March 28, 2021

*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

# 2   Methods

## 2.1   Tools Used

Surveys of major social science journals routinely fail to reproduce the findings of a plurality or majority of published articles, even when using data provided with the paper (Nuijten, Hartgerink, et al. 2015; Nuijten and Polanin 2020; Eubank 2016). In many cases, failures to replicate are due to coding errors or mistakes in transcribing the results of a calculation into the published manuscript (Eubank 2016).

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). I use Kintr to integrate statistical calculations into the paper, eliminating the possibility of transcription errors (Xie 2014).

Finally, I provide a Docker image (Merkel 2014) with the reproducibility materials to ensure others can replicate the figures in this paper on their own systems (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 on their own systems from the reproducibility materials.

Figure 1: Sinclair News Anchors Reading a "Must-Run" Script (May 2018)



(a) Images from 30 of 210 Stations Reading the Script (Burke 2018)
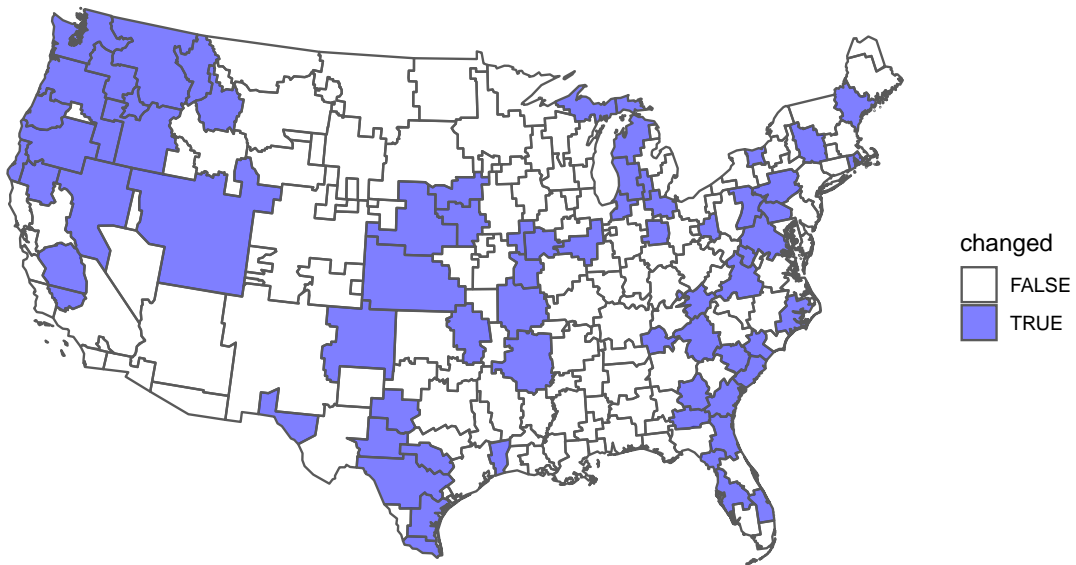
**(A):** But we're concerned about the troubling trend of irresponsible, one sided news stories plaguing our country. The sharing of biased and false news has become all too common on social media.

**(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)

Map of Sinclair Stations That Bought or Sold in 2004–2021



changed
☐ FALSE
■ TRUE

# 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).

Eubank, Nicholas (Apr. 2016). "Lessons from a Decade of Replications at the Quarterly Journal of Political Science". In: *PS: Political Science & Politics* 49.02, pp. 273–276.

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.

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

Nuijten, Michèle B., Chris H. J. Hartgerink, et al. (Oct. 2015). "The prevalence of statistical reporting errors in psychology (1985–2013)". In: *Behavior Research Methods* 48.4, pp. 1205–1226.

Nuijten, Michèle B. and Joshua R. Polanin (Apr. 2020). "" statcheck ": Automatically detect statistical reporting inconsistencies to increase reproducibility of meta-analyses". In: *Research Synthesis Methods*.

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.

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.

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

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.

# A Codings for Offensive Words

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

# B Code In This Document

```
1  ## ----echo=FALSE, warn=FALSE, include=FALSE-----------------------------------
2  # Set options to knit this document, load libraries
3  knitr::opts_chunk$set(out.width="100%", message=F, cache=T,echo=F, warn=F, include=F)
4  library(tidyverse)
5  library(sf)
6
7
8  ## ------------------------------------------------------------------------------
9  # Load and create data to show markets Sinclair has expanded to / moved out of
10 sinclair_expansions <- read_csv("../data/clean_sinclair_data.csv") %>%
11     nest(-code) %>%
12     mutate(
13         any_true =  data %>% map_lgl(~any(.$sinclair_present)),
14         any_false =  data %>% map_lgl(~any(!.$sinclair_present)),
15         changed = any_true & any_false
16     )
17
18 dma_boundaries <- st_read("../data/dma_boundaries/dma_boundary.shp")
19 dma_boundaries <- merge(dma_boundaries, sinclair_expansions, by.x="dma0", by.y="code")
20
21
22 ## ----include = T---------------------------------------------------------------
23 # Plot markets Sinclair has expanded to / moved out of
24 ggplot() +
25     geom_sf(data=dma_boundaries, aes(fill=changed), alpha=.5) +
26     scale_fill_manual(values  = c(NA, "blue")) +
27     ggtitle("Map of Sinclair Stations That Bought or Sold in 2004-2021") +
28     theme_void()
29
30
31 ## ----echo=FALSE, include=T-----------------------------------------------------
32 # Pull R code out of this document to put into the appendix
33 code <- knitr::purl("diss.Rnw", quiet=T)
```

# C Web Scraping Code

## C.1 Master Web Scraping Script

```python
1  #!/usr/bin/python
2  import glob
3  import os
4  import pickle as pkl
5  from between_regions import between_reigion_many
6  from in_region import get_region_trend
7  from utils import itr_split_overlap
8  import tqdm
9
10 def extract_keywords(filename):
11     with open(filename, "r") as f:
12         keywords = f.read().splitlines()
13     return keywords
14
15
16 def get_basename(filename):
17     name_w_ext = os.path.basename(filename)
18     basename = os.path.splitext(name_w_ext)[0]
19     return (basename)
20
21
22 def run_keywords(name, keywords):
23     regions = pkl.load(open("dma_abbreviations.pkl", "rb"))
24     filename = f"../data/google_trends_data/{name}_time_serires.csv"
```

```
25        with open(filename, "w") as f:
26            f.write("row,date,score,ispartial,code,term\n")
27            for keyword in keywords:
28                for region in tqdm.tqdm(regions):
29                    df = get_region_trend(keyword, region, True)
30
31                    if df is not None:
32                        df.to_csv(f, header=False)
33
34        filename = f"../data/google_trends_data/{name}_between_regions.csv"
35        with open(filename, "w") as f:
36            df = between_reigion_many(itr_split_overlap(keywords, 5, 1),
37                                      censor=True,
38                                      timeframe="all",
39                                      geo="US",
40                                      gprop="")
41            df.to_csv(f, header=True)
42
43
44 if __name__ == "__main__":
45     keyword_files = sorted(glob.glob("keywords/*.csv"))
46     base_names = [get_basename(filename) for filename in keyword_files]
47     keywords = [extract_keywords(filename) for filename in keyword_files]
48     first = (list(zip(base_names, keywords)))[1]
49     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
5
6 def between_region(query, censor, **kwargs):
7     pytrends = TrendReq(hl='en-US', tz=360)
8     pytrends.build_payload(query, cat=0, **kwargs)
9     df = pytrends.interest_by_region(resolution='DMA',
10                                      inc_low_vol=True,
11                                      inc_geo_code=True)
12     df = df.set_index("geoCode")
13     df.index.name = 'code'
14     if censor:
15         df = df.rename(censor_string, axis="columns")
16     return (df)
17
18
19 def between_reigion_many(iterable, censor, **kwargs):
20     iterable = iter(iterable)
21     chunk1 = next(iterable)
22     df1 = between_region(chunk1, censor, **kwargs)
23
24     for chunk in iterable:
25         chunk2 = chunk
26         shared = list(set(chunk1) & set(chunk2))[0]
27         if censor:
28             shared = censor_string(shared)
29         df2 = between_region(chunk2, censor, **kwargs)
30
31         mean1 = df1[shared].mean()
32         mean2 = df2[shared].mean()
33         normaliation_factor = mean1 / mean2
34
35         df2 = df2 * normaliation_factor
36         df2 = df2.drop(columns=[shared])
37
38         df1 = df1.join(df2)
39
40         chunk1 = chunk2
41
42     return (df1)
```

## 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
6  def get_region_trend(query, region, censor):
7      pytrends = TrendReq(hl='en-US', tz=360)
8      try:
9          pytrends.build_payload(kw_list=[query],
10                                 geo=f"{region}",
11                                 timeframe="all")
12         df = pytrends.interest_over_time()
13         if not df.empty:
14             df.columns = ["n", "ispartial"]
15             df.index.name = 'date'
16             df.reset_index(inplace=True)
17             df["region"] = region
18
19             df["query"] = query
20
21             if censor:
22                 df["query"] = df["query"].apply(censor_string)
23
24             return df
25     except:
26         if censor:
27             print(f"Rate error: {censor_string(query)} in {region}")
28         else:
29             print(f"Rate error: {query} in {region}")
30
31         time.sleep(60)
32         get_region_trend(query, region, censor)
```

## C.4   Utility Functions

```
1  import itertools as it
2  import requests
3  import random
4
5
6  # Returns list of N random words from MIT dictionary
7  def random_words(n):
8      word_site = "https://www.mit.edu/~ecprice/wordlist.10000"
9
10     response = requests.get(word_site)
11     words = response.text.splitlines()
12
13     keywords = random.sample(words, n)
14
15     return (keywords)
16
17
18 # Tests if computer is connected to internet (used in tests)
19 def connected():
20     url = "http://google.com"
21     timeout = 5
22     try:
23         requests.get(url, timeout=timeout)
24         return (True)
25     except (requests.ConnectionError, requests.Timeout):
26         return (False)
27
28
29 # Censors strings so that senstive words aren't uploaded to github / used in scripts
30 def censor_string(string):
31     return (string[0] + "_" * (len(string) - 2) + string[-1])
32
33
34 # credit to Ilja Everila for this implimentation
35 # https://stackoverflow.com/questions/48381870/a-better-way-to-split-a-sequence-in-chunks-
        with-overlaps
36 def itr_split_overlap(iterable, size, overlap):
37
38     if overlap >= size:
39         raise ValueError("overlap must be smaller than size")
40
41     itr = iter(iterable)
42
```

```
43      next_ = tuple(it.islice(itr, size))
44
45      yield next_
46
47      prev = next_[-overlap:] if overlap else ()
48
49      while True:
50          chunk = tuple(it.islice(itr, size - overlap))
51
52          if not chunk:
53              break
54
55          next_ = (*prev, *chunk)
56          yield next_
57
58          if overlap:
59              prev = next_[-overlap:]
```

# D   Analysis Code

```
1  # Load Libraries for Analysis
2  library(tidyverse)
3  library(lubridate)
4  library(broom)
5
6  # Load utility functions
7  source("utils.R")
8
9  # Creat tidy dataframe of years / markets in which sinclair was present
10 sinclair_present <- tibble(filename = list.files(path = "../data/sinclair_data/", full.names
       = T, pattern = "*.csv")) %>%
11   mutate(
12     year = str_extract(filename, "[0-9]+") %>% as.double(),
13     data = filename %>% map(~ read_csv(.x, locale = locale(encoding = "UTF-8"))),
14     data = map(data, ~ rename_with(.x, ~ gsub("\n", "", .x))),
15     data = map(data, ~ .x %>% rename_with(tolower)),
16     data = map(data, ~ .x %>% dplyr::select(market))
17   ) %>%
18   unnest(data) %>%
19   select(-filename) %>%
20   mutate(
21     sinclair_present = T,
22     market = gsub("\n", "", market)
23   )
24
25 # Loads in Traslation dictionary of market names in the Sinclair dataset to DMA
26   # codes
27 sinclair_codes <- read_csv("../data/sinclair_names.csv")
28
29 # Loads in Traslation dictionary of DMA codes to Standardized names
30 dma_names <- read_csv("../data/dma_list.csv")
31
32 # Adds codes to sinclair present datset
33 sinclair_present <- sinclair_present %>%
34   full_join(sinclair_codes) %>%
35   dplyr::select(-market)
36
37 # Fills in data to include stations where sinclair was not present
38 sinclair_data <- sinclair_present %>%
39   complete(code = dma_names$code, year = 2005:2020) %>%
40   arrange(code) %>%
41   replace_na(list(sinclair_present = F)) %>%
42   left_join(dma_names)
43
44 write_csv(sinclair_data, "../data/clean_sinclair_data.csv")
45
46 sinclair_data <- read_csv("../data/clean_sinclair_data.csv")
47
48 # Not done annotating / updating after here
49 ################################################################################
50 between_data <- read_csv("../data/google_trends_data/set_1_between_regions.csv") %>%
51   pivot_longer(-code, names_to = "term", values_to = "overall") %>%
52   mutate(term = censor_string(term))
53
54 search_data <- read_csv("../data/google_trends_data/set_1_time_serires.csv") %>%
```

```r
55     mutate(year = year(date)) %>%
56     mutate(term = term) %>%
57     group_by(year, code, term) %>%
58     nest() %>%
59     mutate(
60       score = data %>% map_dbl(~ mean(.x$score, na.rm = T)),
61       code = str_extract(code, "[0-9]+") %>% as.numeric()
62     )
63
64  search_data <- search_data %>%
65     dplyr::select(-data) %>%
66     distinct(term, code, year, .keep_all = TRUE) %>%
67     pivot_wider(names_from = term, values_from = score) %>%
68     mutate(overall_score = rowSums(across(everything()), na.rm = T))
69
70
71  full_data <- search_data %>%
72     right_join(sinclair_data) %>%
73     filter(year != 2021) %>%
74     group_by(market) %>%
75     mutate(years_before = years_before(sinclair_present))
76
77
78  full_data %>%
79     ggplot(aes(x = years_before, y = n____r, group = market)) +
80     geom_line()
81
82
83  full_data %>%
84      nest(-code) %>%
85      mutate(
86            any_true =  data %>% map_lgl(~any(.$sinclair_present)),
87            any_false =  data %>% map_lgl(~any(!.$sinclair_present)),
88            changed = any_true & any_false
89      ) %>%
90      unnest() %>%
91     ggplot(aes(x = year, y = overall_score, col = sinclair_present, group = market)) +
92     scale_color_manual(values = c("grey", "black")) +
93     geom_line() +
94     facet_wrap(~market) +
95     theme(
96       plot.background = element_rect(fill = "transparent", colour = NA),
97       panel.border = element_rect(colour = "black", fill = NA, size = 0.1),
98       panel.background = element_blank(),
99       legend.key = element_rect(fill = "transparent", colour = NA),
100      legend.background = element_rect(fill = "transparent", colour = NA),
101      legend.position = "bottom",
102      panel.grid = element_blank(),
103      panel.grid.minor = element_blank(),
104      panel.grid.major = element_blank()
105    )
106
107 lm(overall_score ~ as.factor(year) + as.factor(code) + sinclair_present, data = full_data)
        %>%
108    summary()
109
110 lm(overall_score ~ as.factor(year) + as.factor(code) + sinclair_present + year:as.factor(code
        ), data = full_data) %>%
111    summary()
112
113 lm(overall_score ~ as.factor(year) + as.factor(code) + as.factor(years_before), data = full_
        data) %>%
114    tidy() %>%
115    filter(grepl("years_before", term)) %>%
116    mutate(years_before = as.numeric(str_extract(term, "-?[0-9]+"))) %>%
117    ggplot(aes(x = years_before, y = estimate)) +
118    geom_point() +
119    geom_line() +
120    geom_errorbar(aes(x = years_before, ymin = estimate - 1.96 * std.error, ymax = estimate +
        1.96 * std.error)) +
121    geom_hline(yintercept = 0, linetype = 2, col = "grey") +
122    geom_vline(xintercept = 0, linetype = 1)
```

## D.1   Utility Functions

```r
1 censor_string <- function(string) {
```

```r
 2    substr(string, 2, nchar(string) - 1) <- paste0(rep("_", nchar(string) - 2), collapse = "")
 3    names(string) <- NULL
 4    return(string)
 5  }
 6  censor_string <- Vectorize(censor_string, USE.NAMES=F)
 7
 8  years_before <- function(bool) {
 9    if (any(bool) & any(!bool)) {
10      out <- numeric(length = length(bool))
11      start <- min(which(bool))
12
13      before <- (seq(start, 1) - 1) * -1
14      if (start != length(bool)) {
15        after <- seq(1, length(bool) - start)
16      } else {
17        after <- c()
18      }
19
20
21      out <- c(before, after)
22    } else {
23      out <- rep(-99, length(bool))
24    }
25
26    return(out)
27  }
```

# E   Unit tests

## E.1   For Python Code

```python
 1  import unittest
 2  import pandas as pd
 3  from between_regions import between_region
 4  from utils import connected
 5
 6
 7  @unittest.skipIf(not connected(), "not connected to the internet")
 8  class TestBetweenRegion(unittest.TestCase):
 9      def test_between_region(self):
10
11          result_1 = between_region(["socks"],
12                                    timeframe="2016-12-14 2017-01-25",
13                                    geo="US",
14                                    gprop="")
15          result_2 = between_region(["socks", "shoes", "fish"],
16                                    timeframe="2016-12-14 2017-01-25",
17                                    geo="US",
18                                    gprop="")
19
20          expected_1 = pd.read_parquet("tests/test_data/between_region_1.parquet")
21          expected_2 = pd.read_parquet("tests/test_data/between_region_2.parquet")
22
23          self.assertTrue(expected_1.equals(result_1))
24          self.assertTrue(expected_2.equals(result_2))
```

```python
 1  import unittest
 2  from utils import itr_split_overlap, censor_string
 3
 4  class TestItrSplitOverlap(unittest.TestCase):
 5      def test_itr(self):
 6          test_list = ["one", "two", "three", "four", "five"]
 7
 8          expected_result = [('one', 'two'), ('two', 'three'), ('three', 'four'),
 9                             ('four', 'five')]
10          self.assertEqual(expected_result,
11                           list(itr_split_overlap(test_list, 2, 1)))
12
13          expected_result = [('one', 'two', 'three', 'four'),
14                             ('two', 'three', 'four', 'five')]
15          self.assertEqual(expected_result,
16                           list(itr_split_overlap(test_list, 4, 3)))
17
18      def test_exceptions(self):
```

```
19          with self.assertRaises(ValueError):
20              test_list = ["one", "two", "three", "four", "five"]
21
22              list(itr_split_overlap(test_list, 2, 3))
23
24
25 class CensorString(unittest.TestCase):
26     def test_censor(self):
27         self.assertEqual("h_y", censor_string("hey"))
28         self.assertEqual("f____r", censor_string("fender"))
29         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 })
```