fox News Stations for two reasons. First, Fox News' expansion strategy has involved purchasing a larger stations. As television companies in the US can only expand until they broadcast to 39% of U.S. households (Scherer, 2018), Fox News has been able to buy fewer stations than Sinclair, which entails a smaller sample size of stations that changed ownership. Second, Fox News' expansion primarily happened before 2004, the first year for which Google Trends data is available, which further limits the sample size of stations which changed ownership for which there is data on racial animus.

## 3.3   Google Trends Data

I use the concentration of Google searches containing racial epithets as a proxy measurement of the level of racial animus in an area. Data from Google searches are ideal are searches not subject to the social-desirability biases that confound traditional, survey-based measures of racial resentment, and provide a large, regularly-sampled source of data which represents a great deal of the population (at present, Google as an over 85% market share in the US).

The core assumption of this strategy is that Google searches for the word [Word 1] reflect underlying racial animus in an area. If readers find this assumption unconvincing, then they will have little reason to accept the conclusions of this analysis. Accordingly, I offer several reasons to suggest that searches for racial epithets well proxy racial animus in an area, namely: searches are conducted in private, and are not subject to the same social-desirability biases as traditional, survey-based measures; and measurements of area-level racism obtained from Google trends data correspond well to traditional survey-based and non-survey-based measures of racial animus.
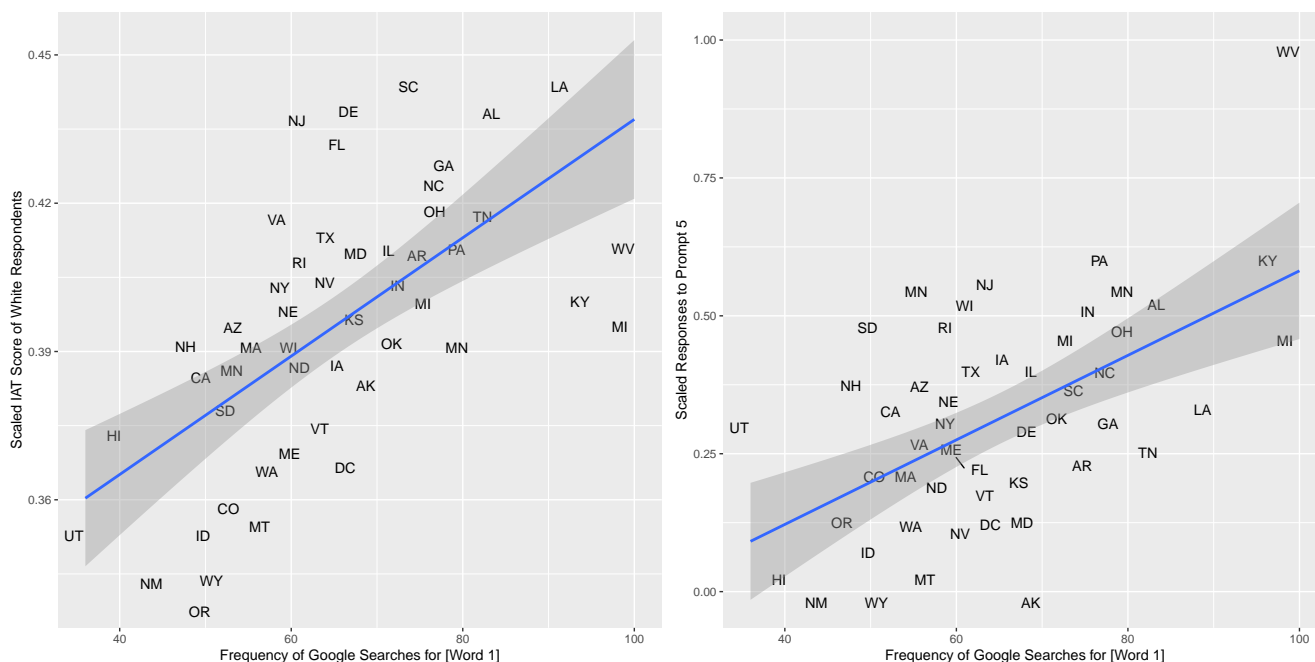
First, Google searches do not suffer from the same social censoring as traditional measures of public opinion that rely on face-to-face or over-the-phone interviews. Overt expressions of racism are no longer socially palatable in the U.S., Social desirability bias is a pressing concern when it comes to measuring racial animus; research has repeatedly demonstrated that (Kuklinski, Cobb, and Gilens, 1997)

This measurement approach has an advantage over traditional survey-based measures of racial animus in that it is less subject to a social desirability bias; "Google searchers are online and likely alone, both of which make it easier to express socially taboo thoughts (Kreuter et al., 2009)" (Stephens-Davidowitz, 2014, p. 26). Further, it provides a high-resolution set of data that would be prohibitively expensive to collect from a traditional survey, especially given that the difference-in-differences approach requires a repeated survey comparable across multiple time periods. [1]

As such, I give several reasons to suggest that Google trends data can be used to measure racial animus.

---

[1]Stephens-Davidowitz, 2014 evidences this claim by reporting statistics for pornography searches. Over the past 16 years, the number of searches for "porn" and "news" are commensurate, yet only 14% of GSS respondents tell the GSS they have visited a pornographic website in the past 30 days.

A pressing concern is that Google searches for "[Word 1](s) " may not actually capture the extent of racial bias in an area, but simply reflect users learning about the term. In fact, "definition of [Word 1] " and "what does [Word 1] mean" are both among the top 5 search queries related to the "[Word 1] ". These queries suggest that many who search for the term are searching out of curiosity to investigate the term.

This is not to suggest that searches for the term do not capture any variation in racial animus: among the top 10 most searched related queries to the term [Word 1]'s are "I don't like [Word 1]'s ," "fuck the [Word 1]'s ," and "ship those [Word 1]'s back."

The concern that Google Searches for racial slurs may largely reflect curiosity about the term is valid. Subject to data availability, I suggest controlling for the frequency of searches for the definitions of these terms in each area, to try and isolate the effect of Sinclair media entering a market on searches expressing "hardcore" racial animus rather than curiosity.

## 3.4   Scaling Google Trends Data:

Google trends data measures the popularity of a search on an idiosyncratic scale: a term's "search score" in a given time is given as the volume of searches over that time divided by the volume of searches when the term was most popular. In this section, I describe the scaling process used to back-out an interval-level measure of search popularity from these search scores.

First, a full explanation of the problem. Search popularity can be obtained from Google in two flavors: a measure that compares the popularity of a search across all regions at a given time, and a measure that gives the popularity of a search in a given region across all times.

The search score between regions at a given time is defined as the following:

$$\text{Between Regions Search Score} = \frac{\text{Popularity in Region } i}{\text{Popurality in the Most Popular Region}}$$

The search score that can be used to compare the popularity across different times in a given

region is defined as the following.

$$\text{Between Times Search Score} = \frac{\text{Popularity at Time } i}{\text{Popurality at the Most Popular Time}}$$

The first measure allows for comparisons between regions but not between times, while the second permits comparisons between times but only within one region.

# 4 Methods

## 4.1 Tools Used

I use Kintr to integrate statistical calculations into the paper, eliminating the possibility of transcription errors (Xie, 2014). To ensure that the methods of this paper have been properly implemented and the finding are reproducible, I tested the analysis routines using the *testthat* package in R (Wickham, 2011) and the *unittest* module in Python (Van Rossum and Drake, 2009). But I won't make you take my word for it – I provide a Docker image with the reproducibility materials to ensure others can replicate the calculations on their own systems (Merkel, 2014; Boettiger, 2015). The net result is "one-click reproducibility" (Nüst et al., 2020); readers can reproduce this exact paper with the push of a button from the linked materials. [2]

## 4.2 Preregistration

To avoid the possibility of fitting hypotheses to the data after results are known, I created a preregistration plan of my analysis using the Google trends data. The plan can be seen in section **??**. As I decided to conduct the analysis using the project implicit data to confirm the Google trends results, I did not preregister this analysis plan.

I have made just one deviation from the preregistration. In my preregistration, I describe a different strategy to scale the Google trends data than the one I actually employ. This original strategy is based on a misunderstanding of the format of Google trends data, and does not actually produce the desired measure.

In the analysis I perform, I correct this mistake. I describe the correct scaling procedure in section **??**.

# 5 Results

Racially Charged Search Rate $= \beta_1(\text{ Sinclair Present })+\beta_2(\text{ DMA fixed effects })+\beta_3(\text{ year fixed effects })$

---

[2]Replication materials available here. By default, the webscraping and models using IAT data are not re-run as they are memory and time-intensive.

Figure 2: Sinclair News Anchors Reading a "Must-Run" Script (May 2018)



(a) Images of 30 of 210 Sinclair Station Hosts Reading a Must-Run Script (Burke, 2018)

**(A):** But we're concerned about the troubling trend of irresponsible, one sided news stories plaguing our country. The sharing of biased and false news has become all too common on social media.

**(B):** More alarming, some media outlets publish these same fake stories... stories that just aren't true, without checking facts first.

**(A):** Unfortunately, some members of the media use their platforms to push their own personal bias and agenda to control 'exactly what people think'...This is extremely dangerous to a democracy.

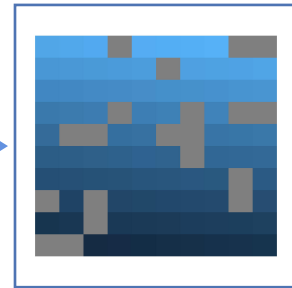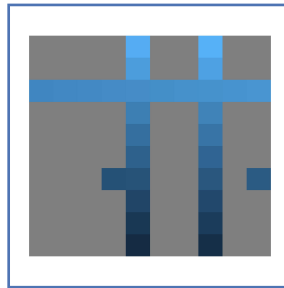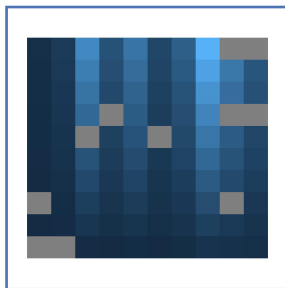(b) Transcript of Segment (Cohen, 2018)
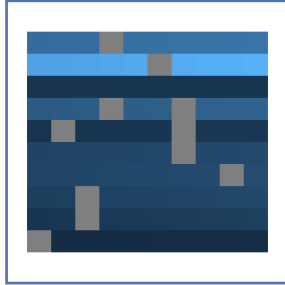
Table 1: Fixed-Effect Model Results

|  | *Dependent variable:* | |
| --- | --- | --- |
|  | Frequency of Google Searches for [Word 1] | |
|  | (1) | (2) |
| Sinclair Present | −0.122* | −0.068 |
|  | (0.065) | (0.107) |
| Constant | −0.212 | 41.959 |
|  | (0.234) | (126.128) |
| Year Fixed Effects | Yes | Yes |
| Region Fixed Effects | Yes | Yes |
| Region Time Trends | No | Yes |
| Observations | 3,570 | 3,570 |
| $R^2$ | 0.240 | 0.309 |
| Adjusted $R^2$ | 0.188 | 0.214 |
| Residual Std. Error | 0.901 (df = 3343) | 0.887 (df = 3134) |
| F Statistic | 4.668*** (df = 226; 3343) | 3.229*** (df = 435; 3134) |

*Note:*           *p<0.1; **p<0.05; ***p<0.01

Figure 3: Illustration of Scaling Algorithm

Horizontally Comparable Matrix



Scaled

Repaired Dataset

Vertically Comparable Matrix

(a) Media Markets Sinclair Bought or Sold a Station in between 2004-2020. DMA Boundaries From Hill, 2015

(b) Media Markets by Sinclair Ownership Status, Time

## 5.1 Indentification Asssumption

In this section, I test the indefication assumption, the assumption that the treated and control units would have the outcomes if the treatment were absent.

Racially Charged Search Rate = $\beta_1$( Sinclair Present )+$\beta_2$( DMA fixed effects )+$\beta_3$( year fixed effects )+$\beta_4$(Year / D

Figure 5: Fixed-Effects Estimates of the effect of Sinclair Acquisition on Rate of Racially Charged Google Searches

| Code   | Word   |
|--------|--------|
| Word 1 | nigger |

# A  Codings for Offensive Words

# References

Boettiger, Carl (Jan. 2015). "An introduction to Docker for reproducible research". In: *ACM SIGOPS Operating Systems Review* 49.1, pp. 71–79.

Burke, Timothy (2018). *How I Made A Dumb Video Making Fun Of Sinclair Broadcasting And Somehow Started A Media War*. URL: https://deadspin.com/how-i-made-a-dumb-video-making-fun-of-sinclair-broadcas-1825106452 (visited on 03/23/2020).

Chae, David H., Sean Clouston, Mark L. Hatzenbuehler, et al. (Apr. 2015). "Association between an Internet-Based Measure of Area Racism and Black Mortality". In: *PLOS ONE* 10.4. Ed. by Hajo Zeeb, e0122963.

Chae, David H., Sean Clouston, Connor D. Martz, et al. (Feb. 2018). "Area racism and birth outcomes among Blacks in the United States". In: *Social Science & Medicine* 199, pp. 49–55.

Cohen, Stephen (2018). *KOMO attacks 'biased and false news' in Sinclair-written promos*. URL: https://www.seattlepi.com/seattlenews/article/KOMO-fake-news-Sinclair-promos-12792032.php (visited on 03/23/2020).

Connor, Paul et al. (Jan. 2019). "Income Inequality and White-on-Black Racial Bias in the United States: Evidence From Project Implicit and Google Trends". In: *Psychological Science* 30.2, pp. 205–222.

Hill, Andrew (2015). *DMA Boundaries*. URL: https://team.carto.com/u/andrew/tables/dma_boundary/public.

Isoya, Hidehito and Hiroyuki Yamada (Mar. 2021). "Hidden internet-based measure of discrimination and mental health of Black people". In: *Social Science & Medicine* 273, p. 113776.

Kuklinski, James H., Michael D. Cobb, and Martin Gilens (1997). "Racial Attitudes and the "New South"". In: *The Journal of Politics* 59.2, pp. 323–349.

Mark Jurkowitz Amy Mitchell, Elisa Shearer and Mason Walker (Jan. 2020). *U.S. Media Polarization and the 2020 Election: A Nation Divided*. URL: https://www.journalism.org/2020/01/24/u-s-media-polarization-and-the-2020-election-a-nation-divided/.

Martin, Gregory J. and Joshua McCrain (Feb. 2019). "Local News and National Politics". In: *American Political Science Review* 113.2, pp. 372–384.

Merkel, Dirk (2014). "Docker: lightweight linux containers for consistent development and deployment". In: *Linux journal* 2014.239, p. 2.

Miho, Antonela (2018). "Small screen, big echo? Estimating the political persuasion of local television news bias using the Sinclair Broadcasting Group as a natural experiment". In: *SSRN Electronic Journal*.

Nüst, Daniel et al. (Nov. 2020). "Ten simple rules for writing Dockerfiles for reproducible data science". In: *PLOS Computational Biology* 16.11. Ed. by Scott Markel, e1008316.

Pleat, Zachary (2020). *These seven national Sinclair Broadcast news segments on the Floyd protests didn't mention police violence against protesters*. URL: https://www.mediamatters.org/black-lives-matter/these-seven-national-sinclair-broadcast-news-segments-floyd-protests-didnt.

Pleat, Zachary and Rob Savillo (2020). *Sinclair reporters pushed debunked myths in national segments on police violence against Black Americans*. URL: https://www.mediamatters.org/black-lives-matter/sinclair-reporters-pushed-debunked-myths-national-segments-police-violence.

Scherer, Dana A. (Oct. 2018). "Federal Communications Comission (FCC) Media Ownership Rules". In: *Congressional Research Service Reports*.

Stephens-Davidowitz, Seth (Oct. 2014). "The cost of racial animus on a black candidate: Evidence using Google search data". In: *Journal of Public Economics* 118, pp. 26–40.

Van Rossum, Guido and Fred L. Drake (2009). *Python 3 Reference Manual*. Scotts Valley, CA: CreateSpace.

Wetts, Rachel and Robb Willer (May 2018). "Privilege on the Precipice: Perceived Racial Status Threats Lead White Americans to Oppose Welfare Programs". In: *Social Forces* 97.2, pp. 793–822.

Wickham, Hadley (2011). "testthat: Get Started with Testing". In: *The R Journal* 3, pp. 5–10.

Willer, Robb, Matthew Feinberg, and Rachel Wetts (2016). "Threats to Racial Status Promote Tea Party Support Among White Americans". In: *SSRN Electronic Journal*.

Xie, Yihui (2014). "knitr: A Comprehensive Tool for Reproducible Research in R". In: *Implementing Reproducible Computational Research*. Ed. by Victoria Stodden, Friedrich Leisch, and Roger D. Peng. ISBN 978-1466561595. Chapman and Hall/CRC.

# B    Code In This Document

# C    Web Scraping Code

## C.1    Master Web Scraping Script

```python
#!/usr/bin/python
import glob
import os
import pickle as pkl
from between_regions import between_reigion_many
from in_region import in_region
from utils import itr_split_overlap
import tqdm


def extract_keywords(filename):
    with open(filename, "r") as f:
        keywords = f.read().splitlines()
    return keywords


def get_basename(filename):
    name_w_ext = os.path.basename(filename)
    basename = os.path.splitext(name_w_ext)[0]
    return (basename)


def run_keywords(name, keywords):
    regions = pkl.load(open("data/dma_abbreviations.pkl", "rb"))
    filename = f"data/google_trends_data/{name}_time_serires.csv"
    with open(filename, "w") as f:
        f.write("row,date,score,ispartial,code,term\n")
        for keyword in keywords:
            for region in tqdm.tqdm(regions):
                df = in_region(keyword, region, True, timeframe="all")

                if df is not None:
                    df.to_csv(f, header=False)

    filename = f"data/google_trends_data/{name}_between_regions.csv"
    with open(filename, "w") as f:
        df = between_reigion_many(itr_split_overlap(keywords, 5, 1),
                                  censor=True,
                                  timeframe="all",
                                  geo="US",
                                  gprop="")
        df.to_csv(f, header=True)
```

```
44
45  if __name__ == "__main__":
46      keyword_files = sorted(glob.glob("data/keywords/*.csv"))
47      base_names = [get_basename(filename) for filename in keyword_files]
48      keywords = [extract_keywords(filename) for filename in keyword_files]
49      first = (list(zip(base_names, keywords)))[0]
50      run_keywords(first[0], first[1])
```

## C.2    Collecting Data for Between-Region Comparisons

```
1   #!/usr/bin/env python
2   from pytrends.request import TrendReq
3   from utils import censor_string
4   import pandas as pd
5
6
7   def between_region(query, censor, **kwargs):
8       pytrends = TrendReq(hl='en-US', tz=360)
9       pytrends.build_payload(query, cat=0, **kwargs)
10      df = pytrends.interest_by_region(resolution='DMA',
11                                       inc_low_vol=True,
12                                       inc_geo_code=True)
13      df = df.set_index("geoCode")
14      df.index.name = 'code'
15      if censor:
16          df = df.rename(censor_string, axis="columns")
17      return (df)
18
19
20  def between_reigion_many(iterable, censor, **kwargs):
21      iterable = iter(iterable)
22      chunk1 = next(iterable)
23      df1 = between_region(chunk1, censor, **kwargs)
24
25      for chunk in iterable:
26          chunk2 = chunk
27          shared = list(set(chunk1) & set(chunk2))[0]
28          if censor:
29              shared = censor_string(shared)
30          df2 = between_region(chunk2, censor, **kwargs)
31
32          mean1 = df1[shared].mean()
33          mean2 = df2[shared].mean()
34          normaliation_factor = mean1 / mean2
35
36          df2 = df2 * normaliation_factor
37          df2 = df2.drop(columns=[shared])
38
39          df1 = df1.join(df2)
40
41          chunk1 = chunk2
42
43      return (df1)
44
45
46  def create_v_df(term, year):
47      timeframe = f"{year}-01-01 {year}-12-31"
48      df = between_region([term], False, timeframe=timeframe, geo="US")
49      df = df.rename(columns={term: str(year)})
50      return df
51
52
53  if __name__ == "__main__":
54      dfs = tuple(create_v_df("hello", year) for year in range(2004, 2021))
55      df_c = pd.concat(dfs, axis=1)
56
57      df_c.to_parquet("v_matrix.parquet")
```

## C.3    Collecting Within-Region Data

```
1   #!/usr/bin/python
2   from pytrends.request import TrendReq
3   from utils import censor_string
4   import time
5   import pandas as pd
```

```python
 6  import re
 7  import pickle as pkl
 8
 9
10  def in_region(query, region, censor, **kwargs):
11      """ Returns a set of data showing the popularity of a search term over time
12
13      :query: The search term that time-series data is collected for
14      :region: The geographic region to retrieve time-series data for.
15      :censor: boolean should the search terms be censored?
16      :**kwargs: kwargs to be passed to pytrends.build_payload()
17      :returns: DataFrame giving time-series data for the popularity of a
18          search term in a given region
19
20      """
21      pytrends = TrendReq(hl='en-US', tz=360)
22      try:
23          pytrends.build_payload(kw_list=[query], geo=f"{region}", **kwargs)
24          df = pytrends.interest_over_time()
25          if df.empty:
26
27              df = pd.DataFrame()
28              df['date'] = pd.period_range(start='2004-01-01',
29                                           end='2021-01-01',
30                                           freq='M').to_timestamp()
31              df['n'] = 0
32              df['ispartial'] = pd.Series([True]).bool()
33
34          else:
35
36              df.columns = ["n", "ispartial"]
37              df.index.name = 'date'
38              df.reset_index(inplace=True)
39
40          df["query"] = query
41          df['code'] = re.findall("\d+", region)[0]
42          if censor:
43              df["query"] = df["query"].apply(censor_string)
44
45          return df
46      except:
47          if censor:
48              print(f"Rate error: {censor_string(query)} in {region}")
49          else:
50              print(f"Rate error: {query} in {region}")
51
52          time.sleep(60)
53          return in_region(query, region, censor, **kwargs)
54
55
56  def to_wide(df):
57      """TODO: Turns time-series search popularity data into a 'wide' dataframe to be used in
58      scaling
59
60      :df: 'long' dataframe of search data, as from in_region()
61      :returns: 'wide' DataFrame of search data, averaged by year
62
63      """
64      print(df['date'])
65      df['year'] = pd.DatetimeIndex(df['date']).year
66      df['year'] = df['year'].apply(str)
67      df = df.groupby(['year', 'code'])["n"].mean()
68      df = df.unstack(level=0)
69      return (df)
70
71
72  # with open("data/dma_abbreviations.pkl", "rb") as f:
73  #     dmas = pkl.load(f)
74
75  # dmas = dmas
76
77  # in_region_dfs = tuple(
78  #     in_region("economist", dma, True, timeframe="all") for dma in dmas)
79  # wide_dfs = map(to_wide, in_region_dfs)
80  # h_df = pd.concat(wide_dfs).sort_index()
81  # print(h_df)
```

## C.4 Utility Functions

```python
#!/usr/bin/python
import itertools as it
import requests
import random


# returns mean of nonzero values in iterator
def mean_nonzero(iterator):
    nonzero = tuple(filter(lambda x: x != 0, iterator))
    if len(nonzero) == 0:
        return None
    elif None in nonzero:
        return None
    else:
        return sum(nonzero) / len(nonzero)


# Returns list of N random words from MIT dictionary
def random_words(n):
    word_site = "https://www.mit.edu/~ecprice/wordlist.10000"

    response = requests.get(word_site)
    words = response.text.splitlines()

    keywords = random.sample(words, n)

    return (keywords)


# Tests if computer is connected to internet (used in tests)
def connected():
    url = "http://google.com"
    timeout = 5
    try:
        requests.get(url, timeout=timeout)
        return (True)
    except (requests.ConnectionError, requests.Timeout):
        return (False)


# Censors strings so that senstive words aren't uploaded to github / used in
# scripts
def censor_string(string):
    return (string[0] + "_" * (len(string) - 2) + string[-1])


# credit to Ilja Everila for this implimentation
# https://stackoverflow.com/questions/48381870/a-better-way-to-split-a-sequence-in-chunks-
    with-overlaps
def itr_split_overlap(iterable, size, overlap):

    if overlap >= size:
        raise ValueError("overlap must be smaller than size")

    itr = iter(iterable)

    next_ = tuple(it.islice(itr, size))

    yield next_

    prev = next_[-overlap:] if overlap else ()

    while True:
        chunk = tuple(it.islice(itr, size - overlap))

        if not chunk:
            break

        next_ = (*prev, *chunk)
        yield next_

        if overlap:
            prev = next_[-overlap:]
```

# D Analysis Code

```
1  #!/usr/bin/Rscript
2  # Load Libraries for Analysis
3  library(tidyverse)
4  library(lubridate)
5  library(broom)
6  # Load utility functions
7  source("utils.R")
8
9  search_data <- read_csv("data/google_trends_data/word_1.csv") %>%
10   pivot_longer(-code, names_to = "year", values_to = "word1") %>%
11   mutate_all(as.numeric)
12
13 stopifnot(nrow(search_data) == 3570)
14 stopifnot(all(!is.na(search_data)))
15
16 write_csv(search_data, "../data/google_trends/word1.csv")
17
18 sinclair_data <- read_csv("../data/clean_sinclair_data.csv")
19 stopifnot(nrow(sinclair_data) == 3570)
20 stopifnot(all(!is.na(sinclair_data)))
21
22 # search_data <- search_data %>%
23 #   dplyr::select(-data) %>%
24 #   distinct(term, code, year, .keep_all = TRUE) %>%
25 #   pivot_wider(names_from = term, values_from = score) %>%
26 #   mutate(overall_score = rowSums(across(everything()), na.rm = T))
27 dma_names <- read_csv("data/dma_list.csv")
28 stopifnot(nrow(dma_names) == 210)
29
30 full_data <- search_data %>%
31   right_join(sinclair_data) %>%
32   full_join(dma_names) %>%
33   filter(year != 2021) %>%
34   group_by(code) %>%
35   mutate(years_before = years_before(sinclair_present)) %>%
36   ungroup() %>%
37   mutate(sword1 = (word1-mean(word1))/sd(word1)) %>%
38   mutate(years_before = relevel(as.factor(years_before),"-99"))
39
40 stopifnot(nrow(full_data) == 3570)
41 stopifnot(all(!is.na(full_data)))
42
43 write_csv(full_data, "../data/full_data.csv")
44
45 model_1 <- lm(sword1 ~ as.factor(year) + as.factor(code) + sinclair_present, data = full_data
      )
46 model_1 %>% summary()
47
48 model_2 <- lm(sword1 ~ as.factor(year) + as.factor(code) + year:as.factor(code)+ sinclair_
      present, data = full_data)
49 model_2 %>% summary()
50
51 model_3 <- lm(sword1 ~ as.factor(year) + as.factor(code) + as.factor(years_before), data =
      full_data)
52 model_3 %>% summary()
53
54 model_3 %>%
55   tidy() %>%
56   filter(grepl("years_before", term)) %>%
57   mutate(term = as.numeric(gsub("[^0-9\\-]+", "", term))) %>%
58   ggplot(aes(x = term, y = estimate)) +
59   geom_point() +
60   geom_errorbar(aes(ymin = estimate - 1.96 * std.error, ymax = estimate + 1.96 * std.error))
      +
61   geom_hline(aes(yintercept=0), linetype=2) +
62   geom_vline(aes(xintercept=0))
```

## D.1 Utility Functions

```
1  #!/usr/bin/Rscript
2  censor_string <- function(string) {
3    substr(string, 2, nchar(string) - 1) <- paste0(rep("_", nchar(string) - 2), collapse = "")
4    names(string) <- NULL
```

```r
 5    return(string)
 6  }
 7  censor_string <- Vectorize(censor_string, USE.NAMES=F)
 8
 9  years_before <- function(bool) {
10    if (any(bool)) {
11      out <- numeric(length = length(bool))
12      start <- min(which(bool))
13
14      before <- (seq(start, 1) - 1) * -1
15      if (start != length(bool)) {
16        after <- seq(1, length(bool) - start)
17      } else {
18        after <- c()
19      }
20
21
22      out <- c(before, after)
23    } else {
24      out <- rep(-99, length(bool))
25    }
26
27    return(out)
28  }
```

# E  Unit tests

## E.1  For Python Code

```python
 1  #!/usr/bin/python
 2  import unittest
 3  import pandas as pd
 4  from between_regions import between_region
 5  from utils import connected
 6
 7
 8  @unittest.skipIf(not connected(), "not connected to the internet")
 9  class TestBetweenRegion(unittest.TestCase):
10      def test_between_region_uncensored(self):
11
12          result_1 = between_region(["socks"],
13                                     censor=False,
14                                     timeframe="2016-12-14 2017-01-25",
15                                     geo="US",
16                                     gprop="")
17
18          result_2 = between_region(["socks", "shoe", "fish"],
19                                     censor=False,
20                                     timeframe="2016-12-14 2017-01-25",
21                                     geo="US",
22                                     gprop="")
23
24          expected_1 = pd.read_parquet(
25              "tests/test_data/between_region_1_uc.parquet")
26          expected_2 = pd.read_parquet(
27              "tests/test_data/between_region_2_uc.parquet")
28
29          self.assertTrue(expected_1.equals(result_1))
30          self.assertTrue(expected_2.equals(result_2))
31
32      def test_between_region_censored(self):
33          result_1 = between_region(["socks"],
34                                     censor=True,
35                                     timeframe="2016-12-14 2017-01-25",
36                                     geo="US",
37                                     gprop="")
38
39          result_2 = between_region(["socks", "shoe", "fish"],
40                                     censor=True,
41                                     timeframe="2016-12-14 2017-01-25",
42                                     geo="US",
43                                     gprop="")
44
45          expected_1 = pd.read_parquet(
46              "tests/test_data/between_region_1_c.parquet")
```

```
47          expected_2 = pd.read_parquet(
48              "tests/test_data/between_region_2_c.parquet")
49
50          self.assertTrue(expected_1.equals(result_1))
51          self.assertTrue(expected_2.equals(result_2))
```

```python
1  #!/usr/bin/python
2  import unittest
3  import pandas as pd
4  from in_region import in_region
5  from utils import connected
6
7
8  @unittest.skipIf(not connected(), "not connected to the internet")
9  class TestBetweenRegion(unittest.TestCase):
10     def test_in_region_uncensored(self):
11
12         result_1 = in_region("hello",
13                              "US_AL_630",
14                              censor=False,
15                              timeframe="2016-12-14 2017-01-25")
16
17         result_2 = in_region("France",
18                              "US_GA_522",
19                              censor=False,
20                              timeframe="2016-12-14 2017-01-25")
21
22         expected_1 = pd.read_parquet("tests/test_data/in_region_1_uc.parquet")
23         expected_2 = pd.read_parquet("tests/test_data/in_region_2_uc.parquet")
24
25         self.assertTrue(expected_1.equals(result_1))
26         self.assertTrue(expected_2.equals(result_2))
27
28     def test_in_region_censored(self):
29         result_1 = in_region(["socks"],
30                              censor=True,
31                              timeframe="2016-12-14 2017-01-25",
32                              geo="US",
33                              gprop="")
34
35         result_2 = in_region(["socks", "shoe", "fish"],
36                              censor=True,
37                              timeframe="2016-12-14 2017-01-25",
38                              geo="US",
39                              gprop="")
40
41         expected_1 = pd.read_parquet("tests/test_data/in_region_1_c.parquet")
42         expected_2 = pd.read_parquet("tests/test_data/in_region_2_c.parquet")
43
44         self.assertTrue(expected_1.equals(result_1))
45         self.assertTrue(expected_2.equals(result_2))
```

```python
1  #!/usr/bin/python
2  import unittest
3  import numpy as np
4  import random
5  from scaling import compare, compare_horizontal, compare_vertical
6  from scaling import recover_scores, repair_matrix
7
8
9  class TestCompareHorizontal(unittest.TestCase):
10     def test_compare_horizontal(self):
11         h_matrix = np.array([[2, 4, 6], [12, 15, 18], [3.5, 4, 4.5]])
12
13         self.assertEqual(1 / 3, compare_horizontal(h_matrix, 0, 0, 2))
14         self.assertEqual(2 / 3, compare_horizontal(h_matrix, 1, 0, 2))
15         self.assertEqual(5 / 4, compare_horizontal(h_matrix, 1, 1, 0))
16
17
18 class TestCompareVertical(unittest.TestCase):
19     def test_compare_vertical(self):
20         v_matrix = np.array([[3, 1, 1], [12, 2.5, 2], [21, 4, 3]])
21
22         self.assertEqual(1 / 4, compare_vertical(v_matrix, 0, 0, 1))
23         self.assertEqual(1 / 4, compare_vertical(v_matrix, 1, 0, 2))
24         self.assertEqual(4, compare_vertical(v_matrix, 0, 1, 0))
25
```

```
26
27  class TestCompare(unittest.TestCase):
28      def test_compare(self):
29          h_matrix = np.array([[2, 4, 6], [12, 15, 18], [3.5, 4, 4.5]])
30          v_matrix = np.array([[3, 1, 1], [12, 2.5, 2], [21, 4, 3]])
31
32          self.assertEqual(1 / 8, compare(h_matrix, v_matrix, 0, 0, 1, 2))
33          self.assertEqual(1 / 9, compare(h_matrix, v_matrix, 0, 0, 2, 2))
34          self.assertEqual(1 / 4, compare(h_matrix, v_matrix, 1, 0, 1, 2))
35
36
37  class TestRecoverScores(unittest.TestCase):
38      def test_recover_scores(self):
39          for _ in range(100):
40              arr = np.array([[x + 8 * y for x in range(8)] for y in range(8)])
41
42              rowmod = np.array(tuple(random.random() for _ in range(8)))
43              colmod = np.array(tuple(random.random() for _ in range(8)))
44
45              h_matrix = arr * rowmod[:, np.newaxis]
46              v_matrix = arr * colmod
47
48              self.assertTrue(
49                  np.allclose(arr,
50                              recover_scores(h_matrix, v_matrix) * 63))
51
52
53  class TestRepairScores(unittest.TestCase):
54      def test_repair_scores(self):
55          for _ in range(100):
56              arr = np.array([[x + 10 * y for x in range(10)]
57                              for y in range(10)])
58
59              n = 15
60              index = np.random.choice(arr.size, n, replace=False)
61              arr.ravel()[index] = 0
62
63              rowmod = np.array(tuple(random.random() for _ in range(10)))
64              colmod = np.array(tuple(random.random() for _ in range(10)))
65
66              h_matrix = arr * rowmod[:, np.newaxis]
67              v_matrix = arr * colmod
68
69              gen1 = recover_scores(h_matrix, v_matrix)
70              if not gen1 is None:
71
72                  repaired = repair_matrix(h_matrix, v_matrix, gen1)
73
74                  repaired_scaled = repaired / np.max(repaired)
75                  arr_scaled = arr / np.max(arr)
76
77                  self.assertTrue(np.allclose(repaired_scaled, arr_scaled))
```

```
1   #!/usr/bin/python
2   import unittest
3   from utils import itr_split_overlap, censor_string
4
5
6   class TestItrSplitOverlap(unittest.TestCase):
7       def test_itr(self):
8           test_list = ["one", "two", "three", "four", "five"]
9
10          expected_result = [('one', 'two'), ('two', 'three'), ('three', 'four'),
11                             ('four', 'five')]
12          self.assertEqual(expected_result,
13                           list(itr_split_overlap(test_list, 2, 1)))
14
15          expected_result = [('one', 'two', 'three', 'four'),
16                             ('two', 'three', 'four', 'five')]
17          self.assertEqual(expected_result,
18                           list(itr_split_overlap(test_list, 4, 3)))
19
20      def test_exceptions(self):
21          with self.assertRaises(ValueError):
22              test_list = ["one", "two", "three", "four", "five"]
23
24              list(itr_split_overlap(test_list, 2, 3))
```

```
25
26
27 class CensorString(unittest.TestCase):
28     def test_censor(self):
29         self.assertEqual("h_y", censor_string("hey"))
30         self.assertEqual("f____r", censor_string("fender"))
31         self.assertEqual("fr", censor_string("fr"))
```

## E.2   For R Code

```
1 library(testthat)
2 source("../../utils.R")
3
4 test_that("censor string works", {
5       expect_equal(censor_string("hey"),"h_y")
6       expect_equal(censor_string("hy"),"hy")
7       expect_equal(censor_string("watermelon"),"w_____n")
8       expect_equal(censor_string("fantastic"),"f_____c")
9 })
10
11 test_that("censor string vectorized correctly", {
12      expect_equal(censor_string(c("hello","there")),c("h___o","t___e"))
13      expect_equal(censor_string(c("watermelon","paper")),c("w_____n","p___r"))
14 })
15
16 test_that("years before works", {
17   expect_equal(years_before(c(F, F, F, T, T)), c(-3, -2, -1, 0, 1))
18   expect_equal(years_before(c(F, F, F, F, T)), c(-4, -3, -2, -1, 0))
19   expect_equal(years_before(c(T, T, T, T, T)), c(0, 1, 2, 3, 4))
20 })
```