

NHH



# Explaining Election Outcomes Using Web Search Data:

Evidence from the U.S. Presidential Elections 2008 - 2016

**Hannah Skaar Hauge s135897, Therese Borge Lied s136139**

**Supervisor: Po Yin Wong**

Master thesis, MSc in Economics and Business Administration,  
Economics and Finance

NORWEGIAN SCHOOL OF ECONOMICS

This thesis was written as a part of the Master of Science in Economics and Business Administration at NHH. Please note that neither the institution nor the examiners are responsible – through the approval of this thesis – for the theories and methods used, or results and conclusions drawn in this work.



## Abstract

Can Google Trends data be used to proxy socially sensitive sentiments, and can such proxies explain election outcomes? This thesis analyzes the effects of different social sentiments, proxied by Google search data, on outcomes for the Republican Party in the 2008, 2012 and 2016 U.S. presidential elections.

We assess the relationship between three socially sensitive sentiments and the outcome of presidential elections. The sentiments we examine are racial animus, immigration skepticism and far-right sentiment. We use data on the search terms “nigger” or “niggers”, “Breitbart News”, “Stormfront” and “Drudge Report”, and “Illegal immigration and residence” ahead of the three latest U.S. presidential elections to assess the prevalence of these sentiments. We look at the sentiments in both a long-term context, over a year, and short-term context, over two weeks.

Using a fixed effects model controlling for both state and time fixed effects, we find a positive effect of long-term immigration skepticism, and both long-term and short-term far-right sentiment, on the Republican election outcome. The estimated effects are small in magnitude. However, in the context of presidential elections, where a change of 1 percentage point can alter the election outcome, even small effects have potentially large consequences. Thus, our findings should be of value to both opinion pollsters and campaign strategists. Also, our analysis shows that higher increases than 1 index point in the proxied social sentiments should be regarded when interpreting the estimated effects, suggesting that the actual effect on the election outcome is likely larger in magnitude.

The findings presented are especially interesting in two regards. Firstly, they contribute to the existing literature on the use of Internet search data in predicting and explaining election outcomes, as well as the literature on determinants of voting. Secondly, they bolster the argument for consideration of web search data in future election predictions and analyses. Further, the size, variation and availability of search data increases constantly due to continuing and substantial growth in online searches.



# Table of Contents

<b>1</b>	<b>Introduction</b>	<b>7</b>
1.1	Motivation and Purpose	7
1.2	Research Question	8
1.3	Literature Review	9
<b>2</b>	<b>Historical Background</b>	<b>12</b>
2.1	U.S. Presidential Election in 2008	12
2.2	U.S. Presidential Election in 2012	13
2.3	U.S. Presidential Election in 2016	14
<b>3</b>	<b>Voting Mechanisms and Theoretical Frameworks</b>	<b>17</b>
3.1	The Determinants of Voting	17
3.2	Opinion Polls	18
3.3	Social Desirability Bias	19
3.4	Determinants of Information Seeking	21
<b>4</b>	<b>The Socially Sensitive Sentiments</b>	<b>24</b>
4.1	Far-Right Sentiment	24
4.2	Immigration Skepticism	25
4.3	Racial Animus	26
<b>5</b>	<b>Google Trends Data</b>	<b>28</b>
5.1	Google Trends	28
5.2	Measuring Social Sentiments Using Google Trends Data	30
<b>6</b>	<b>Data</b>	<b>33</b>
6.1	The Google Trends Proxies	33
6.2	Polling Data	43
6.3	Economic and Demographic Control Variables	43
<b>7</b>	<b>Empirical Strategy</b>	<b>49</b>
7.1	Estimation with Panel Data	49
7.2	The Estimation Model	50
<b>8</b>	<b>Empirical Analysis</b>	<b>53</b>
8.1	Main Results	53
8.2	Testing for Regions	56
8.3	Testing for Swing States	59
8.4	Testing for each Presidential Election	60
8.5	Limitations to the Estimation Strategy	62
<b>9</b>	<b>Discussion</b>	<b>63</b>
9.1	Discussion of the Results	63
9.2	Limitations to the Data Set	67
9.3	External Validity of the Study	69
<b>10</b>	<b>Conclusion</b>	<b>71</b>
	<b>References</b>	<b>73</b>
	<b>Appendix</b>	<b>81</b>



# 1 Introduction

## 1.1 Motivation and Purpose

On several occasions over the past few years, opinion polling prior to elections has failed in predicting the correct election outcome (Zukin, 2015). During the 2014 midterm election in the U.S., the pollsters did not capture the Republican support that led to strong Republican majorities in both the Senate and the House of Representatives. In the 2015 Israeli Legislative election, polls wrongly predicted the failure of prime minister Benjamin Netanyahu. In 2016, the polls failed in anticipating the outcome of the U.S. presidential election and the EU referendum vote in the UK (Zukin, 2015). Understanding and anticipating voting behavior contributes to a stable and foreseeable economy. Unexpected election results can give large fluctuations in corporate stock prices, yields and currencies (Scott (2017); Kiersz (2012)). From an economic aspect, it will thus be desirable to predict election outcomes more precisely.

The support for far-right populist parties and leaders<sup>1</sup> has increased in Europe, Canada, New Zealand and the U.S. over the past two decades (Rydgren, 2008). Across Europe, the average share of vote in national and European parliament elections for what can be defined as populist parties, has more than doubled since the 1960's, from an average of 3.8 percent of the vote share, to 12.8 percent (Inglehart & Norris, 2016). Populist leaders like Marine Le Pen, Norbert Hoffer, Geert Wilders and Donald Trump have changed the rules of political competition in several modern Western Societies (Inglehart & Norris, 2016). This suggests that substantial change is occurring in public sentiment towards the far-right political agenda, a change that has not entirely been captured by opinion polls and public surveys (Zukin, 2015).

Research on the topic of social desirability bias and socially sensitive sentiments has identified several topics that tend to yield high measurement errors in public surveys and opinion polls (Kreuter, Tourangeau, & Presser, 2008). Among such topics are voting intention and political affiliation (Brownback & Novotny, 2017). In addition, anticipating the emergence of far-right or far-left political sentiments is often challenging as the majority of

---

<sup>1</sup> Inglehart & Norris identify populist parties and leaders that share three core features: Authoritarianism, nativism and anti-establishment.

people will hesitate to admit to views that violate the social norm in the society, as far-right or far-left opinions tend to do (Krumpal, 2013). This complicates the use of opinion polls and surveys as they might not capture underlying sentiments that could be an important factor in predicting election outcomes. In situations where traditional survey based methods like public opinion polls yield high measurement errors, the use of a non-survey based measure can be a helpful supplement.

Google searches provide a lot of information on public opinion and the sentiment in a population. It has for instance been successfully used as a tool for measuring racial animus (Stephens-Davidowitz, 2013) and for predicting anti-muslim hate crime (Stephens-Davidowitz & Soltas, 2015). Google search data has been shown to have predictive power in forecasting consumer behavior and the election turnout for specific minority groups, and it has been used as a non-survey based measure of public opinion in election prediction models (Choi & Varian (2012); Stephens-Davidowitz (2013); Shimshoni et al. (2009); Goel et al. (2010); Chen et al. (2012)).

This thesis is a case study of the past three presidential elections in the U.S., where we further examine the use of Google searches as a non-survey based measure of public sentiment. We have chosen the U.S. presidential elections due to the large volume of Google searches in the U.S. and because of the failure of the polls to predict the correct election outcome during the 2016 presidential election (Ad Hoc Committee on 2016 Election Polling, 2016). The purpose of the study is to explore the more general question of whether social sentiments that are difficult to capture by survey based measures are reflected in election outcomes, and to examine how Google searches can be used to proxy social sentiments in a population.

## **1.2 Research Question**

Based on the motivation and purpose of this thesis, we attempt to answer the following question:

*“Can Google Trends data be used as a proxy for socially sensitive sentiments, and can such proxies be used in models explaining election outcomes?”*

where we define socially sensitive sentiments as sentiments that conflict with the social norm in a society. The term is defined in detail in section 3.2.

In the following section, we review previous literature on the subject of social sentiment affecting election outcome and on the use of web search data in research on social sentiments and elections. We present some historical background on the U.S. presidential elections in 2008, 2012 and 2016 and highlight key points from each of the elections in chapter 2. In chapter 3 we present some mechanisms behind voting decisions, a theory explaining why survey based methods often fail to capture public sensitive sentiments correctly, and two approaches for explaining internet search behavior. In chapter 4 we present the social sentiments we proxy using Google search data. Google search data and its applications are described in chapter 5. In chapter 6 we describe our data and elaborate on how the sentiment proxies are created. We then present our empirical strategy in chapter 7. Our results are described in chapter 8. In chapter 9 we discuss the results retrieved. In chapter 10 we will make a conclusion based on our research question.

### **1.3 Literature Review**

One of the first papers published on the topic of the use of web search data in forecasting economic statistics, was Ettredge et al. (2005), which applied web search data to the U.S. unemployment rate. After this, several empirical studies have used search data for forecasting in different fields, for instance in the consumer market (Choi and Varian (2012); Goel et al. (2010)) and regarding health issues (Cooper et al. (2005); Polgreen et al. (2008); Ginsberg et al. (2009); Brownstein et al. (2009)). Both Hal Varian and Seth Stephens-Davidowitz have contributed largely to the field of Google Trends analysis, how to use the data and its predictability powers, see especially Choi and Varian (2012) and Stephens-Davidowitz and Varian (2015). Shimshoni et al. (2009) has also contributed on the matter of predictability of search trends, focusing on how seasonal decomposition methods can give predictive power to a large amount of search terms.

Several researches have looked at how internet behavior can be used in explaining election outcomes affected by social sentiments. Stephens-Davidowitz (2013) examines the relationship between racial animus and election outcome in the 2008 and 2012 U.S. presidential election. Using data on Google search queries containing racially charged language he finds that racial animus negatively affected the share of votes received by Barack Obama, the first African American U.S. presidential candidate. He concludes that racial

animus cost Obama around 4 percentage points of the national popular vote in both the 2008 and the 2012 election. Research has also shown that social sentiments, such as feelings towards a specific gender or a religion, and xenophobia, influences voting decisions (King & Leight (2010); Berggren et al. (2010); Rydgren (2008)).

Stephens-Davidowitz (2013) investigates how Google searches prior to an election can be used to proxy voting intention. By comparing the Google search rates for [vote] or [voting] before the 2010 U.S. midterm election to the same search rates prior to the 2006 election, he finds that the search rate explained 20-40 percent of state-level changes in turnout rates. C. Douglas Swearingen and Joseph T. Ripberger (2014) proposes a new indicator of public attention to electoral candidates in the U.S. Senate elections prior to elections. Their index is based on the relative pattern of internet queries for the different candidates. They find that their proposed index behaves in a manner consistent with a credible measure of public attention. This finding holds when they include the index in a model explaining U.S. Senate election outcomes. Chen et al. (2012) examines how to predict the results of the U.S. presidential elections in 2012 using Google Trends data. They describe how different search terms related to the U.S. economy and candidate policies can be classified, using both supervised and unsupervised learning methods, where they found a support vector machine model to be the most efficient. Search data have also been used to measure social sentiments directly. In his book on how web search data can be used to observe and measure social sentiments, Stephens-Davidowitz (2017) explains how social biases and preferences tend to be eliminated online. By using internet searches, one can measure social preferences, behavior and sentiments that are hard to measure using survey based methods.

A fixed effects model for predicting and explaining election outcomes was proposed by Strumpf and Philippe (1999). They show that state partisan predisposition was the most important explanatory variable for election outcome in the period 1972-1992, and with that highlight the importance of using fixed effects in election models, due to the bias caused by time-invariant factors within a state. Strumpf and Philippe provide a utility model for explaining voter's choice, which they use in explaining how economic and demographic factors affect election outcomes.

Our research is closely related to the empirical work done by Stephens-Davidowitz (2013) and Choi and Varian (2009) on how Google search data can provide information about the social sentiment in a population, and on how web search data can be a tool for proxying socially sensitive sentiments. This thesis applies the relationship between social sentiments and voting decision provided in the literature. We use the index for racial animus proposed by Stephens-Davidowitz (2013), and contribute further to the research by proposing new measures for immigration skepticism and the far-right sentiment in a population. We apply the use of such proxies in a general model for explaining election outcome. Using a similar approach as Stephens-Davidowitz (2013) and C. Douglas Swearingen and Joseph T. Ripberger (2014), we propose a method for measuring socially sensitive sentiments in a population by indexing search rates. Our empirical method is based on the fixed effects model proposed by Strumpf and Philippe, but we extend the model by adding measures of socially sensitive sentiments. We also include polls data, arguing that this is the most used tool for predicting election outcomes.

## 2 Historical Background

In the following chapter, we will provide some background for the empirical analysis from the U.S. presidential elections in 2008, 2012 and 2016. We discuss the election results and how accurate the polls were in predicting them, and we briefly present some of the major issues during the respective presidential campaigns. See Figure A1 in the appendix for a graphical display of presidential election outcomes in 2008 to 2016, segmented on regions and states.

In the United States, the president is chosen through a process called the Electoral College (The United States Government, 2017). Each state has a certain number of electors based on how many members of Congress the state has. Each elector has one electoral vote. There are 538 electors in all. In order to become president, a candidate needs the vote of at least 270 electors. The political parties in each state choose their potential electors. During the presidential election people vote for either a Democratic or a Republican elector in their state<sup>2</sup>. In 48 out of 51 states, the electoral candidate who receives the highest total amount of votes, gets all the electoral votes. Thus, it is the election outcome in each state that determines who becomes president. The implications of this is that the candidate that gets the highest total number of votes, i.e. wins the popular vote, does not necessarily win the the presidency.

### 2.1 U.S. Presidential Election in 2008

The Democratic nominee in the 2008 U.S. presidential election was Barack Obama, a U.S. senator from Illinois. He ran against the Republican U.S. senator from Arizona, John McCain (Federal Election Commission, 2009). In 2008, the incumbent president was the Republican George W. Bush.

Barack Obama won the popular vote with 52.93 percent of the national vote and secured a total of 365 electoral votes, while John McCain received 173. He won the national vote by 7.28 percentage points (Federal Election Commission, 2009). With that, Obama received the largest percentage of the popular vote for a Democrat since 1964 (Nagourney, 2008).

---

<sup>2</sup> It is also possible to vote for candidates from other parties, e.g. the Libertarian Party, but we choose to disregard this from this deliberation on the presidential elections in chapter 2.

Many of the nationwide polls prior to the 2008 election pointed towards the election victory of Obama in November (Ejara, Nag, & Upadhyaya, 2008). In the final week before the election, every major polling company predicted that Obama would win with the popular vote with between 2 to 13 percentage points<sup>3</sup>. On average the polls predicted a win of 7.6 percentage points, which was only 0.3 percentage points away from the actual results (Real Clear Politics, 2008).

Among the most debated issues during the 2008 presidential campaign, were the financial crisis, health care and the war in Iraq. In 2008, it was clear that the world was facing the biggest financial crisis since the Great Depression. This became a big issue during the presidential campaign, with the candidates proposing different measures to limit the crisis (CNN, 2008). Obama pledged to create a national health insurance program for individuals who did not have health care provided through their employer, or did not qualify for other existing health care programs (CNN, 2008). This later became known as Obamacare and was unpopular amongst conservatives (BBC, 2017). Barack Obama also pledged to remove all troops from Iraq by the summer of 2010, while McCain did not believe in setting a withdrawal timetable (CNN, 2008).

Barack Obama was the first African American to run for presidency in the U.S., and it is of interest in this election to know whether or not racism was an issue during the election campaign, since we in our thesis examine if racial animus to some extent can explain the election outcome. With some exceptions, explicit racism was not a part of the 2008 campaign. It has however been argued by social psychologists that implicit biases towards black people did, through affecting how people evaluate each other, play a role in the 2008 presidential election (Parks, Rachlinski, & Epstein, 2009).

## **2.2 U.S. Presidential Election in 2012**

In the 2012 U.S. presidential election, Barack Obama ran for his second term against the former Republican Massachusetts Governor and businessman, Mitt Romney (Federal Election Commission, 2013).

---

<sup>3</sup> In the week before the election the following companies made predictions on the outcome of national vote: Marist, Battleground (Lake and Tarrance), Rasmussen Reports, Reuters, IBD, FOX News, Wall St. Jnl, Gallup, Diageo, CBS News, Ipsos, ABC News, CNN/Opinion Research, Pew Research.

On the election day November 6, the result of the national vote was 51.06 percent to Obama and 47.20 percent to Romney, giving Obama a victory of 3.86 percentage points (Federal Election Commission, 2013). In the electoral vote, Obama secured 332 votes while Romney got 206 electoral votes.

The average of national polls taken in the week before the election had predicted Obama to win by a margin of 0.7 percentage points. Two of the nine national polls conducted in the week before the election predicted a Republican victory, while three predicted a tie (Real Clear Politics, 2012). Most of the national polls were somewhat biased towards the Republican candidate during the entire presidential campaign (Enten, 2012). One of the national polling firms that largely overestimated the support for Romney, was Gallup. In a final survey, Gallup gave Romney a 1 percentage point lead on Obama, when Obama in reality won by nearly 4 percentage points (Blumenthal & Edwards-Levy, 2013).

The economy was a highly debated issue during the 2012 presidential campaign. While Obama advocated for government interference to stimulate economic growth, Romney argued that short-term stimulus does not work, and that it only increases government debt (Lauter, 2012). Among the other central issues was healthcare, foreign policy and immigration. A great part of the immigration debate was focused on what to do with the estimated 11 million illegal immigrants, mostly latinos, in the United States.<sup>4</sup> Moral value issues such as gay rights, abortion and stem cell research also played a role in the 2012 election (Lauter, 2012).

### **2.3 U.S. Presidential Election in 2016**

In the 2016 presidential election the Republican candidate was the New York based businessman Donald Trump. The Democratic candidate was Hillary Clinton, former First Lady and Secretary of State (State Election Office, 2017).

Hillary Clinton won the popular vote by 2.1 percentage point. She got 48.18 percent of the national vote while Donald Trump got 46.09 percent. Despite winning the popular vote, Hillary did not secure the 270 electoral votes necessary to secure the presidency. Trump got 304 electoral votes against Clinton's 227 votes (State Election Office, 2017). The election

---

<sup>4</sup> 11 million illegal immigrants are 2012 estimations (Lauter, 2012).

was extremely tight. 77,744 votes combined in the states of Pennsylvania, Michigan and Wisconsin gave Trump the 46 electoral votes he needed to win (Sabato, Kondik, & Skelley, 2017). Losing the national popular vote but winning the electoral vote, like Donald Trump did, has rarely happened in American history.

During the presidential race, pre-election polls stimulated high-profile predictions that Hillary Clinton's probability of winning the election was about 90 percent. When it became clear on the election day that Donald Trump was the winner, it surprised even his own pollsters (Ad Hoc Committee on 2016 Election Polling, 2016). In the week before the election, most pollster predicted that Hillary Clinton would win the popular vote by between 1 and 5 percentage points (Real Clear Politics, 2016). On average they predicted the Democrats to win with 3.2 percentage points, 1.1 percentage points higher than the actual result <sup>5</sup>.

In the contest for electoral votes, statewide polls showed Hillary Clinton leading, but with a smaller margin than what the national polls predicted. The polls indicated on average that Donald Trump was one state away from winning the election. In hindsight, the polls overestimated the Democratic vote in what was assumed to be Hillary's "blue wall": Pennsylvania, Michigan and Wisconsin. These states had voted Democratic in the past elections, and pollsters saw no sight of this election being any different. Donald Trump also did better than expected in battleground states like Florida, Ohio, North Carolina and Iowa, states that the pollsters predicted would vote Democratic (Ad Hoc Committee on 2016 Election Polling, 2016).

There are several explanations as to how the polls missed this. One explanation is that the turnout for Trump voters were higher than expected. Another is that last minute voters chose Trump instead of Hillary (Ad Hoc Committee on 2016 Election Polling, 2016). The turnout grew more in heavily Republican areas than in the Democratic ones relative to 2012, and a number of polls were adjusted to align with turnout patterns from 2012. The polls might also have underestimated the turnout among rural whites while overestimating the turnout among for example African Americans (Ad Hoc Committee on 2016 Election Polling, 2016). All explanations as to why the polls failed suggest that there has been some underlying sentiment

---

<sup>5</sup> The average is based on the polling numbers published within a week before the election (Real Clear Politics, 2016).

in certain states that prompted people to vote for Trump instead of Hillary, that the opinion-polls did not capture.

According to a survey conducted by Pew Research Center, the top five voting issues during the 2016 election was the economy, terrorism, foreign policy, health care and gun policy (Pew Research Center, 2016)<sup>6</sup>. Among the registered Republicans, i.e Trump voters, immigration and foreign policy is given higher priority than among Clinton voters. Registered Democrats are more concerned about the treatment of racial, ethnic minorities and the environment compared to Republican voters (Pew Research Center, 2016)<sup>7</sup>.

A lot of attention in the 2016 election was on Donald Trump and his personal image (Sabato, Kondik, & Skelley, 2017). His unpolished speaking style and populist approach separates him from most of America's previous presidential candidates. His support among white nationalist and other far-right groups have fueled the perception that racial resentment and hostility towards globalization and undocumented immigrants were strong forces benefiting Trump (Sabato, Kondik, & Skelley, 2017).

---

<sup>6</sup> Based on registered voters, Q40, survey conducted June 15-26, 2016, % of registered voters saying each is "very important" to their vote in 2016" (Pew Research Center, 2016).

<sup>7</sup> Registered Republicans vs. registered Democrats: Immigration; 79 vs 65, Foreign Policy: 79 vs 71, Treatment of racial, ethnic minorities 42 vs 79, Environment: 32 vs 69, based on registered voters saying each is very important to their vote in 2016 (Pew Research Center, 2016).

## 3 Voting Mechanisms and Theoretical Frameworks

In the following chapter, we present a theoretical framework for the empirical analysis. We discuss determinants of voting and explain how opinion polls prior to elections usually are conducted, as well as looking at some of the implications of using opinion polls. We suggest a theory for explaining why opinion polls in some cases fail in predicting the correct outcome, called the social desirability bias. We present two theories for understanding what drives people to search for certain things online: confirmation bias and information utility.

### 3.1 The Determinants of Voting

Identifying the determinants of voting are of importance in an economic aspect. In a functioning democracy, it is the aggregated preferences of the voters that decide the political agenda (Blais, 2000). Whether the voters value the environment more than infrastructure, or if they prefer public welfare over low taxes, has economical consequences on both country and business level. Unexpected election results can also give large fluctuations in corporate stock prices, yields and currencies (Scott (2017); Kiersz (2012)). Understanding and predicting voting behavior thus contributes to a stable and foreseeable economy.

Rational choice theory says that personal gains fully determines voting and that voters will re-elect candidates who deliver personal gains for them, i.e they will vote rationally out of their economic self-interest (Gelman & Kaplan, 2008). Several researches have found that not to be true. One example of non-rational voting behavior is that national economic factors seem to matter more to the voter than regional factors, which in reality affects the voter more (Wolfers, 2002). Several non-economic factors also matter to the voters. Among them are physical appearance of the candidates, gender, party affiliation and race (Rosar, Klein & Beckers (2008); Streb et. al (2008); Kever (2017); Stephens-Davidowitz (2013)). The notion that people do not vote rationally, makes predicting election outcome harder than if people voted solely based on personal expected economic outcome.

Rational choice theory relies on people having perfect information about their situation (Gelman & Kaplan, 2008). One problem with this assumption is that voters increasingly can choose which information they wish to be exposed to, and thus their worldview will be biased in regards of the information they are exposed to. This has been called the echo-chamber effect, and is a consequence of the exponential growth of online news sites and forums, as

well as the spreading of information on social media (Pariser, 2011). A large source to the echo-chamber effect is algorithms employed by companies like Google and Facebook. These companies aim to serve tailored content to their users, i.e provide information that is in line with the interest of the user. Personalized algorithms allows people to quickly obtain the information they want without having to shift through irrelevant content (Pariser, 2011). The implication of this is that people to a less degree is exposed to information that nuances their views. Thus, two voters can have a very different picture of their situation, while in reality their situations are identical. This further complicates the use of actual economic and demographic situations when identifying the determinants of voting, as the actual situation might not reflect the voters experienced situation. Another factor that affects the experienced situation of the voters is the confirmation bias, which we will discuss further in section 3.4.

Based on the theory provided on determinants of voting, we can conclude that analyzing voting behavior and election outcomes are interesting from an economic perspective. Research on determinants of voting shows that one cannot rely on the rationality of the voters when predicting and explaining election outcomes. It is therefore of interest to examine methods for identifying determinants of voting that could help us better predict and explain voting behavior.

### **3.2 Opinion Polls**

It is common to use opinion polls and public surveys to predict voting behavior (Rogers & Aida, 2012). Opinion polls prior to elections play a significant role in deciding winners of political televised debates. They influence electoral turnout and also affects how candidates advance with their political campaigns (Brownback & Novotny, 2017). In the 2016 election, the polls showed a significant edge for the the Democratic candidate, Hillary Clinton, in some upper Midwest states, causing her to forgo campaigning in these states that the Republicans initially won. Market prices also fluctuate in response to opinion polls, since they can be perceived as forecasts on future business environment (Kantchev & Whittall, 2017). Election results viewed as unlikely by opinion polls can therefore lead to market shocks <sup>8</sup>.

---

<sup>8</sup> The European stock market indices were in a slight upward movement in the first half of 2016 until the beginning of June. In response to the Brexit vote stock prices dropped by around 10 percent (Raddant, 2016).

An opinion poll is a scientific survey designed to identify and measure the views of a group of people (BBC , 2016). Election polls are usually conducted by polling companies. Among the major polling companies in the U.S. are The Gallup Poll, Mason-Dixon, Rasmussen, PPP and SurveyU.S.A (Electoral Vote, 2017). The polling company typically get a request from a client to conduct a poll, and then agrees with the client on polling questions and method and on how many people to include in the survey. Traditionally, polls are conducted through phone interviews with the help of computers. Computers cannot call cell phones, which has to be done manually, and as more and more people drop their landlines this is becoming a problem for pollsters (Electoral Vote, 2017). Some polling companies have started to conduct polls over the internet by asking people to sign up. This complicates randomization, but the companies often use careful normalization to remove the sampling bias, by for example treating each response from a woman as four if the sample is largely made up by men. State polls in the U.S. are usually conducted on rather small samples. The margin of error is usually between 3 and 6 percent for the sample sizes used in state polling (Electoral Vote, 2017).

Different polling companies use different formulations in their surveys, which makes comparing them problematic. One polling company might ask: If the presidential election was held today, would you vote for the Democratic or Republican candidate? While another company asks: If the presidential election was held today, for whom would you vote? These formulations can lead to different responses (Electoral Vote, 2017).

We have presented some of the implications of using public opinion polls prior to elections and pointed out some of the limitations to survey based, traditional opinion polling. This provides an understanding of the importance of accurate opinion polling, and of why depending solely on opinion polls when explaining election outcomes is likely to yield measurement errors.

### **3.3 Social Desirability Bias**

Over the past few years, opinion polling prior to election has on several occasions failed in predicting the correct election outcome, see section 1.1. The reliability of polls and surveys is not a new problem and the issue has been widely researched. The problems regarding sampling bias, methodology and questioning method are natural sources of errors and

misreporting, see section 3.2. Social scientists present an additional explanation to why opinion polls generate errors: preferences surveys can be subject to social desirability bias (SDB, hereafter) (Brownback & Novotny, 2017). SDB causes respondents to conceal preferences that are not perceived to be socially desirable. This can lead to misreporting of public opinions in polls and surveys prior to elections. Identifying questions affected by SDB is important when it comes to understanding the non-economic determinants of voting (Blais, 2000). It can also be used to improve election polls and thus make elections more predictable (Rogers & Aida, 2012).

Social desirability refers to making oneself look good in terms of prevailing the social norms defined within a society (Krumpal, 2013). A society can be defined as a grouping of individuals that share common interest and may have distinctive institutions and culture (New World Encyclopedia, 2017). A society might refer to an ethnic group, to a nation or to a broader cultural group, like the Western society. Personal interactions within a society create norms that translate into socially and undesirable behavior (Brownback & Novotny, 2017). Views that violate the social norm in the society are considered to be unsocial or socially unacceptable and are thus affected by SDB (Krumpal, 2013).

Self-reported intention to vote is often used as a dependent variable in research regarding political campaigns, but several researchers have found intention to vote to be a sensitive question yielding high measurement errors and non-response rates (Rogers & Aida, 2012). Not voting violates the social norm in the modern Western society that everyone should perform their civil duty and vote. Thus, people will be biased towards the socially correct answer: “Yes, I intend to vote” due to SDB (Rogers & Aida, 2012). In their research paper on the accuracy of voting, Belli, Traugott and Beckmann (2001) found that 20% of actual non-voters reported in a survey that they voted. Rogers and Aida (2012) points out the fact that people tend to overestimate the fact that they will perform a socially desirable behavior in the future, meaning that a significant fraction of people that say they will vote in an election does not vote. Researchers have also found that respondents more often claim to be indecisive when asked a question subject to SDB (Brownback & Novotny, 2017).

Researchers have found SDB in several political, economic and social contexts. Feelings toward as African American presidential candidates, female politicians and Jewish politicians are among the sentiments affected by SDB (Stehens-Davidowitz (2013); Steb et. al (2008);

Kane et. al (2004)). People tend to hide political preferences that are discriminatory when discriminating against the group in question is socially unacceptable, i.e. violates the social norm defined in a given society (Brownback & Novotny, 2017). Kane et. al (2004, s. 1) writes: "Although national surveys indicate that Americans have become more accepting of the prospect of a Jewish presidential candidate, this could reflect some voter's desire to be seen as having socially correct opinions...". Questions on topics such as immigration, abortion, gay marriage, sex and drug abuse have also been proven to be affected by SDB (Fisher, 1993). In the recent U.S. presidential election, researchers found marginally significant evidence that that SDB caused people to understate their agreement with Donald Trump in the pre-election polls, while they tended to overstate their agreement with Hillary Clinton (Brownback & Novotny, 2017).

In this section we have provided additional theory on why survey based opinion polling prior to elections can be inaccurate. The theory of SDB explains what can cause people to lie in surveys. Thus, opinion polls based on public surveys might not reflect the true opinion in a population. We have focused on how SDB affects the way people respond to socially sensitive questions, which provides an understanding for why obtaining additional ways of identifying and measuring socially sensitive sentiments in a population, is of interest in the context of explaining and predicting voting behavior.

### **3.4 Determinants of Information Seeking**

In this section we will present the theoretical frameworks of two different motivations for seeking information: confirmation bias and information utility.

#### *3.4.1 Confirmation bias*

Confirmation bias can be defined as "the seeking or interpreting of evidence in ways that are partial to existing beliefs, expectations, or a hypothesis in hand" (Nickerson, 1998, s. 175). The theory of confirmation bias has evolved from Festinger's (1957) work on cognitive dissonance, and seeks to explain the motivation for why people selectively exposes themselves for certain information over other. Festinger's theory is built upon the idea that people need cognitive consistency, or else a dissonance, or a mental unbalance, will emerge because the different cognitions deviates.

If a person is motivated by confirmation bias when searching for information, he will go to places where the chance of being exposed to contradicting information in regards to his hypothesis, prejudice or belief is minimized (Nickerson, 1998). Further, people driven by confirmation bias will tend to only seek information that they believe will confirm and strengthen their beliefs. If a person strongly believes in a hypothesis, and is exposed to places containing contradicting information, or other people having other hypotheses or opinions regarding the same subject, the person might fail to consider the relevance of this opinion or information. This phenomena is called restriction of attention, and represents the failure to taking likelihood ratios into account, according to Bayesian theory (Nickerson, 1998).

Even though cognitive dissonance and confirmation bias theory is widely researched and has strong empirical evidence, research also suggest that there exist other motivations for seeking information where the person is responsive to information which also contradicts his existing beliefs. Further, this research suggests that a person can even be motivated to actively seek information which deviates from existing beliefs. One of such motivations is called information utility (Knobloch-Westerwick & Kleinman, 2011).

#### *3.4.2 Information utility*

Atkin (1973) and Knobloch-Westerwick (2008) identifies four primary functions for the need for information: surveillance, performance, guidance and reinforcement. The surveillance factor implies that people need information to know about changes in the environment and to monitor potential threats. The performance factor implies that people need information in order to know how to execute different tasks. The guidance factor explains that people will need information in order to monitor their emotions, and know how to relate to and feel about different things. Reinforcement is a factor which lie closely to the confirmation bias, where the objective is that the need for information is to be able to confirm attitudes.

Knobloch-Westerwick and Kleinman (2011) demonstrates that people will be more willing to seek out information regardless of whether it will be consistent with existing beliefs or not, as long as the information is expected to be of beneficial value. An example of this is when information can help people make a more informed decision on who to vote for. In these situations, the information utility motivation will overrule the confirmation bias.

We will use the theory on confirmation bias, information utility and SDB to validate each of the proxies on socially sensitive sentiments used in the empirical analysis. Due to confirmation bias and information utility it would be likely to assume that information seeking online will grow as the election day approaches. This is because people are likely to search for information to help them make a decision.

Based on the theory of SDB, confirmation bias and information utility, we assume that search data will give us a good measure of sensitive sentiments in a population. We believe that an election model including proxies on socially sensitive sentiments in a population, will have a higher explanatory power compared to a model including only opinion polls. This is due to the theory provided on determinants of voting, opinion polls, SDB and determinants of information seeking.

## 4 The Socially Sensitive Sentiments

The purpose of this thesis is to build a proxy for the far-right sentiment in the different U.S. states, a proxy for the level of immigration skepticism and a proxy for racial animus, and examine how well these proxies fit in a model explaining election outcome. These sentiments are hard to measure due to the SDB, since questions on the topics are associated with views that violate the social norm in the modern western society (Krumpal, 2013). In the following chapter we explain the reasoning behind the choice of sentiments, and how these might generate SDB. It is important to note that we do not believe that the measurement errors in election polls are solely due to the incapability of measuring socially sensitive sentiments, or that the only sentiments that are hard to measure are associated with the radical far-right political side. However, for the purpose of examining the use of non-survey based measures in election research, we have limited our research to proxying the three sentiments mentioned.

### 4.1 Far-Right Sentiment

In the majority of the elections in recent years, where the polls were insufficient in predicting the election outcome, it was largely the far-right political side that was underestimated (Inglehart & Norris, 2016). Radical far-right parties have witnessed a markedly increase in popularity in Europe over the last three decades, and have re-emerged as an electoral force in Western Europe, Australia, Canada and New Zealand (Rydgren, 2008). The political situation in the United States stands out, as the majority of the political spectre consist of the Democratic and the Republican Party. Even though there are several parties on both the extreme right and the extreme left political side in the U.S. as well, these parties rarely compete for seats in Congress. Since World War II, only two out of the 535 member of Congress have been anything other than Republicans or Democrats (Blake, 2016). This makes it harder to track the progress of radical far-right political parties in the U.S. compared to other similar countries. However, the Republican candidate in the 2016 presidential election, Donald Trump, does share many similarities with radical far-right politicians in Europe. He is viewed by many as an anti-immigration, anti-globalization populist, like many of his radical far-right political colleges in Europe (Müller, 2017). Based on these observations, we assume that the radical far-right trend visible in Europa, New Zealand and Canada is applicable in the U.S..

The far-right or far-left wing of politics is defined as the extreme right or left wing of a political party or group (Carlisle, 2005). Far-right ideologies typically advocate the freedom of the individual, question the legitimacy of the democracy, reject social equality and the social integration of marginalized groups, and are associated with anti-immigration views, racism and anti-Semitism (Betz, 1994). In the U.S., the far-right wing consists of several marginal groupings, for instance white supremacist, white nationalists, neo-Nazis, the alternative-right and the Ku Klux Klan (Ford, 2017). A political opinion is characterized as far-right or far-left when it rejects the established socio-cultural and socio-political systems (Krumpal, 2013). In his paper on political correctness, Stephen Morris (2001) explains why admitting to views that support far-right ideologies would be socially undesirable, and thus subject to SDB. He explains that far-right ideologies are considered extreme and are marginal, and thus most people do not sympathize with such ideologies. This means that sympathizing with far-right ideologies would be viewed as socially incorrect and in conflict with the social norm (Morris, 2001). In general, people wish to remain in good standing with their society, and we would therefore suspect misreporting about the public opinion on far-right issues in polls and surveys due to the SDB (Yan & Tourangeau, 2007). We can therefore conclude that the far-right sentiment in a population might be better measured using a non-survey based method.

## **4.2 Immigration Skepticism**

Since 2011, we have witnessed history's largest refugee crisis since World War II (Egeland, 2014). This has led to a substantial increase in immigration, which has affected especially the countries and societies in Europe, but also the U.S. (Ostrand, 2015). Immigration was a hot topic during both the 2012 and 2016 presidential election in the U.S. (Agbafe (2016); Lauter (2012)). The immigration issue has also been central for radical far-right parties in Europe, like Front National in France, the Freedom Party of Austria and the Law and Justice party in Poland (Rydgren, 2008).

The desire to reduce immigration has been shown to be among the principal factors for predicting who will vote for a radical far-right party. Jens Rydgren (2008) has conducted a study on the importance of immigration on voters choice. The study uses election outcome for six radical far-right parties in Europe and self-reported immigration skepticism, and shows that people who wanted to allow only a few immigrants into their country, were

significantly more likely to vote for the far-right (Rydgren, 2008). Thus, measuring attitudes towards immigration is important when it comes to identifying the determinants of voting. Most research on the topic provides limited guidance as to which attitudes behind immigration skepticism who drives voters to vote for the far-right. Immigration skepticism has often uncritically been equated with racism, fascism and xenophobia, behaviors that are perceived as anti-social in the modern western society, and have been highly unpopular since World War II (Rydgren, 2005). In his paper, Rydgren highlights this issue. He points out the fact that ethno-nationalism<sup>9</sup> and opposition to the multicultural society seems to be of less importance to radical far-right voters compared to consequences of immigration like social unrest, unemployment and higher crime-rates. This generates an issue for voters who are immigration skeptical, but who do not want to be associated with xenophobic, racist or fascist views. Admitting to having immigration-skeptical views and thus risk being perceived xenophobic or racist, is socially undesirable since it violates the social norm in the society. Therefore, non-survey based methods for measuring immigration skepticism might give more accurate measures compared to survey based methods.

### **4.3 Racial Animus**

As defined by the Anti-Defamation League (2017), racism is the belief that a particular race is superior or inferior to another, and that a person's social and moral traits are predetermined by his or her inborn biological characteristics. In the United States, racism has mostly been targeted at the black population, with laws segregating the black and the white population and reducing African Americans to lower caste status (Fredrickson, 2002).

Like the question of radical far-right sympathies and immigration skepticism, survey based measures of racial animus are problematic due to the SDB. Negative feelings towards blacks are socially unacceptable, and individuals tend to withhold such feelings (Berinsky, 1999). Stephens-Davidowitz (2013) suggest a non-survey based measurement of racial animus in the United States using Google searches, and he finds evidence that racial animus cost Barack Obama 4.2 percentage points of the national popular vote in 2008 and 4.0 percentage points in 2012. This finding suggest that race is a factor in voter's choice in the U.S., and that racial animus is of interest when it comes to understanding the determinants of voting. A non-

---

<sup>9</sup> Ethno-nationalists believe that nations are defined by a shared heritage, which usually includes a common language, a common faith, and a common ethnic ancestry (Muller, 2008).

survey based measure of racial animus can also help with grasping the extent of contemporary prejudice (Stephens-Davidowitz, 2013). We wish to include a racial animus proxy in our model to see if Stephens-Davidowitz' findings holds when we include the 2016 election and the other social sentiments proxies, and do not control for the 2004 election outcome.

Even though the concepts of far-right sentiment, immigration skepticism and racial animus overlap, we believe that they separately might have explanatory power when it comes to explaining election outcome in the U.S.. While racism specifically cover people with negative feelings towards people of a different race, having radical far-right sympathies can mean everything from actively being a member of a Neo-Nazi groups to being anti-Government and opposing the establishment (Rydgren, 2008). People that are racist or sympathise with traditional radical far-right ideologies are prone to be immigration skeptical, but that does not mean that every person that is skeptical towards immigration is racist or a radical far-right sympathizer. Based on this argument, we believe that each of these sentiments represents an independent source of information on determinants of voting.

## 5 Google Trends Data

Google is by far the most popular search engine in the world, with more than 74 percent of the Global market share in 2017, according to Net Market Share (2017). In 2006, Google launched a new tool for downloading and analyzing Google searches, called Google Trends (Google, 2017). In the following chapter we explain the features of Google Trends and how one can use Google Trends to analyze search data. We further motivate the use of Google search data to measure social sentiments in a population.

### 5.1 Google Trends

#### 5.1.1 *The query index*

Google Trends provides a time series index of search frequency on specific terms and categories of terms across different geographic areas (Varian & Stephens-Davidowitz, 2015). The query index is given by the total query volume for the specific search term within the chosen geographical area, divided by the total number of search queries in that area during the given time period. The highest query share for a given time period is normalized to be 100 (Choi & Varian, 2012). A query share of for example 30, means that at that time, the query was 30 percent as popular as it was at the peak of the search frequency.

Stephens-Davidowitz and Varian (2015) explains that this normalization can lead to a negative trend in one search term over time, without this meaning that the overall searches for this query has decreased. It can mean that the search term has been less popular over time relatively to other search terms. The same applies for comparing regions. If a search term gets a higher query index for Rhode Island compared to California, this does not necessarily mean that in absolute numbers there are more searches for the query in Rhode Island. But relatively to other searches, it is a more popular search term in Rhode Island than in California (Varian & Stephens-Davidowitz, 2015)

#### 5.1.2 *The search query*

In Stephens-Davidowitz and Varian's paper (2015), they explain how to formulate the search query for different purposes:

- "+" means "or". If you type Trump+Hillary, the results will be searches that include either the word Trump or the word Hillary.

- "-" means to exclude a word. If you type Clinton - Bill, results will be searches that include Clinton but do not include Bill
- A space means "and". If you type Donald Trump, the results will be searches that include both the word Donald and the word Trump. The order does not matter.
- Quotes force a phrase match. If you type "Make America Great Again", results will be searches that include the exact phrase Make America Great Again.

Google Trends presents several alternatives when typing a search term in the query field (Google, 2017). When typing [guardian], Google Trends present the following suggestions: [guardian] as an independent search term, [Guardian] as a topic and [The Guardian] as in the newspaper. For topic searches, Google uses an algorithm which collect all searches that are related to the topic, but does not necessarily contain the exact query, i.e. [guardian] (Google, 2017), e.g. searches on the query [guard minor].

#### *5.1.3 Sampling method*

The search data is computed using a sampling method (Choi & Varian, 2012). Google Trends analyzes a percentage of all Google searches to determine how many searches have been conducted for the given search term compared to the number of total Google searches in the given time period (Google, 2017). This sampling method can lead to a few per cent variation in results from day to day (Choi & Varian, 2012).

#### *5.1.4 Segmentation and comparison opportunities*

Stephens-Davidowitz and Varian (2015) describes the possibilities in segmenting the different queries by different geographical levels and over different time periods. The query index is available at country, state/region/county and city level for several countries. For different time periods, the query index presents different scales on the data. A time period of 3 months or shorter will present daily data, or else weekly data. For a time period longer than or equal to 3 years, the query index presents monthly data.

Google Trends gives the opportunity to compare up to five search terms or categories at the same time (Varian & Stephens-Davidowitz, 2015). It is possible to compare different queries over the same time period, compare queries over different regions or different time periods (Google, 2017). For comparisons over different time periods, ie. queries such as [Election] in [2011], [2012], [2013] and [2014] the query index will normalize the index over region and

time (Google, 2017). Thus, the query index will differ from downloading [Election] for [2011] and [Election] for [2012] separately rather than when compared in the query tool in Google Trends.

#### *5.1.5 Limitations due to privacy considerations*

Google Trends has limitations due to privacy considerations (Choi & Varian, 2012). If the frequency of a search term is below an unreported privacy threshold, then the index will show zero. This threshold is measured in absolute numbers, such that smaller geographic areas and shorter time periods will more often generate zeros compared to larger areas or longer periods. As will searches conducted closer to the beginning of Google searches in 2004 (Varian & Stephens-Davidowitz, 2015).

## **5.2 Measuring Social Sentiments Using Google Trends Data**

In 2004, 64 percent of all American adults had access to the internet from their homes. In 2015 this number was 84 percent (Perrin & Duggan, 2015). As of April 2017, 4,464,000,000 searches were made on Google each day (Internet Live Stats, 2017). Through the large number of all demographics using the service, Google searches are likely to provide information about a significant part of the American population (Stephens-Davidowitz, 2013). Google search data, aggregating millions of searches, systematically correlates with the demographics of those who conduct the searches (Stephens-Davidowitz, 2013). Stephens-Davidowitz provides two examples in his paper: the search rate for the word “God” explains 65 percent of the variation in a state’s share of residents believing in God, and the search rate for “gun” explains 62 percent of the variation in a state’s gun ownership rate.

Furthermore, there is reason to believe that Google draws out socially sensitive attitudes. Alone and online, the limit for sharing personal information becomes lower than if one is asked in a survey or in an opinion poll, as the use of Google limit the concern of social censoring (Conti and Sobiesk, 2009). The large number of pornographic searches and sensitive health information that is shared on Google, substantiates the assumption that people are more forthcoming online than otherwise (Stephens-Davidowitz, 2013). The effect of the SDB is reduced as people no longer worry about what the pollster or survey maker believes is the “right” answer. Thus, search queries provides a non-survey source for examining sentiments towards social sensitive topics (Stephens-Davidowitz, 2013).

In his research on islamophobic internet searches and anti-Muslim hate crimes, Stephens-Davidowitz (2014) found a correlation between anti-Muslim Google searches and hate crimes, using 2004-2013 weekly data on negatively loaded Google searches containing the word Muslim. One of the search terms he used was [kill muslims], which after the San Bernardino attack in 2015 was the most popular search term containing the word Muslim in the U.S. (Stephens-Davidowitz, 2014). This is illustrated in Figure 1.

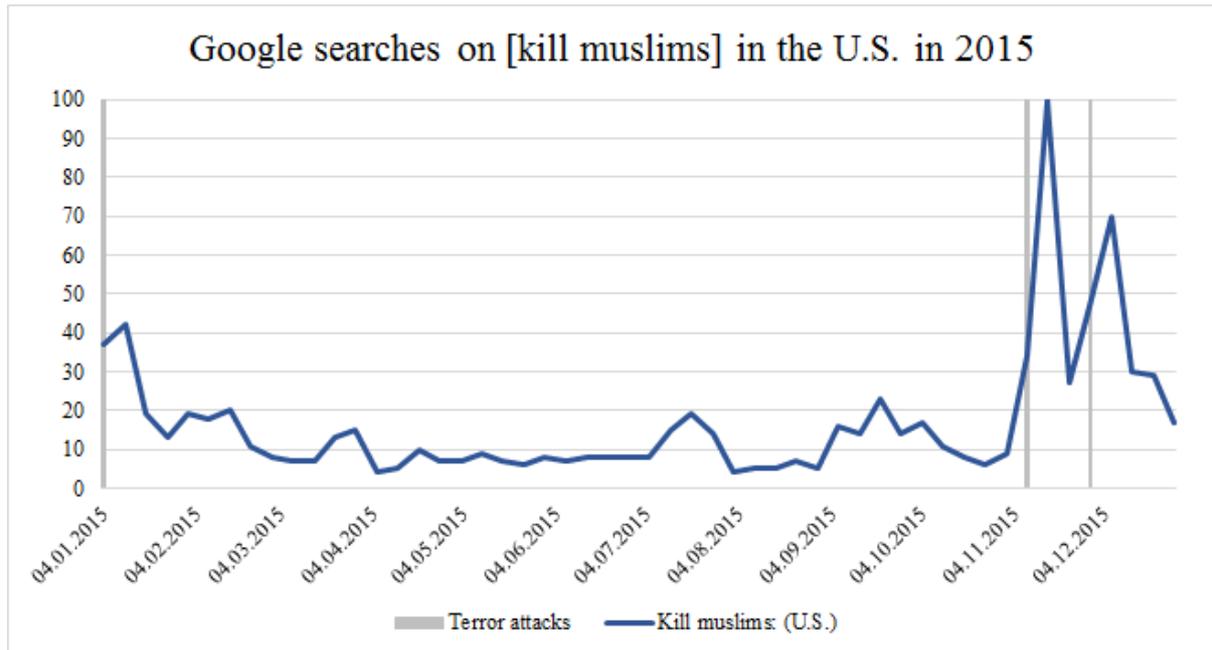


Figure 1: Query index for [kill muslims] in 2015. The following terror attacks are marked in grey: January 2, 2015: Charlie Hebdo Paris; November 13, 2015: Bataclan, Paris; December 2, 2015: San Bernadino attack (Google, 2017).

Stephens-Davidowitz explain the use of such a search term as a result of people typing their uncensored thoughts into Google, without the intention of getting relevant information back. Other examples of such searches are [I hate my boss], [I am drunk] and [people are annoying]. There are about 1600 Google searches for “I hate my boss” every month in the U.S.. Being unhappy with one’s boss is a common reason for why people leave their jobs (Arnold-Smeets, 2013). One can assume that the searches for [I hate my boss] represents a tiny fraction of those that actually leave their job because of their boss, the same way the number of searches for [kill muslims] represents a tiny fraction of people that actually resent Muslims (Stephens-Davidowitz, 2014).

Google searches can also be used to measure racism. In his paper, Stephens-Davidowitz uses the percent of Google search queries from 2004-2007 containing the word [nigger] or [niggers] as a proxy for the racial animus in a state. He compares the proxy to Barack Obama's vote share in 2008 and 2012, controlling for John Kerry's 2004 vote share. Studies using survey data to measure racial animus found little evidence of racial animus affecting Obama negatively in the 2008 election, while Stephens-Davidowitz using the non-survey proxy found evidence that racial animus did cost Obama popular votes in both the 2008 and 2012 election.

While surveys and polls are conducted on a representative sample and can tell us something about the average American's thoughts on a topic, Google searches tells us what someone excited enough over a topic to search for it thinks about the issue. People searching for [muslims] are not a representative group of Americans, and thus the Google searches will suffer from selection bias. For the purpose examining the change in sentiments over time in an area, as well as comparing the extent of certain sentiments across states, Google searches can still be used as a proxy. The relative popularity of a search term in an area in one period compared to another can tell us something about how the size of the population that is excited enough about a topic to search for it has changed. If we further assume that the fraction of the population that searches on a topic represents a small percentage of the people that feel the same way, then comparing the popularity of a search word can tell us something about relatively how many people that felt a certain way in one period compared to another.

## 6 Data

In the following chapter we present the search terms we have used to construct proxies for immigration skepticism, far-right sentiment and racial animus. We further explain the time frame set for the Google Trends proxies. We also explain how the polling data is collected and present the control variables used in the estimation model.

### 6.1 The Google Trends Proxies

We have created two proxies for each sentiment chosen; a one year proxy (long-term) and a two weeks proxy (short-term). In total we have constructed six proxies.

The long-term proxy is constructed to capture the underlying effect of the general sentiment in the population. This proxy uses search data from the month of November in the year ahead of the presidential election, to the last day of October in the year of the presidential election, and compares the frequencies in the three election years: 2008, 2012 and 2016. The objective of this thesis is to test the proxies' explanatory power in advance of the election to see whether they can be used for predictive purposes in the future, thus we have excluded the election month<sup>10</sup>. The periods defined in the query tool are [11.1.2007 - 10.31.2008], [11.1.2011 - 10.31.2012] and [11.1.2015 - 10.31.2016].

The short-term proxy is constructed to capture the effect of higher search activity close to the election. This proxy includes search data from the two weeks before the election date. The periods defined in the query tool are [10.21.2008 - 11.4.2008], [10.23.2012 - 11.6.2012] and [10.25.2016 - 11.8.2016].

The region chosen for all three periods is the United States. We have further segmented the data for the different U.S. states. When downloading data on search frequency in the U.S. and further filtering on states, each state gets an index score between 0 and 100 that reflects the relative popularity for the search query in that state on average over the time period specified. The state with the highest score is the state where the search query on average was relatively most popular during the specific time period. We use the index score for each state in the different time periods as the proxy of the three sentiments of interest: the far-right sentiment,

---

<sup>10</sup> For the U.S. Presidential Elections, the election date is always the Tuesday after the first Monday in November (The United States Government, 2017).

immigration skepticism and racial animus. In order to investigate if search activity on Google can be used as a measure of sensitive social sentiments in a population, we will apply the theories of confirmation bias and information utility to validate the specific queries we use to build the proxies.

### *6.1.1 Far-right proxy*

The condition for using a Google search to build a proxy for the level of far-right sentiment in a state, is that using the search query makes having far-right sympathies more likely.

One way to examine far-right sentiment using Google search data is to look at queries that express a hateful opinion associated with the far-right ideology, for example [kill muslims] or [I hate jews], like the methodology used by Seth Stephens-Davidowitz (2014). Due to the privacy threshold, it is hard to collect statewide data on certain search words. This applies especially for search words associated with the far-right sentiment, as it regards a marginal group of people in the U.S. in terms of absolute numbers. When examining the frequency of the query [kill muslims] on Google Trends, only Texas and California, the two most populous states in the U.S., have a large enough search frequency in absolute numbers for the results to show in our chosen time frame. Therefore, in order to proxy the level of far-right sentiment in a state, other types of Google queries must be used.

We argue that the search frequency for specific news sites and online forums can reflect attitudes towards certain social sensitive topics in a population as a result of confirmation bias. Studies have shown that due to confirmation bias, people tend to selectively search for information that support their ideas and values (Nickerson, 1998), see section 3.4. A person with liberal values would thus rely on liberal news sources for information, while an ultra-conservative American would feel like those news sites contradicted his or hers values and believes and thus look for information elsewhere.

This can also be explained by Atkin's (1973) guidance and reinforcement factors, where he argues that people search for information in order to know how to feel about things or search for information to confirm their attitudes. From this, it is fair to argue that if knowing the sites would provide this guidance or reinforcement of a belief, the information utility argument holds. Examining how the relative popularity of searches on far-right websites in

the U.S. has developed can possibly tell us something about how the far-right support has developed.

There are numerous online forums and news sites with far-right content, many of them with different geographical reach and different life span. Therefore we find it necessary to compute a weighted index of three different queries, ensuring that we minimize the risk of choosing a site for instance only present in the South, or only active between 2008 and 2010, which would issue spurious results. We use the search frequency for three right-oriented conservative web pages: Stormfront.org [Stormfront (Web page)], Breitbart.com [Breitbart (News site)] and Drudgereport.com [Drudge report (Web page)]. We use the filters for web pages and news sites, which means that when using for example the query [Stormfront] and select the alternative “web page”, Google Trends applies an algorithm that filter out all searches that are irrelevant for the webpage, like weather queries.

Stormfront is the oldest and most popular online hate forum in the U.S. (Stephens-Davidowitz, 2014). The site was formed in 1995 by a former Ku Klux Klan leader, and members of the site has been linked to around 100 hate-murders in the past five years according to a report by the Southern Poverty Law Center. Stormfront.org reaches over 113,000 people every month in the U.S. (Quantcast, 2017). Simply searching for the website does not prove that someone has far-right sympathies: a person could for instance search for the website out of curiosity or in the purpose of research. The condition for using searches for Stormfront in the proxy is that searching for the website makes it more likely that you have far-right sympathies, which we argue it does due to the confirmation bias. Also, being a member of the site makes having far-right sympathies more likely than simply searching for the site, since that would increase the probability of being active on the forum. In 2014 the states with the highest number of Stormfront members per capita was Montana, Alaska and Idaho. Montana and Alaska were also the states with the highest query index for the website in 2014 (Google, 2017). This indicates that searches for the page can work as a good proxy.

The Drudge Report is a conservative, alt-right news aggregation site which was first launched in 1995 (Jones & Salter, 2011). The site has had between one and one and a half million visits every month the past ten years (Quantcast, 2017). The site is considered to promote mainly Republican interests, but it also promotes populist anti-government attitudes which are associated with far-right political ideology (Quantcast, 2017) According to Quantcast, on

average 73.8 percent of the visits to the site are made by what they define as fanatics, which means people that visit a site more than 30 times in a given month. This implies that the majority of people who search for [drudge report] do so on a regular basis. Based on this, we assume that searches for the website can work as a good proxy for far-right sentiment.

Breitbart is an American news site founded in 2007. The website has on multiple occasions been accused for being racist, xenophobic and misogynist (Piggott, 2016). Breitbart claims to have three main opponents: the mainstream media, the Democratic Party and the Republican establishment in Washington. The target group of the website is ultra conservative Americans (Piggott, 2016). The site has among other things posted several stories denying global warming, conspiracy theories about Barack Obama, as well as a false story about a Muslim mob burning down a church in Dortmund on New Year's Eve 2017 (The Guardian, 2017). Whether Breitbart News is part of the far-right movement or if it just is a conservative news site is debatable, but we choose to include the site in our proxy due to the site's reputation of publishing stories that sympathize with far-right ideology. From 2007 to 2012, Breitbart News Network consisted of three websites: Big Government, Big Journalism and Big Hollywood. In 2012 a new consolidated website called just "Breitbart" was launched (Bromwich, 2016). Using the query [Breitbart (News site)], we ensure that all searches for [Big Government], [Big Journalism], and [Big Hollywood] are included due to the filter algorithm.

The connection between Breitbart News and the former chief strategist in the Trump administration, Steve Bannon, could possibly interfere with the search frequency of the news site<sup>11</sup>. If a lot of people google Breitbart News because of the attention Steve Bannon received in the media, the assumption for using the query in the proxy does not hold. Since we only look at the search frequency for the website prior to the election in November 2016 and the focus on Steve Bannon in the media mainly began after Trump was elected, we believe that the connection does not generate a problem.

It can be argued that people who visit a website frequently will have the site listed as a favorite or have it bookmarked, and thus does not have to type the site into the search field. For the purpose of this thesis, we assume that a significant part of the people visiting a site

---

<sup>11</sup> Steve Bannon was the executive chair for Breitbart News from 2012 until he joined Donald Trump's campaign as chief strategist (Rahn, 2016). He received major attention in the media after the election, due to his connections to Breitbart, and was accused of being a white nationalist (Anti-Defamation League, 2017).

frequently does so via Google. We believe this is likely as most computers and smartphones use Google as the standard search engine which makes it easy for people to google.

In Figure 2 we see the query index for [Breitbart], [Stormfront] and [Drudge Report] using monthly data for the time period 1.1.2008 to 31.12.2016 in the U.S.. Each election month is marked in grey. We see that the query index spikes before every election month for all of the queries, indicating a systematic spike in the query before presidential elections. We cannot draw any inference from this, but argue that the chosen search queries might capture some valuable variation worth analyzing in regards of predicting election results.

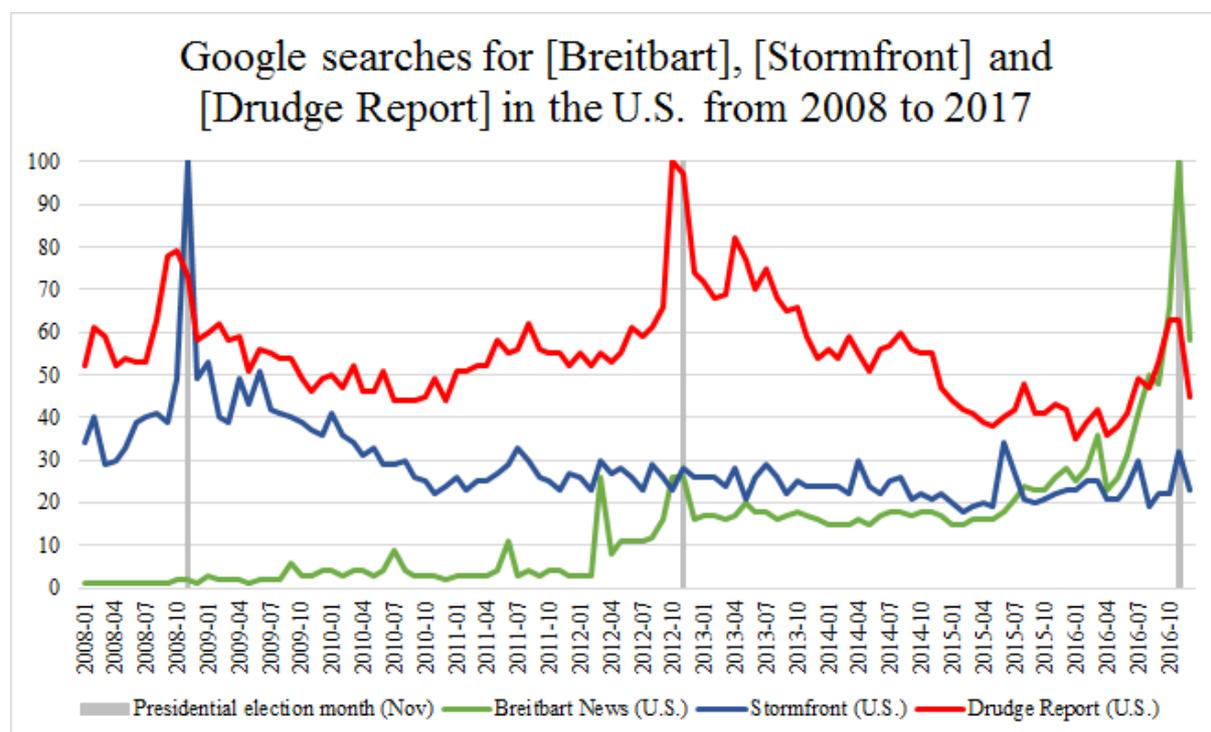


Figure 2: Query index for [Breitbart News], [Stormfront] and [Drudge Report] from 1.1.2008 - 31.12.2016. Election months (November) 2008 - 2016 are marked in grey (Google, 2017).

From the argumentation above we conclude that the average of the query index for the three web sites works as a good proxy for the level of far-right sentiment in a state.

### 6.1.2 Immigration skepticism proxy

In order to create a proxy for immigration skepticism using Google Trends data, we have used the search query [Illegal immigration and residence (Topic)]. The condition for using

this query to build a proxy for the level of immigration skepticism in a state, is that the queries captured by the topic algorithm makes being an immigration skeptic more likely.

Common perception and prejudice in many Western societies is that immigration leads to higher crime rates, that immigrants take away jobs from the native residents and that they contribute to social unrest in the communities (Rydgren, 2005). A lot of research has been conducted on this field, showing that this has little statistical support. Some research has shown that the two in fact might be negatively correlated (Reid et. al (2005), Olson et. al (2009)). Although facts and statistics to a large extent reject the common perception many people have towards immigration, asymmetric information and strong prejudice makes it difficult to change such views. From this, it is fair to assume that the confirmation bias is applicable in the case of immigration skepticism, in the sense that people will search for information that will confirm their prejudices.

The prejudices mentioned above are to a large degree based on uncertainty and fear (Rydgren, 2008). Due to information utility and Atkin's surveillance factor (1973), we can argue that in the matter of monitoring threats, people will search for specific terms to collect more information on topics they fear will affect them. Driven by fear, people will expect that the information they seek will benefit them either way, regardless of being consistent or not with the existing hypothesis (Knobloch-Westerwick & Kleinman, 2011). On one hand, they can confirm their rationale behind the fear, and act accordingly, or the information sought could discard the hypothesis, giving them less a reason to fear. Either way, we can argue that fear of immigration can cause people, motivated by information utility, to seek further knowledge on the topic.

Preferably we would have been able to include search queries like [immigration crime], [effects of immigration], [illegal immigration] and [immigration jobs]. Unfortunately for our purpose, the privacy threshold on Google Trends gave us a query index for only a few states when we segmented down on state level. Therefore, we chose to use a topic search, where Google uses an algorithm to collect all the queries belonging to that subject, which includes more than one search term, and accordingly brings us above the privacy threshold for all U.S. states. The topic must have such a relevance for immigration skepticism that it is still possible to argue that a search within the topic makes being an immigration skeptic more likely. Using solely [Immigration (Topic)] gives a too wide range on the subject. It is impossible to argue,

due to confirmation bias or information utility, that searching for immigration in general would incline a likelihood of being immigration skeptic.

There are different types of immigrants. Among these are legal immigrants, documented refugees and undocumented illegal immigrants. In 2011 there was 11 million illegal immigrants in the U.S., which account for 5 % of the civilian U.S. workforce (Krogstad, Passel, & Cohn, 2017). Rising illegal immigration has been used by politicians on several occasions as an argument for a strict immigration policy, no later than in the U.S. presidential election 2016 by Donald Trump (Swanson, 2016). Based on the rhetoric many politicians use on this political issue, it is fair to assume that the uncertainty and fear many people feel towards immigration, is due to the fear of illegal immigration.

From an information utility perspective, we argue that people whose immigration skepticism is based on fear, will search on the topic of illegal immigration to gather information regardless of whether the information will confirm or reject their hypothesis. We can also argue from a confirmation bias perspective, that people who are skeptic towards immigrants will be likely to search for the topic of illegal immigration to confirm their beliefs, regardless of whether their beliefs is driven by fear or not.

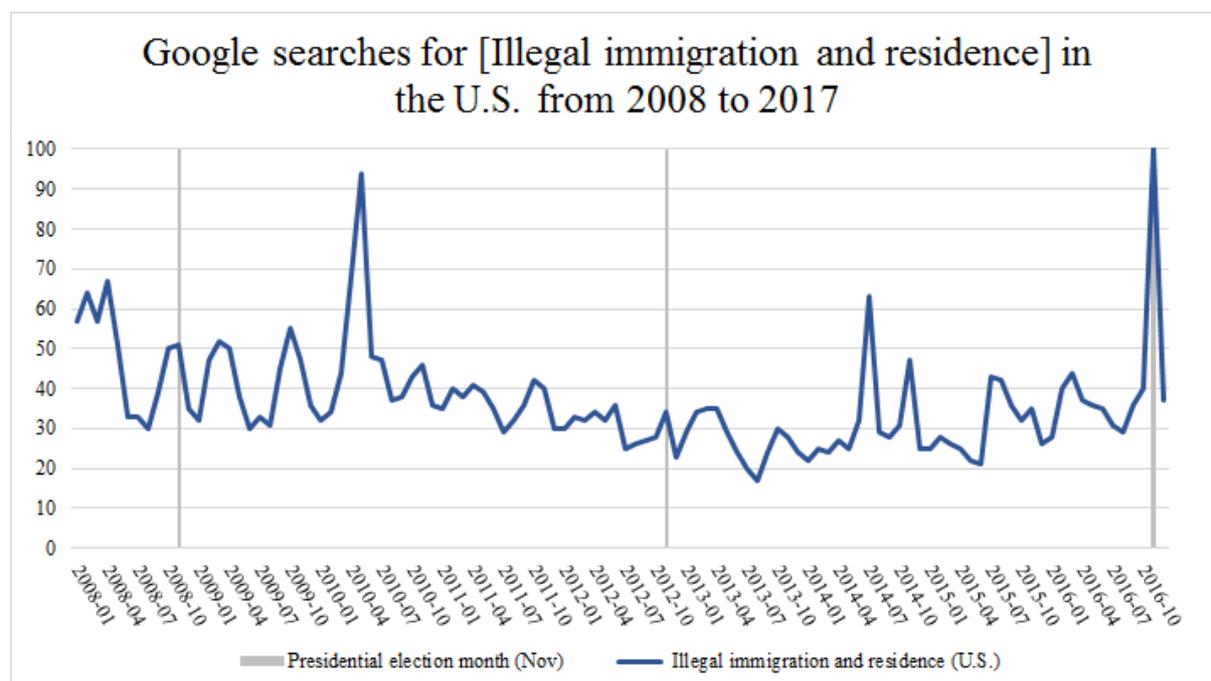


Figure 3: Query index for [Illegal immigration and residence] from 1.1.2008 - 31.12.2016. Election months (November) 2008 - 2016 are marked in grey (Google, 2017).

In Figure 3 we see the query index for [Illegal immigration and residence (Topic)] using monthly data for the time period 1.1.2008 to 31.12.2016 in the U.S.. Election month is marked in grey. As with the previous proxy, we see that the query spikes before every election month, indicating some interesting variation worth analyzing further.

From this we conclude that searches on the query [Illegal immigration and residence (Topic)] can serve as a good proxy for immigration skepticism due to both information utility and confirmation bias, and that a state's query index will tell us something about the level of immigration skepticism in that state.

### *6.1.3 Racial animus proxy*

In order to build a proxy for the level of racial animus in a state, we use the search frequency for queries including the word [nigger] or [niggers]. These are the same queries used by Stephens-Davidowitz (2013) in his research on how racial animus affected the election outcome for a black presidential candidate. In 2013 the racial epithet was included in more than 7 million searches annually on a world basis (Stephens-Davidowitz, 2013). Search queries including the word “nigger” is on average 83 percent as popular as queries including the word “migraine”, about half as popular as queries including one of the most common word in Google searches: “weather”, and more than twice as popular as searches for stormfront.org (Google, 2017). These comparisons are arbitrary chosen to provide an understanding of how common the racial epithet is in Google searches. The numbers are based on the popularity of the Google queries from 2004 until November 2017.

In Figure 4 we see the query index for [nigger + niggers] using monthly data for the time period 1.1.2008 to 31.12.2016 in the U.S., with each election month marked in grey. Also here, we see that the query spikes before every election month, although not as much as for the previous two proxies. We argue that the variation is interesting and valuable for our research purposes.

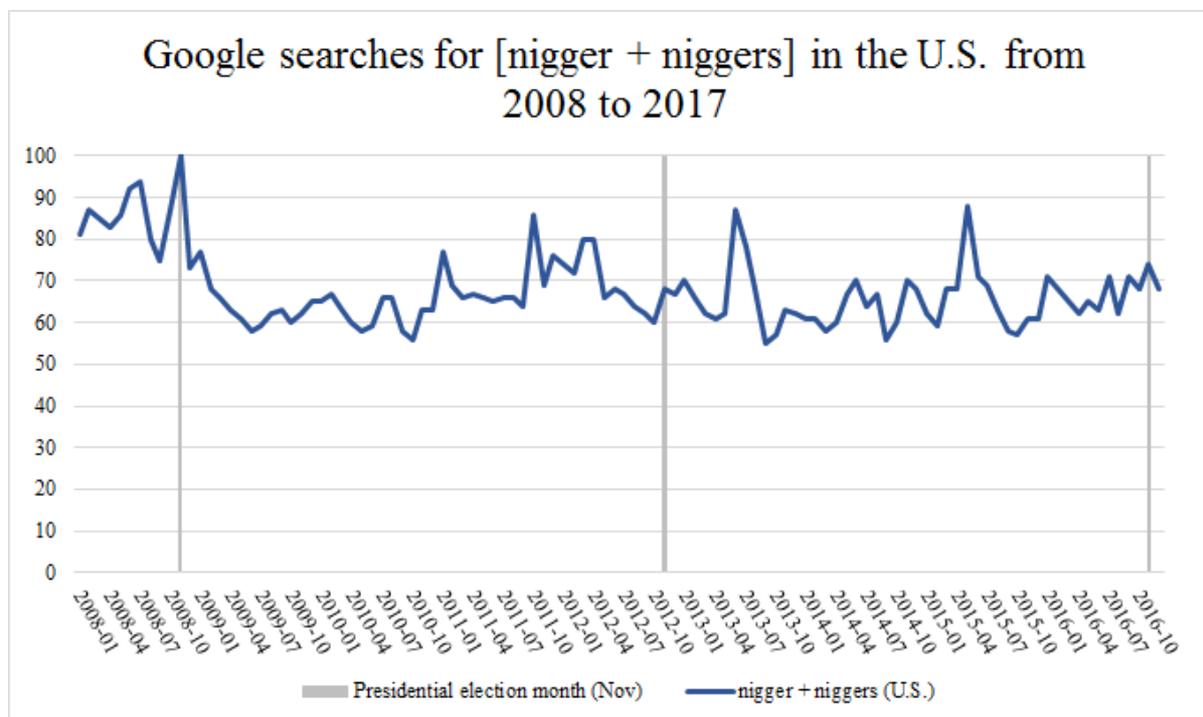


Figure 4: Query index for [nigger + niggers] from 1.1.2008 - 31.12.2016. Election months (November) 2008 - 2016 are marked in grey (Google, 2017)

The condition for using the racial epithet query [nigger] to proxy racial animus, is that searching for the term makes it more likely that a person harbours racial animus. When looking at searches related with “nigger”, the most common search is “nigger jokes”, which returns websites with degrading material about African Americans. “Nigger” is the most salient of racial slurs in America<sup>12</sup> and it is considered to be extremely offensive (Kennedy, 2003). Among top five states ranked by the popularity of Google queries containing the racial epithet are Mississippi, Kentucky, Louisiana and Alabama, which all are ranked on a list of the most racist states based on hate-crime, with Mississippi and Alabama taking first and second place (Durankiev, 2015). Stephens-Davidowitz (2013) did detect a small positive correlation between racially charged search rate and the percent of black people in the area. However, he disregards that the popularity of the racial epithet is due to African Americans using the term, which would limit the value of the proxy. Stephens-Davidowitz points out the difference between the word “nigger” and “nigga”, with the last being a commonly used term

<sup>12</sup> Kennedy (2003, s. 22) says this is “the best known of the American language’s many racial insults”.

in African American culture, while the first is associated with black slavery and the suppression of African Americans (Fredrickson, 2002). Based on this evidence, it is fair to assume that searching for the query [nigger] or [niggers] increases the probability that a person harbours racial animus.

#### 6.1.4 Summary Statistics by year

In Table 1 we present some descriptive statistics of our proxies.

Table 1: Summary Statistics: Social Sentiment Proxies

Year		Long-term Racial Animus	Long-term Immigration Skepticism	Long-term Far-right sentiment	Short-term Racial Animus	Short-term Immigration Skepticism	Short-term Far-right sentiment
2008	Max	100	100	94	100	100	59
	Min	36	27	35	5	13	16
	Range	64	73	59	95	87	42
	Mean	63	56	53	30	40	34
2012	Max	100	100	87	100	100	71
	Min	30	42	47	16	19	28
	Range	70	58	40	84	81	43
	Mean	59	64	62	50	44	46
2016	Max	100	100	96	100	100	83
	Min	27	43	53	13	34	32
	Range	73	57	43	87	66	51
	Mean	59	75	67	39	70	55
Total	Max	100	100	96	100	100	83
	Min	27	27	35	5	13	16
	Range	73	73	61	95	87	67
	Mean	60	65	61	40	52	45

Table 1: Summary Statistics of the Social Sentiment Proxies

*Notes:* For short-term immigration skepticism and short-term far-right sentiment some of the states failed in showing index results for the queries used to proxy the sentiments. We have discussed the implications of this in section 9.2.

For both the racial animus proxy and the immigration skepticism proxy, the maximum value is 100, while for the far-right sentiment the maximum value ranges from 87 to 96 in the long-

term proxy and 59 to 83 in the short-term proxy. This is due to the fact that this proxy is an average of three queries. Further, we observe that the mean has increased over the years for immigration skepticism, and for the far-right sentiment proxy both on long-term and short-term. The mean of short-term immigration skepticism has increased with 30 index points from 2008 to 2016, while for the long-term proxy it has increased with 19 index points. For short-term far-right sentiment, the increase in the mean is 21 index points from 2008 to 2016. This implies a growth in frequency for these search queries, which indicates that both immigration skepticism and far-right sentiment are growing sentiments in the American society. Descriptive statistics segmented for each state and region, as defined in section 8.2, can be found in Table A2 and A1 in the appendix.

## **6.2 Polling Data**

Polling data is the main predictor of election outcome, see section 3.2, and we would thus expect that this data has a significant explanatory power when it comes to explaining election outcome. We therefore include polling data in our estimation model.

The statewide polling data used in this paper is collected from electoral-vote.com. The website analyzes the most recent polls on state level to predict the winner of the electoral vote (Electoral Vote, 2017). Collecting statewide polling data from one single polling company was not possible, since different companies conduct polls in different states and different times. We have used the most recent polls data available on electoral-vote.com, and mainly used data from surveyU.S.A, surveymonkey, PPP and youGov. Most of the polls was conducted within two weeks before the election, but for some of the states we did manage to find the newest polls. For Alaska in 2012 we used polls data from August 16 and for Tennessee in 2012 we did not find any polls data later than February. We discuss the implications of this in section 9.2.

## **6.3 Economic and Demographic Control Variables**

In order to select the regressors belonging in the final model specification, we use forward induction sequentially adding the economic and demographic variables<sup>13</sup>. We include the regressors that yield the best estimation model. In this section we present the control variables included, explain why they are relevant for our estimation model and how they are

---

<sup>13</sup> This is the method used by Strumpf and Philippe (1999) in order to select control variables. They use forward induction and include the variables that yield the best model, i.e. yield the highest R-squared value.

collected. In Table 3 we present some descriptive statistics for each of the control variables in 2008, 2012 and 2016.

All of our control variables are collected from The United States Census Bureau. We have collected raw, statewide data from 2008, 2012 and 2016. The statewide control variables we use in our final estimation model is: Percentage of male population (% Male), median age, level of black or African Americans (% Black or African American), level of Hispanics and Latinos (% Hispanic or Latino), level of the population with a bachelor degree or more (% Bachelor or more), level of the population with a High School degree or less (% Less than High School), unemployment rate (UE) and mean income in U.S. dollars per household, inflation adjusted (Mean income U.S.D). We will further explain the reasoning behind including each of the control variables.

### *6.3.1 Median Age*

In the presidential election in 2016, the support for the Democratic candidate was much higher among young voters compared to older voters (Castillo & Schramm, 2016). According to a CNN Politics exit poll<sup>14</sup> based on 24 537 respondents, 55 percent of 18-29 year olds voted for Hillary Clinton, while 38 percent voted for Donald Trump (CNN, 2016). In general, young people tend to have a stronger Democratic orientation than the rest of the population. According to Pew Research Center survey data, about 36 to 57 percent more Millennial voters, i.e people born between 1980 and 1994, identify as Democrats or lean towards the Democratic Party compared to older voters (Pew Research Center, 2016). Based on this, one can assume that states with a younger population are more likely to vote for the Democratic candidate.

### *6.3.2 Male ratio*

Since 1980, women in America has moved towards the Democratic side of politics. This has happened parallel to women in general becoming more engaged in politics (Inglehart & Norris, 2000). Men has, on a long and consistent basis, moved towards the Republican side. While the Democratic Party holds a large advantage among women, the Republican Party has a remarkable advantage in party affiliation among men (Pew Research Center, 2016). This modern gender gap in American politics has widened in recent years, and it leads us to assume that if the male ratio in a state changes significantly, it could have an effect on

---

<sup>14</sup> Exit polls are surveys of a small percentage of voters taken after they leave their voting place.

election outcome. This assumption is discussed in section 9.2. Table 2 shows how women and men have voted in the last twelve presidential elections.

Election	Men		Women	
	% Dem	% Rep	% Dem	% Rep
1972	36 %	<b>62 %</b>	37 %	<b>61 %</b>
1976	<b>50 %</b>	48 %	<b>50 %</b>	48 %
1980	36 %	<b>55 %</b>	45 %	<b>47 %</b>
1984	37 %	<b>62 %</b>	44 %	<b>56 %</b>
1988	41 %	<b>57 %</b>	49 %	<b>50 %</b>
1992	<b>41 %</b>	38 %	<b>45 %</b>	37 %
1996	43 %	<b>44 %</b>	<b>54 %</b>	38 %
2000	42 %	<b>53 %</b>	<b>54 %</b>	43 %
2004	44 %	<b>55 %</b>	<b>51 %</b>	48 %
2008	<b>49 %</b>	48 %	<b>56 %</b>	43 %
2012	45 %	<b>52 %</b>	<b>55 %</b>	44 %
2016	41 %	<b>52 %</b>	<b>54 %</b>	41 %

Table 2: Exit poll data in each presidential election from 1972-2016 (*Sabato, Kondik, & Skelley, 2017*).  
*Notes:* The highlighted cells show which party got the majority of the vote from male and female voters.

### 6.3.3 African American and Hispanic share of population

In 2012, Barack Obama won the support of nine out of ten non-white voters, and 19 out of 20 black voters (Kirk & Scott, 2016). This was most likely partly due to the attraction of voting for a black president. Race was also a major topic of debate in the 2016 election, with Donald Trump advocating for a wall between the U.S. and Mexico, and giving comments about mexicans having “lots of problems” and “bringing drugs to the U.S.” (Kopan, 2016). According to a poll conducted by ABC news in september 2016, just 17 percent of Hispanic and Latino people intended to vote for Trump, and only 3 percent of black people (Kirk & Scott, 2016). The trend in party affiliation among African American voters has in general been stable over the past years. Based on surveys conducted by Pew Research Center, 87 percent of black voters identify with the Democratic Party, while only 26 percent of white people identify as Democratic. Among Hispanic voters, the trend is similar, with 63 percent of hispanics identifying as Democrats (Pew Research Center, 2016)<sup>15</sup>. Based on this, we

<sup>15</sup> Annual totals of Pew Research Center survey data; 2016 data based of surveys conducted January-August.

believe that if the level of black or African Americans, or Hispanic or Latinos, in a state has changed in the period 2008-2016, this can have an effect on election outcome.

#### *6.3.4 Education*

Over the last two decades, less educated voters, i.e. voters with a High School degree or less, has gone from being mostly Democrats to being mostly Republicans (Pew Research Center, 2016). In 2008, most voters with a High School degree or less were Democrats, while in 2016 the gap between the two parties regarding party affiliation among these voters was as good as closed. For college graduates, over half of the voters today support the Democratic Party. The Republican Party began to lose ground among people with a college degree in the second half of George W. Bush's first term. By the 2008 election, the Democrats held a 10-point edge among college graduates and the gap has increased some since then (Pew Research Center, 2016). This inclines us to believe that a growth in people with a bachelor's degree or more would increase the support for the Democratic candidate.

#### *6.3.5 Unemployment and mean income*

Previous research on election outcome have found that economic factors explain a considerable portion of variation in vote outcomes (Strumpf & Philippe, 1999). We include the statewide unemployment rate in our model, which could serve as a proxy for the regional economic performance and indicate voters' job security. It is fair to assume that a state experiencing a growth in the unemployment rate from 2008 to 2016 would be eager for a change in government, and thus be less likely to vote for the incumbent party, i.e the Democratic Party. A rise in income reflect increasing employment and wages and can thus be used as a measure of the economic health in a state (Iceland, 2003). Changes in income can tell something about how the economic situation for people has changed, and median income has traditionally been used as a measure of the economic health of the middle class (Dorfman, 2014). For the purpose of measuring the economic situation for the middle class it would be preferable to use median income and not mean income. We explain this further in section 9.2. Income can also be viewed as an indirect measure of the level of education. As people become higher educated, they are more likely to vote for the Democratic Party, see subsection 6.3.4. Given the positive correlation between income and education, we can assume that as people earn more money, they are more likely to vote for the Democratic Party (Porter, 2014).

### 6.3.6 Summary statistics by year

In Table 3 we have presented some descriptive statistics for the control variables used in the estimation model.

Table 3: Summary Statistics by Year, Control Variables

Year		% Male	Median Age	% African American	% Hispanic	Mean Income U.S.D	% Bachelor or more	% UE	% Less than High School
2008	Max	0,521	41,9	0,534	0,449	96 572	0,482	0,061	0,204
	Min	0,473	28,7	0,005	0,01	52 642	0,171	0,022	0,083
	Mean	0,493	37,210	0,110	0,098	69 199	0,274	0,039	0,134
	Sd	0,008	2,240	0,112	0,098	11 359	0,056	0,008	0,035
2012	Max	0,521	43,5	0,495	0,47	108 168	0,53	0,08	0,19
	Min	0,473	29,9	0,004	0,013	55 371	0,19	0,02	0,07
	Mean	0,493	37,778	0,112	0,110	70 686	0,289	0,055	0,120
	Sd	0,008	2,397	0,109	0,100	12 156	0,059	0,013	0,032
2016	Max	0,526	44,5	0,471	0,485	119 777	0,568	0,054	0,176
	Min	0,475	30,7	0,004	0,015	61 753	0,208	0,02	0,068
	Mean	0,494	38,247	0,113	0,117	80 143	0,310	0,035	0,109
	Sd	0,009	2,475	0,108	0,102	13 359	0,062	0,006	0,029
Total	Max	0,526	44,5	0,534	0,485	119 777	0,568	0,08	0,204
	Min	0,473	28,7	0,004	0,01	55 371	0,171	0,02	0,068
	Mean	0,493	37,745	0,111	0,108	49 804	0,291	0,043	0,121
	Sd	0,008	2,395	0,109	0,100	36 960	0,061	0,013	0,034

Table 3: Summary Statistics, Control Variables (*United States Census Bureau, 2017*).

We observe from Table 3 that the mean level of education has increased. It is therefore reasonable to assume that the general level of education in most states has increased. According to the theory on party affiliation and level of education, see subsection 6.3.4, this alone could indicate that more people would vote Democratic in 2016 compared to 2008. The mean unemployment rate increased from 2008 to 2012, which could imply that people in general would be more dissatisfied with the incumbent party, i.e. the Democratic Party, in 2012, and that the Republican Party would get a higher percentage of the votes, see subsection 6.3.5. From 2012 to 2016 the mean unemployment rate decreased, implying higher satisfaction with the incumbent party and thus a higher percentage of the votes to the

incumbent party, i.e. the Democratic Party. Mean income has also increased from 2008 to 2016. From the relationship between income and education discussed above, a higher income level could be a reflection of a higher education level, which implies that there might be a positive correlation between income and Democratic votes. However, change in mean income is not a good measure on whether income has increased or decreased for the average American. It could be a reflection of an increase only for the richest Americans. One should therefore be careful when making inference about the change in mean income.

## 7 Empirical Strategy

In this thesis, we examine the intrastate variation in socially sensitive sentiments. Our primary interest is to see whether the socially sensitive sentiment proxies can explain some of the variation in the election outcome within a state. This can be helpful in identifying the determinants of voting in the future and contribute to building better models for predicting election outcomes. In this chapter, we explain how we can estimate causal effects with panel data using a fixed effects model and how we can use this model to explain election outcomes, before we present our estimation model.

### 7.1 Estimation with Panel Data

Consider the model presented in equation 1,

$$y_{it} = \alpha + \lambda t + \beta x_{it} + u_{it} \quad (1)$$

where  $i$  = entity (state),  $n$  = number of entities, so  $i = 1, \dots, n$ ,  $t$  = time period (year) and  $T$  = number of time periods, so  $t = 1, \dots, T$ .  $y_{it}$  then denotes the dependent variable at time  $t$  for state  $i$ , while  $x_{it}$  denotes the explanatory variable at time  $t$  for state  $i$ .  $u_{it}$  is the error term.  $\lambda t$  is the year fixed effect. In order to estimate the causal effect  $\beta$ , we must assume that  $cov(u_{it}, x_{it}) = 0$ , i.e. that there is no unobserved effect in the error term that correlates with our explanatory variable and causes our estimated effect  $\beta$  to be biased. When we have several observations for the same unit, it is very unlikely that the error terms for the different periods are independent of each other. This implies that there will likely be an unobserved fixed effect in the error term causing omitted variable bias. The fixed effect does not vary over time but between units. We denote the component of the error term that is specific for a state as  $A_i$ , while  $v_{it}$  denotes the component that is idiosyncratic to observations. Thus, we can write the model as in equation 2,

$$y_{it} = \alpha_i + \lambda t + \beta x_{it} + v_{it} \quad (2)$$

where  $\alpha_i = \alpha + A_i$  and can be interpreted as each state's unique intercept for the regression line. This is illustrated in Figure 5 with the regression lines for California (CA), Texas (TX) and Alaska (AL). The intercept  $\alpha_i$  is unique for each state, while the slope,  $\beta$ , is the same for each state.

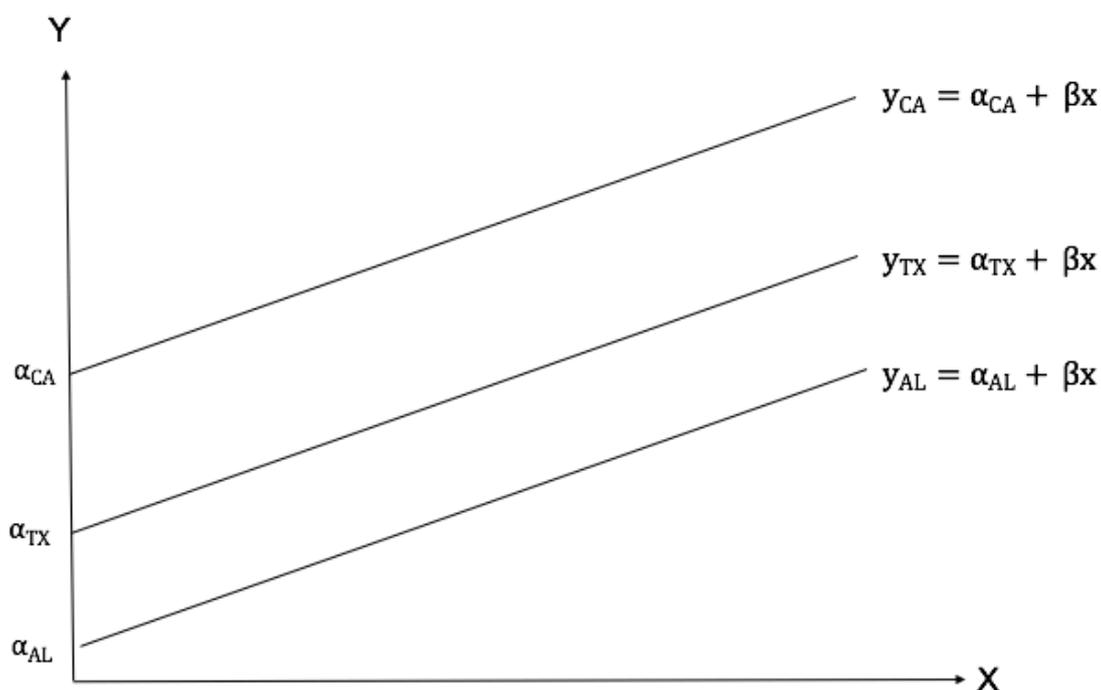


Figure 5: State fixed effects

We can estimate the causal effect  $\beta$  by treating each fixed effect as a parameter to be estimated. Treating each state fixed effect as a parameter to be estimated is algebraically the same as estimation in deviations from means. This method is called the within-group estimator. We also treat the year effects as a parameter to be measured. The estimation model will then look like presented in equation 3.

$$y_{it} - \bar{y}_i = \beta (x_{it} - \bar{x}_i) + (v_{it} - \bar{v}_i) \quad (3)$$

Using the fixed effects model controlling for both unobserved time fixed and state fixed effects allows us to measure the causal effect of interest,  $\beta$  (Torres-Reyna, 2007).

## 7.2 The Estimation Model

A fixed effects model for predicting election outcomes was introduced in 1998 by Strumpf and Phillippe. They used a model of voter choice where the voters difference in expected welfare for the two candidates was measured as the difference between a voter best forecast of the incumbent party's performance, given available information about economic factors, and the challenger's relative ideological advantage (Strumpf & Philippe, 1999). The use of a

fixed effects model eliminates the omitted variable bias that is caused by intrastate time-invariant factors, such as unobservable across state party preferences. This bias would cause the effect on election outcome to be over or underestimated in states where for example the Republicans already are more inclined to win the election. In our estimation model we also control for year fixed effects, ensuring that factors that are specific for one of the election years does not interfere with the coefficients.

The basis for our model is the same as the model of voter choice used by Strumpf and Philippe (1999). The dependent variable is the election outcome for the Republican Party in 2008, 2012 and 2016. As regressors we use statewide polling data on the predicted election outcome for the Republican Party prior to the elections in 2008, 2012 and 2016, proxies on immigration skepticism, far-right sentiment and racial animus, as well as some economic and demographic control variables. Strumpf and Phillippe omits polling data in their model as they believe that it only reflects historical economic information, and does not provide an independent source of information. We choose to include polling data as it is the main predictor of election outcome, see section 3.2, and we would thus expect that this data has a significant explanatory power when it comes to explaining election outcome. In addition to the polls data, we include statewide unemployment rate, level of education, level of black or African Americans, level of Latinos and Hispanics, level of male voters and mean income as control variables. Including these variables in the model specification controls for observed state fixed effects, and help avoid omitted variable bias by contributing to the conditional mean independence assumption that the error term is independent of the variables of interest. Our fixed effect estimation model is shown in equation (4),

$$y_{it} = \beta_1 \gamma_{LTit} + \beta_2 \delta_{LTit} + \beta_3 \theta_{LTit} + \beta_4 \gamma_{STit} + \beta_5 \delta_{ST,it} \quad (4)$$

$$+ \beta_6 \theta_{STit} + \beta_7(\text{Control}) + \lambda t + \alpha_i + v_{it}$$

where  $\gamma_{it}$  is the proxy for immigration skepticism in state  $i$  at time  $t$ ,  $\delta_{it}$  is the proxy for far-right sentiment in state  $i$  at time  $t$  and  $\theta_{it}$  is the proxy for racial animus in state  $i$  at time  $t$ . LT denotes the long-term proxies while ST denotes the short-term proxies. Unemployment, median age, percent of population with a bachelor's degree or higher, percent of population with High School or less, mean income, Hispanic or Latino percentage of population and

black or African American percentage of population is included as control variables, denoted by (*Control*) in equation 4.  $\lambda t$  is the year fixed effect,  $\alpha_i$  is the state fixed effect and  $v_{it}$  is the idiosyncratic error term. For the purpose of the model, we assume that there were no major changes within a state from 2008 to 2016 that is not reflected in the polling data or in our control variables.

#### *District of Columbia*

District of Columbia is an outlier in the dataset, in terms of having a stable Democratic vote share above 90 % in all three elections. The state also has special features compared to many of the other U.S. states. According to Census (2017), 681,170 people live in District of Columbia, while 831,531 people over the age of 16 work there. Stephens-Davidowitz and Varian (2015) mentions the same in “A Guide to Google Data”, implying that a lot of the search activity might come from the commuters who cast their vote in a different state. For these reasons, we choose to exclude District of Columbia from the empirical analysis.

## 8 Empirical Analysis

This chapter presents the results of five regressions conducted using a fixed effects model, controlling for different variables. All regressions are reported in Table 4. Our main results are from the estimation including all proxies, controlling for year effects and using standard errors clustered at state level. This model has the highest explanatory power, see column (4) in Table 4<sup>16</sup>. We further present the estimation results from analyses where we have estimated regional effects, swing state effects and the separate effect for the presidential election in 2008, 2012 and 2016. In the end of chapter 8, we point out limitations to the estimation strategy.

### 8.1 Main Results

Our main results are shown in column (4) in Table 4.

Table 4: Fixed Effect Estimates on Republican U.S. Presidential Election Results, 2008 – 2016

	(1) Only proxies	(2) Adding year fixed effects	(3) Adding robust st. errors	(4) Main results, all variables, clustered st. errors	(5) Only control variables
Long-term Racial Animus	-0.000333 (0.000488)	-0.000246 (0.000501)	-0.000246 (0.000636)	0.0000621 (0.000326)	
Short-term Racial Animus	-0.0000964 (0.000245)	-0.000333 (0.000296)	-0.000333 (0.000323)	0.000238 (0.000181)	
Long-term Immigration Skepticism	0.000288 (0.000439)	0.000520 (0.000443)	0.000520 (0.000751)	0.000798* (0.000344)	
Short-term Immigration Skepticism	-0.000208 (0.000227)	0.000105 (0.000253)	0.000105 (0.000254)	0.000199 (0.000156)	

<sup>16</sup> Regression (4) has the highest within R-squared value

Short-term far-right sentiment	0.000224 (0.000422)	0.000552 (0.000453)	0.000552 (0.000444)	0.000474 (0.000260)	
Long-term far-right sentiment	0.00107* (0.000534)	0.00154* (0.000588)	0.00154* (0.000715)	0.000676 (0.000363)	
2012.year		-0.00185 (0.0124)	-0.00185 (0.0109)	0.0235 (0.0129)	0.0397** (0.0116)
2016.year		-0.0314 (0.0169)	-0.0314 (0.0234)	0.0674** (0.0203)	0.0933*** (0.0169)
Election Polls				0.303* (0.139)	0.356 (0.181)
% Male				3.488* (1.562)	3.062 (1.750)
Median age				-0.0165** (0.00572)	-0.00836 (0.00638)
% Hispanic or Latino				-1.215** (0.394)	-0.680 (0.388)
% African American				0.422* (0.204)	0.577* (0.245)
Mean income U.S.D				-0.00000557* (0.00000219)	-0.00000430* (0.00000214)
<i>State fixed effects</i>	Yes	Yes	Yes	Yes	Yes
<i>Time fixed effects</i>	No	Yes	Yes	Yes	Yes
<i>Robust st. errors</i>	No	No	Yes	Yes	Yes
<i>Economic control variables</i>	No	No	No	Yes	Yes
<i>Google trends proxies</i>	Yes	No	Yes	Yes	No
<i>N</i>	150	150	150	150	150
<i>R</i> <sup>2</sup>	0.137	0.196	0.196	0.748	0.680

Standard errors in parentheses

\*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$

Table 4: Main Results

*Notes:* We have conducted five regressions using a fixed effects model. We have controlled for state fixed effects in all of the regressions. In the first regression, we did not include the economic control variables and did not control for year fixed effects or state-level clustered standard errors. In the second one we included the year fixed effects and in the third regression we added state-level clustered standard errors as well, in order to control for heteroscedasticity and autocorrelation in the idiosyncratic error term. In the fourth regression, we included all variables and controlled for all factors mentioned. Only significant control variables are presented here. At last we conducted one regression where we left out of all the proxies to see if this could indicate that the model including the proxies is the better model.

We find that the proxy for the level of long-term immigration skepticism in a state has a small, but significant, effect on the Republican election outcome in the U.S. presidential election. A 1 index point increase in the measured immigration skepticism sentiment yields a 0.000798 percentage point increase in the Republican election outcome, significant at a 5 percent level <sup>17</sup>. This indicates that when the level of immigration skepticism rises, measured by higher search frequency for the query topic [illegal immigration and residence], people are slightly more likely to vote for the Republican Party. We also find that both the short-term and long-term proxy for far-right sentiment is significant on a 10 percent level, and they have a small, but positive, effect on the Republican election outcome. A 1 index point increase in the measured short-term far-right sentiment yields a 0.000474 percentage point increase in the Republican election outcome, while the long-term sentiment yields a 0.000676 percentage point increase. This implies that when the search frequency for the selected right-oriented webpages increases in a state, people are more likely to vote for the Republican candidate.

Neither the short-term proxy for immigration skepticism, nor any of the proxies for racial animus in a state, seem to have any significant effect on the Republican election outcome in the past three U.S. presidential elections.

An increase of 1 percentage point in the election polls' prediction of the Republican vote-share, results in a 0.303 increase in the actual election outcome. The effect is significant at a 5 percent level, see Table 4 column (4).

---

<sup>17</sup> The socially sensitive sentiment has an index score between 0 and 100 for the different states. A 1-point increase in the index score is therefore equal to a change of 1.

### 8.1.1 Effect of control variables

We observe significant effects on the share of Republican votes for male ratio, median age, the share of Hispanic residents, the share of black or African American residents, and mean income, see Table 4 column (4). The male ratio and the level of black or African American people in the population have a positive effect on the Republican election outcome, while the other control variables have a negative effect. Median age and the share of Hispanics in the population are significant at a 1 percent level, while the other effects are significant at a 5 percent level. The effect of an increase in mean income is very small, so the statistical importance might be debated. An increase in mean income of 10 000 U.S. dollars decreases the Republican election outcome by 0.0557 percentage points. The male ratio has a relatively large effect on the Republican election outcome: a 1 percent increase in the male ratio in a state, increases the Republican election outcome by 3.5 percentage points.

## 8.2 Testing for Regions

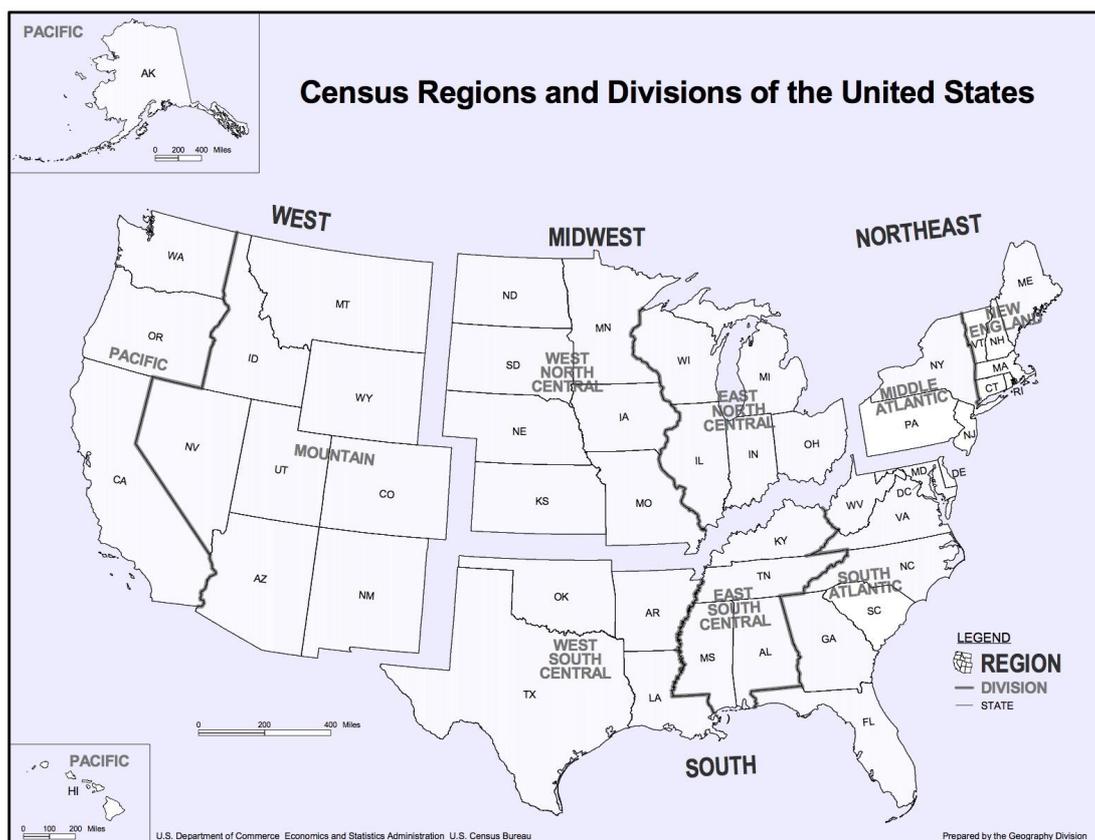


Figure 6: Census Regions and Divisions of the United States (*United States Census Bureau, 2017*).

We test if the overall effect found in Table 4 is different for various geographical areas in the U.S.. We used the four U.S. regions Midwest, Northeast, South and West as configured by the United States Census Bureau, see Figure 6.

We perform one regression for each region using a fixed effects model, controlling for year fixed effects. We do not control for serial correlation by clustering standard error at state level here, due to a low number of clusters in each of the regional regressions. All regressions are reported in Table 5. In the following section, we present the key results from the four different regressions.

Table 5: Fixed Effect Estimates on Republican U.S. Presidential Election Results, 2008 – 2016, clustered at U.S. Census Regions

	(1) Midwest	(2) Northeast	(3) South	(4) West
Long-term Racial Animus	0.000107 (0.00121)	0.000797 (0.000268)	0.00117* (0.000656)	0.00114* (0.00143)
Short-term Racial Animus	0.000248 (0.000676)	-0.00163 (0.000385)	0.000178 (0.000250)	0.000952 (0.000708)
Long-term Immigration Skepticism	0.000594 (0.000702)	0.000106 (0.000243)	0.000139 (0.000378)	-0.000191 (0.00147)
Short-term Immigration Skepticism	0.000607 (0.000583)	0.00175* (0.000254)	0.0000132 (0.000260)	0.000463 (0.000673)
Short-term far-right sentiment	0.00105 (0.000594)	0.00183* (0.000247)	-0.0000146 (0.000389)	-0.000441 (0.000928)
Long-term far-right sentiment	0.000602 (0.000639)	-0.00186 (0.000421)	-0.000201 (0.000595)	0.00324 (0.00144)

Election Polls	0.359 (0.178)	1.138 (0.122)	-0.0209 (0.0635)	0.695*** (0.184)
<i>State fixed effects</i>	Yes	Yes	Yes	Yes
<i>Time fixed effects</i>	Yes	Yes	Yes	Yes
<i>Robust st. errors</i>	No	No	No	No
<i>Economic control variables</i>	Yes	Yes	Yes	Yes
<i>Google Trends proxies</i>	Yes	Yes	Yes	Yes
<i>N</i>	36	27	48	39
$R^2$	0.956	0.998	0.930	0.940

Standard errors in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table 5: Results Clustered at U.S. Regions

In the Northeast region, see column (2) in Table 5, the short-term proxies for immigration skepticism and far-right support are both significant at a 10 percent significant level. The effects are small, estimating an increase of 0.00175 percentage points in the Republican election outcome, given a 1 index point increase in the short-term immigration skepticism proxy. An increase of 0.00183 percentage points is estimated for a 1 index point increase in the short-term far-right proxy. This effect is higher than for the U.S. in total.

We observe that in the South, see column (3) in Table 5, the long-term proxy for racial animus is significant at a 10 percent level, indicating a positive relationship between the long-term racial animus and Republican votes. If the proxied long-term racial animus increases with 1 index point, the Republican election outcome will increase with 0.00117 percentage points. The same applies for the West, see column (4) in Table 5. The effect is small, estimating an increase of 0.00114 percentage points in the Republican election outcome, given a 1 index point increase in the long-term racial animus.

None of the control variables came out significant in the regression, except the election polls. We see that the effect of election polls varies between the regions, and that the effect is only

significant in West. This will be further discussed in section 9.1. For the West, the effect of election polls indicates a significant positive relationship between the polls and the election results, which is to be expected. If the Republican election outcome predicted by the polls increase by 1 percentage point, the actual Republican election outcome will increase by 0,695 percentage points. This is significant at a 1 % level.

### **8.3 Testing for Swing States**

Due to both confirmation bias and information utility, it is fair to assume that people will search for information to get a better decision base if they experience uncertainty ahead of an election. Given this, we find it interesting to test our model on swing states. Our hypothesis is that the effect of the proxies might be stronger in swing states, especially for the short-term proxies. Nate Silver (2016) defines the following states as traditional swing states, i.e. states that have regularly experienced close races in the last U.S. presidential elections, see List 1.

#### Traditional Swing States

---

Colorado  
Florida  
Iowa  
Michigan  
Minnesota  
Nevada  
New Hampshire  
North Carolina  
Ohio  
Pennsylvania  
Virginia  
Wisconsin

---

List 1, (Silver, 2016)

The results of the swing state regressions can be found in Table 6. The estimation results provide no evidence to support our hypothesis that the effect of the proxies is stronger in swing states.

Table 6: Fixed Effect Estimates on Republican U.S. Presidential Election Results, for Swing States in 2008 – 2016

	(1) Proxies and polls	(2) Controlling for time-effects	(3) Controlling for economic variables
Short-term Racial Animus	-0.000698* (0.000293)	-0.000696 (0.000438)	-0.000350 (0.000624)
Short-term Immigration Skepticism	0.000673* (0.000242)	0.000555 (0.000293)	-0.0000627 (0.000397)
2012.year		0.0120 (0.0142)	0.0776* (0.0314)
2016.year		0.0173 (0.0173)	0.138 (0.0628)
Election Polls	0.580*** (0.114)	0.543*** (0.126)	0.256 (0.146)
<i>State fixed effects</i>	Yes	Yes	Yes
<i>Time fixed effects</i>	No	Yes	Yes
<i>Robust st. errors</i>	No	No	No
<i>Economic control variables</i>	No	No	Yes
<i>Google trends proxies</i>	Yes	Yes	Yes
N	36	36	36
R <sup>2</sup>	0.780	0.795	0.932

Standard errors in parentheses

\*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$

Table 6: Swing states estimates

## 8.4 Testing for each Presidential Election

Based on the technological development from 2008 to 2016, especially in terms of personal technical gadgets, which has contributed to a significant increase in Google searches, we find it interesting to do a cross-sectional analysis for the different elections. An important note to this exercise is awareness of the limitations of doing a cross-sectional analysis. There are several issues with using a cross-sectional method to estimate a model where the parameters

vary across states and over time. Because the parameters are assumed constant across states and time in such a model, a cross-sectional analysis often fail to identify the true parameters, leading to biased inferences. This will limit the scope of measuring the true relationship investigated, and we will not be able to determine any causal effects (Bowen & Wiersema, 1999). We still include the additional cross-sectional result to add some insight to the discussion around our main results. Significant variables are presented in Table 7.

Table 7: Cross-Sectional Estimates on Republican U.S. Presidential Election Results, 2008 – 2016

	(1) 2016	(2) 2012	(3) 2008
Short-term Racial Animus	0.00163*** (0.000449)	-0.000110 (0.000368)	0.000300 (0.000345)
Short-term far-right sentiment	0.00154* (0.000687)	0.00106 (0.000842)	0.000134 (0.000419)
Election Polls	0.807*** (0.0909)	0.850*** (0.0909)	0.851*** (0.0657)
Median age	-0.00939* (0.00358)	-0.00744 (0.00385)	-0.00642* (0.00286)
% African American	-0.0136 (0.113)	-0.228* (0.104)	-0.0831 (0.0661)
% Bachelor or more	-0.597* (0.289)	-0.143 (0.250)	-0.539** (0.184)
<i>N</i>	50	50	50
$R^2$	0.923	0.938	0.960

Standard errors in parentheses

\*  $p < 0.05$ , \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$

Table 7: Cross Sectional Estimates

The short-term racial animus proxy and the short-term immigration skepticism proxy are significant in 2016. Both proxies indicate a positive relationship with the Republican election

outcome, though the effect is very small for a 1 index point increase in the proxies. We note that the effect of election polls appears to have fallen since 2008, and that the standard error has increased with the years. We discuss these elements further in section 9.1.

### **8.5 Limitations to the Estimation Strategy**

In the main results for the U.S. in total, and for the regressions segmented on region level and on swing states, we have controlled for time fixed effects using year-specific dummies. This may result in overfitting. An overfitted regression model creates estimates which fit all the little quirks and noise in the data, meaning the model will not be applicable for a new sample of data (Frost, 2015). Such a model is not representative for the population, and is therefore not a good estimation model. Time fixed effects should however be included if we assume that time-dependent exogenous shocks that affect the outcome may have occurred. We deem presidential elections to be highly exposed for time-dependent outcome shocks due to the history of heated campaigns, special candidate characteristics, scandals and a media industry on edge. We therefore include year-specific dummies in our main results, despite the chance of overfitted results.

## 9 Discussion

In the following chapter we will discuss our results, point out limitations to the dataset and look at the external validity of the study.

### 9.1 Discussion of the Results

#### *9.1.1 Main results*

An increase of 1 index point in the proxied long-term immigration skepticism in a state increases the Republican election outcome by 0.000798 percentage points, see Table 4 column (4). This effect is too small to be of any statistical importance. However, we do observe an increase in the mean index score of the long-term immigration skepticism proxy, see Table 1 in section 6.1. The mean index score increased by 20 points between 2008 and 2016, which indicates that some states have witnessed a high increase in measured long-term immigration skepticism. An increase of 20 points in the measured immigration skepticism proxy yields a 0.016 percentage point increase in the Republican election outcome, which is an effect of more statistical importance.

The positive sign of the long-term immigration skepticism coefficient implies that as more people in a state searches for queries within the topic of illegal immigration, the state is more likely to vote for the Republican Party. This is consistent with the theory that as people grow more skeptical towards immigration, they are more likely to vote for a far-right party, see section 4.2. The Republican candidate that won the 2016 election, i.e. the only election in our estimation model where the Republicans won, has advocated for stricter immigration policies associated with far-right parties, see section 4.2. It is also consistent with the perception that the Democratic Party typically is more liberal towards immigration, while the Republican Party typically advocates stricter policies towards immigration (Diffen, 2017). Our estimation results indicates that the immigration skepticism proxy catch some underlying sentiment not captured in election polls.

We observe an effect significant at a 10 percent level for both the short-term and long-term proxy for far-right sentiment. The mean index value of the long-term far-right sentiment across all U.S. states has increased by 14 points, see Table 1 in section 6.1. An increase of 14 index points in the proxied long-term far-right sentiment, yields a 0.009 percentage point increase in in the Republican election outcome. The positive sign of both coefficients is

consistent with our theory on the radical far-right sentiment, and that it propagates into election results, see section 4.1. The rise of radical far-right sentiment in a population could be hard to measure correctly due to SDB, and thus election polls would not be able to completely capture the emergence of the sentiment. It is difficult to say whether the measured sentiment in this estimation model captures something not observed in the election polls. The estimated effects are quite small and the statical importance can be debated. The estimation results still indicates that the proxy captures an underlying sentiment in the population, that the opinion polls does not reflect.

The observed effect of election polls on the actual election outcome is surprisingly low. Our findings are not consistent with the fact that the margin of error is usually between 3 and 6 percent for the sample sizes used in state polling, see section 3.2. A reason why we observe such a small effect, could be that not all polls data is collected within two weeks before the election, which would increase the margin of error. This is further discussed in section 9.2.

The effect of the polls data is only significant when we include our proxies on sensitive sentiments, see Table 4 column (4) and (5). This suggest that supplementing the election model with estimates on socially sensitive behavior could in fact give better outcome estimates, although the effect of our measured sentiments is either insignificant or very small. The within explanatory value of the model also improves when we include the proxies, which support the previous argument<sup>18</sup>.

An increase in mean income seem to have a negative effect on the vote outcome for the Republican Party. One could argue that this is consistent with the assumption that an increase in welfare increases the satisfaction with the incumbent government and decreases the likelihood of people voting for the opposing candidate, see section 6.3. The incumbent president in 2012 and 2016 was Democratic, and thus a vote for the opposing party would reflect a dissatisfaction with the economic situation. We have also argued that an increase in income can reflect an increase in education, see section 6.3. Given this correlation, the estimated effect is consistent with the belief that income is higher among Democrats. However, the estimated effect of mean income is very small, so we disregard its statistical

---

<sup>18</sup> We assume that the model is improves as the models within explanatory power increases. However, the R-square value of a model tend to increase with the number of variables, so a higher R-square value does not necessarily mean that the model has a good fit.

importance in this case. We discuss the implications of using mean income instead of median income in section 9.2.

The positive effect of an increase in male ratio on the election outcome for the Republican Party is consistent with the voting pattern for men and women in the U.S, see section 6.3. The effect we observe is quite high: a one percent increase in the male ratio increases the election outcome for the Republican Party by 3.5 percentage points. When we look at how the male ratio has developed, the mean ratio hardly changed between 2008 and 2016, see Table 3 section 6.3, implying that a 1 percentage point change in the male ratio within a state is rare.

The negative effect we observe for the level of the Hispanic or Latino population is also consistent with theory. Hispanics are in general much less likely to vote for the Republican candidate, and thus a state that experiences an increase in the Hispanic population (relative to other races) would be more likely to vote Democratic, see section 6.3. The observed effect is however quite large. An increase of 1 percentage point in the level of the Hispanic population decreases the election outcome for the Republican Party by 1.2 percentage points. An effect this large seem unlikely, as the Hispanic population in most states is not very high compared to other races, see Table 3 section 6.3. When discussing the relevance of this result one must also look at the turnout among Hispanic population to see how much this in reality would affect an election. We discuss this further in section 9.2.

An increase in the median age in a state seem to decrease the vote outcome for the Republican Party, which is also consistent with theory, see section 6.3. The size of the effect also seems plausible: an increase of the median age by one year decreases the election outcome for the Republican Party by 0.017 percentage points. Something that is problematic with median age is that as people live longer, median age tends to increase, which we observe in Table 3 in section 6.3. A major decrease in the median age in a state would thus be unlikely.

For the level of African American people in the population we observe that an increase in the level actually has a positive effect on the election outcome for the Republican Party. This is inconsistent with theory, which suggest that black or African American voters to a large extent identify with the Democratic Party, see section 6.3. The reason why we observe this effect could be a low turnout among African Americans compared to the rest of the

population, but this only holds if we assume that an increase of the level of African Americans in a state would lead the rest of the population to increase their likelihood of voting for the Republican Party. This assumption has little scientific evidence, so we disregard the statistical importance of this effect.

It is important to note that since the Democrats had the incumbent president in the majority of the time period, the positive effects we have detected could be due to a general wish for change in government, and not necessarily due to a growing support for the Republican Party or candidate.

### *9.1.2 Regional results*

Compared to the U.S. in total, the effect of short-term far-right sentiment on the Republican election outcome is 0.0014 percentage point higher in the Northeast. For immigration skepticism, the effect is approximately 0.0016 percentage point higher. This might seem minimal, but can be of more importance if we regard larger changes in the proxies than 1 index point. In the Northeastern region the short-term immigration skepticism index has increased by 30 from 2012 to 2016, see Table A1 in the appendix. This gives an increase of 0.056 percentage points in the Republican election outcome, thus 0.042 percentage points larger than the U.S. in total, which is a highly relevant effect.

We observe significant positive effects of long-term racial animus on Republican election outcome for both the South and the West. The relationship is not surprising in the South, as the top 5 U.S. states with the highest racial animus index, see Table A2 in the appendix, are Southern states, and 14 of 17 the Southern states voted Republican in the 2016 presidential election, see Figure A1-d) in the appendix. The positive significant relationship is also present in the West, which on average has the lowest share of racial animus of the regions, see Table A1 in the appendix. This can tell us that the effect is prominent in states that are not traditionally perceived to be racist, suggesting that there is a correlation between an increase in racial animus and the Republican election outcome.

### *9.1.3 Cross-sectional results*

The results in Table 6 indicated that both the short-term racial animus and short-term far-right sentiment were significant only in 2016. This may have many different explanations, and the

following reflection will not give us any statistical grounds to make inference from the analysis, but is included to supplement the discussion regarding the main results.

One hypothesis regarding the cross-sectional regression is that the proxy variables are more likely to have a significant effect on the election outcome in 2016 rather than in 2008 and 2012. In 2008, 8 billion google searches were conducted in the U.S. (Sullivan, 2009), almost all on desktop. Approximately 80 billion searches were conducted in the U.S. in 2012, where 60 billion came from a desktop and 20 billion from a mobile device. In 2016 almost 200 billion searches were conducted in the U.S., where over 110 billion of the searches came from a mobile platform (Allen, 2017). Thus, we believe that the Google Trends samples pulled from the population in 2016 compared to 2008, will contain a better distributed variance within the sample. As more and more people have constant access to Google through their smartphones, it is fair to assume that the samples will become more representative for the population in the future. Given this, one can assume that Google Trends data pulled ahead of the U.S. presidential election 2016, would provide better explanatory power for the 2016 election, compared to Google Trends data pulled in 2008 and 2012 explaining the 2008 and 2012 election.

From this we can partially conclude that the proxies will have better predictive power in the future compared to what our main results have indicated, given that Google searches continue to grow in usage.

## **9.2 Limitations to the Data Set**

In this section we will point out limitations to our data set. Among them are the limitations in the statewide polling data which is collected at different times, and the fact that we only have three time periods, since Google Trends data can only be collected after 2004.

We did not manage to collect statewide data provided by one polling company. This means that some of our polling data is based on different sampling methods, and is collected at different times. This is potentially a source of error in our estimation and a weakness in our model. The low explanatory power provided by statewide election polls in our estimation model is inconsistent with earlier research on the predictive power of election polls.

An important aspect to note regarding the search data, is that when we segment the U.S. data for the different states, the query index is the average index score over the weekly data for the given time period, relatively to the other U.S. states (Google, 2017). Thus, the only time variation in the proxy is the three time periods defined. The proxies is constructed this way to capture the within-state variation for all the U.S. states. For cases where the main objective is to look at the time variation over a given period, e.g. applying the seasonal decomposition method proposed by Shimshoni et al. (2009) to different search terms, it is possible to do a comparison of maximum five states at the time in the query tool (Google, 2017).

Due to the privacy considerations in Google Trends, we could not collect data on explicit search words that could better reflect the social sentiments we attempted to capture. We believe that the far-right sentiment proxies and the immigration skepticism proxies would reflect reality better if we had we been able to use more explicit search terms. This is also discussed in section 6.1.

When we constructed the short-term proxies, we did not get an index value for some of the search queries in certain states. This is because of the privacy considerations in Google Trends. In order to avoid very large variations in the proxy values for the states, we normalized the index value for these states to be equal to the lowest index value among the other states for that query. This issue regards the queries [Breitbart News] in 2008, [Stormfront] for all the three time periods and [Illegal immigration and residence] in 2012. For the states that did not show any results for [Breitbart News], the index value was normalized to 5. For the states that did not show any results for [Stormfront], the index value was normalized to 3. For the query topic [illegal immigration and residence], the index value was normalized to 19. This is a problem that often occurs when one looks at short time periods in Google Trends, as earlier argued. This was impossible to avoid without using very general search terms. For the short-term far-right sentiment proxy, two of the search terms used have been altered, but since the normalized index values are relatively small, it is likely that it does not affect the estimation results significantly. For the short-term immigration skepticism proxy, the normalized index value is higher, which means that the normalization to a small degree can have affected the estimation results.

There are some limitations regarding the control variables used in the model. For one, research has shown that economic variables on national level typically are more important

determinants of voting (Strumpf & Philippe, 1999). This implies that the model would improve if we controlled for national unemployment rates and income data, and that our regional data is of limited importance. The use of mean income is also problematic, as the income data within a state is not symmetrically distributed. We would have preferred to use median income, but this data was not yet available on a state level for 2016.

Further, it would have been preferable to look at the share of Hispanic or Latino and black or African American *voters* in a state instead of the share of *population*. The problem of using the population is that it does not necessarily reflect the vote turnout. An increase in the African American population does not necessarily increase the African American share of voters, and thus the increase in the population is of limited importance to election research. In addition, Stephens-Davidowitz (2013) points out that election turnout amongst black or African Americans is typically low.

### **9.3 External Validity of the Study**

The external validity of our research is varying. The long-term proxy for immigration skepticism is built on a topic search, and is directly applicable for most countries. We tested the query topic [Illegal immigration and residence] and segmented on regions for randomly chosen countries, using the U.S. presidential election long-term time frame for 2016. This yielded incomplete results for most of the countries in Google Trends, due to the privacy threshold. We got the following region turnout in Google Trends for [Illegal immigration and residence]: 8 of 13 regions in Canada, 6 of 19 regions in Norway, 10 of 22 regions in France, 18 of 27 in Brazil and 19 of 36 regions in India. Belgium showed 3 of 3 regions, but this is too few cross-sectional elements to build a solid fixed effects model. We tested the query on the German federal election in September 2017, using the appropriate long-term time frame for this election, since immigration was a hot topic prior to that election as well (Mudde, 2017). [Illegal immigration and residence] showed query results for 14 of 16 regions. This shows that even for large countries where one would assume there would be high enough search frequency in all of the regions, the privacy threshold makes it difficult to get the complete region-wide data needed to build an appropriate fixed effect model. Thus, it can be difficult to use the methodology of this thesis to create a good supplement for opinion polls in other countries, given today's search data. For campaign strategists in other countries than the

U.S. the methodology can still be of interest, because the regions above the privacy threshold would still provide some interesting data.

The validity of the research in areas with strict governmental internet regulations will be minimal due to extremely low search frequency on sensitive topics. For China, we only got query results for Shanghai out of 31 regions, and for Russia we got for Moscow and St. Petersburg out of 89 regions using [Illegal immigration and residence]. The reason for this can also be the fact that they use other national search engines in a higher scale, or that illegal immigration is not a large issue in these countries. These areas should be followed closely by election scientists in the future, if the time comes for the governmental internet regulations to open up, as these countries might have some interesting features of social desirability bias and opinion polling bias in general.

For developing countries with limited access to the Internet, the research in this thesis is less valid as well. These countries would most likely only show query results for the region containing the capital, if they show regional results at all. In addition, many developing countries suffer from governmental regulations and non-democratic institutions. In such countries, the preferences of the people would not necessarily be reflected in the election outcome, and therefore it would be difficult to detect a relationship between the social sentiment in the population and the election outcome. A fixed effects model trying to measure a linear relationship would be of low statistical importance if the underlying process is by no chance to be assumed as linear, due to the amount of noise in the model. This means that opinion polling alone will still be a better measure for many of these countries.

## 10 Conclusion

This thesis aims to answer the research question “*Can Google Trends data be used as a proxy for socially sensitive sentiments, and can such proxies be used in models explaining election outcomes?*”. Using theory on social desirability bias, we explain why opinion polls prior to elections sometimes fail in predicting the correct outcome. Further, we use confirmation bias and information utility to explain how internet search data can provide valid measures of socially sensitive sentiments in a population. We also provide theory on how Google search data can be applied in empirical research and on how this data is collected. We argue that the three following socially sensitive sentiments: *far-right sentiment, immigration skepticism and racial animus*, are sentiments of interest to the U.S. presidential election, why they might not be captured in opinion polls and we explain how we can use Google search data to build proxies for these sentiments. Using a one-year period and a two-weeks period ahead of each election, we have constructed long-term and short-term proxies for each social sentiment.

Using a fixed effects model controlling for both state fixed and year fixed effects, we estimate the effect of the socially sensitive sentiment proxies on the election outcome for the Republican Party in the 2008, 2012 and 2016 U.S. presidential election. Our results suggest that some of our social sentiment proxies have an explanatory power for the election outcome in the past three U.S. presidential elections. This implies that proxies based on web search data can be used in models explaining election outcomes. We find significant effects for the proxy for long-term immigration skepticism, as well as both proxies for far-right sentiment. The estimated effects are, however, small for a 1 index point increase. For the long-term immigration skepticism proxy, we observe an estimated effect of 0.000789 percentage points on the election outcome for the Republican Party, if the immigration skepticism proxy increases by 1 index point. For the far-right sentiment proxies the effect is 0.000676 for the long-term proxy and 0.00474 for the short-term proxy.

Even though the effects are small, the results can still serve as a complement for polling agencies. We argue that one should assume changes larger than 1 index point due to the nature of our proxies, which suggests that the estimated effect could be around 0,01 percentage points or higher. This might still seem small, but the effect is not negligible. In a

presidential election, a 1 percentage point change in a party's election outcome could be decisive in which party wins the state.

Further, we find heterogeneity across regions for all proxies. This insight is valuable for campaign strategists in order to choose which states to visit and arrange speeches. Strategists can also learn about the underlying sentiments within a region, and tailor the candidate's message accordingly.

We do not find any statistical evidence to support the hypothesis of higher effects of the proxies in swing states. The lack of statistical evidence for this hypothesis could indicate that for these states, the analyzed sentiments are not decisive in which party who wins the election in the end.

The external validity of our research is varying, due to three main reasons. First, in most countries, Google's privacy threshold makes it unfeasible to extract search data and examine the proxies on a regional level. Second, for countries where the Internet is subject to strict governmental regulations, the search frequency of sensitive topics will be minuscule, rendering this analysis unviable. Third, in developing countries, both the lack of Internet access and democratic institutions are detrimental to the ability to draw inference from this model.

This thesis contributes to research on how social sentiments in a population affects election outcomes, by investigating how hard-to-measure sentiments can be used in models explaining election outcome. It also contributes to research on how web search data can be utilized in economic and social research. We conclude that the results are applicable in the U.S. due to its large statewide variation, and should be included in the work by opinion pollsters and campaign strategists. We further conclude that the results are of less validity outside the U.S. due to Google's privacy threshold, given the search data available today.

## References

- Ad Hoc Committee on 2016 Election Polling. (2016). *An Evaluation of 2016 Election Polls in the U.S.* From American Association for Public Opinion Research: <http://www.aapor.org/Education-Resources/Reports/An-Evaluation-of-2016-Election-Polls-in-the-U-S.aspx>).
- Agbafé, V. (2016, January 18). Immigration and the 2016 Election. *Harvard Political Review*.
- Allen, R. (2017, April 13). *Search Engine Statistics 2017*. From Smart Insights: <https://www.smartinsights.com/search-engine-marketing/search-engine-statistics/>
- Anti-defamation League. (2017). *Racism*. From ADL: <https://www.adl.org/racism>
- Arnold-Smeets, L. (2013). *Survey Says: 87 Percent of Employees Do not Trust Their Bosses*. From Payscale Human Capital: <https://www.payscale.com/career-news/2013/12/survey-says-87-percent-of-employees-don-t-trust-their-bosses>
- Atkin, C. (1973). Instrumental Utilities and Information Seeking. *SAGE Journal*.
- BBC . (2016). *Q&A: How Do Opinion Polls Work*. From BBC News: <http://www.bbc.com/news/uk-politics-35350361>
- BBC . (2017, August 13). *White supremacy: are U.S. far-right groups on the rise?* From BBC World: <http://www.bbc.com/news/world-us-canada-40915356>
- BBC. (2017). *Why Is Obamacare so Controversial*. From BBC World: <http://www.bbc.com/news/world-us-canada-24370967>
- Belli , R. F., Traugott, M. W., & Beckmann, M. N. (2001). What Leads to Voting Overreports? Contrasts of Overreporters to Validate Voters and Admitted Nonvoters in the American National Election Studies. *Journal of Official Statistics*, 17, 479-498.
- Berggren, N., Jordahl, H., & Poutvaara, P. (2010, December 20). The Right Look: Conservative Politicians Look Better and Voters Reward It. *Journal of Public Economics*, 79-86.
- Betz, H.-G. (1994). *Radical right Wing Populism in Western Europe*. The Macmillan Press LMT.
- Blais, A. (2000). *To Vote or Not to Vote* . University of Pittsburg Press.
- Blake, A. (2016, April 27). *Why are There Only Two Parties In American Politics*. From Washington Post: [https://www.washingtonpost.com/news/the-fix/wp/2016/04/27/why-are-there-only-two-parties-in-american-politics/?utm\\_term=.3ec62756b153](https://www.washingtonpost.com/news/the-fix/wp/2016/04/27/why-are-there-only-two-parties-in-american-politics/?utm_term=.3ec62756b153)
- Blumenthal, M., & Edwards-Levy, A. (2013, June 5). *Gallup Poll Reveals 4 Reasons It Got The 2012 Election Wrong*. From Huffington Post: [https://www.huffingtonpost.com/2013/06/04/gallup-poll-2012\\_n\\_3384882.html](https://www.huffingtonpost.com/2013/06/04/gallup-poll-2012_n_3384882.html)
- Bowen, H. P., & Wiersema, M. F. (1999, July 20). Matching Method to Paradigm in Strategy Research: Limitations of Cross-Sectional Analysis and Some Methodological Alternatives. *Strategic Management Journal*, 20(7), pp. 625-636.
- Bromwich, J. E. (2016, August 17). *What is Breitbart News*. From The New York Times: <https://www.nytimes.com/2016/08/18/business/media/what-is-breitbart-news.html?ribbon-ad-idx=4&rref=politics>
- Brownback, A., & Novotny, A. (2017). Social Desirability Bias and Polling Errors in the 2016 Presidential Election. *Arkansas Poll*, 2-5, 14.
- Brownstein, J. S., Freifeld, C. C., & Madoff, L. C. (2009). Digital Disease Detection — Harnessing the Web for Public Health Surveillance. *The New England Journal of Medicine*(360), pp. 2153-2157.

- 
- Castillo, W., & Schramm, M. (2016, November 9). *How we voted — by age, education, race and sexual orientation*. From USA Today :  
<http://college.usatoday.com/2016/11/09/how-we-voted-by-age-education-race-and-sexual-orientation/>
- Chen, Y., Zhang, F., & Yue, Y. (2012, December 14). *Predicting U.S. Presidential Election result based on Google Insights*. Stanford University.
- Choi, H., & Varian, H. (2012, June). Predicting the Present With Google Trends. *Economic Record*, 88, 2-9.
- Choi, H., & Varian, H. (2012, June). Predicting the Present With Google Trends. *Economic Record*, 88(s1), pp. 2-9.
- CNN. (2008). *Election Center 2008 - Campaign Issues*. From CNN Election:  
<http://edition.cnn.com/ELECTION/2008/issues/>
- CNN. (2016, November 23). *Exit Polls, 2016 Election*. From CNN Politics:  
<http://edition.cnn.com/election/results/exit-polls/>
- Cooper, C. P., Mallon, K. P., Leadbetter, S., Pollack, L. A., & Peipins, L. A. (2005, Jul-Aug 1). Cancer Internet Search Activity on a Major Search Engine, United States 2001-2003. *Journal of Medical Internet Research*, 7(3).
- Diffen. (2017). *Democrat vs. Republican, Stand on Immigration*. From Diffen, Compare Anything : [https://www.diffen.com/difference/Democrat\\_vs\\_Republican](https://www.diffen.com/difference/Democrat_vs_Republican)
- Dorfman, J. (2014, December 2014). *The Obvious Reason for the Decline in Median Income*. From Real Clear Markets:  
[http://www.realclearmarkets.com/articles/2014/12/02/the\\_obvious\\_reason\\_for\\_the\\_decline\\_in\\_median\\_income\\_101409.html](http://www.realclearmarkets.com/articles/2014/12/02/the_obvious_reason_for_the_decline_in_median_income_101409.html)
- Durankiev, G. (2015, June 17). *11 Most Racist States Ranked By Hate Crimes*. From Insider Monkey: <https://www.insidermonkey.com/blog/11-most-racist-states-ranked-by-hate-crimes-353960/10/>
- Egeland, J. (2014). *This Is the Worst Refugee Crisis Since WWII. It's Time for Us to Rethink Our Response*. From Huffington Post: [https://www.huffingtonpost.com/jan-egeland/refugee-crisis-wwii-aid-\\_b\\_5791776.html](https://www.huffingtonpost.com/jan-egeland/refugee-crisis-wwii-aid-_b_5791776.html)
- Ejara, D., Nag, R., & Upadhyaya, K. (2008). Opinion polls and the stock market: evidence from the 2008 US presidential election. *Applied Financial Economics*, 22(6), pp. 437-443.
- Electoral Vote. (2017). *Data Galore*. From Electoral Vote: <http://www.electoral-vote.com/evp2017/Info/datagalore.html>
- Electoral Vote. (2017). *Polling FAQ*. From Electoral Vote: <http://electoral-vote.com/evp2017/Info/polling-faq.html>
- Enten, H. J. (2012, November 21). *How the 2012 Election Polling Really Was Skewed - for Mitt Romney*. From The Guardian:  
<https://www.theguardian.com/commentisfree/2012/nov/21/2012-election-polling-skewed-for-mitt-romney>
- Ettredge, M., Gerdes, J., & Karuga, G. (2005). Using web-based search data to predict macroeconomic statistics. *Communications of the ACM*, 48(11), 87-92.
- Federal Election Commission. (2009). *Federal Elections 2008, Election Result for the U.S. President, the U.S. Senate and the U.S. House of Representatives*. Washington D.C.
- Federal Election Commission. (2013). *Election Results for the U.S. Presidential, the U.S. Senate and the U.S. House of Representatives*. Washington D.C.
- Festinger, L. (1957). *A Theory of Cognitive Dissonance*. Stanford University Press.

- 
- Fisher, R. J. (1993). Social Desirability Bias and the Validity of Indirect Questioning. *Journal of Consumer Research*, 20, 303-315.
- Ford, M. (2017, January 22). *The Far Right in America: A Brief Taxonomy*. From The Atlantic: <https://www.theatlantic.com/politics/archive/2017/01/far-right-taxonomy/509282/>
- Fredrickson, G. M. (2002). *Racism: A Short History*. Princeton University Press.
- Frost, J. (2015, September 3). *The Danger of Overfitting Regression Models*. From The Minitab Blog: <http://blog.minitab.com/blog/adventures-in-statistics-2/the-danger-of-overfitting-regression-models>
- Gelman, A., & Kaplan, N. (2008, April 5). Voting As a Rational Decision. *CERS Policy Portal*.
- Ginsberg, J., Mohebbi, M. H., Patel, R. S., Brammer, L., Smolinski, M. S., & Brilliant, L. (2009, February 19). Detecting influenza epidemics using search engine query data. *Nature*(457), pp. 1012-14.
- Goel, S., Hofman, J. M., Lahaie, S., Pennock, D. M., & Watts, D. J. (2010). Predicting Consumer Behaviour With Web Search. *Proceedings of the National Academy of Sciences of the United States of America* .
- Google. (2017). *Google Trends*. From Google Trends: <https://trends.google.com/trends/>
- Iceland, J. (2003, August). Why Poverty Remains High: The Role of Income Growth, Inequality and Changes in Family Structure, 1949-1999. *Demography* , 40.
- Inglehart, R. F., & Norris, P. (2016, August). Trump, Brexit, and the Rise of Populism: Economic Have-Nots and Cultural Backlash. *HKS Faculty Research Working Paper Series*.
- Inglehart, R., & Norris, P. (2000, October 1). The Developmental Theory of the Gender Gap: Women's and Men's Voting Behavior in Global Perspective. *International Political Science Review*, 21(4).
- Internet Live Stats. (2017). *Internet live stats*. From Internet live stats: <http://www.internetlivestats.com/>
- Jones, J., & Salter, L. (2011, October 11). *Digital Journalism*. SAGE.
- Kane, J. G., Craig, S. C., & Wald, K. D. (2004). Religion and Presidential Politics in Florida: A List Experiment. *Social Science Quarterly*.
- Kantchev, G., & Whittall, C. (2017, June 10). *After Another Surprise Election Result, Investors Grow Sceptical of Pollsters*. Retrieved November, 2017 from The Wall Street Journal: [https://www.ted.com/talks/eli\\_pariser\\_beware\\_online\\_filter\\_bubbles/transcript](https://www.ted.com/talks/eli_pariser_beware_online_filter_bubbles/transcript)
- Kennedy, R. (2003). *Nigger: The Strange Career of a Troublesom Word*. New York: Vintage Books.
- Kever, J. (2017). *Does a Candidate's Religion Matter to Voters*. From University of Houston: <https://www.uh.edu/news-events/stories/2017/June/07112017Religion-in-Politics.php>
- Kiersz, A. (2012, December 12). *Here is why stocks have been on a tear since Trump's election*. Retrieved November, 2017 from Business Insider: <http://www.businessinsider.com/stock-markets-after-trump-election-2016-12?r=US&IR=T&IR=T>
- King, A., & Leigh, A. (2010). Bias at the Ballot Box? Testing Whether Candidates' Gender Affects Their Vote. *Social Science Quarterly*.
- Kirk, A., & Scott, P. (2016, November 7). *US election: How age, race and education are deciding factors in the race for President*. From The Telegraph:

- <http://www.telegraph.co.uk/news/0/us-election-how-age-race-and-education-are-deciding-factors-in-t/>
- Knoblick-Westerwick, S. (2008, June 5). Information Utility. *The International Encyclopedia of Communication*.
- Knobloch-Westerwick, S., & Kleinman, S. B. (2011). Confirmation Bias Versus Information Utility. *SAGE Journals*.
- Kreuter, F., Tourangeau, R., & Presser, S. (2008, December). Social Desirability Bias in CATI, IVR, and Web Surveys: The Effects of Mode and Question Sensitivity. *Public Opinion Quarterly*, 72(5), pp. 847-865.
- Krogstad, J., Passel, J., & Cohn, D. (2017, April 27). *5 facts about illegal immigration in the U.S.* From Pew Research Center: <http://www.pewresearch.org/fact-tank/2017/04/27/5-facts-about-illegal-immigration-in-the-u-s/>
- Krumpal, I. (2013, June). Determinants of Social Desirability Bias in Sensitive Surveys: A Literature Review. *Quality & Quantity, International Journal of Methodology*, 47(4), pp. 2025-2047.
- Lauter, D. (2012, October 13). *Central Issues of Election 2012*. Retrieved October, 2017 from L.A Times: <http://timelines.latimes.com/central-issues-election-2012/>
- Müller, J.-W. (2017). *What Is Populism*. University of Pennsylvania Press.
- Morris, S. (2001). Political Correctness. *Journal of Political Economy*, 109(2), pp. 231-265.
- Muller, J. Z. (2008, April/May). Us and Them: The Enduring Power of Ethnic Nationalism. *Foreign Affairs*, 87(2), pp. 9-14.
- Nagourney, A. (2008, November 4). *Obama Wins Election*. From The New York Times: <http://www.nytimes.com/2008/11/05/us/politics/05campaign.html>
- Net Market Share. (2017). *Market Share Statistics for Internet Technologies*. From Net Market Share: <https://www.netmarketshare.com/search-engine-market-share.aspx?options=%7B%22filter%22%3A%7B%22%24and%22%3A%5B%7B%22deviceType%22%3A%7B%22%24in%22%3A%5B%22Desktop%22Flaptop%22%5D%7D%7D%5D%7D%2C%22dateLabel%22%3A%22Trend%22%2C%22attributes%22%3A%22share%22%2C%22group%22%3A%22searchEngine%22%2C%22sort%22%3A%7B%22share%22%3A-1%7D%2C%22id%22%3A%22searchEnginesDesktop%22%2C%22dateInterval%22%3A%22Monthly%22%2C%22dateStart%22%3A%222016-12%22%2C%22dateEnd%22%3A%222017-11%22%2C%22segments%22%3A%22-1000%22%7D>
- New World Encyclopedia. (2017). *New World Encyclopedia*. From Society - Definition: <http://www.newworldencyclopedia.org/entry/Society>
- Nickerson, R. S. (1998). Confirmation Bias, A Ubiquitous Phenomenon in Many guises. *Review of General Psychology*, 2(2), pp. 175-220.
- Olson, C. P., Laurikkala, M. K., Huff-Corzine, L., & Corzine, J. (2009, June 10). Immigration and Violent Crime: Citizenship Status and Social Disorganization. *SAGE*, 13(3), pp. 227-241.
- Pariser, E. (2011). Beware Online "Filter Bubbles". *Eli Pariser at TED2011*.
- Parks, G. S., Rachlinski, J. J., & Epstein, R. A. (2009). Implicit Race Bias and the 2008 Presidential Election: Much Ado About Nothing. *University of Pennsylvania Law Review*, 157, p. 210.
- Perrin, A., & Duggan, M. (2015). *American's internet access 2000-2015*. Pew Research Center.

- 
- Pew Research Center. (2016, July 7). *Pew Research Center, U.S. Politics and Policy*. From Top Voting Issues in the 2016 Election : <http://www.people-press.org/2016/07/07/4-top-voting-issues-in-2016-election/>
- Pew Research Center. (2016). *The Parties on the Eve of the 2016 Election: Two Coalitions Moving Further Apart*. Pew Research Center.
- Piggott, S. (2016, April 28). *Is Breitbart.com Becoming the Media Arm of the "Alt-right"?* From The Southern Poverty Law Center: <https://www.splcenter.org/hatewatch/2016/04/28/breitbartcom-becoming-media-arm-alt-right>
- Polgreen, P. M., Chen, Y., Pennock, D. M., Nelson, F. D., & Weinstein, R. A. (2008, December 1). Using Internet Searches for Influenza Surveillance. *Clinical Infectious Diseases*, 47(11), pp. 1443-1448.
- Porter, E. (2014, September 10). *A Simple Equation: More Education = More Income*. From The New York Times: <https://www.nytimes.com/2014/09/11/business/economy/a-simple-equation-more-education-more-income.html>
- Quantcast. (2017). *Drudge Report*. From Quantcast: <https://www.quantcast.com/measure/drudgereport.com#trafficCard>
- Quantcast. (2017). *Stormfront*. From Quantcast: [https://www.quantcast.com/stormfront.org?qcLocale=en\\_US](https://www.quantcast.com/stormfront.org?qcLocale=en_US)
- Raddant, M. (2016). *The Response of European Stock Markets To The Brexit*. Kiel Policy Brief, Vol. 100. Kiel: Kiel Institute for the World Economy.
- Real Clear Politics. (2008). *McCain vs. Obama*. From Real Clear Politics: [https://realclearpolitics.com/epolls/2008/president/us/general\\_election\\_mccain\\_vs\\_obama-225.html](https://realclearpolitics.com/epolls/2008/president/us/general_election_mccain_vs_obama-225.html)
- Real Clear Politics. (2012). *General Election: Romney vs. Obama*. From Real Clear Politics: [https://www.realclearpolitics.com/epolls/2012/president/us/general\\_election\\_romney\\_vs\\_obama-1171.html](https://www.realclearpolitics.com/epolls/2012/president/us/general_election_romney_vs_obama-1171.html)
- Real Clear Politics. (2016). *General Election: Trump vs. Clinton, Latest Polls*. From Real Clear Politics: [https://www.realclearpolitics.com/epolls/2016/president/us/general\\_election\\_trump\\_vs\\_clinton-5491.html%23polls](https://www.realclearpolitics.com/epolls/2016/president/us/general_election_trump_vs_clinton-5491.html%23polls)
- Reid, L. W., Weiss, H. E., Adelman, R. M., & Jaret, C. (2005, December). The Immigration-crime Relationship: Evidence across U.S. Metropolitan Areas. *Social Science Research*, 34(4), pp. 757-780.
- Rogers, T., & Aida, M. (2012). *Why Bother Asking? The Limited Value of Self-Reported Vote Intention*. Harvard Kennedy School. John F. Kennedy School of Government.
- Rosar, U., Klein, M., & Beckers, T. (2008, January). The Frog Pound Beauty Contest: physical Attractiveness and Electoral Success of the Constituency Candidates at the North Rhine Westphalia State Election of 2005. *European Journal of Political Research*, 47(1), pp. 64-79.
- Rydgren, J. (2005, May 17). Is extreme right-wing populism contagious? Explaining the emergence of a new party family. *European journal of political research*, 44(3), pp. 413-437.
- Rydgren, J. (2008, October). Immigration Skeptics, Xenophobes or Racist? Radical Right Wing Voting in Six West European Countries. *European Journal of Political Research*, 47(6), pp. 737-765.
- Sabato, L. J., Kondik, K., & Skelley, G. (2017). *Trumped: The 2016 Election That Broke All The Rules*. The Rowman and Littlefield.

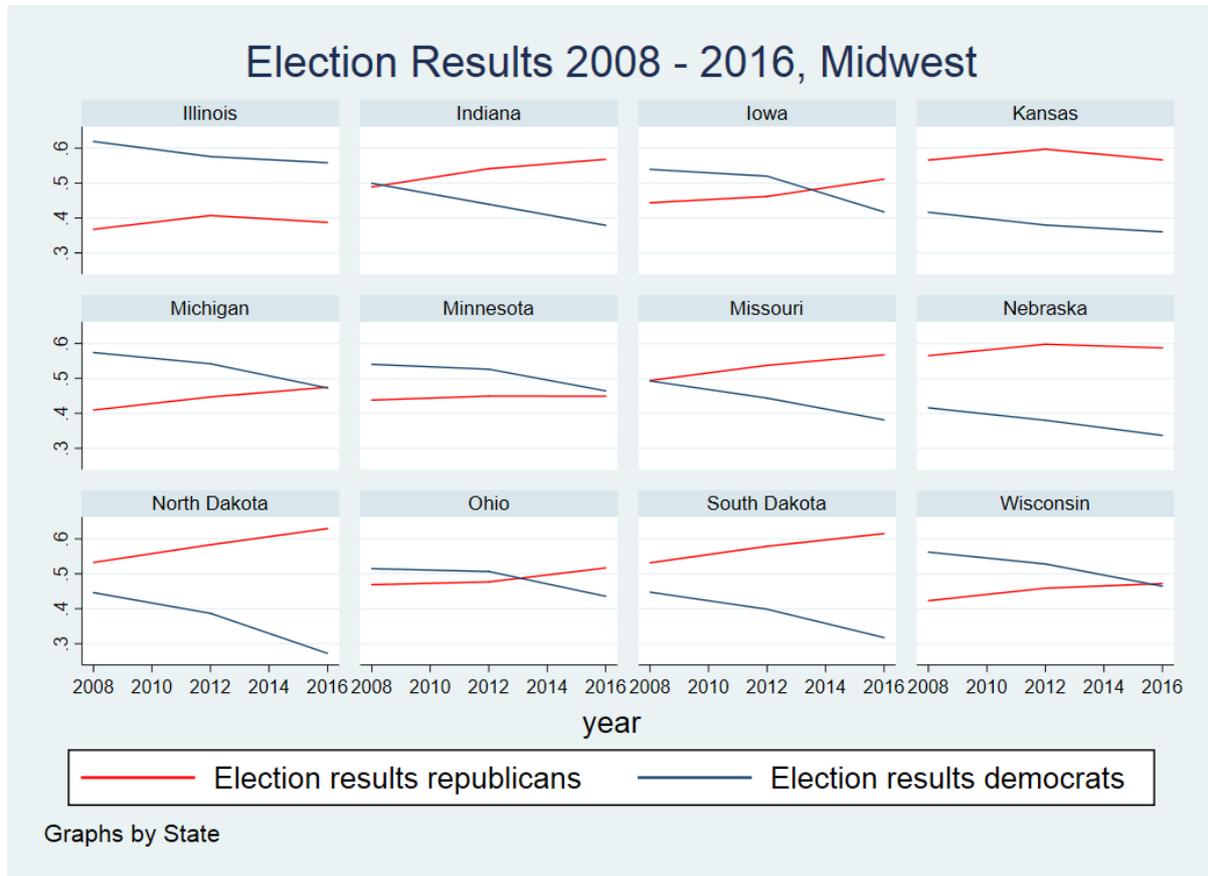
- SAGE Reference Publications. (2005). *Encyclopedia of Politics: The Left and The Right*. SAGE.
- Scott, P. (2017, March). *From the EU Referendum to Article 50: the charts that show how the UK's Economy has performed*. From The Telegraph: <http://www.telegraph.co.uk/business/2017/03/28/eu-referendum-article-50-charts-show-uks-economy-has-performed/>,
- Shimshoni, Y., Efron, N., & Matias, Y. (2009, August). *On the Predictability of Search Trends*. Google, Israel Labs.
- Silver, N. (2016, October 31). *The Odds Of An Electoral College-Popular Vote Split Are Increasing*. From FiveThirtyEight: <https://fivethirtyeight.com/features/the-odds-of-an-electoral-college-popular-vote-split-are-increasing/>
- State Election Office. (2017). *Official 2016 Presidential General Election Results*. State Election Office. From <https://transition.fec.gov/pubrec/fe2016/2016presgeresults.pdf>
- Stephens-Davidowitz, S. (2013, March 24). *The Cost of Racial Animus on a Black Presidential Candidate: Using Google Search Data to Find What Surveys Miss*.
- Stephens-Davidowitz, S. (2013). *Who Will Vote? Ask Google* .
- Stephens-Davidowitz, S. (2014, July 12). *The Data of Hate*. From The New York Times: <https://www.nytimes.com/2014/07/13/opinion/sunday/seth-stephens-davidowitz-the-data-of-hate.html>
- Stephens-Davidowitz, S. (2017). *Everybody Lies: Big Data, New Data, and What the Internet Can Tell Us About Who We Really Are*. Dey Street Books.
- Stephens-Davidowitz, S., & Soltas, E. (2015, December 12). *The Rise of Hate Search*. From The New York Times: [https://www.nytimes.com/2015/12/13/opinion/sunday/the-rise-of-hate-search.html?\\_r=0](https://www.nytimes.com/2015/12/13/opinion/sunday/the-rise-of-hate-search.html?_r=0)
- Streb, M. J., Burell, B., Frederick, B., & Genovese, M. A. (2008, January). *Social Desirability Effects and Support for a Female American President*. *Public Opinion Quarterly*, 72(1), pp. 76-89.
- Strumpf, K. S., & Philippe, J. R. (1999, March). *Estimating Presidential Elections: The Importance of State Fixed Effects and the Role of National versus Local Information*. *Economics & Politics*, 11(1), pp. 33-50.
- Sullivan, D. (2009, January 26). *Search Market Share 2008: Google Grew, Yahoo & Microsoft Dropped & Stabilized*. From Search Engine Land: <https://searchengineland.com/search-market-share-2008-google-grew-yahoo-microsoft-dropped-stabilized-16310>
- Swanson, A. (2016, November 7). *Rising illegal immigration in the US is one of the 2016 election's biggest myths*. From Independent: <http://www.independent.co.uk/news/world/americas/us-elections/us-illegal-immigration-donald-trump-us-presidential-election-2016-myth-a7403371.html>
- Swearingen, C., & Ripberger, J. (2014, January 20). *Does Search Traffic Provide a Valid Measure of Public Attention to Political Candidates*. *Social Science Quarterly*, 95(3), pp. 882-893.
- The Guardian. (2017, January 7). *German police quash Breitbart story of mob setting fire to Dortmund church*. From The Guardian: <https://www.theguardian.com/world/2017/jan/07/german-police-quash-breitbart-story-of-mob-setting-fire-to-dortmund-church>
- Torres-Reyna, O. (2007). *Panel Data Analysis, Fixed and Random Effects using Stata*. Princeton University.

- United States Census Bureau. (2017). *American Fact Finder, Selected Economic Characteristics, Table S0102*. From Fact Finder, United States Census Bureau: [https://factfinder.census.gov/faces/tableservices/jsf/pages/productview.xhtml?pid=ACS\\_15\\_5YR\\_DP03&src=pt](https://factfinder.census.gov/faces/tableservices/jsf/pages/productview.xhtml?pid=ACS_15_5YR_DP03&src=pt)
- United States Census Bureau. (2017). *Census Regions and Divisions of the United States*. From [https://www2.census.gov/geo/pdfs/maps-data/maps/reference/us\\_regdiv.pdf](https://www2.census.gov/geo/pdfs/maps-data/maps/reference/us_regdiv.pdf)
- usa.gov. (2017, August 7). *usa.gov*. From Presidential Election Process: <https://www.usa.gov/election#item-36072>
- Varian, H., & Stephens-Davidowitz, S. (2015). *A Hands-on Guide to Google Data*. Google Inc.
- Wolfers, J. (2002). Are Voters Rational? Evidence from Gubernatorial Elections. *Stanford GSB Working Paper, 1730*.
- Zukin, C. (2015, June 20). *The New York Times*. From What's The Matter With Polling: <https://www.nytimes.com/2015/06/21/opinion/sunday/whats-the-matter-with-polling.html>

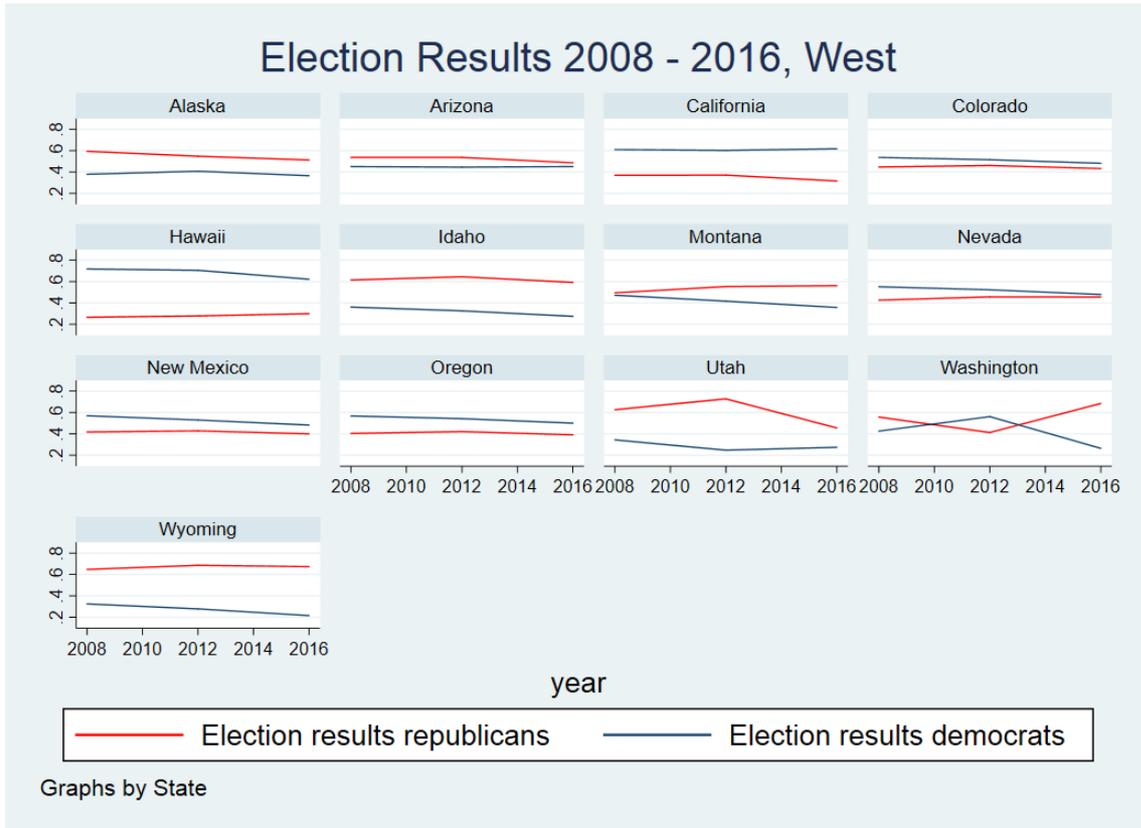


# Appendix

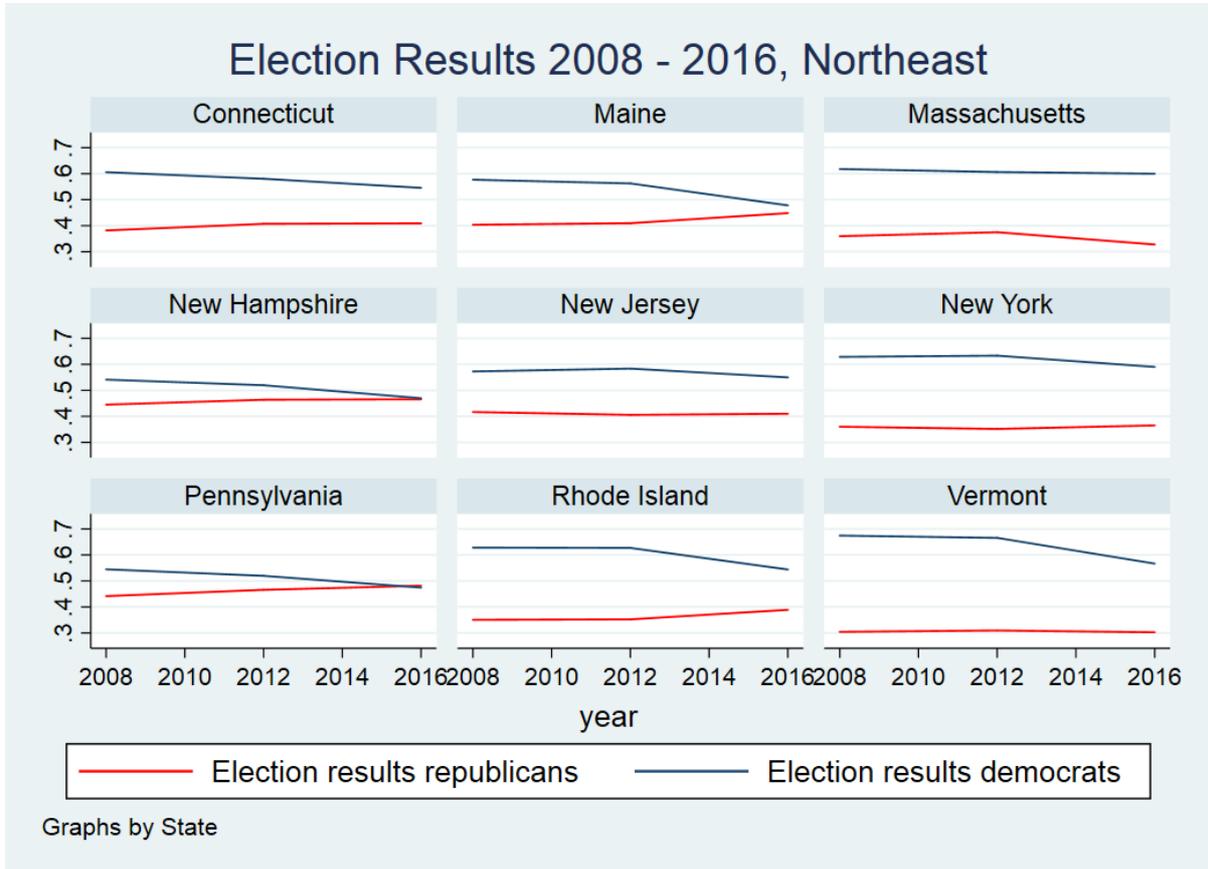
A1: Presidential Election results, 2008-2016, segmented on regions and on states.



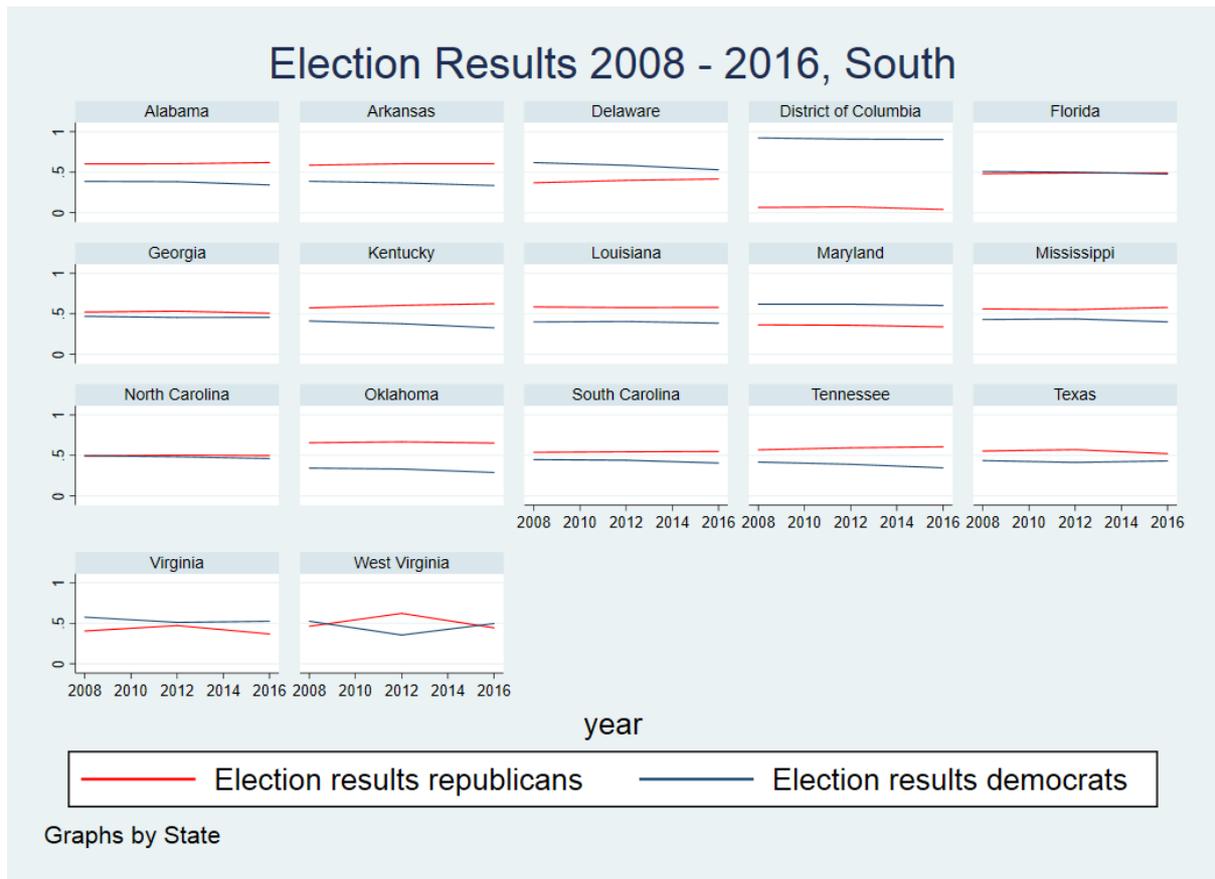
a) The Midwest



b) The West



c) Northeast



## d) The South

*Notes:* These figures show how the presidential election results in the different states in 2008, 2012 and 2016, clustered by region.

Table A1: Social sentiment proxies showing the mean value for each region in the three periods 2008, 2012 and 2016.

Table A1: Social sentiment proxies							
Region	Year	Long-term Racial Animus	Short-term Racial Animus	Long-term Immigration Skepticism	Short-term Immigration Skepticism	Long-term Far-right sentiment	Short-term Far-right sentiment
Midwest	2008	60	30	58	48	48	35
	2012	57	47	61	40	59	44
	2016	58	39	71	65	63	51
Northeast	2008	65	30	50	37	48	34
	2012	57	44	58	36	60	43
	2016	54	37	84	67	67	57
South	2008	74	36	55	39	58	35
	2012	71	63	59	44	64	48
	2016	71	48	72	71	67	54
West	2008	49	24	59	38	54	33
	2012	48	41	76	55	65	49
	2016	49	28	78	75	72	58

Table A1: Summary statistics segmented on regions

Table A2: Social sentiment proxies segmented on states, showing the maximum and minimum value, range, as well as the mean value for each state, over the three periods 2008, 2012 and 2016.

Table A2: State-wide Social Sentiment Proxies

State		Long-term Racial Animus	Short-term Racial Animus	Long-term immigration Skepticism	Short-term Immigration Skepticism	Long-term Far-Right Sentiment	Short-term Far-Right Sentiment
Alabama	Max	76	54	69	63	80	64
	Min	71	30	45	41	67	30
	Range	5	24	24	22	13	34
	Mean	74	45	61	49	75	50
Alaska	Max	68	62	49	57	96	71
	Min	59	18	38	23	83	58
	Range	9	44	11	34	13	14
	Mean	65	33	45	41	91	65
Arizona	Max	54	38	100	95	76	62
	Min	51	28	100	56	51	31
	Range	3	10	0	39	25	31
	Mean	52	32	100	76	67	51
Arkansas	Max	77	100	65	100	74	44
	Min	69	45	57	42	60	24
	Range	8	55	8	58	14	20
	Mean	72	64	61	64	67	34
California	Max	48	44	97	82	58	46
	Min	46	19	57	36	43	25
	Range	2	25	40	46	15	21
	Mean	47	31	79	56	51	35
Colorado	Max	48	43	85	99	65	54
	Min	39	25	65	41	47	32
	Range	9	18	20	58	19	22
	Mean	43	36	75	63	59	46
Connecticut	Max	66	41	100	73	64	52
	Min	55	18	53	25	46	43
	Range	11	23	47	48	17	9
	Mean	59	28	74	47	57	48
Delaware	Max	72	100	83	99	70	70
	Min	64	40	51	21	58	16
	Range	8	60	32	78	13	53
	Mean	68	76	64	51	65	49
District of Columbia	Max	68	44	86	57	70	72
	Min	55	28	67	38	55	34
	Range	13	16	19	19	15	39
	Mean	61	36	74	46	62	51

Florida	Max	68	46	78	69	65	53
	Min	59	27	44	30	51	37
	Range	9	19	34	39	14	16
	Mean	63	37	60	49	58	44
Georgia	Max	77	55	69	63	60	48
	Min	71	27	52	35	56	39
	Range	6	28	17	28	5	9
	Mean	74	40	60	47	58	45
Hawaii	Max	36	24	60	59	64	43
	Min	27	8	27	13	44	30
	Range	9	16	33	46	20	13
	Mean	31	17	44	33	52	35
Idaho	Max	48	39	95	100	86	62
	Min	36	27	68	42	68	53
	Range	12	12	27	58	18	9
	Mean	43	34	81	78	75	58
Illinois	Max	69	60	79	77	58	42
	Min	58	28	57	40	46	32
	Range	11	32	22	37	13	10
	Mean	63	41	65	57	51	39
Indiana	Max	72	60	85	89	66	53
	Min	59	36	67	42	47	23
	Range	13	24	18	47	18	30
	Mean	65	45	74	61	58	41
Iowa	Max	60	64	96	81	58	46
	Min	48	32	83	28	41	33
	Range	12	32	13	53	17	13
	Mean	53	44	87	61	52	40
Kansas	Max	54	33	79	66	70	59
	Min	41	28	61	34	53	40
	Range	13	5	18	32	17	19
	Mean	50	31	69	54	63	47
Kentucky	Max	89	83	77	54	64	52
	Min	78	41	49	35	58	34
	Range	11	42	28	19	6	17
	Mean	82	66	61	43	61	42
Louisiana	Max	100	64	61	75	67	62
	Min	79	33	42	37	58	41
	Range	21	31	19	38	9	22
	Mean	89	44	51	50	63	52
Maine	Max	54	50	73	78	82	69
	Min	48	36	32	31	57	49
	Range	6	14	41	47	25	21
	Mean	51	41	50	50	71	56
Maryland	Max	63	43	79	82	71	58
	Min	59	24	62	29	54	42
	Range	4	19	17	53	16	16

	Mean	61	36	72	60	62	48
Massachusetts	Max	52	53	80	75	63	52
	Min	48	20	47	24	47	29
	Range	4	33	33	51	15	23
	Mean	50	33	61	45	56	41
Michigan	Max	71	53	69	63	61	51
	Min	67	34	47	34	44	30
	Range	4	19	22	29	18	21
	Mean	69	44	55	45	52	38
Minnesota	Max	55	36	60	72	62	48
	Min	45	21	48	23	44	38
	Range	10	15	12	49	18	10
	Mean	50	28	53	46	53	42
Mississippi	Max	93	67	61	83	63	68
	Min	82	33	43	26	52	34
	Range	11	34	18	57	11	34
	Mean	88	46	50	46	58	52
Missouri	Max	72	62	70	65	76	63
	Min	67	27	51	33	48	39
	Range	5	35	19	32	28	25
	Mean	69	44	59	46	62	48
Montana	Max	61	75	74	78	91	83
	Min	41	28	53	26	72	29
	Range	20	47	21	52	19	54
	Mean	49	52	65	49	83	57
Nebraska	Max	56	42	84	57	69	57
	Min	49	34	67	19	48	40
	Range	7	8	17	38	22	17
	Mean	52	38	73	41	58	50
Nevada	Max	69	47	93	78	71	62
	Min	57	20	52	30	42	23
	Range	12	27	41	48	29	39
	Mean	62	36	77	50	60	44
New Hampshire	Max	56	55	85	56	75	58
	Min	46	16	57	43	55	44
	Range	10	39	28	13	20	14
	Mean	50	31	67	49	65	52
New Jersey	Max	71	49	92	81	66	61
	Min	60	36	52	22	46	31
	Range	11	13	40	59	20	29
	Mean	64	42	71	49	58	44
New Mexico	Max	50	33	88	84	73	70
	Min	43	9	59	15	64	37
	Range	7	24	29	69	9	33
	Mean	46	24	78	55	67	52
New York	Max	65	45	86	70	61	46

	Min	53	32	48	31	47	28
	Range	12	13	38	39	14	18
	Mean	59	37	66	46	52	35
North Carolina	Max	69	65	72	63	64	47
	Min	64	33	58	42	56	27
	Range	5	32	14	21	8	19
	Mean	66	49	65	50	60	39
North Dakota	Max	56	51	53	66	87	51
	Min	40	5	51	25	41	37
	Range	16	46	2	41	46	14
	Mean	49	31	52	44	65	44
Ohio	Max	79	59	62	67	65	60
	Min	70	37	49	30	51	35
	Range	9	22	13	37	14	25
	Mean	73	45	54	45	60	48
Oklahoma	Max	69	50	84	84	72	57
	Min	61	28	53	30	69	33
	Range	8	22	31	54	3	24
	Mean	66	37	67	55	70	48
Oregon	Max	47	31	69	61	70	53
	Min	45	24	55	49	39	32
	Range	2	7	14	12	31	21
	Mean	46	26	64	54	55	40
Pennsylvania	Max	86	60	78	77	70	60
	Min	70	33	45	29	48	34
	Range	16	27	33	48	22	26
	Mean	76	48	58	48	60	45
Rhode Island	Max	73	59	83	65	60	43
	Min	49	40	65	21	35	18
	Range	24	19	18	44	25	25
	Mean	58	49	74	38	51	32
South Carolina	Max	82	66	69	75	70	60
	Min	62	31	53	50	52	57
	Range	20	35	16	25	17	3
	Mean	70	46	61	65	62	59
South Dakota	Max	68	50	64	100	55	56
	Min	34	19	46	41	42	28
	Range	34	31	18	59	13	28
	Mean	49	33	57	65	50	42
Tennessee	Max	75	68	67	66	67	57
	Min	72	31	57	21	53	45
	Range	3	37	10	45	14	11
	Mean	74	51	62	43	62	52
Texas	Max	63	47	82	84	58	48
	Min	59	27	62	48	55	33
	Range	4	20	20	36	3	14
	Mean	60	37	72	61	56	41

---

Utah	Max	37	29	84	80	65	56
	Min	36	12	64	52	35	16
	Range	1	17	20	28	31	40
	Mean	36	22	77	67	53	39
Vermont	Max	74	36	76	62	60	74
	Min	52	16	42	34	50	24
	Range	22	20	34	28	10	50
	Mean	60	23	55	47	56	46
Virginia	Max	100	97	79	59	77	57
	Min	47	17	40	26	45	23
	Range	53	80	39	33	32	34
	Mean	68	53	57	48	64	44
Washington	Max	56	55	72	65	76	66
	Min	49	20	52	39	47	34
	Range	7	35	20	26	28	32
	Mean	52	35	64	49	59	48
West Virginia	Max	100	100	63	56	74	38
	Min	51	33	39	38	63	21
	Range	49	67	24	18	11	17
	Mean	83	66	52	45	68	30
Wisconsin	Max	61	44	78	64	58	48
	Min	52	27	54	34	51	32
	Range	9	17	24	30	7	16
	Mean	56	33	64	48	55	43
Wyoming	Max	65	30	76	63	60	51
	Min	51	13	61	40	46	17
	Range	14	17	15	23	14	34
	Mean	58	22	70	55	53	37

---

Table A2: Summary statistics segmented on states