

# From Hashtag to Hate Crime: Twitter and Anti-Minority Sentiment \*

Karsten Müller<sup>†</sup> and Carlo Schwarz<sup>‡</sup>

July 24, 2020

(First version: March 2018)

## Abstract

We study whether social media can contribute to hatred against minorities with a focus on Donald Trump’s political rise. To establish causality, we construct an instrument for Twitter usage based on the platform’s early adopters at the South by Southwest (SXSW) festival in 2007, who were crucial for Twitter’s diffusion across US counties. Instrumenting with the home counties of SXSW followers who joined in March 2007, while controlling for the counties of SXSW followers who joined before the festival, we find that a one standard deviation increase in Twitter usage is associated with a 32% larger increase in anti-Muslim hate crimes since the 2016 presidential primaries. Further, Trump’s tweets about Islam-related topics predict increases in xenophobic tweets by his followers, cable news attention paid to Muslims, and hate crimes on the following days. These correlations persist in an instrumental variable framework exploiting that Trump is more likely to tweet about Muslims on days he plays golf.

---

\*A previous version of this paper was circulated under the title “Making America Hate Again? Twitter and Hate Crime Under Trump.” We are grateful to Roland Bénabou, Dan Bernhardt, Leonardo Bursztyn, Rafael di Tella, Mirko Draca, Ruben Durante, Didi Egerton-Warburton, Ruben Enikolopov, James Fenske, Thomas Fujiwara, Scott Gehlbach, Matthew Gentzkow, Andy Guess, Atif Mian, Jonathan Nagler, Maria Petrova, Joshua Tucker, Alessandra Voena, Hans-Joachim Voth, Fabian Waldinger, David Yanagizawa-Drott, Ekaterina Zhuravskaya, and seminar participants at the Princeton University, University of Warwick, New York University, RES Conference 2019, Young Economist Symposium 2019, EMCON Conference 2019, Barcelona Political Economy Workshop 2019, Galatina Summer Meetings 2019, EEA Conference 2019, CEP Policing & Crime Workshop, ENS de Lyon, Toulouse School of Economics, Bocconi University, Nottingham University, Queen Mary University, IE Business School, University of Mannheim, University of Munich, University of Illinois, Cornell University, and University of St. Andrews for their helpful suggestions. Christian Kontz provided excellent research assistance. Schwarz was supported by a Doctoral Scholarship from the Leverhulme Trust as part of the *Bridges* program.

<sup>†</sup>Princeton University, Julis-Rabinowitz Center for Public Policy and Finance, karstenm@princeton.edu

<sup>‡</sup>Department of Economics, Bocconi University, www.carloschwarz.eu, carlo.schwarz@unibocconi.it

# 1 Introduction

Social media platforms have been widely accused of enabling hatred of minorities (e.g. New York Times, 2019a; United Nations, 2020). An influential line of argument posits that social media may be particularly effective in reinforcing extreme beliefs in what has often been described as “echo chambers” (e.g. Sunstein, 2002, 2017). Despite much public debate, there is limited empirical evidence that social media can spur anti-minority sentiments.

We investigate this question in the context of a particularly notable case study: the political rise of Donald Trump. With more than 80 million followers, Trump is one of Twitter’s most prominent users and has used Twitter almost every day since joining in March 2009, racking up a total number of more than 32,000 tweets until the end of 2017. Trump’s rhetoric on Twitter has been widely criticized as inflammatory and is frequently cited as an example of how social media might increase anti-minority sentiments (New York Times, 2017).<sup>1</sup> Both Twitter and Facebook have recently flagged or deleted posts by Trump or his campaign that were considered hateful (e.g. Wall Street Journal, 2020; Financial Times, 2020). Such steps have further fueled discussions about how platform providers and governments should moderate content on social media (e.g. CNN, 2020).

We start by documenting that the frequency of anti-Muslim hate crimes has doubled since the 2016 presidential primaries. We investigate the potential role of social media in enabling such hate crimes using a difference-in-differences approach and find that their increase predominantly originates from counties with high Twitter usage. These regressions, however, may not isolate a pure “social media effect” because counties with many Twitter users likely also differ in many unobservable dimensions. This may bias our estimates upwards or downwards, depending on how individuals select into social media usage. For example, areas where many people use relatively new technologies such as Twitter may be more liberal (Pew Research Center, 2019b, 2020), which could bias our estimates downwards. On the other hand, such areas may have larger minority communities and thus more potential targets for perpetrators of hate crimes, leading to an upward bias.

To overcome these concerns, we develop an instrument for county-level Twitter usage in the United States based on the home towns of the platform’s early adopters at the South by Southwest (SXSW) festival in March 2007.<sup>2</sup> SXSW is widely regarded as the tipping point

---

<sup>1</sup>Minnesota congresswoman Ilhan Omar, for example, has linked tweets by Trump targeting her Muslim faith to “an increase in direct threats on my life—many directly referring or replying to the president’s video” (BBC, 2019).

<sup>2</sup>SXSW is an annual event, held since 1987, that comprises a number of festivals, conferences, trade shows, and exhibitions. In 2019, more than 230,000 people attended the festivals, where almost 2,000 acts from all over the world performed. More than 70,000 people attended the SXSW conference, which featured almost 4,800 speakers. Around 30,000 people attended SXSW Interactive, which focuses on emerging technology. For simplicity, we refer to the event as “SXSW festival” or similar short forms throughout the paper.

for Twitter's popularity and an important early catalyst for the site's diffusion. The number of daily tweets *tripled* during the festival and increased by a factor of 60 between 2007 and 2008 (Twitter, 2010). We show that 60% of early Twitter adopters were connected to SXSW, and the platform's growth accelerated disproportionately in counties with SXSW followers who joined Twitter during the 2007 festival.

Building on the literature on path dependence in technology adoption (e.g. Arthur, 1989, 1994; Liebowitz & Margolis, 1999; Arrow, 2000), our identification strategy exploits that the locations of Twitter's early adopters at SXSW are a strong predictor of county-level Twitter usage today. Using data on the profiles of more than four million Twitter users, we document an S-shaped adoption impact of the SXSW festival over time, consistent with theories of innovation diffusion (Griliches, 1957; Rogers, 2010; Bass, 1969; Geroski, 2000; Fagerberg et al., 2009). Similar to the empirical strategy in Enikolopov et al. (2016), we control for the locations of SXSW followers that signed up before the festival to mitigate concerns that our findings could be driven by selection into attending SXSW. The identifying assumption underlying our approach is that differences in the locations of SXSW followers who joined Twitter in March 2007 relative to earlier months are not related to unobserved county characteristics that explain the rise in expressed hatred of minorities with the 2016 presidential campaign. In support of this assumption, we show that hate crimes did not increase in counties with SXSW followers who signed up before the festival. This holds true even though the Twitter profiles and counties of residence of these pre-period users are indistinguishable from those of SXSW followers who signed up in March 2007.

Instrumenting for Twitter usage with SXSW followers who joined in March 2007, we show that higher exposure to social media increased anti-Muslim hate crimes around the time of Donald Trump's political rise. We find that a one standard deviation higher exposure to Twitter is associated with a 32% larger increase in hate crimes with the 2016 presidential campaign period. These findings suggest that social media platforms have played a role in the recent spread of expressed xenophobic hatred in the United States. We also find a similar but slightly weaker pattern for hate crimes targeting Hispanics, the second minority group often targeted by Trump. Using data from the National Crime Victimization Survey, we find no evidence for changes in the propensity of victims to report hate crimes they experienced, indicating that these patterns are indeed driven by an increased incidence of such crimes.

To understand the potential channels between social media and hate crimes, we analyze Trump's Twitter feed in the second part of the paper. In particular, we test whether incendiary tweets by Trump may have contributed to anti-Muslim sentiment, building on existing evidence that celebrities can have a disproportionate effect on public opinion (e.g. Beaman et al., 2009; Bursztyn et al., 2017; Alatas et al., 2019). We find a strong time series correlation between

Trump's tweets on Islam-related topics and the number of anti-Muslim hate crimes after the start of his presidential campaign, even after controlling for general attention paid to topics associated with Muslims. We find no such link for the period before his campaign.

To establish causality, we leverage Trump's well-documented golf habit. In 2017 alone, Trump played golf on more than 90 days, and many commentators have argued that golfing shifts Trump's state of mind. In the data, we find a clear pattern: Trump's golf days coincide strongly with changes in the content, but not the number of his tweets. In particular, Trump is more likely to send messages aimed at Muslims and the media on his golf days, and fewer about policy, a fact we exploit in an instrumental variable framework. One intuitive explanation of this finding is that day-to-day politics may be less salient to the President when outside of Washington, DC. There is also anecdotal evidence that Trump may be influenced by his social media director Dan Scavino—former manager of Trump National Golf Club Westchester and Trump's former caddie—who is the likely source of many particularly inflammatory tweets (New York Times, 2018; Reilly, 2019; CNN, 2020).

Using golf days as an instrument, we find evidence consistent with the idea that Trump's tweets about Muslims trigger waves of anti-Muslim sentiment. In particular, we find that his instrumented tweets not only predict the frequency of hate crimes, but also measures of media attention paid to Muslim-related topics. Using transcript data on the reporting of the major cable news networks Fox News, CNN, and MSNBC, we show that Trump's golf-induced tweets are associated with more mentions of Muslims. This link seems to be particularly pronounced for Fox News, which tends to support rather than oppose Trump's rhetoric. Based on a sample of more than 100 million tweets, we also find that Trump's anti-Muslim tweets are widely shared by his followers, who produce further xenophobic content in response, such as messages containing the hashtags “#StopIslam” and “#BanIslam”. Additionally, we investigate whether the transmission of Donald Trump's tweets is stronger in counties with more Twitter users in a panel regression setting. Interacting county-level Twitter usage and Trump's Twitter activity, we document that the spike in anti-Muslim hate crime in the days after Donald Trump's tweets is driven by counties with higher Twitter penetration.

Taken together, our evidence is consistent with the interpretation that social media can play a role in rising anti-minority sentiment. We investigate three possible mechanisms that could broadly explain our findings: changes in the costs of coordinating hate crimes between potential perpetrators (a *coordination channel*), changes in people's beliefs (a *persuasion channel*), and changes in the perceived societal acceptance of xenophobic beliefs or actions (a *social norms channel*). A range of additional results is most consistent with the idea that social media might enable anti-minority sentiments by changing social norms.

Three pieces of evidence support this interpretation. First, contrary to what a coordination channel would suggest, the effects we find are largely accounted for by hate crimes committed by individual perpetrators, rather than groups. Second, we find no evidence that social media has persuaded people to become more xenophobic. Indeed, data from nationally representative surveys and implicit association tests (IAT) suggest that Americans' views about Muslims and immigrants have become more positive since Trump's political rise, particularly among social media users. We also show that social media has a precisely estimated zero effect on people's implicit bias against Muslims. Third, in contrast to the predictions of Bayesian persuasion models (e.g. Kamenica & Gentzkow, 2011), we find that our effects are largely driven by areas that are more likely to harbor pre-existing hatred of minorities. We argue that the sum of these findings is most easily explained by social media changing people's perception of which actions toward minorities are socially acceptable.

Our paper contributes to the literature on the relationship between media consumption and violence. Yanagizawa-Drott (2014), Adena et al. (2015), and DellaVigna et al. (2014) find that traditional media can contribute to ethnic hatred and violence. Other research has linked media such as television (Card & Dahl, 2011) and movies (Dahl & DellaVigna, 2009) to short-lived spikes (or decreases) in violence. Bhuller et al. (2013) document increases in sex crime associated with the roll-out of broadband internet in Norway; Chan et al. (2016) find a correlation between broadband availability and hate crimes in the US. Our findings speak to the role of social media in the spread of violence against minority groups.

We most directly contribute to a growing literature on the influence of social media on real life outcomes. Enikolopov et al. (2016) show that social media can increase participation in protests in Russia by reducing coordination costs. Petrova et al. (2017) study whether adopting Twitter helps politicians attract donations.<sup>3</sup> In previous work, we found evidence that social media affects the propagation of anti-refugee incidents in Germany, using Facebook and internet disruptions as a source of short-lived exogenous variation (Müller & Schwarz, 2018). Bursztyn et al. (2019) study social media and xenophobia in Russia. In contrast to this existing work, our paper suggests that social media may shift societal norms in the medium-term, based on the particularly salient case study of Donald Trump's presidency.

A separate related literature studies political polarization. While there is evidence that polarization has increased over the past decades (Fiorina & Abrams, 2008; Gentzkow, 2016; Draca & Schwarz, 2018), existing studies have found no or even a negative correlation with

---

<sup>3</sup>A growing body of experimental studies on the effects of social media include Bond et al. (2012), Jones et al. (2017), Bail et al. (2018), Mosquera et al. (2020), Chen & Yang (2019), Levy (2019), and Allcott et al. (2020).

social media use (Boxell et al., 2017; Barberá, 2014).<sup>4</sup> Our findings suggest that social media may enable those with extreme viewpoints to find sources of social legitimacy. A widely shared discriminatory tweet by the President, for example, could signal to potential perpetrators of hate crimes that their actions are more widely accepted than they really are.

The paper proceeds as follows. In Section 2, we introduce the data sources and present descriptive evidence on hate crimes since 2010. In Section 3, we discuss our empirical strategy and introduce our instrument for Twitter usage based on the SXSW festival. Section 4 presents the main empirical results. In Section 5 we discuss evidence for the link between Trump’s tweets and anti-Muslim sentiment. Section 6 discusses plausible mechanisms behind our results. Section 7 concludes.

## 2 Data and Background

We create two datasets for our analysis. First, we build a county-level dataset containing information on hate crimes, Twitter usage, and numerous other variables. Second, we construct a daily time series dataset that combines Donald Trump’s Twitter activity, the number of total hate crime incidents in the US, data on TV news coverage, and time series control variables. The key sources we draw on are (1) hate crime data reported by the FBI’s Uniform Crime Reporting (UCR) program; (2) a county-level measure of Twitter usage based on almost 500 million tweets collected by Kinder-Kurlanda et al. (2017); (3) hand-collected county-level data on the locations of early adopters of Twitter in 2006 and 2007; (4) data on the universe of Donald Trump’s tweets; and (5) information on Trump’s golf activity from his inauguration in early 2017 until the end of that year. We describe these and all other data sources in more detail in the following subsections. Table A.8 and Table A.22 in the online appendix present the full descriptive statistics.

### 2.1 FBI Hate Crime Data

We use data on hate crime in the US from the FBI for the years 1990 to 2017. The data set contains all hate crimes reported to the FBI as part of the Uniform Crime Reporting (UCR) program. The FBI defines hate crimes as:

“[...] criminal offenses that are motivated, in whole or in part, by an offenders bias against a race, religion, disability, sexual orientation, ethnicity, gender, or gender identity.” (FBI, 2015, p. 4)

---

<sup>4</sup>A separate literature has analyzed the effects of the media on elections and other political outcomes. See, among others, the work by Adena et al. (2015), DellaVigna et al. (2014), Stephens-Davidowitz (2014), Gavazza et al. (2015), Gentzkow (2016), Manacorda & Tesei (2020), and Martin & Yurukoglu (2017).

To classify hate crimes, the FBI uses a two-tier decision making process. First, the law enforcement officer recording an incident decides whether it might constitute a hate crime. Second, potential hate crime cases are evaluated by officers with special training in hate crime matters. The FBI (2015) states (p. 35): “For an incident to be reported as a hate crime, sufficient objective facts must be present to lead a reasonable and prudent person to conclude that the offenders actions were motivated, in whole or in part, by bias.” For more information on the FBI classification procedure, see appendix A.1.

Because considerable evidence needs to be available for an offense to be classified as a hate crime, the numbers reported by the FBI have been criticized as dramatic underestimates (e.g. ProPublica, 2017; NBC News, 2017).<sup>5</sup> Nonetheless, the FBI data constitute the most complete record of hate crimes committed in the United States for which incident details are available. Among others, they include information on the exact date of the crime, the type of crime (e.g. vandalism, theft, assault), the number of victims, and the number of perpetrators. The median and mode incident has a single perpetrator. We map these data to counties using the location of the more than 24,000 original reporting agencies based on their Originating Agency Identifier (ORI). Figure 1a plots the geographic distribution of hate crimes across the mainland USA.<sup>6</sup> The counties in grey never report any hate crime to the FBI.

The FBI differentiates hate crimes by motivating bias (e.g. anti-Muslim). Overall, they report 34 bias motivations for the broad categories race, religion, sexual orientation, disability, and gender/gender identity. We report all codes for the motivating bias in Table A.2. We use this classification to identify hate crimes against Muslims. The other categories used in the paper are defined according to the codes listed in Table A.1.

## 2.2 Measuring County-Level Twitter Usage

Twitter does not publish statistics on the number of active users per US county. We create an approximate measure of local Twitter usage using 475 million geo-located tweets collected by Kinder-Kurlanda et al. (2017) made available through the Gesis Datorium. The data were collected between June and November in 2014 and 2015 by repeatedly calling the Twitter

---

<sup>5</sup>Note that time-invariant reporting bias across counties is unlikely to drive our results. We accommodate potential geographical reporting differences in our cross-sectional tests by estimating our model in first-differences or including county fixed effects. In further robustness checks we restrict the sample to counties where at least one hate crime is reported. We discuss why changes in reporting over time are unlikely to explain our results below.

<sup>6</sup>The FBI hate crime data do not contain information on the US territories of Virgin Island, Puerto Rico, Northern Mariana Islands, American Samoa, and Guam.

streaming API, restricted to US tweets.<sup>7</sup> These tweets were then assigned to counties based on the geo-location of each tweet.

To create a measure based on the users in each county, we scraped the underlying user profiles for each tweet and assigned users to the county from which they tweet most frequently. Overall, our data contain over four million users, around 7% of the US Twitter population in 2015. Figure 1b visualizes the distribution of Twitter users per capita across the continental United States.<sup>8</sup> The user profiles also provide us with information about names, join dates, and profile descriptions (“bios”).

For robustness, we also consider three other measures of Twitter use. The first is simply the number of tweets sent from each country. The other two are based on the Survey of the American Consumer, conducted by GfK Mediamark Research & Intelligence, and capture the number or share of households who used Twitter in the past 30 days (in 2015).

### 2.3 Twitter Data for South by Southwest and Other Outcomes

We collect additional Twitter data using the platform’s application programming interface (API). In particular, we collect the universe of users following the Twitter account of SXSW Conference & Festivals (SXSW). This yields 658,240 unique user IDs. For each of these users, we collect information on their location and the date their account was created. In line with the findings of Takhteyev et al. (2012), around 75% of Twitter users in the sample report their geographical location. Previous research suggests that these user locations yield valid proxies for Twitter usage (e.g. Takhteyev et al., 2012; Haustein & Costas, 2014). To compare Twitter activity around the 2007 SXSW festival to other festivals in the same year, we additionally collect the tweets and user data for the Austin City Limited Festival, Burning Man, Coachella, Electric Daisy Festival, New Orleans Jazz and Heritage Festival, Lollapalooza, Pitchfork Music Festival, and the West by Southwest Festival. The full list of search terms for these festivals can be found in Table A.5.

Since we are also interested in the impact of the SXSW festival on overall Twitter activity, we create a proxy for the total number of tweets using the 100 most common English words for January through March 2007 (the full list of words is reported in Table A.6). The tweets are collected by calling the Twitter search for each word one day at a time. While this approach does not give us the universe of tweets in this time window, it should serve as a valid proxy for how many people are using Twitter in a given county over time.

---

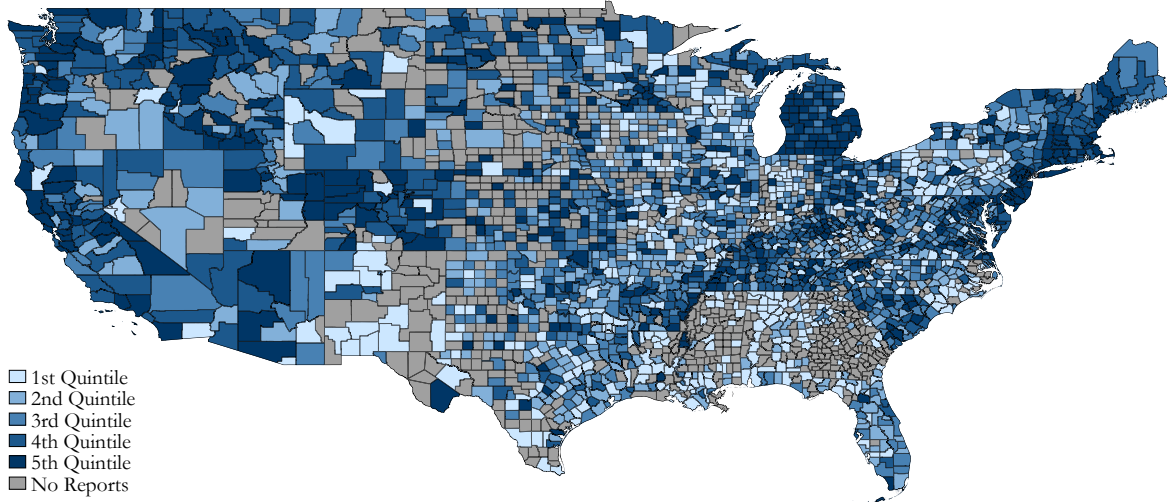
<sup>7</sup>The streaming API provides a 1% sub-sample of public tweets each time it is called. While the exact underlying sampling procedure is unknown, this process should result in a good approximation of overall Twitter activity.

<sup>8</sup>The data do not contain information for Alaska and Hawaii; we thus focus on the continental US.

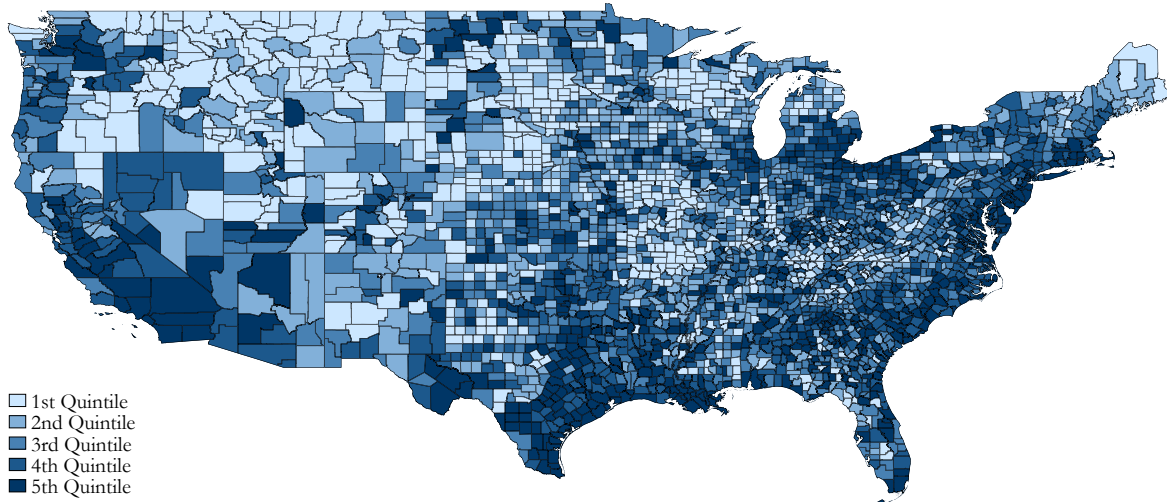


Figure 1: Hate Crimes and Twitter Usage by US County

(a) Hate Crimes per Capita



(b) Twitter Users per Capita



*Notes:* These maps plot the geographical distribution of hate crimes and Twitter usage across the counties of the mainland United States. Panel (a) plots quintiles of the total number of hate crimes per capita between 1990 and 2017 as reported by the FBI. Counties in grey never reported any hate crime. Panel (b) plots the number of Twitter users per capita based on the tweets collected by Kinder-Kurlanda et al. (2017).

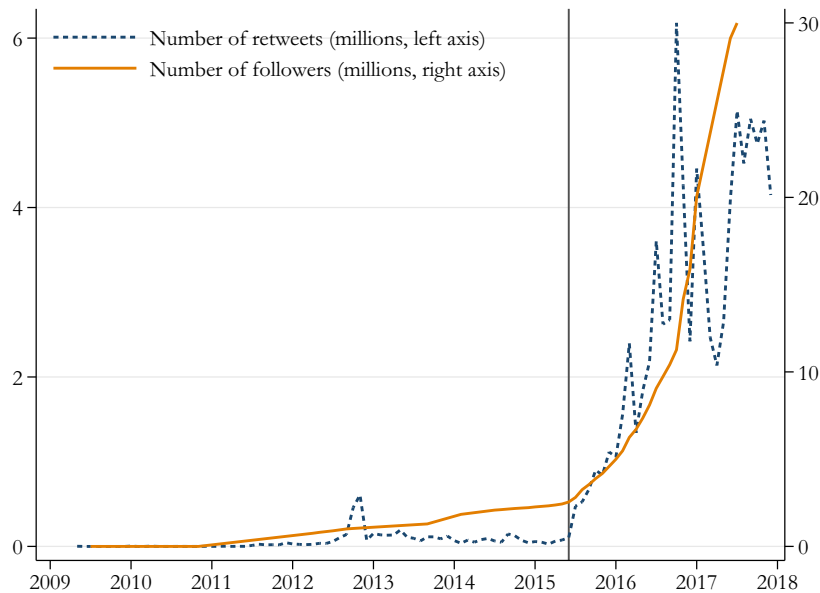
We create proxies for anti-Muslim Twitter content by collecting tweets containing the hashtags “#BanIslam” or “#StopIslam” from 2010 to 2017. We selected these hashtags because they are both clearly anti-Muslim and commonly used on Twitter (Miller & Smith, 2017). Following the same procedure as for the SXSU tweets, we assign these tweets to counties based on the location of the users.

## 2.4 Measuring Trump’s Twitter Activity

To understand Trump’s Twitter activity, we collect the universe of his tweets from the Trump Twitter Archive (Brown, 2018). Our version of this data set contains 32,794 tweets for the time period of April 2009 to December 2017. The data contain the date, time, and text of each tweet and the number of retweets a tweet received.

We analyze Trump’s Twitter reach and shows that he has the potential to influence a considerable fraction of Americans. Figure 2 plots the monthly number of retweets he received since joining Twitter. The number of retweets increased distinctly with his presidential run (indicated by the vertical line). The same also holds true for the number of his followers.

**Figure 2: Trump’s Twitter Reach Over Time**



*Notes:* The figure plots the number of monthly retweets (in millions) Donald Trump’s Twitter account received since he joined the platform in 2009. The grey vertical line marks the start of Trump’s presidential campaign in June 2015.

**Identifying Trump’s negative tweets about Muslims.** We classify tweets in a four-step process. First, we use the text of Trump’s tweets to identify messages about Muslims or

Islam-related topics and hand-code negative tweets about Muslims from a random subsample of 5,000 tweets. Second, we use this subsample as the training sample for a machine learning classifier that we apply to the entire body of tweets.<sup>9</sup> We train a classifier based on a logistic regression model with L1 regularization. We decide the optimal regularization strength using 5-fold cross-validation. The final model achieves an out-of-sample F1 score of 0.98. In the total sample of Trump’s tweets the classifier tags 265 anti-Muslim tweets.

Third, we also add any tweets containing the words “muslim”, “islam”, “terror”, “mosque”, “refugee”, and ‘sharia”, because we use these terms to identify Google searches and news reports on Muslims. This process tags an additional 58 potential tweets about Muslims. Finally, to rule out that we are picking up unrelated topics by mistake, we hand-check all selected tweets. We list examples of negative tweets about Muslims in Table A.4 in the online appendix; Table A.3 shows 24 tweets we manually removed in the hand-coding step.<sup>10</sup>

To further understand the topics of Trump’s tweets during his presidency, we code Trump’s tweets in 2017 into the following categories: media, islam and terrorism, party politics, immigration, foreign policy, domestic policy, and other topics. We also code the sentiment of each tweet. More specifically, we code the sentiment of each tweet either as “very negative”, “negative”, “neutral”, “positive”, or “very positive”. We recode these categories into a scale from -2 (very negative) to 2 (very positive).

To provide direct evidence for spillovers of Trump’s negative tweets about Muslims on his followers, we collect a random sample of tweets by 630,000 of his followers. This yields a dataset with over 115 million tweets.

## 2.5 Information on Trump’s Golf Trips

Information on Trump’s golf outings is from the New York Times (New York Times, 2019b). These data cover Trump’s travels and identify sources indicating that he was in fact golfing on any given trip. We cross-check the information from the New York Times using information from *trumpgolfcount.com* and the official Presidential schedule from the White House, and add a few additional days of golf in the process. Table A.7 in the online appendix describes these sources in more detail; Figure A.7a graphs the days in 2017 Trump spent golfing, where the darker shade of orange indicates golf outings longer than three days. More than two

---

<sup>9</sup>We remove stopwords from and reduce all words to their morphological roots, so called lemmas. We then extract all unigram, bigrams, and trigrams that appear in at least three tweets. The extracted n-grams are reweighted using term frequency-inverse document frequency (tf-idf). In this step, the frequency of a n-gram  $v$  in document  $d$  is replaced by  $tfidf(f_{d,v}) = (1 + \ln(f_{d,v}) \cdot (\ln(\frac{1+D}{1+d_v}) + 1))$ , where  $d_v$  is the number of documents n-gram  $v$  appears in.

<sup>10</sup>Our results are not driven by excluding these tweets.

thirds of golf days are on the weekend. However, Table A.23 shows that Trump has golfed multiple times on all days of the week.

## 2.6 Additional Data Sources

We construct a large number of additional variables, which mostly serve as controls. A more detailed variable description and the relevant data sources can be found in Table A.9.

**County-level variables.** We collect demographic control variables at the county level from the United States Census and the American Community Survey. In particular, we use information on yearly population, the share of the population by age group, the ethnic composition of the population, the poverty rate, and education levels. Information on a county’s unemployment rate and industry level employment shares were obtained from the Bureau of Labor Statistics. County-level election results are available from the MIT Election lab. The number of Muslims in each US county is derived from the 2010 US Religious Census. Additionally, we make use of county-level crime statistics based on the FBI’s UCR data. Information on TV viewership patterns was collected from Simply Analytics.

Lastly, we proxy for potential preexisting xenophobic sentiments at the county level using data on hate groups from the Southern Poverty Law Center (SPLC). We assign hate groups to counties based on the reported state and city information. While the classification of hate groups is subjective and subject to controversy, the information gathered by the SPLC is widely used as a proxy for where hate groups are located.<sup>11</sup>

**Time series variables.** To study the content of cable news, we collect news mentions of Muslims from the TV News Archive (part of the Internet Archive). We scrape news mentions for Fox News, CNN, and MSNBC based on the search terms “sharia”, “refugee”, “mosque”, “muslim”, and “islam”, consistent with those used to classify Trump’s tweets. We collect a total of 75,193 mentions from the start of Trump’s presidential campaign to the end of 2017.

We are also interested in the overall salience of Islam-related topics on the internet. We use Google Trends to obtain daily trends for the above search terms for the US. Unfortunately, Google trends only allows us to collect the daily search interest for a 90 day period. We therefore separately collect the Google trends in 90 day intervals for the period since Trump started his presidential campaign. Since Google normalizes the search interest between 0-100 for each 90 day period, we use weekly search interest—which is available for the period as a

---

<sup>11</sup>Note that, as long as the geography of potential misclassification of hate groups by SPLC is random, this will bias our estimates towards zero.

whole—to bring the daily search results to the same scale. We describe this process in more detail in Appendix A.1.4.

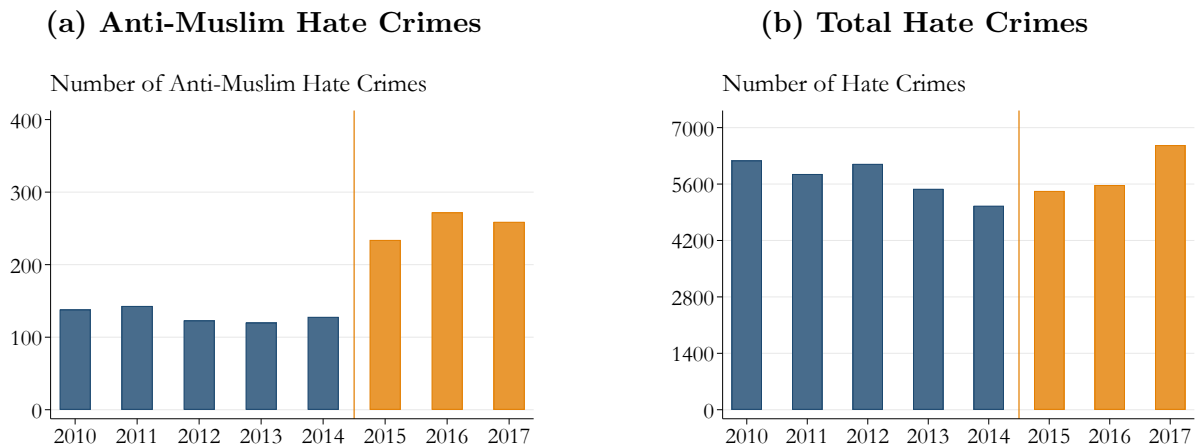
Lastly, we compile information on the daily number of Islamist terror attacks from the Global Terrorism Database, where we focus on terror attacks in the US and Europe. For the years 2015-2017 our data contain 37 terror attacks.

### 3 Social Media and Anti-Muslim Sentiment

#### 3.1 Introductory Findings

To motivate our analysis, we begin by investigating how the number of hate crimes has evolved over time. Panel A of Figure 3 plots the total number anti-Muslim hate crimes for each year from 2010 to 2017. The data suggest that anti-Muslim hate crimes have become considerably more common since 2015, which coincides with the 2016 presidential primaries. In fact, the average number of hate crimes has approximately doubled in the 2015-2017 period compared to 2010-2014.

**Figure 3: Trends in Hate Crimes Since 2010**



*Notes:* This figure plots the number of yearly hate crimes in the United States based on the Uniform Crime Reporting (UCR) data of the Federal Bureau of Investigation (FBI). Panel (a) shows the number of anti-Muslim hate crimes. Panel (b) shows the total number of hate crimes. The years that include Donald Trump’s presidential campaign start and election win are marked orange.

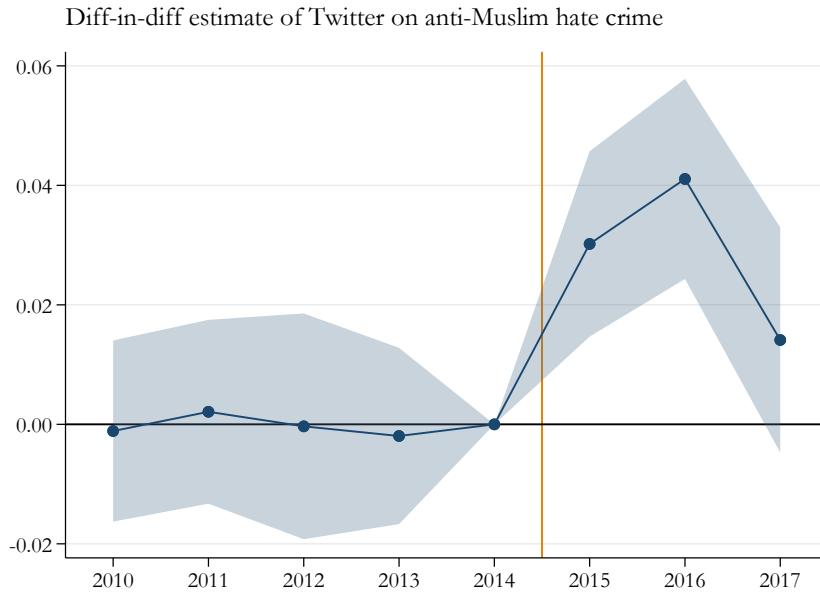
We also plot the number of total number of hate crimes, for which we do not observe a similar increase, in Panel B of of Figure 3. Similarly, we do not find comparable increases in hate crimes when we split them into the other underlying bias categories in Figure A.3. We conclude that the period of the presidential primaries in 2015-2016 coincided with a clear rise in measured anti-Muslim sentiment in the United States.

Could Twitter play a role in this spread of xenophobic sentiments starting around the time of the 2016 presidential campaign? If that were the case, we would expect the increase in hate crimes documented in Figure 3 to be concentrated in areas where many people use Twitter. To get a first pass at this question, we estimate panel regressions of the form:

$$Hate\ Crimes_{it} = \sum_{t=2010}^{2017} \beta_t \cdot Twitter\ Usage_i + \mathbf{X}'_{it}\gamma + County\ FE + Year\ FE + \epsilon_{it} \quad (1)$$

where the outcome variable is the natural logarithm of anti-Muslim hate crimes in county  $i$  and year  $t$  (with one added inside). *Twitter Usage* is the natural logarithm of the total number of Twitter users in a county (also with one added inside). The county fixed effects in the regression control for underlying differences in the number of hate crimes per county. Year fixed effects absorb changes in such crimes that affect all counties to the same extent. The main regressors of interest are  $\beta_t$ , which measure the differential change in anti-Muslim hate crimes in counties with higher Twitter usage in year  $t$ .

**Figure 4: Twitter Usage and the Increase in Anti-Muslim Hate Crimes**



*Notes:* This figure plots the coefficients from running an event study regression as in Equation (1). The dependent variable is the natural logarithm of anti-Muslim hate crimes (with 1 added inside). The omitted category is 2014, the year leading up to the 2016 presidential primaries (indicated with the vertical line). The shaded area indicates 95% confidence intervals based on standard errors clustered by state.

Figure 4 plots the estimated coefficients of Equation (1). The figure reveals that the increase in anti-Muslim hate crimes starting in 2015 appears to be concentrated in areas with

high Twitter usage. The magnitude of the coefficients indicate that a one standard deviation increase in Twitter usage is associated with 7% increase in anti-Muslim hate crimes per year. The coefficients for previous years are close to zero and not statistically significant, which suggests county-level social media use did not matter for the incidence of hate crimes before the 2016 presidential primaries.

The evidence suggests a potential connection between anti-Muslim sentiment and Twitter usage. However, our proxy for Twitter usage is likely correlated with a host of observable and unobservable factors that might also affect hate crimes. To overcome this challenge, the next section develops an identification strategy to isolate the effect of social media.

### 3.2 Identification Strategy

The results in the previous section suggests that the increase in anti-Muslim hate crimes around the 2016 presidential campaign has been concentrated in areas with high Twitter usage. In this section, we address the concern that Twitter usage may be correlated with other factors by developing an instrumental variable strategy based on the early diffusion of Twitter. The starting point is a county-level first-difference model relating the shift in anti-Muslim hate crimes around 2015 to a measure of Twitter usage:

$$\Delta Hate Crimes_i = \alpha + \beta \cdot Twitter Usage_i + \mathbf{X}'_i \gamma + State FE + \epsilon_i. \quad (2)$$

As a baseline,  $\Delta Hate Crimes$  will refer to the log-change of average hate crime incidents aimed at Muslims or other groups (with one added inside) between 2010-2014 and 2015-2017.<sup>12</sup> *Twitter Usage* is the natural logarithm of Twitter users in a given county, our measure of social media use. All regressions control for state fixed effects and dummies for each decile of the population distribution.

$\mathbf{X}_i$  is a vector of control variables that includes demographic controls for population growth and the share of the population in five-year age buckets; the linear distance of each county centroid from Austin Texas, the location of the SXSW festival (for more details see below); controls for ethnic composition and the share of Muslims; socioeconomic controls including the share of high school graduates or people with a graduate degree, the poverty rate, the unemployment rate, local GINI index, the share of uninsured individuals, the log median household income, the employment shares in eight sectors; media controls for the viewership share of Fox News, the cable TV spending to population ratio, and the prime time

---

<sup>12</sup>In robustness checks, we show that our results neither depend on the precise pre-period we use in the first-difference, functional form, or estimation method. The results also hold for the *level* of hate crimes after 2015 or the precise start date of Trump's presidential run.

TV viewership to population ratio; and the county-level vote share of the Republican party in 2012. Standard errors in all specifications are clustered at the state level.<sup>13</sup>

When estimating Equation (2) using OLS, the point estimates for  $\beta$  in Equation (2) are likely biased because Twitter usage is not exogenous. In particular, one may be concerned that the factors driving people to commit hate crimes are correlated with the decision to adopt social media. This could give rise to alternative interpretations of the graph in Figure 4 and the  $\beta$  estimate in Equation (2). To give one example, perhaps the potential perpetrators of hate crimes live predominantly in areas with a sizable presence of minority groups, and those areas are also more likely to use Twitter. In that case, the period around the 2016 presidential primaries and Trump's political rise could still be interpreted as a trigger point for anti-Muslim sentiments, but it is not clear whether or not social media plays a role.

To circumvent this issue, we exploit plausibly exogenous variation in the early adoption of Twitter in the United States. We make use of the fact that Twitter's popularity reached a tipping point at the SXSW conference and festival in 2007. During the festival, Twitter held a launch event with a special option that allowed users to join Twitter by simply sending a text message, and screens in the main hallways were used to show tweets about the festival. These measures proved to be extremely successful in spurring Twitter adoption. The daily volume of tweets increased from around 20,000 to 60,000 (Gawker, 2007). Figure 5a gives a first indication of the impact of SXSW on the success of Twitter: we see a clear spike of tweets about the event during the SXSW conference in mid-March 2007, followed by an upward shift in the growth of the total number of tweets. While total tweets grew by 55% from February to March, this growth accelerated to over 190% from March to April. March 2007 is also a clear outlier in the number of SXSW followers that signed up to Twitter (see A.5 in the online appendix). As another indication, there were more tweets about SXSW than about any other major festivals in 2007 (see Figure 5b). This is particularly noteworthy because of the relatively low number of attendees at SXSW Interactive.

Our identification strategy exploits that the home counties of SXSW attendees were most heavily exposed to this Twitter adoption shock, as these counties received a boost in the early-stage inflow of Twitter users. This pattern is in line with the literature on the path dependence of technology adoption (e.g. Arthur, 1989, 1994; Liebowitz & Margolis, 1999; Arrow, 2000). In Figure 5, we provide three pieces of evidence that these early adopters who were induced to join Twitter by SXSW were key to Twitter's rise in their home counties.

First, we show the short-term impact of SXSW on local tweets in Figure 5c by estimating event study panel regressions to compare Twitter activity in counties with and without new SXSW followers in March 2007. The graph clearly indicates that areas with early adopters at

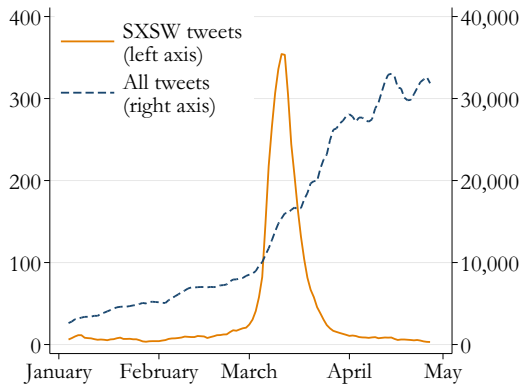
---

<sup>13</sup>In Table A.20 in the online appendix, we show the results with a range of alternative standard errors.

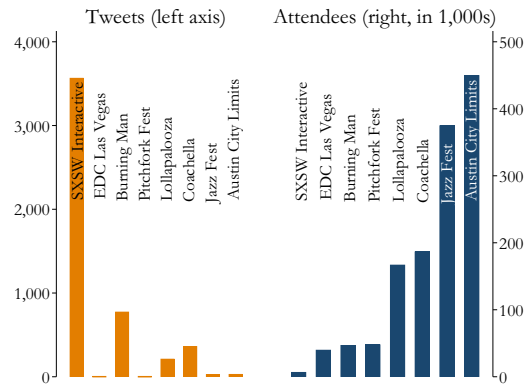


**Figure 5: South by Southwest 2007 and the Diffusion of Twitter**

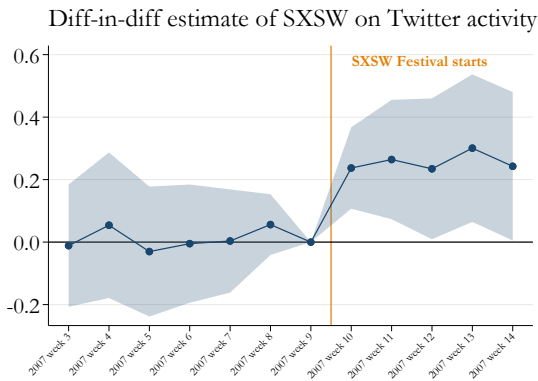
**(a) Twitter Activity SXSXW 2007**



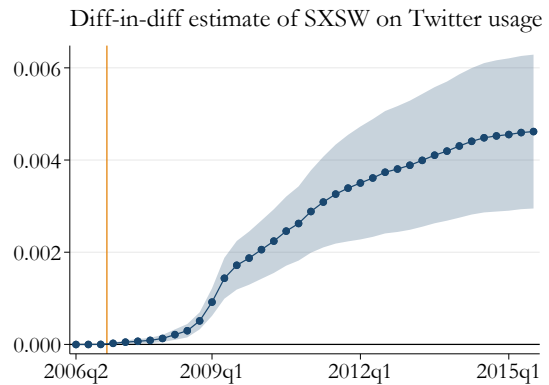
**(b) Other Festivals 2007**



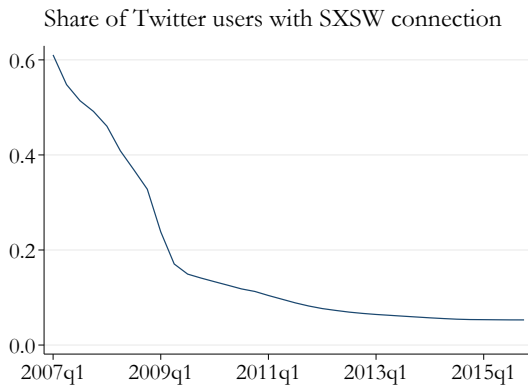
**(c) Short-Term Adoption Effect**



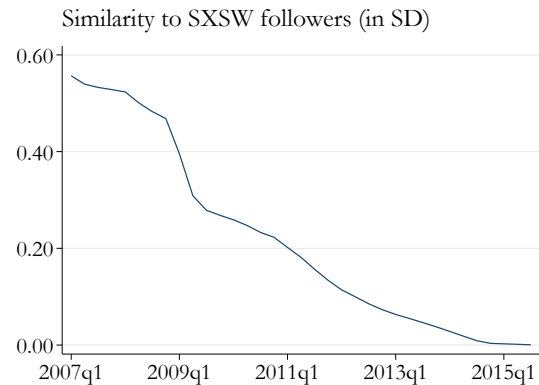
**(d) Long-Term Adoption Effect**



**(e) Connections to SXSXW Followers**



**(f) Similarity to SXSXW Followers**



*Notes:* Panel (a) plots the total number of tweets and those containing the term “SXSXW” over time, smoothed using a 7-day moving average. Panel (b) plots the number of tweets mentioning major festivals in 2007 in a 14-day window before and after the event. Attendee numbers are from various internet sources. Panel (c) and (d) plot the estimates of  $\beta_t$  from panel event study regressions of the type  $Outcome_{it} = \sum \beta_t SXSXW\ followers, March\ 2007_i \times Time_t + \mathbf{X}_{it} + County\ FE + Time\ FE + \varepsilon_{it}$ . In Panel (c), *Time* refers to weeks and *Outcome* to  $Log(1 + \# \text{ of tweets})$ . In Panel (d), *Time* refers to quarters and *Outcome* to *Twitter users/Capita*. Panel (e) plots the fraction of Twitter users that either follow SXSXW or follow a user who follows it over time. Panel (f) plots the similarity of all Twitter users to those that follow SXSXW based on their profile descriptions.

SXSW did not exhibit a higher growth rate of Twitter activity prior to SXSW Interactive 2007 but the growth rate increased in its aftermath. Quantitatively, counties with a one standard deviation higher number of SXSW followers in March (0.32) increased their local Twitter activity by around 10% in April compared to February 2007.

Secondly, Figure 5d shows the long-term adoption impact of the SXSW festival. For this exercise, we exploit the fact that we know when the more than four million Twitter users, we have data on, joined Twitter. We use this information to construct the cumulative number of Twitter users per capita in a county for each quarter between the launch of Twitter until the beginning of 2015. We then estimate an event study panel regressions and compare counties with and without new SXSW followers in March 2007. Again, we observe no pre-existing trends in the adoption of Twitter before the festival and an uptick in Twitter adoption with the beginning of SXSW that persists until the end of 2015. The two pre-SXSW quarters are not statistically significant, in contrast to all coefficients after the event. Consistent with theory, the pattern of Twitter adoption in these counties exhibits an S-shape typical for the diffusion of innovations (Griliches, 1957; Rogers, 2010; Bass, 1969; Geroski, 2000; Fagerberg et al., 2009). The estimates imply that a one standard deviation increase in SXSW followers who signed up in March 2007 increased Twitter adoption by around 22% by the end of 2015.

Third, we provide evidence that early Twitter adopters were indeed largely connected to the SXSW festival. Figure 5e plots the share of Twitter users in our data that either follows the SXSW festival or a SXSW follower who joined in March 2007. In March 2007, as many as 60% of Twitter users had either a first or second degree connection to the SXSW festival. With the diffusion of Twitter over time, this decreased to around 5% today. A similar pattern also holds for a text-based measure that captures the similarity of Twitter users generally with SXSW followers based on their user descriptions (“bios”).<sup>14</sup> Twitter users in March 2007 were close to 0.6 standard deviations more similar to SXSW followers than the average Twitter user today.

Taken together, we conclude that the 2007 SXSW festival led to increased adoption of Twitter in the home counties of attendees. We exploit that this pattern of technology adoption persists until today. The concern with this identification strategy is that, even after controlling for a large number of county characteristics, the home counties of SXSW followers who joined in March 2007 might be selected in a way that could explain increases in hate crime with Donald Trump’s presidential run without an impact of Twitter usage.

To address concerns about inherent differences of counties with SXSW followers, we include the number of SXSW followers that joined before the festival at any point in 2006 as

---

<sup>14</sup>This measure is constructed using Latent Semantic Analysis and cosine similarity. See Appendix Appendix A.1.6 for details.

a control variable in our regressions. In contrast to what one would expect if the persistent effect of Twitter adoption around SXSWS was driven by selection, these “control” counties do not exhibit systematic differences in Twitter usage today. The home counties of SXSWS followers who signed up before the 2007 event also do not systematically differ in observable characteristics from those of users who signed up during the event (see Table A.10). Out of 38 county characteristics, only three exhibit a mean difference that is marginally statistically significant, which vanishes once we apply a Šidák correction to account for multiple hypothesis testing. We also use user-level data to compare the profiles of people signing up for Twitter around SXSWS with those who signed up before. The analysis in Table A.11 again suggests that these user groups are highly similar: their first names and the terms they use to describe themselves in their Twitter “bio” are almost indistinguishable. As one indication, the correlation of words mentioned in the Twitter biographies between these groups is 0.92. Twitter users who reside in counties with SXSWS followers in March 2007 also do not differ systematically from those who live in other US counties (see Table A.12).

The identifying assumption underlying our empirical strategy is similar to Enikolopov et al. (2016). Conditional on a large number of county characteristics, the decision to start following SXSWS in March 2007 rather than before should drive increases in anti-Muslim sentiments with the 2016 presidential campaign only through the diffusion of Twitter usage.<sup>15</sup>

The historical diffusion of Twitter gives rise to a first-difference instrumental variable framework, with the first stage equation given by:

$$\begin{aligned} Twitter\ Usage_i &= \alpha + \delta_1 \cdot SXSWS\ followers,\ March\ 2007_i \\ &+ \delta_2 \cdot SXSWS\ followers,\ Pre_i \\ &+ \mathbf{X}_i' \psi + State\ FE + \xi_t, \end{aligned} \tag{3}$$

where *SXSWS followers, March 2007* is the number of SXSWS followers in county *i* that joined Twitter in March 2007, which serves as the excluded instrument. *SXSWS followers, Pre* are followers that joined before the festival at any point in 2006.

Figure A.1 in the online appendix plots the distribution of our proxy of new SXSWS followers in March 2007 across US counties. 155 counties received an inflow of early adopters of Twitter at the time of SXSWS. Table A.13, also in the online appendix, plots the correlation coefficients between the county-level SXSWS measures and those for the other festivals. Al-

---

<sup>15</sup>As an alternative control, we use tweeting about the much more popular festivals Coachella, Burning Man, and Lollapalooza in the same year. With the alternative festival controls, the assumption is similar in that attending SXSWS rather than other festivals in 2007 should only affect outcomes through this social media adoption channel.

though these variables are strongly correlated, as one would expect, there is enough variation in the locations of SXSW followers we can exploit in our empirical strategy. In robustness exercises, we consider a large range of alternative SXSW metrics, some of which show a considerably lower correlation between “treatment” and “control” group.

### 3.3 South by Southwest and Twitter Adoption: First Stage

Table 1 plots the results of estimating the first stage Equation (3). We can see that, across the board, the number of new Twitter users in March 2007 who followed SXSW is highly predictive of Twitter usage today. The point estimates are always statistically significant at the 1% level. The coefficient for SXSW followers in the months prior to the 2007 event is not statistically significant as soon as we control for observable county characteristics. Indeed, an *F*-test for the equality of coefficients suggests that the March 2007 and pre-period estimates are also statistically different from each other. Importantly, the coefficient estimates for March are highly stable and do not depend on the included covariates. Quantitatively, the estimate of 0.443 in column 8 implies that a one standard deviation increase in the log number of new SXSW followers in March (0.32) is associated with 15% higher Twitter usage today. The estimated effect based on the pre-period estimate implies less than 3% more users, which is not distinguishable from zero.

Based on these estimates and the event study plots in Figure 5, we conclude that county-level differences in the early diffusion of Twitter spread through the 2007 SXSW conference and festival are a reliable predictor of social media usage in the United States today. Because the locations of early adopters in the period before the festival do not predict Twitter usage, it is unlikely that this result is driven by selection into following the SXSW festival’s Twitter page. Put differently, the inflow of early adopters prompted by SXSW put some counties on a higher growth path of Twitter adoption than predicted based on observable county characteristics. In contrast, the otherwise highly similar counties with SXSW followers before this key event did not receive additional early adopters and their level of Twitter usage is accounted for by observable characteristics. In the next sections, we will employ the strong first stage result to estimate the effect of social media propagation on the recent rise in anti-minority sentiments.

## 4 Main Results

### 4.1 Reduced-Form and IV Estimates

The results in the previous section can be interpreted as evidence that social media plays a role in the recent increase in hate crimes in the United States. In this section, we use the new SXSU followers in March 2007 as an instrument for Twitter usage across the US today, while holding interest in SXSU prior to the key event constant to alleviate selection concerns.

Table 2 provides three sets of results. In panel A, we plot the OLS results from regressions of the change in hate crimes against Muslims on our measure of Twitter usage. Panel B shows the reduced-form results using new SXSU followers in March 2007 as instrument for Twitter usage. In panel C, we report the 2SLS results and a number of diagnostic tests. The results suggest that social media penetration, measured by Twitter usage, is positively associated with the increase in hate crimes against Muslims. The 2SLS estimates in column 8 imply that a one standard deviation increase in Twitter usage (1.76) is associated with a 32% ( $0.159 \times 1.76 \approx 0.28$  log points) larger increase in hate crimes after the start of the 2016 presidential primaries and Trump’s campaign launch. In Table A.19 in the online appendix suggests that these results are largely accounted for by a rise in assaults.

Since our baseline outcome variable is differenced over time, we also require that the parallel trends assumption holds. We already saw in Figure 4 above that hate crimes against Muslims disproportionately increased in areas with high Twitter usage only in 2015, *after* Trump’s presidential campaign started. Figure A.4 in the online appendix provides additional reduced form evidence in support of parallel trends when comparing areas with and without users that attended SXSU in March 2007.

A well-known concern with IV estimation is the weak instruments problem, which can lead to biased point estimates. We believe that our estimation does not suffer from a weak first stage. The robust  $F$ -statistic for the excluded regressor ranges between 57 and 98 in columns 1 through 8.<sup>16</sup> The values of the  $F$ -statistic are also above the critical values to reject the null hypothesis of a 5% potential bias with 5% statistical significance derived in Olea & Pflueger (2013), which is 37.42.<sup>17</sup>

We also assess the significance of our main estimates using confidence sets based on test inversion that are valid whether or not instruments are weak. For the case of a single instrument we study here, Andrews et al. (2019) recommend reporting Anderson-Rubin (AR) confidence sets that are efficient and robust to weak identification (Anderson et al.,

---

<sup>16</sup>Note that because the model is just-identified, the robust  $F$ -statistic (also called Kleibergen-Paap) is equivalent to the effective  $F$ -statistic derived in Olea & Pflueger (2013).

<sup>17</sup>These authors extend the well-known thresholds of Stock & Yogo (2005) to the case of heteroskedasticity-robust and, relevant in our case, clustered standard errors.

1949). Andrews (2018) develops a two-step approach to construct these confidence sets that is implemented in Stata by Sun (2018). Basing inference on this two-step approach sidesteps the issue that the usually reported (Wald) confidence intervals for 2SLS estimates can exhibit large coverage distortions. This is because AR confidence sets allow for inference without assessing the strength of first-stage results in a separate initial step. As such, we can determine whether our second stage coefficients are likely to be non-zero even if our instrument was indeed weak. Reassuringly, the AR confidence sets reported below the (instrumented) Twitter usage in panel C always exclude zero.

Across all specifications in Table 2, the OLS estimates are highly statistically significant, but smaller than those obtained using 2SLS. There are a number of potential reasons for the difference in magnitudes. The first possibility is that the selection of individuals into social media adoption biases the OLS estimates downward. To give one example, if people in particularly xenophobic areas commit more hate crimes but are less likely to use Twitter, the OLS estimate would be downward biased. Second, counties with more SXSW followers that signed up in March 2007 may have a higher local average treatment effect (LATE). Third, the endogenous variable (our proxy for Twitter usage) is likely subject to measurement error. This measurement error could also bias the OLS estimate towards zero.

Our findings are unlikely to be affected by social media changing people's propensity to report hate crimes rather than actual incidents. Data from the Bureau of Justice Statistics National Crime Victimization Survey (NCVS) shows that the likelihood of hate crime victims to file a report with the police has, if anything, slightly dropped since 2015 compared to previous years. This can be seen in Figure A.6 in the online appendix. Our empirical strategy also rules out many potential sources of reporting changes as an alternative explanation. The first-difference regressions with state fixed effects mean we consider changes within counties over time and abstract from potential changes in reporting across states. An increase in reporting is thus an unlikely explanation for the increase in anti-Muslim hate crimes.<sup>18</sup>

**Social Media and Changes in Other Hate Crimes.** So far, we have focused on changes in anti-Muslim hate crimes, motivated by the fact we found little change in the frequency of other types of hate crimes around the start of Trump's presidential campaign in the FBI data. However, one might expect that Trump's presidential run could also affect other hate crimes, in particular anti-Hispanic incidents.<sup>19</sup> If social media plays a role, such incidents may have

---

<sup>18</sup>Hobbs & Lajevardi (2019) find that the 2016 presidential election was associated with a partial withdrawal of Muslims from public life. This suggests that we might underestimate the effect on anti-minority sentiment.

<sup>19</sup>In his presidential campaign announcement speech, Trump famously singled out Hispanics and Arab Muslims: "When Mexico sends its people, they're not sending their best. ... They're bringing drugs. They're bringing crime. They're rapists. And some, I assume, are good people. ... They're sending us not the right people. It's coming from more than Mexico. It's coming from all over South and Latin America, and it's coming

become more common in areas with high Twitter usage even if their total number remained unchanged. In Table 3, we consider this possibility empirically by replacing the dependent variable with the log change in hate crimes targeting Hispanic ethnicity, other ethnicities, race, sexual orientation, or religion (excluding anti-Muslim). We also consider hate crime data from the Anti-Defamation League (ADL) as an alternative data source in column 7.<sup>20</sup>

Overall, we also find a role for social media in explaining increases in the total number of hate crimes and those targeting Hispanics, the other minority group frequently singled out by Donald Trump. However, only anti-Muslim hate crimes show a clearly consistent pattern across the OLS and 2SLS estimates. There is less evidence for a reallocation of other hate crimes towards areas with higher Twitter usage.<sup>21</sup> In the 2SLS estimation, a one standard deviation increase in Twitter usage is associated with a 31% larger increase in total hate crimes, and a 26% larger increase for incidents targeting Hispanics.

## 4.2 Robustness

We consider a range of sensitivity checks to validate the robustness of our main findings. We begin by considering alternative ways to account for the selection of users into events such as SXSW in Table A.15 in the online appendix. These alternatives vary widely in the number of counties in the “control group”—counties for which we observe other users but no new SXSW followers in March 2007—and their correlation with our instrument. In column 1, we begin by showing that the results also hold when dropping the SXSW control, which makes the results somewhat stronger. In columns 3 through 6, we consider alternative time periods for the pre-period variable or alternatively control for the individual months. Column 7 uses users tweeting about *other* festivals in 2007 to account for the selection of users into events such as SXSW. We consider tweets about three of the most popular festivals in the United States that are, in many respects, very similar to SXSW: Coachella, Burning Man, and Lollapalooza. Our results are robust throughout, which provides additional evidence that our results are not driven by unobserved selection into SXSW in a particular month.

We also use alternative metrics of Twitter usage in Table A.16 in the online appendix. We consider two survey measures of Twitter usage provided by GfK Mediamark Research & Intelligence (via SimplyAnalytics), as well as two alternative transformations of the GESIS Twitter data (only Twitter users who joined before June 2015 or the number of tweets rather than Twitter *users*). All of these measures yield similar estimates.

---

probably—probably—from the Middle East.”

<sup>20</sup>For most counties, the ADL report hate crimes from 2016 onward, so we focus on the *level* rather than the change in hate crimes. In unreported results, we find similar results using changes in ADL hate crimes.

<sup>21</sup>The weaker evidence for bias against other religions is driven by anti-Jewish incidents (unreported).

In Table A.17, in the online appendix, we present additional robustness checks. In column 1, we weight by a county's population, which decreases the difference between OLS and 2SLS estimates. In column 2, we consider the change in anti-Muslim hate crimes since 1990 (rather than 2010); this yields somewhat larger estimates. In column 3, we replace the change in hate crimes with the log number of hate crimes after Trump's presidential run as dependent variable. In columns 4 through 6 of Table A.17, we address the concern that anti-Muslim hate crimes reported by the FBI mainly occur in a relatively small fraction of all counties. In column 4, we begin by dropping all counties that report a zero change in anti-Muslim hate crimes between 2010 and 2017. Because this applies to the majority of counties, the sample size shrinks considerably. One way to think about this estimation is that it captures the intensive margin of hate crimes. Despite the drop in observations, our estimates remain statistically significant. In column 5, we drop counties for which the FBI always reports zero hate crimes, which likely reflects a lack of reporting. We drop all counties for which the (rounded) estimated share of Muslims in the total population is zero from the sample in column 6.<sup>22</sup> Again, these changes leave our results intact.

In column 7, we restrict the sample to neighbouring counties where one has no new SXSX followers in March 2007 and the other one has at least one. This is to purge the estimates of potential unobserved local confounders. In column 8, we restrict the sample to the counties where we have variation in SXSX followers (either in March 2007 or before), i.e. the intensive margin of SXSX Twitter users. This rules out potential concerns about the limited geographical variation of our instrument. Reassuringly, this yields quantitatively similar estimates to our baseline results. At last, in column 9, we show that constructing the dependent variable as the difference in hate crimes around the exact start of Donald Trump's presidential campaign in June 2015—rather than as the period 2015-2017 compared to 2010-2014—leaves our findings unchanged.

Table A.18 considers other estimation techniques: IV probit (with a dummy for increases in hate crimes in a county as dependent variable); IV poisson (with the number of hate crimes after Trump's presidential campaign as dependent variable); OLS regressions where, instead of natural logarithms, we use inverse hyperbolic sine transformations; and OLS regressions where the dependent variable is an index equal to 1 for increases in anti-Muslim hate crimes, 0 for no change, and  $-1$  for decreases. In all of these exercises, the results are highly similar to our baseline findings.

---

<sup>22</sup>Although the Religious Census reports no Muslims living in these counties, this might be the artifact of a very small number, rather than an actual zero.



**Table 1: First Stage - South by Southwest 2007 and the Diffusion of Twitter Usage**

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Log(SXSW followers, March 2007)	0.614*** (0.062)	0.581*** (0.062)	0.554*** (0.067)	0.525*** (0.061)	0.481*** (0.055)	0.472*** (0.057)	0.452*** (0.059)	0.443*** (0.059)
Log(SXSW followers, Pre)	0.213*** (0.072)	0.225*** (0.083)	0.171** (0.078)	0.117 (0.080)	0.115 (0.077)	0.108 (0.075)	0.098 (0.074)	0.090 (0.071)
State FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Population controls	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Demographic controls		Yes	Yes	Yes	Yes	Yes	Yes	Yes
Geographical controls			Yes	Yes	Yes	Yes	Yes	Yes
Race and religion controls				Yes	Yes	Yes	Yes	Yes
Socioeconomic controls					Yes	Yes	Yes	Yes
Media controls						Yes	Yes	Yes
Election control							Yes	Yes
Crime controls							Yes	Yes
Observations	3,107	3,107	3,107	3,107	3,106	3,105	3,105	3,105
R <sup>2</sup>	0.927	0.933	0.934	0.936	0.944	0.945	0.946	0.947
Mean of DV	5.289	5.289	5.289	5.289	5.290	5.290	5.290	5.290
p-value: March 2007 = Pre	0.00	0.01	0.01	0.00	0.00	0.00	0.00	0.00

*Notes:* This table presents county-level regressions where the dependent variable is the number of tweets sent (in natural logarithm). *SXSW followers, March 2007* is the number of Twitter users who joined in March 2007 and follow South by Southwest (SXSW) *SXSW followers, Pre* is the number of SXSW followers who registered at some point in 2006. The bottom row reports *p*-values from *F*-tests for the equality of these coefficients. All regressions control for population deciles and state fixed effects (not shown). Demographic controls include population growth between 2000 and 2016 as well as age cohort controls for the share of people aged 20-24, 25-29, 30-34, 35-39, 40-44, 45-49, and those over 50. Race and religion controls contains the share of people identifying as white, African American, Native American or Pacific Islander, Asian, Hispanic, or Muslim. Socioeconomic controls include the poverty rate, unemployment rate, local GINI index, the share of uninsured individuals, log median household income, the share of highschool graduates, the share of people with a graduate degree, as well as the employment shares in agriculture, information technology, manufacturing, nontradables, construction and real estate, utilities, business services, or other sectors. Media controls include the viewership share of Fox News, the cable TV spending to population ratio, and the prime time TV viewership to population ratio. Election control is the county-level vote share of the Republican party in 2012. Crime controls are the rates of violent or property crime from the FBI. Geographical controls include the linear distance from the SXSW festival location (Austin, Texas), population density, and the natural logarithm of county size. Robust standard errors in parentheses are clustered by state. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

**Table 2: 2SLS - Social Media and the Rise in Hate Crimes Against Muslims**

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
$\Delta \text{Log}(\text{Hate crimes against Muslims})$								
<b>Panel A: OLS</b>								
Log(Twitter users)	0.030*** (0.007)	0.030*** (0.007)	0.032*** (0.008)	0.025*** (0.006)	0.025*** (0.006)	0.025*** (0.006)	0.024*** (0.006)	0.025*** (0.006)
<b>Panel B: Reduced form</b>								
Log(SXSW followers, March 2007)	0.070** (0.032)	0.070** (0.032)	0.079** (0.031)	0.073** (0.031)	0.071** (0.031)	0.071** (0.031)	0.070** (0.031)	0.071** (0.031)
<b>Panel C: 2SLS</b>								
Log(Twitter users)	0.114** (0.052)	0.121** (0.054)	0.143** (0.057)	0.139** (0.061)	0.148** (0.065)	0.151** (0.067)	0.156** (0.070)	0.159** (0.071)
Weak IV 95% AR confidence set	[0.018; 0.211]	[0.021; 0.221]	[0.038; 0.260]	[0.025; 0.264]	[0.028; 0.268]	[0.028; 0.288]	[0.027; 0.299]	[0.028; 0.304]
Log(SXSW followers, Pre)	0.034 (0.065)	0.032 (0.067)	0.052 (0.063)	0.039 (0.061)	0.039 (0.061)	0.041 (0.061)	0.041 (0.060)	0.042 (0.060)
State FE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Population controls	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Demographic controls		Yes	Yes	Yes	Yes	Yes	Yes	Yes
Geographical controls			Yes	Yes	Yes	Yes	Yes	Yes
Race and religion controls				Yes	Yes	Yes	Yes	Yes
Socioeconomic controls					Yes	Yes	Yes	Yes
Media controls						Yes	Yes	Yes
Election control							Yes	Yes
Crime controls							Yes	Yes
Observations	3,107	3,107	3,107	3,107	3,106	3,105	3,105	3,105
Mean of DV	0.018	0.018	0.018	0.018	0.018	0.018	0.018	0.018
Robust F-stat.	98.42	86.97	69.03	74.80	76.10	67.77	58.94	56.63

*Notes:* This table presents county-level OLS and IV regressions where the dependent variable is the log change in hate crimes against Muslims between 2010 and 2017. *Log(Twitter usage)* is instrumented using the number of users who started following SXSW in March 2007. *SXSW followers*, *Pre* is the number of SXSW followers who registered at some point in 2006. All regressions control for population deciles and state fixed effects (not shown). Demographic controls include population growth between 2000 and 2016 as well as age cohort controls for the share of people aged 20-24, 25-29, 30-34, 35-39, 40-44, 45-49, and those over 50. Race and religion controls contains the share of people identifying as white, African American, Native American or Pacific Islander, Asian, Hispanic, or Muslim. Socioeconomic controls include the poverty rate, unemployment rate, local GINI index, the share of uninsured individuals, log median household income, the share of highschool graduates, the share of people with a graduate degree, as well as the employment shares in agriculture, information technology, manufacturing, nontradables, construction and real estate, utilities, business services, or other sectors. Media controls include the viewership share of Fox News, the cable TV spending to population ratio, and the prime time TV viewership to population ratio. Election control is the county-level vote share of the Republican party in 2012. Crime controls are the rates of violent or property crime from the FBI. Geographical controls include the linear distance from the SXSW festival location (Austin, Texas), population density, and the natural logarithm of county size. Weak IV 95% Anderson-Rubin (AR) confidence sets are calculated using the two-step approach of Andrews (2018) using the Stata package from Sun (2018). For the just-identified case we study here, the “robust” *F*-stat. is equivalent to the “Kleibergen-Paap” or the “effective” *F*-statistic of Olea & Pflueger (2013). Robust standard errors in parentheses are clustered by state. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

**Table 3: Social Media and Other Hate Crimes**

	FBI Data				ADL Data		
	Total (1)	Hispanic (2)	Other ethnic (3)	Race (4)	Sexual Orientation (5)	Religion (excl. Muslims) (6)	Total (Levels) (7)
<b>Panel A: OLS</b>							
Log(Twitter users)	0.018 (0.012)	0.002 (0.012)	-0.021** (0.009)	0.013 (0.012)	-0.005 (0.007)	0.026** (0.011)	0.291*** (0.049)
<b>Panel B: Reduced form</b>							
Log(SXSW followers, March 2007)	0.090** (0.041)	0.076** (0.035)	0.010 (0.029)	0.051 (0.043)	0.068* (0.040)	0.061* (0.033)	0.510*** (0.097)
<b>Panel C: 2SLS</b>							
Log(Twitter users)	0.154** (0.072)	0.131** (0.057)	0.017 (0.049)	0.088 (0.075)	0.117 (0.071)	0.105* (0.054)	0.878*** (0.113)
Weak IV 95% AR confidence set	[0.021; 0.302]	[0.025; 0.237]	[-0.074; 0.109]	[-0.052; 0.242]	[-0.014; 0.263]	[-0.006; 0.205]	[0.647; 10.087]
Log(SXSW followers, Pre)	-0.043 (0.077)	-0.092 (0.072)	-0.041 (0.059)	-0.014 (0.082)	-0.055 (0.077)	-0.050 (0.061)	0.154 (0.104)
Observations	3,107	3,107	3,107	3,107	3,107	3,107	3,107
Mean of DV	-0.017	-0.012	-0.015	-0.012	-0.025	0.005	0.230
Robust F-stat.	86.97	86.97	86.97	86.97	86.97	86.97	86.97

*Notes:* This table presents county-level OLS, reduced form, and IV regressions where the dependent variable is the log change in hate crimes against the group in the top row between 2010 and 2017. *Log(Twitter usage)* is instrumented using the number of users who started following SXSW in March 2007. All regressions control for population deciles and state fixed effects (not shown). We include the full set of controls, as in column 8 of Table 2. Demographic controls include population growth between 2000 and 2016 as well as age cohort controls for the share of people aged 20-24, 25-29, 30-34, 35-39, 40-44, 45-49, and those over 50. The hate crime data from the Anti-Defamation League (ADL) is sparse prior to 2016, so we use the log-level of hate crimes in column 7. Weak IV 95% Anderson-Rubin (AR) confidence sets are calculated using the two-step approach of Andrews (2018) using the Stata package from Sun (2018). For the just-identified case we study here, the “robust” *F*-stat. is equivalent to the “Kleibergen-Paap” or the “effective” *F*-statistic of Olea & Pfueger (2013). Robust standard errors in parentheses are clustered by state. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

## 5 Trump's Tweets and Anti-Muslim Sentiment

The previous section suggests that social media may have played a role in the spread of anti-Muslim sentiment around 2015, the time Donald Trump started his presidential campaign. An often-voiced hypothesis is that Trump actively contributes to anti-Muslim sentiment through his incendiary comments on Twitter. Indeed, there is some existing evidence that influential individuals can have a disproportionate effect on public opinion (e.g. Beaman et al., 2009; Bursztyjn et al., 2017; Alatas et al., 2019).

One potential channel underlying the patterns we have documented so far could thus be that Trump's rhetoric, broadcasted via social media, has real-life effects. We attempt to shed some light on this channel by analyzing the time series relationship between Trump's tweets about Muslims, anti-Muslim hate crimes, and media attention. We attempt to get at the issue of causality by again leveraging an instrumental variable.

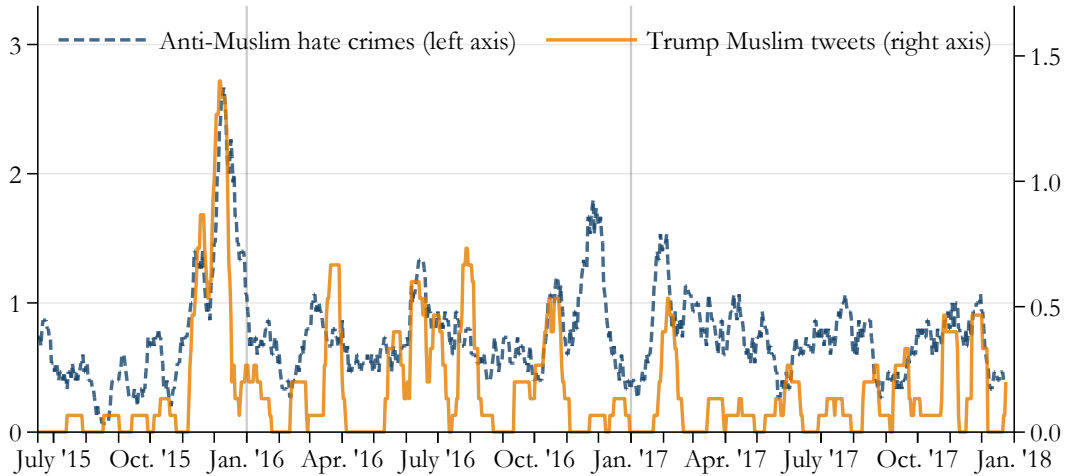
### 5.1 Trump Tweets and Hate Crimes

We begin by plotting the number of Trump's tweets about Islam-related topics and anti-Muslim incidents over time in Figure 6. We define these tweets based on a careful reading of Trump's Twitter feed, combined with a machine learning algorithm; see the data section and online appendix Table A.6 for more details. Since the daily number of tweets is highly volatile, we plot the 14-day moving average of the series; collapsing the data on the weekly level looks very similar (unreported).

It is immediately apparent that Trump's tweets about Muslims and anti-Muslim hate crimes are highly correlated. This correlation could reflect that Trump reacts to US-wide anti-Muslim sentiments driven by observable and unobservable factors, e.g. terrorist attacks. It could also be that Trump's tweets themselves contribute to a climate that enables hate crimes. Clearly, we cannot disentangle these possibilities using the graphical evidence from the data or running a simple OLS regression of hate crimes on tweets.

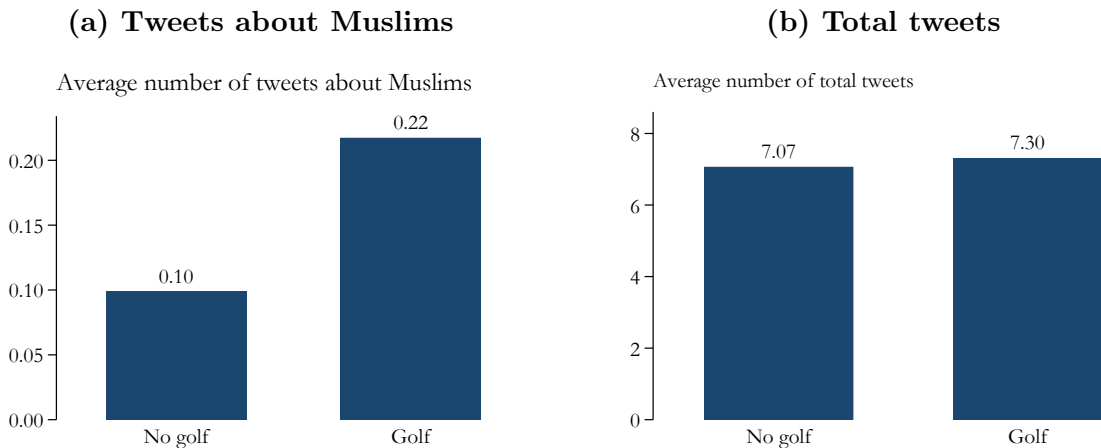
We propose an instrumental variable strategy to get around the most obvious concerns. In particular, we leverage Trump's passion for golf. In 2017 alone, Trump likely golfed on 92 days. As it turns out, the data suggest a strong link between Trump's golf outings and his Twitter feed: Figure 7 shows that while the total number of tweets he sends are unchanged on golf days, the *content* of his tweets sharply tilts towards negative, Muslim-related rhetoric. In 2017, 15 out of the 34 tweets we classify as negatively mentioning Muslims were sent on golf days. In Figure A.8a in the online appendix, we show that the topic shift is explained by a drop in policy-related tweets and more frequent mentions of Muslims and the media. Figure A.7c shows that his tweets also become more negative in sentiment.

**Figure 6: Trump’s Tweets About Muslims and Anti-Muslim Hate Crime**



*Notes:* This figure plots a 14-day moving average of anti-Muslim hate crimes from the FBI and Donald Trump’s tweets about Muslims for the period from Trump’s presidential campaign start in June 2015 until the end of 2017.

**Figure 7: Trump’s Twitter Activity, Split by Golf Days**



*Notes:* These figures plot the daily average number of Trump’s tweets in 2017, split by whether he plays golf on a given day. Panel (a) reports the average number of tweets about Muslims, panel (b) reports the total number of tweets.

One explanation for this pattern is that Trump’s attention shifts away from policy issues once he is away from the White House. Another influence on golf days is his social media manager and former caddie Dan Scavino, who is known to supply Trump with internet content and suggested tweets (Edwards, 2018; Reilly, 2019; CNN, 2020). Figure A.8 in the online appendix provides additional evidence that Trump’s daily schedule influences the content of his tweets. In particular, we show that Trump is more likely to tweet about foreign politics when he is abroad and more likely to tweet about domestic and party politics on days he receives a policy briefing.

The identifying assumption is that Trump’s golf outings are only systematically correlated with anti-Muslim sentiment through their effect in Trump’s tweeting behaviour. As the President’s schedule is to a significant extent predetermined to accommodate security concerns and meetings, it is plausibly exogenous with respect to hate crimes against Muslims. We provide additional evidence supporting the exclusion restriction below.

More formally, we run time series regressions using the following framework:

$$Hate\ Crimes_{t+h} = \alpha + \beta \cdot Muslim\ Trump\ Tweets_t + \mathbf{X}'_t \gamma + \epsilon_{t+h} \quad (4)$$

$$Muslim\ Trump\ Tweets_t = \alpha + \delta \cdot I[Trump\ golfs]_t + \mathbf{X}'_t \psi + \xi_t \quad (5)$$

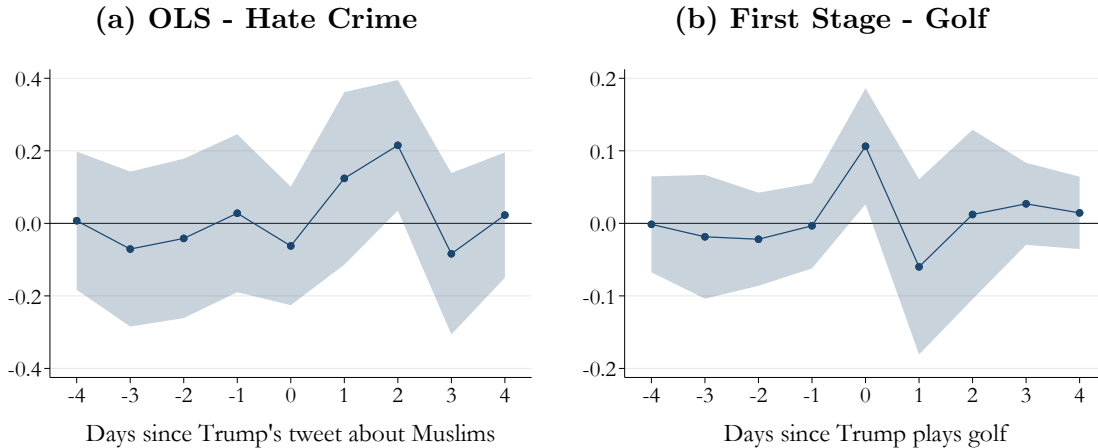
The dependent variable in equation (4) is the natural logarithm of US-wide hate crimes against Muslims on day  $t$  (with one added inside). The main regressor of interest is the natural logarithm of the number of Donald Trump’s Muslim tweets (again with one added inside). In the baseline specification, the vector  $X_t$  includes linear and quadratic time trends and a full set of day-of-week as well as year-month fixed effects. We focus on 2017, for which we have both details about Trump’s schedule and data on hate crimes. We present additional OLS evidence for the full time period since Trump joined Twitter in 2009 below.

Naively estimating equation (4) would not be informative about whether Trump’s Twitter activity might contribute to driving sentiments because his tweets cannot be regarded as random. We will thus instrument for tweets about Muslims in equation (5) using  $I[Trump\ golfs]_t$ , an indicator variable that is 1 for days on which Trump plays golf. We base inference on Newey-West standard errors that allow for heteroscedasticity and autocorrelation.

The appropriate choice of the prediction horizon  $h$  depends on the lead-lag relationship between Trump’s tweets and real-life hate crimes. We plot the result from estimating an event-study OLS specification of equation (4) where we allow for leads and lags of Trump’s tweets about Muslims in panel (a) of Figure 8. As we can see, the log number of anti-Muslim hate crimes is essentially flat prior to Trump’s tweets and subsequently rises to its peak in

$t+2$ . In our baseline regressions, we will thus set  $h$  to 2. Panel (b) also plots the dynamic relationship between Trump’s golf outings and tweets about Muslims. We can see that his tweets only increase on the days he golfs, with no similar spikes before and after.

**Figure 8: Time Series Correlations Trump Tweets about Muslims**



*Notes:* These figures plot the  $\beta_\tau$  coefficients from dynamic versions of equations 4 and 5 of the type  $Y_t = \alpha + \sum_{\tau=-4}^4 \beta_\tau \cdot Z_t + \mathbf{X}'_{t-\tau} + \epsilon_t$ . In Panel (a), the dependent variable is the number of anti-Muslim hate crimes and  $Z_t$  the number of Donald Trump’s tweets about Islam-related topics (both in natural logarithm). In Panel (b),  $Y_t$  is the log number of Trump’s Islam-related tweets and  $Z_t$  a dummy for days when he golfs. 0 indicates the date of tweets about Muslims or golfing ( $\tau = 0$ ). All regressions include linear and quadratic time trends; a full set of day of week and year-month dummies; and four lags of dummies for the incidence of terror attacks in the US and Europe. The sample period is the year 2017. The shaded areas are 95% confidence intervals based on Newey-West standard errors.

Table 4 presents the regression results of this exercise. We plot the first stage coefficients in panel A, OLS coefficients in panel B, reduced form coefficients in panel C, and the 2SLS estimation in panel D. Across the different specifications, the estimates suggest a clear link between Trump’s tweets about Muslims and subsequent real-life hate crimes. To get a sense of the implied magnitudes, consider the estimate in column 7 of panel D in Table 4. The coefficient of 1.648 implies that a one standard deviation increase in the log number of tweets about Muslims (0.25) is associated with a 41 log-point increase in hate crimes. This effect is larger than the OLS estimate of 0.091. A potential explanation is that unobserved third factors lead to a downward bias of the OLS estimates. For example, Trump’s tweets about Muslims might coincide with periods of *low* pre-existing anti-Muslim sentiment. In that case, the OLS estimates would be downward biased because they conflate the true Trump effect with low general anti-Muslim sentiment. This explanation is also consistent with the finding that controlling for general attention paid to Muslims or terror attacks in columns 4 through 6 *increases* the point estimates relative to the baseline specification.

As mentioned above, a concern with instrumental variable estimation is the weak instruments problem. Two pieces of information suggest that our 2SLS estimates do not suffer from this issue. First, the robust  $F$ -statistics we find are consistently above the widely used linear IV rule of thumb of 10. Most of them are above the critical value for a worst case bias of 30% (at 5% statistical significance) using the cutoffs from Olea & Pflueger (2013). Second, the Anderson-Rubin confidence sets constructed using the two-step approach proposed in Andrews (2018) always exclude a zero estimate even if we assume that the instrument is weak. The reduced form and 2SLS results thus suggest that Trump's tweets could indeed be a contributing factor triggering potential perpetrators to commit real-life hate crimes.

Another concern could be that the exclusion restriction of the instrument is violated and Trump's golf visits correlate with anti-Muslim hate crimes for reasons unrelated to his Twitter activity. However, several pieces of evidence are hard to square with this alternative interpretation. First, golf visits only affect the probability of anti-Muslim tweets on the day itself (see Figure 8). The sharp pattern is unlikely to be explained by the news cycle, for which we would expect a smoother pattern that should also affect Trump's tweets on the previous and following days. Second, the reduced form and 2SLS coefficients are almost fully unchanged when we control for measures of the salience of Muslim-related topics based on Google searches and the number of mentions on the big three TV networks (Fox News, CNN, and MSNBC). This suggests that Trump's golf outings do not appear to strongly correlate with pre-existing salience of Muslims. Finally, we present some additional evidence in support of the exclusion restriction in Table A.25. In column 1, we require that no terror attack occurred on the four days before Trump's golf outing (and tweets about Muslims). This exercise results in a slightly larger point estimate, suggesting that our effects are not driven by periods of high salience of Muslims. Columns 2 and 3 next split the sample into periods above and below the median number of reports about Muslims on the previous day on Fox News. Consistent with the idea that Trump's golf trips somewhat coincide with *lower* pre-existing sentiments, we find somewhat stronger (and statistically significant) predictability of hate crimes with Trump's tweets when reporting was low on the previous day.

In Table A.28 in the online appendix, we re-run the OLS estimation for the entire period since Trump's first tweet in 2009 and split the sample into the period before and after the launch of his presidential run on June 16, 2015. We find very similar OLS estimates on his tweets about Muslims, but only after the start of his presidential campaign. For the much longer period from 2009 to mid-2015, his tweets seem to be uncorrelated with anti-Muslim hate crimes. While many factors may explain this finding, it is at least some indication that we are capturing a general pattern.

In Table A.24 in the online appendix, we report more robustness results. Our results



remain largely unchanged when we control for more lags of the dependent variable to capture stronger serial correlation in hate crimes. We further experiment with additional controls for the total length of Trump's golf outings in column 3, a control if Trump golfed in the previous week (column 4), or alternative definitions of the golf dummy in columns 6 and 7. Our results are also robust to using a dummy for days with *any* Islam-related tweet from Trump (column 5).

Given the relatively short sample period, how likely would it be to find an effect if we picked golf days at random? Figure A.7b reports the results of a randomization test for the first stage regression of Trump's tweets about Muslims on a golf dummy, where we randomly pick 92 golf days in 2017 (except the ones used in the actual variable). The distribution of the resulting *t*-statistics of the golf day dummy suggests that none of the placebo coefficients are close to our estimate.

We further investigate which type of anti-Muslim hate crimes drive our results. Based on the most common criteria in the FBI data, we divide anti-Muslim incidents into vandalism, theft, burglary, robbery, and assault. The results of this exercise are presented in Table A.26 in the online appendix. Our high-frequency results appear to be mainly driven by cases of vandalism.<sup>23</sup>

The precise timing in our time series results also go against the idea we are capturing increases in hate crime reporting, rather than actual incidents. If Trump's negative tweets about Muslims make people more willing to report hate crimes, they should also become more likely to report *past* hate crimes. This would lead to a very different time series pattern: increases in reporting should then translate into a larger number of hate crimes not only after but also *before* Trump's tweets. However, the data only shows a spike *after* the tweets.<sup>24</sup>

As a simple validation exercise, we also investigate whether Trump's messages about Muslims are also correlated with hate crimes against other minorities. In particular, we consider incidents motivated by ethnicity, race, sexual orientation, or religions other than Islam. Table A.27 plots the predictive ability of Trump's tweets about Islam-related topics for these different types of hate crimes. We only find clear-cut correlations with crimes against Muslims, not other hate crimes. This suggests that we are not merely capturing anti-minority sentiment but Muslim-specific hatred.

---

<sup>23</sup>Note that this does not stand in contradiction to our cross-sectional results, for which we find the largest role for assault. The daily variation we exploit here likely picks up more spontaneous anti-Muslim incidents relative to the medium-term effects in the cross-section.

<sup>24</sup>It also seems unlikely that the time series findings are driven by changes in the way the FBI classifies hate crimes, because the incident date rarely corresponds to the date a hate crime is reviewed by the FBI as part of the two-tier process. If Trump's tweets change the behavior of FBI analysts, this would again lead to increases in hate crimes before Trump's tweets, which we do not observe in the data.

## 5.2 Trump Tweets and Twitter Spillovers

We next provide evidence that Trump's negative tweets about Muslims have a direct effect on his followers. In particular, we analyze if Trump's followers become more willing to express anti-Muslim content. For this analysis, we use more than 115 million tweets drawn from a random 1% sample of Trump's followers, around 630,000 users. In this dataset, we identify tweets that are retweets of Trump's negative content about Muslims, tweets that refer to Muslim-related topics but are not retweets of Trump, and tweets that contain the hashtags #StopIslam or #BanIslam.

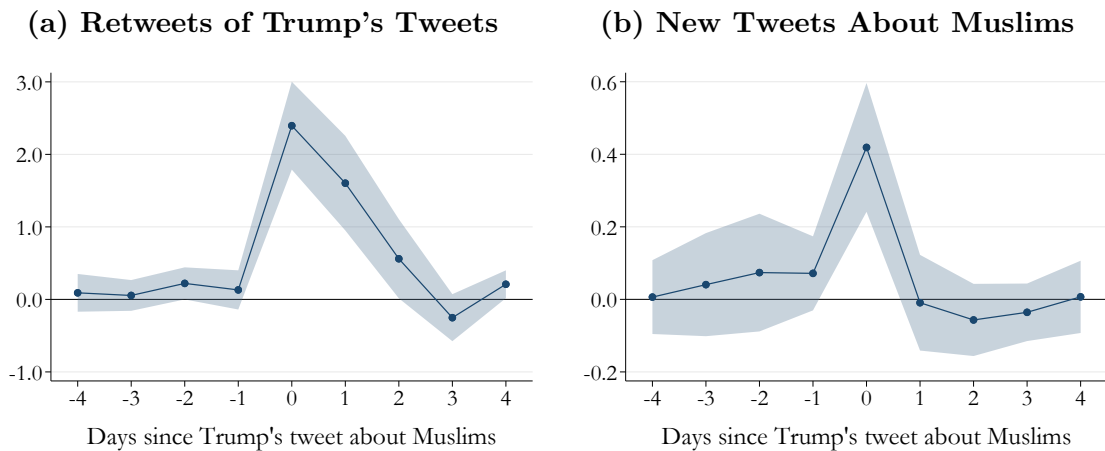
We continue to run time series regressions of the type in equation (4). To start, we plot dynamic correlations in Figure 9, where the dependent variables are different measures of tweets (in natural logarithm). The results show a clear pattern. Trump's negative tweets about Muslims are not only widely shared by his followers over the next days but also systematically followed by a spike in new content about Muslims. The time series pattern suggests no increase of anti-Muslim sentiment before Trump's tweets.

Columns 1 through 3 in Table 5 provide evidence that these patterns also hold when we instrument for the tweets using golf days. We focus on contemporaneous correlations, as suggested by the pattern in Figure 9. The reduced form and 2SLS specifications are highly statistically significant, and the weak IV confidence sets always clearly exclude zero. The 2SLS estimates suggest that a one standard deviation increase in Trump's Muslim tweets (0.25) is followed by more than a doubling of retweets and a 33% increase in new messages about Muslims that do not mention Trump. They are also followed by a 75% increase in the use of the hashtags #StopIslam or #BanIslam by Trump followers.

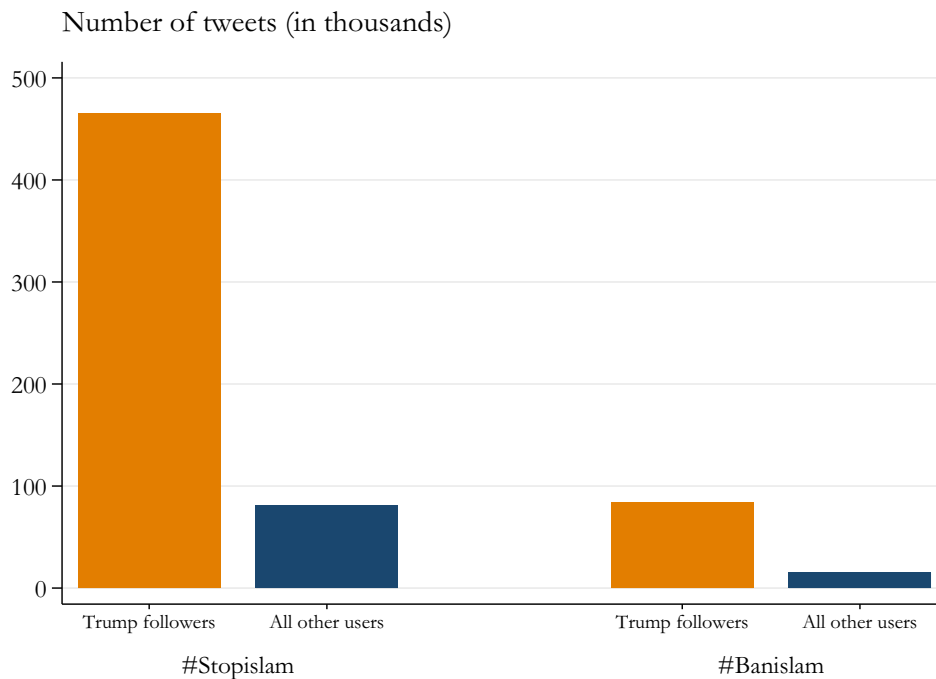
In Figure 9c, we plot the number of tweets using the hashtags #StopIslam and #BanIslam, as well as the number of these tweets coming from Trump's Twitter followers (see section 2.6). To construct these counts, we obtained the IDs of all people who follow Trump on Twitter. The figure shows that the majority of the tweets using these hashtags indeed come from people that also follow Trump.

These results lend credence to the idea that Trump's tweets are trigger points for anti-Muslim sentiment among his followers and that many people who harbor anti-Muslim sentiments self-select into following Donald Trump on Twitter, which exposes them to his tweets. The willingness of Trump's followers to produce their own anti-Muslim messages speaks to changes in the perceived acceptability of such content after a tweet by the president.

**Figure 9: Spillovers of Trump’s Tweets to His Followers**



**(c) Anti-Muslim Tweets**



*Notes:* Panel (a) and (b) plot the  $\beta_\tau$  coefficients from a dynamic version of equation 4 of the type  $Y_t = \alpha + \sum_{\tau=-4}^4 \beta_\tau \cdot Muslim\ Trump\ tweets_t + \mathbf{X}'_{t-\tau} + \epsilon_t$ . In Panel (a), the dependent variable is the number of retweets Donald Trump’s tweets about Muslims receive on a given day (in natural logarithm). In Panel (b), the dependent variable refers to tweets about Muslims by Trump’s followers that are not Trump retweets, and thus new content. All regressions include a full set of day of week and year-month dummies; and four lags of dummies for the incidence of terror attacks in the US and Europe. The sample period is the year 2017. The shaded areas are 95% confidence intervals based on Newey-West standard errors. Panel (c) plots the number of tweets containing the hashtags #StopIslam or #BanIslam between 2010 and 2017, which we interpret as clearly expressing negative sentiment towards Muslims. The orange bars show the number of these tweets posted by followers of Trump’s Twitter account.

### 5.3 Trump Tweets and the News Cycle

As a last time series exercise, we ask whether Trump’s tweets about Muslims affect the news cycle. This is important to understand because, unlike for the social media channel we study here, there is ample evidence that other types of media can persuade people to participate in spontaneous, potentially violent outbursts (see e.g. DellaVigna & Gentzkow, 2010; Yanagizawa-Drott, 2014). As such, one obvious channel through which social media may affect offline outcomes is through influencing what other media report on. Indeed, it has been widely recognized that Twitter has become an important dissemination channel for journalists (Willnat et al., 2019); some estimates suggest that up to a quarter of Twitter users may be working for media outlets (Haje Jan Kamps, 2015).

We investigate the effect of Trump’s tweets on media coverage using transcript data from the *TV News Archive*. In particular, we replace the dependent variable in equation (4) with the log number of mentions of Muslim-related topics on a given day by the three major cable news stations Fox News, CNN, and MSNBC. Columns 4 through 7 in Table 5 present the results of this exercise. Because we find a more immediate correlation between Trump’s Twitter activity and news coverage, we report specifications with  $h = 0$  as the prediction horizon.

Trump’s tweets about Muslims are highly correlated with TV mentions in the OLS, reduced form, and 2SLS regressions. For overall news coverage in column 4, we find that a one standard deviation increase in Muslim Trump tweets (0.25) is associated with a 96% increase in news coverage. The  $F$ -statistics are again almost uniformly above the rule-of-thumb of 10, and mostly above the 12.04 threshold for a maximum 30% coefficient bias with 5% statistical significance derived in Olea & Pflueger (2013). Perhaps more importantly, the Anderson-Rubin confidence sets always clearly exclude zero.

We also consider heterogeneity across news stations. The correlation of instrumented Trump tweets with TV mentions appears to be strongest for Fox News (see column 5). Indeed, for CNN and MSNBC (columns 6 and 7), a zero effect is well within the AR confidence sets. This is interesting because Fox News is well-known to be supportive of Trump, following a longer term move towards more Republican-friendly reporting (Martin & Yurukoglu, 2017). This might allow Trump’s comments to be broadcasted uncritically and even more widely through the channel’s considerable reach. Taken together, these patterns suggests that social media may allow influential individuals—such as the president of the United States—to drive the news cycle. Xenophobic rhetoric that is spread by the media largely unchallenged, in turn, may be one potential trigger point for potential perpetrators of hate crimes.

Table 4: Trump Tweets and Anti-Muslim Hate Crimes

	Baseline (1)	Add lagged dependent variable (2)	Add federal holiday control (3)	Add Google search control (4)	Add TV coverage control (5)	Add terror attack control (6)	Add total tweets control (7)
<b>Panel A: First stage - Log(Trump tweets about Muslims)</b>							
Trump golfs	0.102*** (0.027)	0.101*** (0.027)	0.104*** (0.027)	0.101*** (0.027)	0.090*** (0.026)	0.090*** (0.026)	0.098*** (0.027)
<b>Panel B: OLS - Log(Hate crimes against Muslims) in t+2</b>							
Log(Muslim Trump tweets)	0.109 (0.071)	0.112* (0.068)	0.110 (0.071)	0.100 (0.067)	0.095 (0.062)	0.160** (0.077)	0.091 (0.074)
<b>Panel C: Reduced form - Log(Hate crimes against Muslims) in t+2</b>							
Trump golfs	0.164** (0.069)	0.169** (0.073)	0.159** (0.068)	0.159** (0.070)	0.154** (0.070)	0.168** (0.070)	0.162** (0.069)
<b>Panel D: 2SLS - Log(Hate crimes against Muslims) in t+2</b>							
Log(Muslim Trump tweets)	1.609** (0.791)	1.665** (0.819)	1.534** (0.755)	1.579** (0.799)	1.712* (0.926)	1.859* (0.966)	1.648* (0.852)
Weak IV 95% AR confidence set	[0.278; 40.036]	[0.287; 40.016]	[0.263; 30.850]	[0.234; 40.031]	[0.338; 40.737]	[0.425; 50.014]	[0.384; 40.430]
Fixed effects (month, day of week)	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Time trend	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Observations	363	363	363	362	362	363	363
Robust F-stat.	13.15	13.46	13.55	13.03	11.25	11.63	12.07

Notes: This table presents OLS and IV regressions where the dependent variable is the number of hate crimes against Muslims on any given day based on FBI data. We use a dummy for days on which President Donald Trump golfs used as an instrument for his tweets about Muslims. Column 2 controls for one lag of the dependent variable and column 3 for a dummy that tags federal holidays. Column 4 controls for the first principal component of Google searches for Islam-related terms. Column 5 controls for the number of times Fox News, CNN or MSNBC mention Islam-related words in their reporting on a given day. Column 6 controls for the number of terror attacks in the US, Europe, or other countries. Column 7 controls for the total number of tweets by Donald Trump. The sample year is 2017, for which we have information on Trump's golfing. All regressions include day-of-week and year-month dummies, linear and quadratic time trends as well as a dummy for whether Trump's golfing is the first of a series of golf days. See online appendix for more details on data and variable construction. Newey-West standard errors are reported in parentheses. Weak IV 95% Anderson-Rubin (AR) confidence sets are calculated using the two-step approach of Andrews (2018) with the Stata package from Sun (2018). \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

Table 5: Spillover Effects on Trump’s Followers and Cable News Coverage

	Trump followers’ Muslim tweets			Cable news coverage			
	Trump retweets (1)	New content (2)	#StopIslam or #BanIslam (3)	All stations (4)	Fox News (5)	CNN (6)	MSNBC (7)
<b>Panel A: OLS - Log(Total number of Muslim TV mentions/tweets)</b>							
Log(Muslim Trump tweets)	2.658*** (0.346)	0.680*** (0.105)	0.610*** (0.129)	0.677*** (0.089)	0.607*** (0.117)	0.808*** (0.109)	0.660*** (0.084)
<b>Panel B: Reduced form - Log(Total number of Muslim TV mentions/tweets)</b>							
Trump golfs	0.456** (0.201)	0.117** (0.056)	0.228*** (0.081)	0.273** (0.130)	0.296*** (0.111)	0.285 (0.205)	0.185* (0.106)
<b>Panel C: 2SLS - Log(Total number of Muslim TV mentions/tweets)</b>							
Log(Muslim Trump tweets)	4.508*** (1.305)	1.151** (0.469)	2.250** (0.993)	2.701** (1.114)	2.923*** (0.966)	2.813 (1.891)	1.830** (0.921)
Weak IV 95% AR confidence set	[10.020; 60.962]	[0.177; 20.219]	[0.776; 50.493]	[0.385; 50.237]	[10.107; 50.313]	[-10.493; 70.119]	[-0.267; 30.927]
Fixed effects (month, day of week)	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Time trends	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Observations	364	364	364	364	364	364	364
Robust F-stat.	13.02	13.02	13.02	13.02	13.02	13.02	13.02

Notes: This table presents OLS and IV regressions where the dependent variable is the number of tweets by Trump followers in columns 1 to 3 and the number of times Muslims are mentioned on cable news stations on a given day in columns 4 to 7. We use a dummy for days on which President Donald Trump golfs used as an instrument for his tweets about Muslims. *Trump retweets* are retweets by Trump followers of Trump’s negative tweets about Muslims. *New content* refers to tweets by Trump followers mentioning Muslims that are not Trump retweets and do not mention Trump. *#StopIslam* or *#BanIslam* is the number of tweets by Trump followers containing the hashtags *#StopIslam* or *#BanIslam*. *Cable news coverage* is based on the mentions of Muslim-related words on Fox News, CNN, and MSNBC, which are also reported separately. The sample year is 2017, for which we have information on Trump’s golfing. All regressions include day-of-week and year-month dummies, linear and quadratic time trends as well as a dummy for whether Trump’s golfing is the first of a series of golf days. Newey-West standard errors are reported in parentheses. Weak IV 95% Anderson-Rubin (AR) confidence sets are calculated using the two-step approach of Andrews (2018) with the Stata package from Sun (2018). \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

## 5.4 Panel Evidence: Trump’s Tweets and Twitter Usage

As the last part of our analysis, we combine the cross-sectional and time series evidence. If Trump’s anti-Muslim rhetoric spreads through Twitter, we should observe larger increases in anti-Muslim hate crime in counties with higher Twitter usage. We investigate this hypothesis using the following panel regression specification:

$$\begin{aligned} \text{Hate Crimes}_{it} = & \beta \cdot \text{Twitter Usage}_i \times \text{Muslim Trump Tweets}_t \\ & + \mathbf{X}'_{it}\gamma + \text{County FE} + \text{Day FE} + \epsilon_{it} \end{aligned} \tag{6}$$

where  $\text{Hate Crimes}_{it}$  is an indicator variable for a hate crime in county  $i$  on day  $t$ . The main coefficient of interest  $\beta$  is the interaction of county-level Twitter usage with Trump’s tweets about Muslims. We standardize the independent variables to have mean 0 and standard deviation 1. The coefficient measures if there are disproportionate changes in anti-Muslim hate crimes in counties with high Twitter usage on days Trump tweets about Muslims. The specification additionally controls for a vector of control variables  $\mathbf{X}_{it}$  and includes a full set of county and day fixed effects. We cluster standard errors at the state level.

The setup in equation 6 is akin in spirit to a shift-share design, where  $\text{Twitter Usage}$  measures the local exposure to aggregate shocks  $\text{Muslim Trump Tweets}$ . Because we are interested in estimating the effect of social media, the main concern with this empirical strategy is that the local exposure measure is co-determined with unobserved factors that may also lead to changes in outcomes when Trump tweets (Goldsmith-Pinkham et al., 2017). Apart from estimating equation 6 using OLS, we thus also present results based on 2SLS, where we again instrument for local Twitter usage using temporal fluctuations in when users started following SXSW around the 2007 festival.

We first investigate the timing of Trump’s tweets and hate crime. To do so, we include interactions of local Twitter usage with leads and lags of Trump’s tweets about Muslims. Figure A.9 in the online appendix indicates that we observe differential increases in anti-Muslim hate crime in counties with high Twitter usage one day after Donald Trump’s tweets. Next, we test whether this finding is robust to the inclusion of additional fixed effects and compare the importance of Twitter usage relative to other cross-sectional predictors. In particular we analyze if exposure to Fox News or ideological alignment with Trump (measured by a high Republican vote share) are additional factors.<sup>25</sup>

The results of these exercises can be found in Table 6. The interaction of Trump tweets and social media usage robustly predicts hate crimes on the following day. The magnitude

---

<sup>25</sup>Note that we focus on additional cross-sectional exposure variables because we are interested in the effect of social media per se. As we show above, measures of anti-Muslim sentiment (e.g. Fox News reports) are at least partially *outcomes* of Trump’s tweets.

of the main coefficients remains quantitatively unchanged even when we include state  $\times$  day, county  $\times$  day of week, and county  $\times$  day of month fixed effects in columns 1 to 3. In the following two columns, we show that the inclusion of Fox News exposure and the Republican vote share—both of which we interact with Trump’s tweets—have less robust and quantitatively smaller predictive power for increases in anti-Muslim hate crime. Overall, these findings are again in line with the hypothesis that, when triggered by a shock such as Trump’s tweets about Muslims, social media may contribute to hate crimes against minorities.

**Table 6: Panel Regression Results**

	(1)	(2)	(3)	(4)	(5)
<b>Panel A: OLS</b>					
Muslim Trump tweets $\times$ Twitter usage	0.029** (0.011)	0.028** (0.011)	0.031*** (0.011)	0.033*** (0.012)	0.031*** (0.011)
Muslim Trump tweets $\times$ Fox News viewership				0.003** (0.001)	
Muslim Trump tweets $\times$ Republican vote share 2012					-0.000 (0.001)
<b>Panel B: Reduced form</b>					
Muslim Trump tweets $\times$ Log(SXSW followers, March 2007)	0.013** (0.005)	0.011* (0.006)	0.012** (0.006)	0.013** (0.006)	0.012** (0.006)
Muslim Trump tweets $\times$ Fox News viewership				0.002* (0.001)	
Muslim Trump tweets $\times$ Republican vote share 2012					-0.001 (0.001)
<b>Panel C: 2SLS</b>					
Muslim Trump tweets $\times$ Twitter usage	0.117** (0.044)	0.099** (0.048)	0.112** (0.049)	0.124** (0.054)	0.123** (0.055)
Muslim Trump tweets $\times$ Log(SXSW followers, Pre)	-0.001 (0.008)	-0.001 (0.008)	-0.002 (0.008)	-0.002 (0.008)	-0.002 (0.008)
Muslim Trump tweets $\times$ Fox News viewership				0.011** (0.004)	
Muslim Trump tweets $\times$ Republican vote share 2012					0.009* (0.005)
County FE	Yes	Yes	Yes	Yes	Yes
Day FE	Yes	Yes	Yes	Yes	Yes
Pop. deciles $\times$ Date FE	Yes	Yes	Yes	Yes	Yes
County $\times$ Month FE		Yes	Yes	Yes	Yes
State $\times$ Day FE		Yes	Yes	Yes	Yes
County $\times$ Day of week FE			Yes	Yes	Yes
County $\times$ Day of month FE			Yes	Yes	Yes
Observations	2,887,332	2,886,403	2,886,403	2,885,474	2,886,403
$R^2$	0.01	0.08	0.12	0.12	0.12

*Notes:* This table presents OLS, reduced form and IV regressions where the dependent variable is an indicator of anti-Muslim hate crimes in county  $i$  on day  $t$ . The coefficients are multiplied by 100 for readability. In Panel A, the independent variable is the interaction of Trump’s negative tweets about Muslims with county-level Twitter usage. In Panel B, the interaction is with SXSW followers who signed up in March 2007, while controlling for the interaction with users who joined before the festival (omitted for brevity). Panel C shows interactions where Twitter usage is instrumented with SXSW followers who joined in March 2007. The variables are standardized to have a mean of zero and standard deviation of one. All regressions include population controls, as well as county and date fixed effects. Some regressions include county  $\times$  month, state  $\times$  day, county  $\times$  day-of-week, or county  $\times$  day-of-month fixed effects (as indicated). Robust standard errors in parentheses are clustered by state. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .



## 6 Mechanisms

The evidence provided in the previous sections supports the idea that social media has played a role in the spread of anti-Muslim sentiment since the 2016 presidential primaries. We discuss three potential mechanisms that could drive these results: changes in the ability of perpetrators of hate crimes to coordinate (a *coordination channel*), changes in people's beliefs about minorities (a *persuasion channel*), and changes in the perceived societal acceptance or penalty with regard to xenophobic actions and beliefs (a *social norms channel*). We argue that our overall findings are most consistent with the idea that social media can enable changes in social norms.

To begin, our findings are unlikely to be driven by lower coordination costs due to social media. The vast majority of hate crimes recorded by the FBI are committed by a single perpetrator and the effects we document appear to be exclusively driven by these incidents (see Table A.21 in the online appendix).<sup>26</sup> Further, it is not obvious why the 2016 presidential campaign period or Trump's tweets would sharply and suddenly improve potential perpetrators' coordination capabilities specifically for hate crimes against minorities frequently targeted by Trump. Additionally, most content on Twitter is entirely public. As such, Twitter is not the most likely place for plotting violent crimes, but rather a place to spread ideas.

Another hypothesis is that our findings are driven by the persuasiveness of Twitter content (see DellaVigna & Gentzkow, 2010, for a review of the literature on persuasion), which could make people more xenophobic. However, several pieces of evidence are not easily rationalized by belief-based persuasion models. For one, the idea that Twitter makes people more xenophobic is strongly at odds with evidence from survey data. Figure A.10 and Figure A.11 in the online appendix show, across a range of questions in the American National Election Studies (ANES), Americans' opinions toward immigrants, and Muslims in particular, show little change between 2012 and 2016. In fact, people using the internet or social media have clearly become more tolerant towards Muslims, Hispanics, and illegal aliens, and less likely to agree that immigration should be reduced. These patterns are in line with Hopkins & Washington (2019), who show that white Americans' anti-minority prejudice has, if anything, declined since Trump's political rise; Pew Research Center (2017) show that Americans report considerably warmer feelings towards Muslims in 2017 compared to 2014.

We provide additional evidence suggesting no increase xenophobic attitudes using the results from implicit association tests (IAT) from Project Implicit. We measure changes in implicit bias against Muslims, which are based on the difference in an individual's ability to assign positive or negative words to Muslims and non-Muslims. We follow Chetty et al. (2018)

---

<sup>26</sup>We have information on the number of offenders in around 62% of incidents, which somewhat decreases statistical power.

and calculate mean IAT scores on the county-level.<sup>27</sup> While IAT scores are not a perfect measure of bias, previous research has shown that IAT scores have predictive power for key outcomes: Carlana (2019), for example, shows that implicit gender bias affects the gender gap in students' math performance; Chetty et al. (2018) show that, across counties, higher implicit bias against African Americans correlates with larger black-white wage gaps.

Figure A.12a suggests, consistent with the survey evidence outlined above, a clear decrease in implicit bias towards Muslims between 2010-2014 (before Trump's political rise) and 2017 (afterwards), as shown by a leftward shift in the IAT score distribution. However, it could still be that such biases increased in relative terms in counties with higher social media penetration. To test this hypothesis, we run regressions of the type in equation 2, where the dependent variable is now the change in a measure of implicit bias against Muslims around Trump's presidential campaign start. Table A.29 reports the results. We consider a range of specifications and sub-samples, including test scores restricted to whites or conservatives, and find no evidence that social media increases implicit bias against Muslims. The estimates suggest that we can reject even small increases in implicit bias due to social media. For example, the statistically insignificant 2SLS coefficient of 0.006 in column 1 implies that a one standard deviation increase in Twitter usage shifts IAT scores by 2.9% of a standard deviation ( $(0.006 \times 1.57)/0.32 \approx 0.029$ ). The reduced-form 95% confidence interval of  $[-0.021; 0.027]$  suggests that we can plausibly rule out (reduced-form) effects larger than 3% of a standard deviation in IAT scores.<sup>28</sup> This conclusion is also supported by the event study pattern in Figure A.12b.

Furthermore, Bayesian models of persuasion (e.g. Kamenica & Gentzkow, 2011) would suggest that people with weaker priors adjust their attitudes more strongly. In contrast, we find that the effects of Twitter usage are driven by areas with *higher*, not lower pre-existing prejudice. To show this, we repeat the event study regressions from Section 3.1 and split counties by whether the Southern Poverty Law Center (SPLC) identifies at least one hate group.<sup>29</sup> Figure 10 plots the estimated coefficients from this exercise.<sup>30</sup> We find that higher Twitter usage is only associated with more anti-Muslim hate crime in counties with hate

---

<sup>27</sup>Participation in the IAT is online and largely voluntary, which may give rise to selection bias. While we cannot fully rule out such biases, we also consider a measure of implicit bias based on individuals who were obligated to take these tests, e.g. as part of a work program, and find similar results.

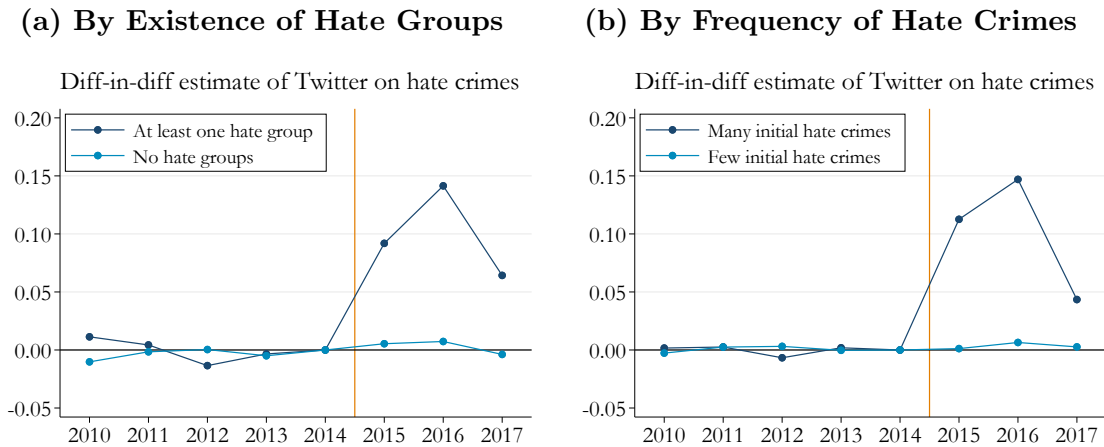
<sup>28</sup>To see this, consider that the standard deviation of *Log(SXSW followers, March 2007)* in this sample is around 0.378. The standard deviation of the change in IAT scores is 0.323. The largest effect of a one standard deviation increase in SXSW followers in the confidence set is thus  $(0.027 \times 0.387)/0.323 \approx 0.03$ . We focus on the reduced form because 2SLS is well-known to be consistent but inefficient.

<sup>29</sup>Note that these sample splits do not estimate whether anti-Muslim hate crimes increased in counties with hate groups; rather, they address the question whether Twitter usage has a different impact in these counties.

<sup>30</sup>To reduce clutter, the figures report the estimated coefficients without confidence bands. We report the full regression results with standard errors in Table A.30 in the online appendix.

groups. In contrast, counties with high Twitter usage but no hate group continue to follow the same trajectory as low Twitter usage counties. In Panel (b), we provide similar evidence for counties that are above the 90th percentile of hate crime per capita in the pre-period. We again observe that the increase in anti-Muslim hate crimes is driven by counties with high Twitter usage and pre-existing biases.

**Figure 10: Heterogenous Effects of Twitter Usage**



*Notes:* These figures plot the coefficients of running panel event study regressions as in Equation (1). Hate crimes and Twitter usage are standardized to have a mean of zero and standard deviation of one. In Panel (a), Equation (1) is estimated separately for counties with and without at least one hate group as defined by the Southern Poverty Law Center (SPLC). In Panel (b), we split counties at the 90th percentile of the average number of hate crimes per capita between 2010 and 2014.

These findings provide suggestive evidence that, at least in our setting, social media does not necessarily change people’s beliefs, but rather triggers real-life action by individuals with existing negative attitudes towards Muslims. This pattern is consistent with the hypothesis that social media facilitated perceived shifts in social norms among people who already harbor extreme viewpoints. In other words, people may infer information about the social acceptability of viewpoints and actions based on what they see online. After observing increased anti-Muslim rhetoric on Twitter, already radicalized individuals may become more willing to commit violent acts against Muslims in real life.

This hypothesis could explain why we observe an effect of social media on public actions (hate crimes and expressed xenophobia on Twitter) but no effect on private beliefs (survey replies and implicit biases). The channel we have in mind is the following. A key feature of social norms is that they are based on people’s *perceptions* of everyone else’s beliefs. These perceptions, in turn, are shaped by the “sample” of beliefs that are most salient to an individual (e.g. Bursztyyn & Jensen, 2015; Perez-Truglia & Cruces, 2017; Enikolopov et al., 2017). However, people are systematically wrong in their perception of what others believe,

particularly when it comes to political topics (e.g. Westfall et al., 2015; Bordalo et al., 2016; Bursztyn et al., 2018).<sup>31</sup>

By enabling relatively few but particularly visible individuals to affect the aggregate discourse, social media could shift beliefs about what is socially acceptable and make people more susceptible to extreme viewpoints.<sup>32</sup> Such effects could be re-enforced by what has often been called “echo chambers” (e.g. Bessi et al., 2015; Del Vicario et al., 2016; Schmidt et al., 2017; Sunstein, 2017). This, in turn, could affect the willingness of a small set of potential perpetrators to take hateful actions online or offline.

This interpretation is in line with the findings of Bursztyn et al. (2017), who show in a range of experiments that Donald Trump’s 2016 election victory increased people’s willingness to publicly express xenophobic views, as well as the tolerance towards expressing such views publicly. It also meshes well with survey evidence suggesting that most Americans believe it has become more acceptable for people to express racist views (Pew Research Center, 2019a). Our findings suggest that social media may contribute to such an unraveling of social norms in a way that affects real-life behavior.

## 7 Conclusion

Social media has come under scrutiny for its oft-alleged potential to increase citizen polarization by creating informational “echo chambers” (Sunstein, 2009, 2017). However, empirical evidence on the real-world effects of social media are limited. Our work suggests that social media usage can enable increases in anti-minority sentiments, particularly when used by powerful individuals such as the president of the United States.

While this paper focused on particularly negative outcomes—hate crimes targeting minorities and other measures of xenophobia— social media may well have a positive impact in other areas. We would also like to caution against using our findings as a basis for policies directed at restricting online communication. History is ripe with cautionary tales of how excessive state power over the media can abet authoritarian rule. The complex trade-offs that policy makers face in this regard thus require nuanced discussion and, above all, more evidence. Notwithstanding, our results suggest that social media can affect offline actions that might endanger minority communities, and should be taken seriously.

---

<sup>31</sup>See Bénabou (2008) for a model of how belief distortions can give rise to a persistence of ideologies in equilibrium; Bénabou (2013) studies “groupthink” more broadly.

<sup>32</sup>This is related to Ali & Bénabou (2016), where the visibility of individuals makes aggregate behavior (*descriptive* norms) less informative about societal preferences (*prescriptive* norms). It is also related to Mukand & Rodrik (2018), where “political entrepreneurs” can change individuals’ perception of who they are, by increasing the salience of particular parts of their identity (e.g. a “true American”). Matz et al. (2017) provide evidence for the effectiveness of social media targeting based on psychological traits.

## References

- Adena, M., Enikolopov, R., Petrova, M., Santarosa, V., & Zhuravskaya, E. (2015). Radio and the Rise of The Nazis in Prewar Germany. *The Quarterly Journal of Economics*, 130(4), 1885–1939.
- Alatas, V., Chandrasekhar, A. G., Mobius, M., Olken, B. A., & Paladines, C. (2019). When Celebrities Speak: A Nationwide Twitter Experiment Promoting Vaccination In Indonesia. Working Paper 25589, National Bureau of Economic Research.
- Ali, S. N. & Bénabou, R. (2016). Image versus information: Changing societal norms and optimal privacy. Working Paper 22203, National Bureau of Economic Research.
- Allcott, H., Braghieri, L., Eichmeyer, S., & Gentzkow, M. (2020). The Welfare Effects of Social Media. *American Economic Review*, 110(3), 629–76.
- Anderson, T. W., Rubin, H., et al. (1949). Estimation of the Parameters of a Single Equation in a Complete System of Stochastic Equations. *The Annals of Mathematical Statistics*, 20(1), 46–63.
- Andrews, I. (2018). Valid Two-Step Identification-Robust Confidence Sets for GMM. *Review of Economics and Statistics*, 100(2), 337–348.
- Andrews, I., Stock, J. H., & Sun, L. (2019). Weak Instruments in IV Regression: Theory and Practice. *Annual Review of Economics*.
- Arrow, K. J. (2000). Increasing Returns: Historiographic Issues and Path Dependence. *The European Journal of the History of Economic Thought*, 7(2), 171–180.
- Arthur, W. B. (1989). Competing Technologies, Increasing Returns, and Lock-in by Historical Events. *The Economic Journal*, 99(394), 116–131.
- Arthur, W. B. (1994). *Increasing Returns and Path Dependence in the Economy*. University of Michigan Press.
- Bail, C. A., Argyle, L. P., Brown, T. W., Bumpus, J. P., Chen, H., Hunzaker, M. B. F., Lee, J., Mann, M., Merhout, F., & Volfovsky, A. (2018). Exposure to Opposing Views on Social Media Can Increase Political Polarization. *Proceedings of the National Academy of Sciences*, 115(37), 9216–9221.
- Barberá, P. (2014). How Social Media Reduces Mass Political Polarization. Evidence from Germany, Spain, and the US.

- Bass, F. M. (1969). A new product growth for model consumer durables. *Management Science*, 15(5), 215–227.
- BBC (2019). Ilhan Omar: Muslim Lawmaker Sees Rise in Death Threats After Trump Tweet.
- Beaman, L., Chattopadhyay, R., Duflo, E., Pande, R., & Topalova, P. (2009). Powerful Women: Does Exposure Reduce Bias? *The Quarterly Journal of Economics*, 124(4), 1497–1540.
- Bénabou, R. (2008). Ideology. *Journal of the European Economic Association*, 6(2-3), 321–352.
- Bénabou, R. (2013). Groupthink: Collective Delusions in Organizations and Markets. *The Review of Economic Studies*, 80(2 (283)), 429–462.
- Bessi, A., Zollo, F., Vicario, M. D., Scala, A., Caldarelli, G., & Quattrociocchi, W. (2015). Trend of Narratives in the Age of Misinformation. *PLOS ONE*, 10(8), 1–16.
- Bhuller, M., Havnes, T., Leuven, E., & Mogstad, M. (2013). Broadband Internet: An Information Superhighway to Sex Crime? *Review of Economic Studies*, 80(4), 1237–1266.
- Bond, R. M., Fariss, C. J., Jones, J. J., Kramer, A. D., Marlow, C., Settle, J. E., & Fowler, J. H. (2012). A 61-Million-Person Experiment in Social Influence and Political Mobilization. *Nature*, 489(7415), 295.
- Bordalo, P., Coffman, K., Gennaioli, N., & Shleifer, A. (2016). Stereotypes. *The Quarterly Journal of Economics*, 131(4), 1753–1794.
- Boxell, L., Gentzkow, M., & Shapiro, J. M. (2017). Greater Internet Use Is Not Associated with Faster Growth in Political Polarization Among US Demographic Groups. *Proceedings of the National Academy of Sciences of the United States of America*, 201706588.
- Brown, B. (2018). The Trump Twitter Archive. <http://www.trumptwitterarchive.com/> (accessed November 2nd, 2018).
- Bursztyn, L., Egorov, G., Enikolopov, R., & Petrova, M. (2019). Social Media and Xenophobia: Evidence from Russia. Working Paper 26567, National Bureau of Economic Research.
- Bursztyn, L., Egorov, G., & Fiorin, S. (2017). From Extreme to Mainstream: How Social Norms Unravel. Working Paper 23415, National Bureau of Economic Research.
- Bursztyn, L., Gonzlez, A. L., & Yanagizawa-Drott, D. (2018). Misperceived Social Norms: Female Labor Force Participation in Saudi Arabia. NBER Working Papers 24736, National Bureau of Economic Research, Inc.

- Bursztyn, L. & Jensen, R. (2015). How does peer pressure affect educational investments? *The Quarterly Journal of Economics*, 130(3), 1329–1367.
- Card, D. & Dahl, G. B. (2011). Family Violence and Football: The Effect of Unexpected Emotional Cues on Violent Behavior. *The Quarterly Journal of Economics*, 126(1), 103–143.
- Carlana, M. (2019). Implicit stereotypes: Evidence from teachers gender bias. *The Quarterly Journal of Economics*, 134(3), 1163–1224.
- Chan, J., Ghose, A., & Seamans, R. (2016). The Internet and Racial Hate Crime: Offline Spillovers from Online Access. *MIS Quarterly*, 40(2), 381–403.
- Chen, Y. & Yang, D. Y. (2019). The Impact of Media Censorship: 1984 or Brave New World? *American Economic Review*, 109(6), 2294–2332.
- Chetty, R., Hendren, N., Jones, M. R., & Porter, S. R. (2018). Race and economic opportunity in the united states: An intergenerational perspective. Working Paper 24441, National Bureau of Economic Research.
- CNN (2020). Trump Signs Executive Order Targeting Social Media Companies.
- Colella, F., Lalive, R., Sakalli, S. O., & Thoenig, M. (2019). Inference with Arbitrary Clustering. IZA Discussion Papers 12584, Institute of Labor Economics (IZA).
- Dahl, G. & DellaVigna, S. (2009). Does Movie Violence Increase Violent Crime? *The Quarterly Journal of Economics*, 677–734.
- Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., & Harshman, R. (1990). Indexing by latent semantic analysis. *Journal of the American society for information science*, 41(6), 391.
- Del Vicario, M., Bessi, A., Zollo, F., Petroni, F., Scala, A., Caldarelli, G., Stanley, H. E., & Quattrocioni, W. (2016). The Spreading of Misinformation Online. *Proceedings of the National Academy of Sciences*, 113(3), 554–559.
- DellaVigna, S., Enikolopov, R., Mironova, V., Petrova, M., & Zhuravskaya, E. (2014). Cross-Border Media and Nationalism: Evidence from Serbian Radio in Croatia. *American Economic Journal: Applied Economics*, 6(3), 103–32.
- DellaVigna, S. & Gentzkow, M. (2010). Persuasion: Empirical Evidence. *Annual Review of Economics*, 2(1), 643–669.
- Draca, M. & Schwarz, C. (2018). How Polarized Are Citizens? Measuring Ideology from the Ground-up.

- Edwards, B. T. (2018). Trump from Reality TV to Twitter, or the Selfie-Determination of Nations. *Arizona Quarterly: A Journal of American Literature, Culture, and Theory*, 74(3), 25–45.
- Enikolopov, R., Makarin, A., & Petrova, M. (2016). Social Media and Protest Participation: Evidence from Russia.
- Enikolopov, R., Makarin, A., Petrova, M., & Polishchuk, L. (2017). Social Image, Networks, and Protest Participation. *Universitat Pompeu Fabra*.
- Fagerberg, J., Mowery, D. C., & Hall, B. H. (2009). Innovation and Diffusion.
- FBI (2015). Hate Crime Data Collection Guidelines And Training Manual. *Criminal Justice Information Services (CJIS) Division Uniform Crime Reporting (UCR) Program*.
- Financial Times (2020). Facebook Takes Down Trump Ads for Violating Organised Hate Policy.
- Fiorina, M. P. & Abrams, S. J. (2008). Political Polarization in the American Public. *Annual Review of Political Science*, 11, 563–588.
- Gavazza, A., Nardotto, M., & Valletti, T. M. (2015). Internet and Politics: Evidence from UK Local Elections and Local Government Policies.
- Gawker (2007). Twitter Blows Up at SXSW Conference. <https://gawker.com/243634/twitter-blows-up-at-sxsw-conference> (accessed March 3rd, 2018).
- Gentzkow, M. (2016). Polarization in 2016. *Toulouse Network of Information Technology white paper*.
- Geroski, P. (2000). Models of Technology Diffusion. *Research Policy*, 29(4), 603 – 625.
- Goldsmith-Pinkham, P., Sorkin, I., & Swift, H. (2017). Bartik Instruments: What, When, Why, and How. *Working Paper*.
- Griliches, Z. (1957). Hybrid Corn: An Exploration in the Economics of Technological Change. *Econometrica*, 25(4), 501–522.
- Haje Jan Kamps (2015). Who Are Twitters Verified Users?
- Haustein, S. & Costas, R. (2014). Determining Twitter Audiences: Geolocation and Number of Followers. *ALM*, 4, 6.
- Hobbs, W. & Lajevardi, N. (2019). Effects of Divisive Political Campaigns on the Day-to-Day Segregation of Arab and Muslim Americans. *American Political Science Review*, 113(1), 270–276.



- Hopkins, D. J. & Washington, S. (2019). The Rise of Trump, the Fall of Prejudice? Tracking White Americans' Racial Attitudes 2008-2018 via a Panel Survey. *Working Paper*.
- Iaria, A., Schwarz, C., & Waldinger, F. (2018). Frontier knowledge and scientific production: evidence from the collapse of international science. *The Quarterly Journal of Economics*, *133*(2), 927–991.
- Jones, J. J., Bond, R. M., Bakshy, E., Eckles, D., & Fowler, J. H. (2017). Social Influence and Political Mobilization: Further Evidence From a Randomized Experiment in the 2012 U.S. Presidential Election. *PLOS ONE*, *12*(4), 1–9.
- Kamenica, E. & Gentzkow, M. (2011). Bayesian persuasion. *American Economic Review*, *101*(6), 2590–2615.
- Kinder-Kurlanda, K., Weller, K., Zenk-Möltgen, W., Pfeffer, J., & Morstatter, F. (2017). Archiving Information from Geotagged Tweets to Promote Reproducibility and Comparability in Social Media Research. *Big Data & Society*, *4*(2), 2053951717736336.
- Landauer, T. K. (2007). *Handbook of latent semantic analysis*. Mahwah, N.J.: Lawrence Erlbaum Associates.
- Levy, R. (2019). Social Media, News Consumption, and Polarization: Evidence from a Field Experiment. *Working Paper*.
- Liebowitz, S. J. & Margolis, S. E. (1999). Path Dependence. *Encyclopedia of law and economics*.
- Manacorda, M. & Tesei, A. (2020). Liberation Technology: Mobile Phones and Political Mobilization in Africa. *Econometrica*, *88*(2), 533–567.
- Martin, G. J. & Yurukoglu, A. (2017). Bias in Cable News: Persuasion and Polarization. *American Economic Review*, *107*(9), 2565–2599.
- Matz, S. C., Kosinski, M., Nave, G., & Stillwell, D. J. (2017). Psychological Targeting as an Effective Approach to Digital Mass Persuasion. *Proceedings of the National Academy of Sciences*, *114*(48), 12714–12719.
- Miller, C. & Smith, J. (2017). Anti-Islamic Content on Twitter. *Centre for the Analysis of Social Media at Demos*.
- Mosquera, R., Odunowo, M., McNamara, T., Guo, X., & Petrie, R. (2020). The Economic Effects of Facebook. *Experimental Economics*, *23*(2), 575–602.
- Mukand, S. & Rodrik, D. (2018). The Political Economy of Ideas: On Ideas Versus Interests in Policymaking. Working Paper 24467, National Bureau of Economic Research.

- Müller, K. & Schwarz, C. (2018). Fanning the Flames of Hate: Social Media and Hate Crime. *Working Paper*.
- NBC News (2017). Advocates Warn of Possible Underreporting in FBI Hate Crime Data, by Chris Fuchs. <https://www.nbcnews.com/news/asian-america/advocates-warn-possible-underreporting-fbi-hate-crime-data-n830711> (accessed March 3rd, 2018).
- New York Times (2017). Trump Shares Inflammatory Anti-Muslim Videos, and Britains Leader Condemns Them, By Peter Baker and Eeileen Sullivan.
- New York Times (2018). The Man Behind the President’s Tweets.
- New York Times (2019a). Free Speech Is Killing Us, By Andrew Marantz.
- New York Times (2019b). Tracking Trump’s Visits to His Branded Properties.
- Olea, J. L. M. & Pflueger, C. (2013). A Robust Test for Weak Instruments. *Journal of Business & Economic Statistics*, 31(3), 358–369.
- Perez-Truglia, R. & Cruces, G. (2017). Partisan Interactions: Evidence from a Field Experiment in the United States. *Journal of Political Economy*, 125(4), 1208–1243.
- Petrova, M., Sen, A., & Yildirim, P. (2017). Social Media and Political Donations: New Technology and Incumbency Advantage in the United States. *Working Paper*.
- Pew Research Center (2017). U.S. Muslims Concerned About Their Place in Society, but Continue to Believe in the American Dream. <https://www.pewforum.org/2017/07/26/findings-from-pew-research-centers-2017-survey-of-us-muslims/>.
- Pew Research Center (2019a). Race in America 2019. <https://www.pewsocialtrends.org/2019/04/09/race-in-america-2019/>.
- Pew Research Center (2019b). Sizing Up Twitter Users. Technical report.
- Pew Research Center (2020). Democrats On Twitter More Liberal, Less Focused On Compromise Than Those Not On The Platform. Technical report.
- ProPublica (2017). Why America Fails at Gathering Hate Crime Statistics, by Ken Schwencke. <https://www.propublica.org/article/why-america-fails-at-gathering-hate-crime-statistics> (accessed March 3rd, 2018).
- Reilly, R. (2019). *Commander in Cheat: How Golf Explains Trump*. Hachette Books.
- Rogers, E. M. (2010). *Diffusion of Innovations*. Simon and Schuster.

- Schmidt, A. L., Zollo, F., Del Vicario, M., Bessi, A., Scala, A., Caldarelli, G., Stanley, H. E., & Quattrocioni, W. (2017). Anatomy of News Consumption on Facebook. *Proceedings of the National Academy of Sciences*, *114*(12), 3035–3039.
- Schwarz, C. (2019). lsemantic: A command for text similarity based on latent semantic analysis. *The Stata Journal*, *19*(1), 129–142.
- Stephens-Davidowitz, S. (2014). The Cost of Racial Animus on a Black Candidate: Evidence using Google Search Data. *Journal of Public Economics*, *118*, 26–40.
- Stock, J. & Yogo, M. (2005). *Testing for Weak Instruments in Linear IV Regression*, (pp. 80–108). New York: Cambridge University Press.
- Sun, L. (2018). Implementing Valid Two-Step Identification-Robust Confidence Sets For Linear Instrumental-Variables Models. *The Stata Journal*, *18*(4), 803–825.
- Sunstein, C. R. (2002). The Law of Group Polarization. *Journal of Political Philosophy*, *10*(2), 175–195.
- Sunstein, C. R. (2009). *Republic.com 2.0*. Princeton University Press.
- Sunstein, C. R. (2017). *# Republic: Divided Democracy in the Age of Social Media*. Princeton University Press.
- Takhteyev, Y., Gruzd, A., & Wellman, B. (2012). Geography of Twitter Networks. *Social networks*, *34*(1), 73–81.
- Twitter (2010). Measuring Tweets.
- United Nations (2020). UN Expert Denounces the Propagation of Hate Speech Through Social Media.
- Wall Street Journal (2020). Inside Twitters Decision to Take Action on Trumps Tweets.
- Westfall, J., Boven, L. V., Chambers, J. R., & Judd, C. M. (2015). Perceiving Political Polarization in the United States: Party Identity Strength and Attitude Extremity Exacerbate the Perceived Partisan Divide. *Perspectives on Psychological Science*, *10*(2), 145–158. PMID: 25910386.
- Willnat, L., Weaver, D. H., & Wilhoit, G. C. (2019). The American Journalist in the Digital Age. *Journalism Studies*, *20*(3), 423–441.
- Yanagizawa-Drott, D. (2014). Propaganda and Conflict: Evidence from the Rwandan Genocide. *The Quarterly Journal of Economics*, *129*(4), 1947–1994.

# A Online Appendix

## A.1. Appendix 1: Additional Details on Data

### A.1.1 FBI Hate Crime Data

As described in the Section 2, the FBI uses a two-tier decision making process for classifying hate crimes. FBI (2015) describes the decision making process in the following way:

“Once the development of this collection was complete, the FBI UCR Program surveyed state UCR Program managers on hate crime collection procedures used at various law enforcement agencies which collected hate crime data employing a two-tier decision-making process. The first level is the law enforcement officer who initially responds to the alleged hate crime incident, i.e., the responding officer (or first-level judgment officer). It is the responsibility of the responding officer to determine whether there is any indication that the offender was motivated by bias. If a bias indicator is identified, the officer designates the incident as a suspected bias-motivated crime and forwards the case file to a second-level judgment officer/unit. (In smaller agencies this is usually a person specially trained in hate crime matters, while in larger agencies it may be a special unit.) It is the task of the second-level judgment officer/unit to review the facts of the incident and make the final determination of whether a hate crime has actually occurred. If so, the incident is to be reported to the FBI UCR Program as a bias-motivated crime.” (FBI, 2015, pp. 2-3)

As indicated, all decisions by the responding officer will be passed on for review to a second examiner. The FBI manual also outlines criteria that have to be full-filled for a crime to be classified as a hate crime:

“An important distinction must be made when reporting a hate crime. The mere fact the offender is biased against the victims actual or perceived race, religion, disability, sexual orientation, ethnicity, gender, and/or gender identity does not mean that a hate crime was involved. Rather, the offenders criminal act must have been motivated, in whole or in part, by his or her bias. Motivation is subjective, therefore, it is difficult to know with certainty whether a crime was the result of the offenders bias. For that reason, before an incident can be reported as a hate crime, sufficient objective facts must be present to lead a reasonable and prudent person to conclude that the offenders actions were motivated, in whole

or in part, by bias. While no single fact may be conclusive, facts such as the following, particularly when combined, are supportive of a finding of bias:

1. The offender and the victim were of a different race, religion, disability, sexual orientation, ethnicity, gender, and/or gender identity. For example, the victim was African American and the offender was white.
2. Bias-related oral comments, written statements, or gestures were made by the offender indicating his or her bias. For example, the offender shouted a racial epithet at the victim.
3. Bias-related drawings, markings, symbols, or graffiti were left at the crime scene. For example, a swastika was painted on the door of a synagogue, mosque, or LGBT center.
4. Certain objects, items, or things which indicate bias were used. For example, the offenders wore white sheets with hoods covering their faces or a burning cross was left in front of the victims residence.
5. The victim is a member of a specific group that is overwhelmingly outnumbered by other residents in the neighborhood where the victim lives and the incident took place.
6. The victim was visiting a neighborhood where previous hate crimes had been committed because of race, religion, disability, sexual orientation, ethnicity, gender, or gender identity and where tensions remained high against the victims group.
7. Several incidents occurred in the same locality, at or about the same time, and the victims were all of the same race, religion, disability, sexual orientation, ethnicity, gender, or gender identity.
8. A substantial portion of the community where the crime occurred perceived that the incident was motivated by bias.
9. The victim was engaged in activities related to his or her race, religion, disability, sexual orientation, ethnicity, gender, or gender identity. For example, the victim was a member of the National Association for the Advancement of Colored People (NAACP) or participated in an LGBT pride celebration.
10. The incident coincided with a holiday or a date of significance relating to a particular race, religion, disability, sexual orientation, ethnicity, gender,

or gender identity, e.g., Martin Luther King Day, Rosh Hashanah, or the Transgender Day of Remembrance.

11. The offender was previously involved in a similar hate crime or is a hate group member.
12. There were indications that a hate group was involved. For example, a hate group claimed responsibility for the crime or was active in the neighborhood.
13. A historically-established animosity existed between the victims and the offenders groups.
14. The victim, although not a member of the targeted racial, religious, disability, sexual orientation, ethnicity, gender, or gender identity group, was a member of an advocacy group supporting the victim group.”

(FBI, 2015, pp. 6-7)

We report the full list of FBI bias motivation categories in Table A.2. The hate crime categories we use in the paper are defined as follows:

**Table A.1: FBI Hate Crimes Codes**

<b>Hate Crime Category</b>	<b>FBI Codes</b>
Muslim	24
Hispanic	32
Other ethnic	33
Racial	11, 12, 13, 14, 15, 16
Sexual orientation	41, 42, 43, 44, 45
Religious (excluding Muslim)	21, 22, 23, 25, 26, 27, 28, 29, 81, 82, 83, 84, 85

**Table A.2: Full List of FBI Bias Motivation Categories**

<b>Bias category</b>	<b>Bias motivation and code</b>
<b>Race/Ethnicity/Ancestry</b>	Anti-American Indian or Alaska Native (13)
	Anti-Arab (31)
	Anti-Asian (14)
	Anti-Black or African American (12)
	Anti-Hispanic or Latino (32)
	Anti-Multiple Races, Group (15)
	Anti-Native Hawaiian or Other Pacific Islander (16)
	Anti-Other Race/Ethnicity/Ancestry (33)
Anti-White (11)	
<b>Religion</b>	Anti-Buddhist (83)
	Anti-Catholic (22)
	Anti-Eastern Orthodox (81)
	Anti-Hindu (84)
	Anti-Islamic (Muslim) (24)
	Anti-Jehovahs Witness (29)
	Anti-Jewish (21)
	Anti-Mormon (28)
	Anti-Multiple Religions, Group (26)
	Anti-Other Christian (82)
	Anti-Other Religion (25)
	Anti-Protestant (23)
	Anti-Sikh (85)
Anti-Atheism/Agnosticism (27)	
<b>Sexual Orientation</b>	Anti-Bisexual (45)
	Anti-Gay (Male) (41)
	Anti-Heterosexual (44)
	Anti-Lesbian (42)
	Anti-Lesbian, Gay, Bisexual, or Transgender (Mixed Group)
<b>Disability</b>	Anti-Mental Disability (52)
	Anti-Physical Disability (51)
<b>Gender</b>	Anti-Female (62)
	Anti-Male (61)
<b>Gender Identity</b>	Anti-Gender Nonconforming (72)
	Anti-Transgender (71)

*Notes:* This table reports the complete list of hate crime bias motivations as classified by the FBI. The table is reproduced from (FBI, 2015, p. 5).

## A.1.2 Trump Twitter Data

**Table A.3: Misclassified Trump’s Anti-Muslim Tweets**

Date	Text	Retweets
12/12/2012	Watching Pyongyang terrorize Asia today is just amazing!	77
26/03/2013	The Scottish windfarm was conceived by the same mind that released terrorist al-Megrahi for humanitarian reasons. ..	101
23/04/2013	Did the Boston terrorists register their guns? No. Another example of why gun control legislation is not the answer!	1192
22/09/2013	”@LebaneseKobe: @realDonaldTrump as a Muslim and as an American, i know for a fact that you Mr. Trump respect all people!	33
22/09/2013	”@mandem3:realDonaldTrump you hate muslims.” Wrong	48
10/10/2013	Obama has called @GOP terrorists during this showdown. Its a shame he really doesnt think it because then he would meet all @GOP demands.	432
29/01/2014	Remember when ”comedian” Bill Maher openly praised the disgusting terrorists who destroyed the World Trade Center-then got canned by ABC?	117
26/01/2015	”tomtumillo: What is worse, Geraldo screaming ’screw the terrorists’ or Kenya feeling she’s ’fabulous’? #CelebrityApprentice	56
15/08/2015	”javonniandjeno:realDonaldTrump AP nbc Donald Trump is Clint Eastwood, the perfect hero not scared of American terrorists. Vote Trump!”	1742
27/08/2015	”jp_sitles:realDonaldTrump HillaryClinton: she compared republicans to terrorist but will not call terrorists , terrorists. #OhMe”	2869
06/09/2015	”jasonusmc2017: blayne_troy @realDonaldTrump: He was right when he called Obama the 5 for 1 president. 5 terrorist for one no good traitor	1016
21/09/2015	”TheBrodyFile: On the Muslim issue: It might help @BarackObama if he actually supported Christians religious liberty rights.	1242
21/09/2015	”TheBrodyFile: On the Muslim issue: It might help @BarackObama if he didn’t take five years to visit Israel”	818
21/11/2015	”WayneDupreeShow: ”Its clear that Donald Trump was NOT even talking about a Muslim Database!” <a href="https://t.co/3tLDZj2WGV">https://t.co/3tLDZj2WGV</a> ”	1020
31/12/2015	”SenSanders: I have a message for Donald Trump: No, were not going to hate Latinos, were not going to hate Muslims.” I fully agree!	1250
23/03/2016	Just watched Hillary deliver a prepackaged speech on terror. Shes been in office fighting terror for 20 years- and look where we are!	11115
23/03/2016	I will be the best by far in fighting terror. Im the only one that was right from the beginning, & now Lyin Ted & others are copying me.	7224
15/06/2016	I will be meeting with the NRA, who has endorsed me, about not allowing people on the terrorist watch list, or the no fly list, to buy guns.	13903
21/05/2017	Speech transcript at Arab Islamic American Summit <a href="https://t.co/eUWxJXJxbe">https://t.co/eUWxJXJxbe</a> nReplay <a href="https://t.co/VtmlSqciXx">https://t.co/VtmlSqciXx</a> #RiyadhSummit #POTUSAbroad	11498
26/05/2017	Getting ready to engage G7 leaders on many issues including economic growth, terrorism, and security.	11322
27/05/2017	Big G7 meetings today. Lots of very important matters under discussion. First on the list, of course, is terrorism. #G7Taormina	9489
18/08/2017	Today, I signed the Global War on Terrorism War Memorial Act (#HR873.) The bill authorizes....cont <a href="https://t.co/c3zIkdtowc">https://t.co/c3zIkdtowc</a> <a href="https://t.co/re6n0MS0cj">https://t.co/re6n0MS0cj</a>	14892
07/09/2017	During my trip to Saudi Arabia, I spoke to the leaders of more than 50 Arab & Muslim nations about the need to confront our shared enemies.[...]	10156
11/11/2017	When will all the haters and fools out there realize that having a good relationship with Russia is a good thing, not a bad thing.[...]	39627

*Notes:* The table lists the tweets we excluded by hand from the set of negative Muslim tweets that were identified by the machine learning model. See text for details.



**Table A.4: Examples of Trump’s Negative Tweets about Muslims**

Date	Text	Retweets
12/10/2015	"mimi_saulino: seanhannity @FoxNews Syrian Muslims escorted into U.S. through Mexico. Now arriving to Oklahoma and Kansas! Congress?"	1223
14/11/2015	Why won't President Obama use the term Islamic Terrorism? Isn't it now, after all of this time and so much death, about time!	6924
15/11/2015	"thewatcher23579: One of Paris terrorist came as Syrian refugee. Donald Trump is right again. BOMB THEIR OIL - TAKE AWAY THEIR FUNDING"	2165
17/11/2015	Refugees from Syria are now pouring into our great country. Who knows who they are - some could be ISIS. Is our president insane?	16285
22/11/2015	We better get tough with RADICAL ISLAMIC TERRORISTS, and get tough now, or the life and safety of our wonderful country will be in jeopardy!	5172
25/11/2015	I LIVE IN NEW JERSEY; @realDonaldTrump IS RIGHT: MUSLIMS DID CELEBRATE ON 9/11 HERE! WE SAW IT! <a href="https://t.co/1SksZU9qlj">https://t.co/1SksZU9qlj</a>	2252
07/12/2015	Obama said in his speech that Muslims are our sports heroes. What sport is he talking about, and who? Is Obama profiling?	9600
07/12/2015	Statement on Preventing Muslim Immigration: <a href="https://t.co/HCWU16z6SR">https://t.co/HCWU16z6SR</a> <a href="https://t.co/d1dhaIs0S7">https://t.co/d1dhaIs0S7</a>	4716
10/12/2015	The United Kingdom is trying hard to disguise their massive Muslim problem. Everybody is wise to what is happening, very sad! Be honest.	6028
10/12/2015	In Britain, more Muslims join ISIS than join the British army. <a href="https://t.co/LQVNz7b2Eb">https://t.co/LQVNz7b2Eb</a>	4325
17/01/2016	Far more killed than anticipated in radical Islamic terror attack yesterday. Get tough and smart U.S., or we won't have a country anymore!	4126
27/03/2016	Another radical Islamic attack, this time in Pakistan, targeting Christian women & children. At least 67 dead,400 injured. I alone can solve	11353
22/05/2016	Crooked Hillary wants a radical 500% increase in Syrian refugees. We cant allow this. Time to get smart and protect America!	9758
12/06/2016	Appreciate the congrats for being right on radical Islamic terrorism, I don't want congrats, I want toughness & vigilance. We must be smart!	27146
13/06/2016	In my speech on protecting America I spoke about a temporary ban, which includes suspending immigration from nations tied to Islamic terror.	13026
25/06/2016	We must suspend immigration from regions linked with terrorism until a proven vetting method is in place.	11726
28/07/2016	Hillary's refusal to mention Radical Islam, as she pushes a 550% increase in refugees, is more proof that she is unfit to lead the country.	20106
18/10/2016	Thank you Colorado Springs. If Im elected President I am going to keep Radical Islamic Terrorists out of our count <a href="https://t.co/N74UK73RLK">https://t.co/N74UK73RLK</a>	12904
19/10/2016	ISIS has infiltrated countries all over Europe by posing as refugees, and @HillaryClinton will allow it to happen h <a href="https://t.co/MmeW2qsTQh">https://t.co/MmeW2qsTQh</a>	16130
11/02/2017	Our legal system is broken! "77% of refugees allowed into U.S. since travel reprieve hail from seven suspect countries." (WT) SO DANGEROUS!	23082
17/08/2017	Study what General Pershing of the United States did to terrorists when caught. There was no more Radical Islamic Terror for 35 years!	30534
18/08/2017	Radical Islamic Terrorism must be stopped by whatever means necessary! The courts must give us back our protective rights. Have to be tough!	37669
15/09/2017	Loser terrorists must be dealt with in a much tougher manner.The internet is their main recruitment tool which we must cut off & use better!	21411
20/10/2017	Just out report: "United Kingdom crime rises 13% annually amid spread of Radical Islamic terror." Not good, we must keep America safe!	29854
01/11/2017	NYC terrorist was happy as he asked to hang ISIS flag in his hospital room. He killed 8 people, badly injured 12. SHOULD GET DEATH PENALTY!	43455

*Notes:* This table reports examples of Trump’s negative tweets about Muslims, including the date of the tweet and the number of retweets the tweet received.

### A.1.3 Geocoded Twitter Data

**Table A.5: Search Terms Used to Identify Users Tweeting about Other Festivals**

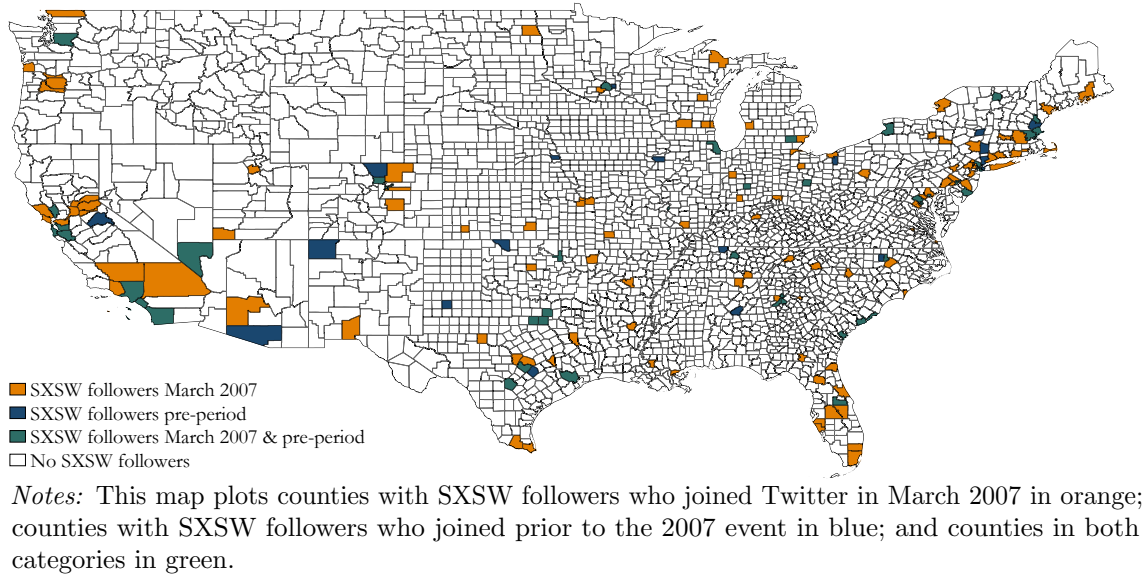
Festival	Search Term
South by Southwest Festival	South by Southwest SXSW
Burning Man	Burningman Burning Man
Coachella	Coachella
Lollapalooza	Lollapalooza
Pitchfork Music Festival	Pitchfork Music Festival Pitchforkfest
West by Southwest Festival	West by Southwest WXSX
Austin City Limited Festival	Austin City Limits Festival
Electric Daisy Festival	EDC Las Vegas Electric Daisy Carnival
New Orleans Jazz and Heritage Festival	New Orleans Jazz and Heritage Festival Jazzfest

**Table A.6: Search Terms Used to Create a Proxy for Total Tweets**

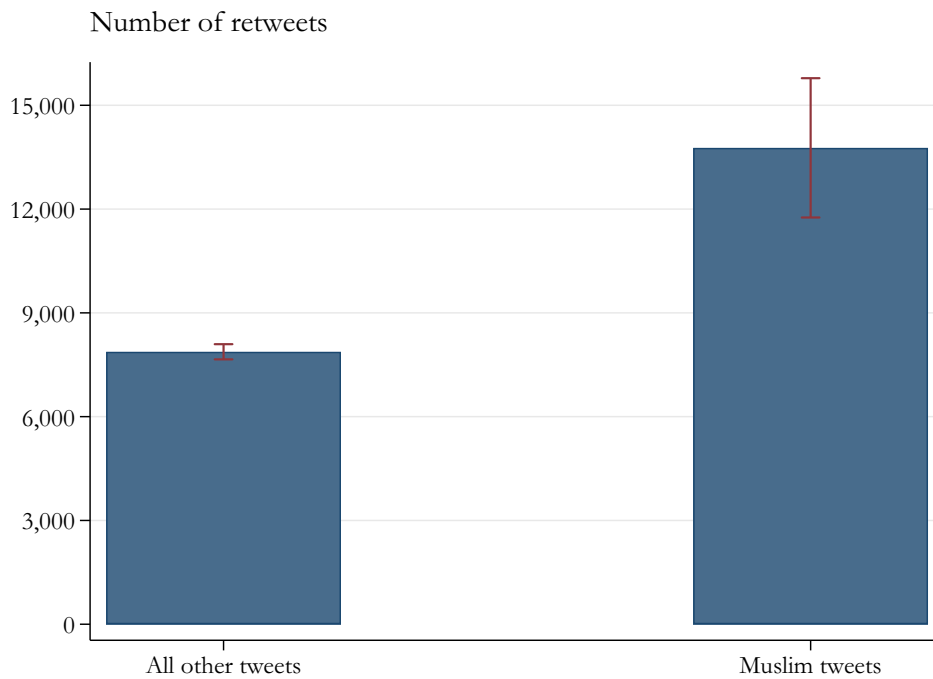
0	I	but	from	his	look	one	she	these	way	would
1	about	by	get	how	make	only	so	they	we	year
2	after	can	give	if	me	or	some	think	well	you
3	all	come	go	in	most	other	take	this	what	your
4	also	could	good	into	my	our	than	time	when	
5	any	day	have	it	new	out	that	two	which	
6	as	do	he	its	no	over	their	up	who	
7	at	even	he	just	not	people	them	us	with	
8	back	first	her	know	now	say	then	use	with	
9	because	for	him	like	on	see	there	want	work	

*Notes:* This table list the search terms we used to collect a proxy of all tweets sent from a given county.

**Figure A.1: Identifying Variation**



**Figure A.2: Average Retweets of Trump’s Tweets, by Muslim Content**



*Notes:* This figure plots the average number of retweets Donald Trump received on his tweets about Muslims compared to all other tweets. We also show 95% confidence intervals.

#### A.1.4 Rescaling of Google trends

As described in Section 2, we use weekly Google trends data to rescale daily values. The daily Google trends data are scaled between 0-100 for each 90 day period, while the weekly Google trends data have a consistent scaling for the entire time period.

To arrive at consistent values, we use the following process. First, we create a scaling factor by dividing the weekly interest by 100. We then multiply the daily data with the scaling factor. If the weekly interest is 100, the scaling factor would be 1, and the daily values would remain the same. On the other hand, if the weekly interest is low, say 10, the daily interest would be scaled down. This way, the adjustment guarantees that daily search interest is on the same scale and thus comparable over time.

As a final step, we divide the rescaled values by their maximum and multiply them by 100. This is to re-normalize the Google trend values to take on values between 0 and 100.

#### A.1.5 Sources for Trump’s golf activity

**Table A.7: Sources for Golf Data**

Source	Description
New York Times	The NYT tracks visits by Trump to his own properties. The data also track how often Trump visited a golf club.
trumpgolfcount.com	This website lists Trump’s visits to golf clubs since his inauguration. It also provides additional analysis during which visits Trump likely played golf.
Presidential Schedule	The presidential schedule lists all past presidential journeys.

#### A.1.6 Calculating the Similarity of SXSW Followers and All Twitter Users

We calculate the similarity of all Twitter user profiles to those of SXSW followers using Latent Semantic Analysis (LSA) (Deerwester et al., 1990; Landauer, 2007). While we could create a similarity measure based on the word count in the Twitter profile bios, this measure would be less reliable at the individual-level as the bio strings are very short and the resulting document-word matrix therefore extremely sparse.

LSA improves on such a measure by reducing the dimensions of the document-word matrix using singular value decomposition. Singular value decomposition derives the components that best describe the semantic space and as a result even profile bios that do not have a single word in common can be similar if they contain words that are used in similar context (e.g. website and webpage). See Iaria et al. (2018) for an example using a similar

approach. For a more extensive description of LSA as well as a Stata implementation see Landauer (2007) and Schwarz (2019).

In our setting, we prepare the data by removing stopwords and reducing all words to their morphological roots, so called lemmas. We then extract all words that appear in at least 5 Twitter bios. This allows us to construct a word-document matrix which is then reweighted using term-frequency inverse document frequency. Afterwards, we use LSA to extract the first 300 principle components of the matrix. The resulting matrix is then used to calculate the cosine similarity between the biography strings of each user in the Kinder-Kurlanda et al. (2017) data with each follower of the SXSU festival. We then normalize the similarity measure to have mean 0 and standard deviation 1 to facilitate the interpretation.

Table A.8: Descriptive Statistics (Main Variables)

	Mean	Std. Dev.	Min.	Median	Max.	N
<b>Hate crime and Twitter variables</b>						
$\Delta$ Log(Hate crimes against Muslims)	0.03	0.14	-0.55	0.00	1.36	3,108
Log(Twitter users)	5.29	1.76	0.00	5.13	12.35	3,108
Log(SXSW followers, March 2007)	0.06	0.32	0.00	0.00	4.98	3,108
Log(SXSW followers, Pre)	0.02	0.18	0.00	0.00	3.61	3,108
<b>Demographic controls</b>						
% aged 20-24	0.06	0.02	0.01	0.06	0.27	3,108
% aged 25-29	0.06	0.01	0.03	0.06	0.15	3,108
% aged 30-34	0.06	0.01	0.03	0.06	0.12	3,108
% aged 35-39	0.06	0.01	0.03	0.06	0.11	3,108
% aged 40-44	0.06	0.01	0.02	0.06	0.10	3,108
% aged 45-49	0.06	0.01	0.02	0.06	0.09	3,108
% aged 50+	0.39	0.07	0.11	0.39	0.75	3,108
Population growth, 2000-2016	0.06	0.18	-0.43	0.03	1.32	3,108
<b>Geographical controls</b>						
Population density	261.27	1733.47	0.10	45.60	69468.40	3,108
Log(County area)	6.53	0.86	0.69	6.47	9.91	3,108
Distance from Austin, TX (in miles)	1450.64	612.61	5.04	1464.66	3098.88	3,108
<b>Race and religion controls</b>						
% white	0.77	0.20	0.03	0.84	0.98	3,108
% black	0.09	0.14	0.00	0.02	0.85	3,108
% native American	0.02	0.06	0.00	0.00	0.90	3,108
% Asian	0.01	0.02	0.00	0.01	0.37	3,108
% Hispanic	0.09	0.14	0.01	0.04	0.96	3,108
% Muslim	0.23	1.08	0.00	0.00	30.35	3,108
<b>Socioeconomic controls</b>						
% below poverty level	16.74	6.58	1.40	16.00	53.30	3,108
% unemployed	5.50	1.94	1.80	5.30	24.10	3,108
Gini index	0.44	0.03	0.33	0.44	0.65	3,108
% uninsured	13.32	5.28	1.80	12.80	49.00	3,108
Log(Median household income)	10.72	0.24	9.87	10.71	11.72	3,107
% employed in agriculture	0.01	0.03	0.00	0.00	0.58	3,108
% employed in IT	0.01	0.01	0.00	0.01	0.21	3,108
% employed in manufacturing	0.16	0.13	0.00	0.13	0.72	3,108
% employed in nontradable sector	0.29	0.11	0.00	0.28	1.00	3,108
% employed in construction/real estate	0.07	0.05	0.00	0.06	1.00	3,108
% employed in utilities	0.04	0.05	0.00	0.03	1.00	3,108
% employed in business services	0.16	0.07	0.00	0.15	0.95	3,108
% employed in other services	0.25	0.10	0.00	0.24	1.00	3,108
% adults with high school degree	34.77	7.07	7.50	35.20	54.80	3,108
% adults with graduate degree	7.05	4.12	0.00	5.80	44.40	3,108

**Table A.8: Descriptive Statistics (Main Variables, Continued)**

	Mean	Std. Dev.	Min.	Median	Max.	N
<b>Media controls</b>						
% watching Fox News	0.26	0.01	0.23	0.26	0.30	3,107
% watching prime time TV	0.43	0.01	0.40	0.43	0.47	3,107
<b>Election control</b>						
Republican vote share, 2012	0.60	0.15	0.06	0.61	0.96	3,108
<b>Crime controls</b>						
Violent crime rate	0.00	0.00	0.00	0.00	0.02	3,108
Property crime rate	0.02	0.01	0.00	0.01	0.10	3,108
<b>Other hate crime variables</b>						
$\Delta \text{Log}(\text{Total hate crimes})$	0.09	0.39	-1.95	0.00	2.34	3,108
$\Delta \text{Log}(\text{Hate crimes against Hispanics})$	0.01	0.17	-1.65	0.00	1.32	3,108
$\Delta \text{Log}(\text{Other ethnicity-based hate crimes})$	-0.00	0.17	-2.60	0.00	1.43	3,108
$\Delta \text{Log}(\text{Racially motivated hate crimes})$	0.06	0.34	-1.69	0.00	2.13	3,108
$\Delta \text{Log}(\text{Hate crimes based on sexual orientation})$	0.01	0.22	-1.32	0.00	1.92	3,108
$\Delta \text{Log}(\text{Hate crimes against other religions})$	0.05	0.24	-1.46	0.00	1.68	3,108
$\text{Log}(\text{Total hate crimes, ADL data})$	0.23	0.64	0.00	0.00	5.38	3,108

Table A.9: Variable Descriptions (Part 1/2)

Variable	Description	Source
<b>Hate crime variables</b>		
Hate crimes	Total number of hate crimes recorded in the FBI hate crime data.	FBI Hate Crime Data
Anti-Muslim hate crimes	Anti-Muslim hate crimes recorded in the FBI hate crime data, based on bias motivation code 24.	FBI Hate Crime Data
Anti-Hispanic hate crimes	Anti-Hispanic hate crimes recorded in the FBI hate crime data, based on the bias motivation codes 32.	FBI Hate Crime Data
Other ethnic-based hate crimes	Anti-ethnic hate crimes recorded in the FBI hate crime data, based on the bias motivation codes 33.	FBI Hate Crime Data
Anti-racial hate crimes	Racial hate crimes recorded in the FBI hate crime data, based on bias motivation codes 11, 12, 13, 14, 15, 16.	FBI Hate Crime Data
Anti-religious hate crimes	Anti-religious hate crimes (except anti-Muslim) recorded in the FBI hate crime data, based on bias motivation codes 21, 22, 23, 25, 26, 27, 28, 29, 81, 82, 83, 84, 85.	FBI Hate Crime Data
Anti-sexual orientation hate crimes	Hate crimes based on sexual orientation recorded in the FBI hate crime data, based on the bias motivation codes 41, 42, 43, 44, 45.	FBI Hate Crime Data
<b>Twitter data</b>		
Trump tweets	The total number of tweets from Donald Trump's Twitter account.	Trump Twitter Archive
Muslim tweets	The number of tweets from Donald Trump's Twitter account about Islam-related topics. We start classifying these tweets by searching for the terms "sharia", "refugee", "mosque", "muslim", "islam" and "terror". We then read all tweets and verify that they indeed mention Muslims in a negative way.	Trump Twitter Archive
Twitter usage	The number of geolocated tweets per county that were collected using the Twitter streaming API in a 12 month period from June to November 2014 and June to November 2015.	Gesis Datatorium
SXSW followers, March 2007	The number of Twitter users following the SXSW account in each county that signed up to Twitter in March 2007.	Twitter Search API
SXSW followers, Pre	The total number of Twitter users following the SXSW account in each county that signed up to Twitter at any point in 2006.	Twitter Search API
Burning Man Twitter Users, August 2007	The number of Twitter users in each county that tweeted about the Burning Man festival in August 2007 and joined Twitter in August 2007.	Twitter Search API
Coachella Twitter Users, April 2007	The number of Twitter users in each county that tweeted about the Coachella festival in April 2007 and joined Twitter in April 2007.	Twitter Search API
Lollapalooza Twitter Users, August 2007	The number of Twitter users in each county that tweeted about the Lollapalooza festival in August 2007 and joined Twitter in August 2007.	Twitter Search API
<b>Trump golf data</b>		
Trump golfs	A dummy variable for each day in 2017 Trump spent on a golf course and likely played golf.	NYT, trumpgolfcount.com and Pres. Schedule
Trump golfs (NYT only)	A dummy variable for each day in 2017 Trump spent on a Golf course and likely golfed, based solely on the information of the New York Times.	NYT
Trump golf (alternative)	A dummy variable for each day in 2017 Trump spent on a golf course and likely golfed, based on the information of trumpgolfcount.com and extended with information from the Pres. Schedule	trumpgolfcount.com and Pres. Schedule
Golf holiday	A dummy for any of Trump's golf outings that lasts longer than 3 days.	NYT and trumpgolfcount.com
Golf at any point in previous week	A dummy variable which is 1 if Trump golfed at any point in the previous week.	NYT and trumpgolfcount.com

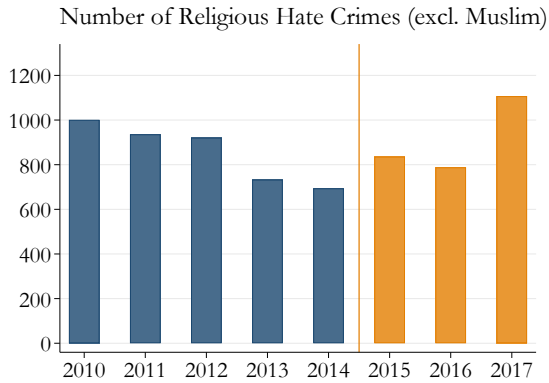


Table A.9: Variable Descriptions (Part 2/2)

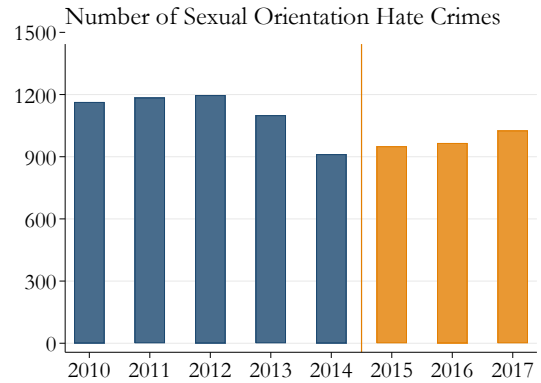
Variable	Description	Source
<b>Other cross sectional controls</b>		
Demographic controls	Contain the share of people in the age buckets 20-24, 25-29, 30-34, 40-44, 45-49 and 50+, and the percentage change in population between 2000 and 2016.	US Census
Education controls	Contains the share of people over 25 with at least a high school degree and the share of people over 25 with at least a graduate degree.	US Census
Race and religion controls	Contains population shares of Muslims, Whites, Blacks, Native Americans, Asians, and Hispanics.	US Census/Religious Census
Socioeconomic controls	Contains a county's poverty rate, unemployment rate, GINI coefficient, share of uninsured, log of median household income, and the share of the population employed in agriculture, manufacturing, accommodation/retail, utilities, information technologies services, and other industries.	US Census/Bureau of Labor Statistics
Media controls	Contains the ratio of prime time TV viewership to population, cable spending to population, and the share of Fox News viewership.	SimplyAnalytics
Election control	Contains the vote share of the Republican party in the 2012 presidential election.	MIT Election Lab
Crime controls	Contains the number of violent crime per capita as well as the number of property crimes per capita based on FBI data.	FBI UCR Data
Distance control	Contains the distance to Austin Texas, the population density, and the logarithm of the land area for each county.	US Census Tigerline File
Change in implicit bias against Muslims	The change in the county-level mean implicit association test score from the Arab-Muslim module between 2015-2017 compared to 2010-2014.	Project Implicit
<b>Other time series variables</b>		
Trump followers' retweets	The number of retweets of Trump's tweets about Muslims by his Twitter followers	Twitter
Trump followers' new content	The number of tweets by Trump followers containing the words "sharia", "refugee", "mosque", "muslim", "islam" or "terror".	Twitter
#StopIslam or #BanIslam	The number of tweets by Trump followers containing the terms "#StopIslam" or "#BanIslam".	Twitter
Muslim mentions (total)	The total number of cable news reports mentioning one of the following terms in their closed captions: "sharia", "refugee", "mosque", "muslim", "islam" and "terror".	Internet Archive
Muslim mentions (Fox News)	The total number of news reports on Fox News mentioning one of the following terms in their closed captions: "sharia", "refugee", "mosque", "muslim", "islam" and "terror".	Internet Archive
Muslim mentions (CNN)	The total number of news reports on CNN mentioning one of the following terms in their closed captions: "sharia", "refugee", "mosque", "muslim", "islam" and "terror".	Internet Archive
Muslim mentions (MSNBC)	The total number of news reports on MSNBC mentioning one of the following terms in their closed captions: "sharia", "refugee", "mosque", "muslim", "islam" and "terror".	Internet Archive
Google searches (PC)	The first principal component of the revealed Google trends for the following terms: "sharia", "refugee", "mosque", "muslim", "islam" and "terror".	Google Trends
Terror attack in the US	The number of Islamist terror attacks committed in the US.	Global Terrorism Database
Terror attack in Europe	The number of Islamist terror attacks committed in the Europe.	Global Terrorism Database
Terror attack elsewhere	The number of Islamist terror attacks committed outside of the US or Europe	Global Terrorism Database

**Figure A.3: Number of Hate Crimes, by Year and Motivating Bias**

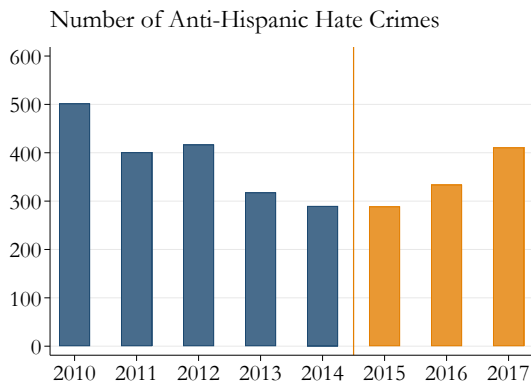
**(a) Religious bias (excl. Muslims)**



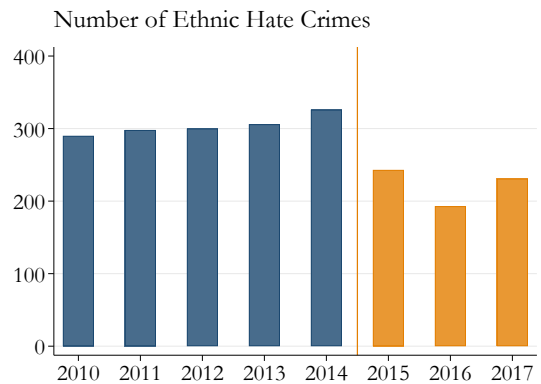
**(b) Sexual orientation bias**



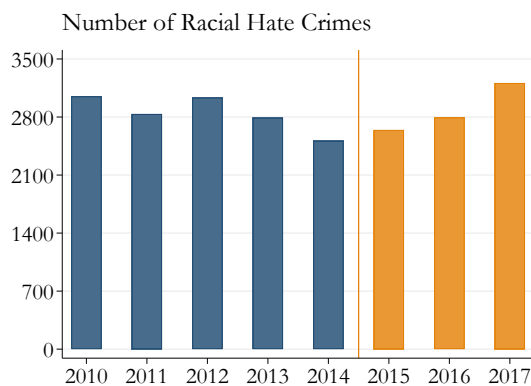
**(c) Anti-Hispanic bias**



**(d) Other ethnic bias**



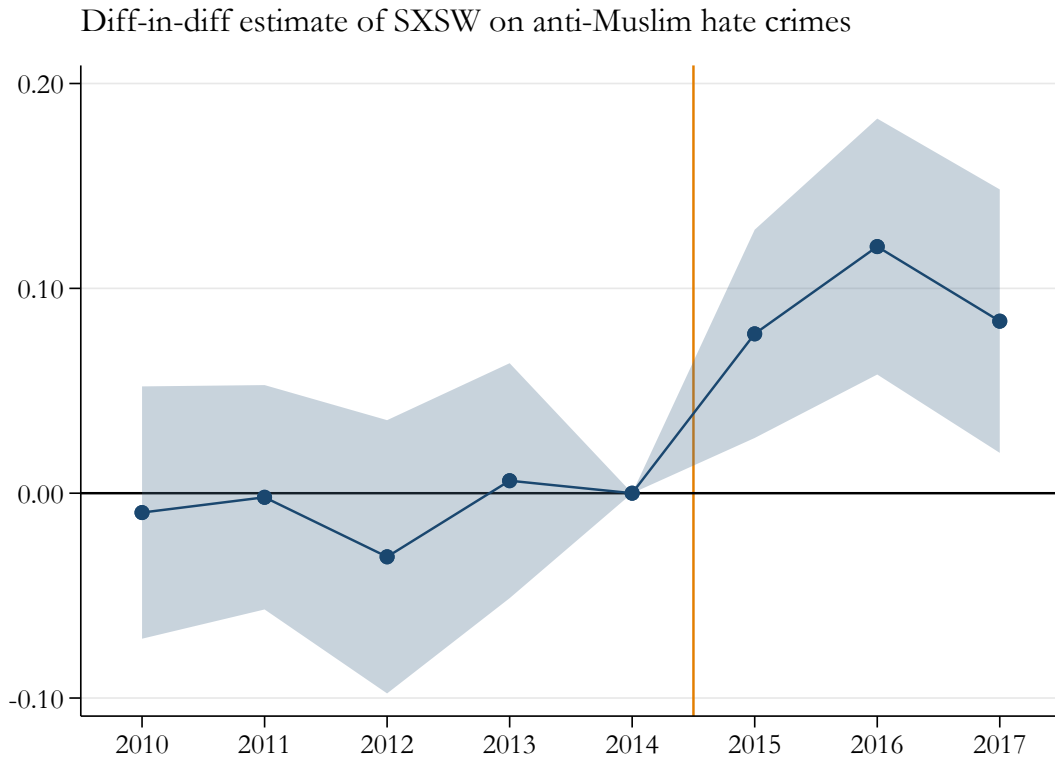
**(e) Racial bias**



*Notes:* These figures plot the number of yearly hate crimes, by year and type of hate crime (as defined by the FBI). The whiskers indicate 95% confidence intervals.

## A.2. Appendix 3: Additional Cross-sectional Evidence

Figure A.4: Anti-Muslim Hate Crimes and Twitter Usage (Reduced Form)



*Notes:* This figure plots the coefficients from running a panel event study regression as in Equation (1), where  $\log(\text{Twitter Usage})$  is replaced by  $\log(\text{SXSW followers, March 2007})$  (with 1 added inside). The dependent variable is the log number of hate crimes. We standardize the variables to have a mean of zero and standard deviation of one. The vertical line indicates the start of Trump's presidential campaign start. The shaded areas are 95% confidence intervals based on standard errors clustered by state.

**Table A.10: Comparing Counties With SXSW Followers, March 2007 vs. Pre**

	March 2007 <i>and</i> Pre (1)	March 2007 <i>only</i> (2)	Pre <i>only</i> (3)	Difference in means (2) - (3)	p-value	Šidàk p-value
<b>Demographic controls</b>						
% aged 20-24	0.07	0.08	0.08	0.00	0.92	1.00
% aged 25-29	0.09	0.07	0.07	-0.00	0.51	1.00
% aged 30-34	0.08	0.07	0.07	-0.00	0.58	1.00
% aged 35-39	0.07	0.06	0.06	-0.00	0.82	1.00
% aged 40-44	0.06	0.06	0.06	0.00	0.82	1.00
% aged 45-49	0.07	0.06	0.06	0.00	0.89	1.00
% aged 50+	0.32	0.35	0.35	-0.00	0.97	1.00
Population growth, 2000-2016	0.18	0.18	0.15	0.03	0.56	1.00
<b>Geographical controls</b>						
Population density	5192.27	1021.39	1998.35	-976.96	0.07*	0.93
Log(County area)	6.30	6.63	6.54	0.09	0.73	1.00
Distance from Austin, TX (in miles)	1775.99	1749.38	1626.64	122.74	0.48	1.00
<b>Race and religion controls</b>						
% white	0.50	0.65	0.67	-0.02	0.62	1.00
% black	0.18	0.12	0.08	0.04	0.20	1.00
% native American	0.01	0.01	0.02	-0.02	0.02**	0.49
% Asian	0.10	0.05	0.05	-0.01	0.55	1.00
% Hispanic	0.20	0.16	0.15	0.01	0.80	1.00
% Muslim	1.31	0.81	0.75	0.05	0.87	1.00
<b>Socioeconomic controls</b>						
% below poverty level	15.71	15.82	13.69	2.14	0.17	1.00
% unemployed	4.86	5.05	4.51	0.54	0.07*	0.93
Gini index	0.48	0.46	0.45	0.01	0.24	1.00
% uninsured	12.87	12.40	11.21	1.19	0.35	1.00
Log(Median household income)	11.00	10.91	10.99	-0.09	0.18	1.00
% employed in agriculture	0.00	0.00	0.00	0.00	0.27	1.00
% employed in IT	0.04	0.02	0.02	-0.00	0.98	1.00
% employed in manufacturing	0.07	0.09	0.09	0.01	0.63	1.00
% employed in nontradable sector	0.23	0.26	0.27	-0.01	0.52	1.00
% employed in construction/real estate	0.06	0.07	0.07	0.01	0.39	1.00
% employed in utilities	0.04	0.04	0.03	0.00	0.56	1.00
% employed in business services	0.29	0.25	0.24	0.01	0.70	1.00
% employed in other services	0.27	0.26	0.28	-0.02	0.27	1.00
% adults with high school degree	21.76	25.99	25.77	0.22	0.88	1.00
% adults with graduate degree	16.15	13.08	14.34	-1.26	0.40	1.00
<b>Media controls</b>						
% watching Fox News	0.25	0.26	0.26	-0.00	0.91	1.00
% watching prime time TV	0.42	0.43	0.43	0.00	0.91	1.00
<b>Election control</b>						
Republican vote share, 2012	0.33	0.46	0.47	-0.02	0.63	1.00
<b>Crime controls</b>						
Violent crime rate	0.01	0.00	0.00	0.00	0.98	1.00
Property crime rate	0.03	0.02	0.02	0.00	0.30	1.00

*Notes:* This table plots the mean values of the control variables for the three types of counties relevant for the cross-sectional results: (1) counties with new SXSW followers in March 2007 *and* the pre-period; (2) counties with new SXSW followers in March 2007 but no new followers in the pre-period; and (3) counties with new SXSW followers in the pre-period but no new followers in March 2007. We report p-values from a two-sided *t*-test for the equality of means between the counties with the key identifying variation, as well as Šidàk-corrected values to account for multiple hypothesis testing. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

**Table A.11: Balancedness - SXSW Twitter Followers' Characteristics**

User first names (Corr. = 0.69)		Terms used in user bios (Corr. = 0.92)	
Pre-SXSW	March 2007	Pre-SXSW	March 2007
michael	michael	http	http
mike	john	founder	com
paul	chris	com	digital
chris	jeff	co	founder
ryan	matt	tech	medium
eric	brian	design	director
david	david	director	tech
matthew	alex	product	music
john	jason	digital	social
jeff	kevin	designer	marketing
robert	paul	medium	design
mark	mike	music	co
andrew	dan	social	writer
daniel	andrew	love	love
james	peter	marketing	lover
kevin	jim	web	dad
jay	tom	geek	creative
jonathan	jennifer	writer	tweet
rob	steve	technology	author
rachel	todd	dad	designer

*Notes:* This table compares the individual characteristics of Twitter users who follow “South by Southwest”, depending on the users’ join date (either in March 2007 or before). We plot the ranking of the most common first names and terms used in a Twitter user’s “bio”.

**Table A.12: Comparison of SXSW Followers and All Twitter Users**

User first names (Corr. = 0.97)		Terms used in user bios (Corr. = 0.94)	
Other counties	SXSW counties	Other counties	SXSW counties
michael	michael	love	co
chris	david	life	love
john	chris	co	life
david	john	http	http
sarah	alex	http co	http co
mike	mike	god	music
emily	matt	ig	lover
ryan	sarah	music	ig
matt	ryan	university	de
alex	andrew	like	like
taylor	emily	fan	fan
ashley	brian	live	world
nick	jessica	lover	instagram
jessica	james	mom	thing
tyler	kevin	husband	la
hannah	daniel	time	live
katie	ashley	follow	time
amanda	jason	one	com
lauren	lauren	wife	artist
brian	mark	thing	one

*Notes:* This table compares the individual characteristics of Twitter users from counties with “South by Southwest” followers who joined in March 2007 (“SXSW counties”) to Twitter users from all other US counties (“Other counties”). We plot the ranking of the most common first names and terms used in a Twitter user’s “bio”.

**Table A.13: Correlation of Log(Twitter Users) Across Events**

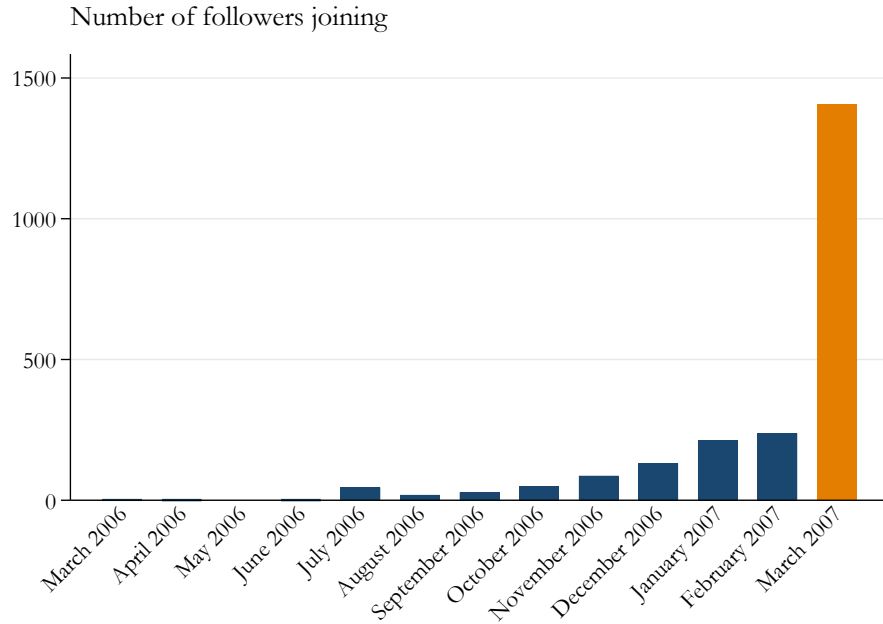
	SXSW March 2007	SXSW Pre	Coachella April 2007	Burning Man August 2007	Lollapalooza August 2007
SXSW followers, March 2007	1				
SXSW followers, Pre	0.77	1			
Coachella users, April 2007	0.44	0.48	1		
Burning Man users, August 2007	0.52	0.56	0.54	1	
Lollapalooza users, August 2007	0.03	0.06	0.00	0.00	1

*Notes:* This table reports the Pearson correlation coefficients between the main measure of interest (*SXSW followers, March 2007*) and different control variables. “Followers” are based on the locations of people who started following SXSW in a given month; “users” are based on people who tweeted at least once about a festival. We take the natural logarithm of these numbers with one added inside.

**Table A.14: Number of Counties With Any Twitter Users at SXSW or Other Festivals**

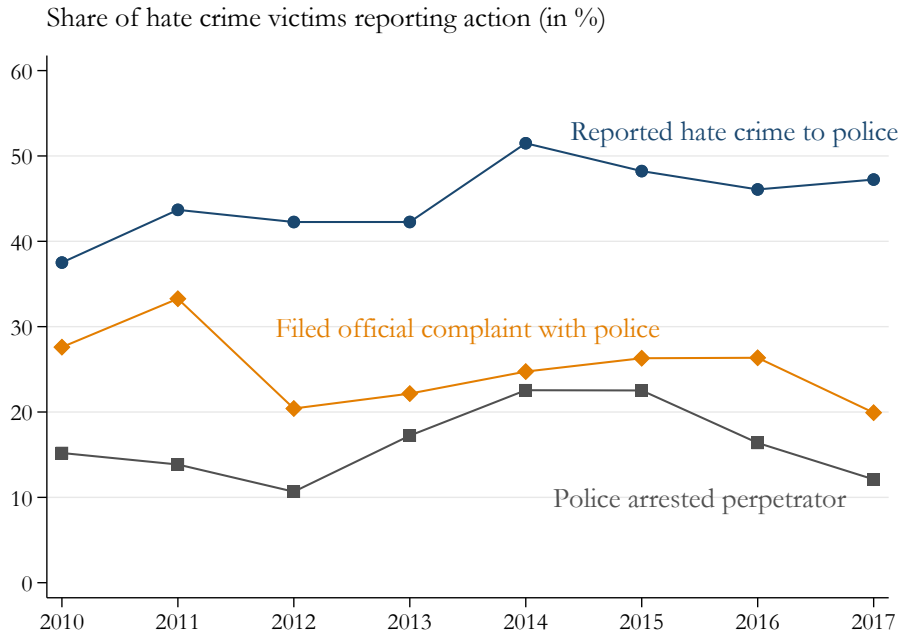
	SXSW March 2007	SXSW Pre	Coachella April 2007	Burning Man August 2007	Lollapalooza August 2007
No followers	2953	2987	3091	3098	3105
At least 1 follower	155	121	17	10	3

**Figure A.5: Number of SXSW Followers Joining Each Month**



*Notes:* This figure plots the number of SXSW followers who joined Twitter each month in the run-up to the 2007 SXSW festival. The orange bar marks the instrument used in the paper.

**Figure A.6: Trends in Hate Crime Reporting**



*Notes:* This figure visualizes time series trends in the reporting of hate crimes and police actions taken in response to them. The source is the Bureau of Justice Statistics National Crime Victimization Survey (NCVS). The sample consists of 1,416 hate crime incidents reported between 2010 and 2017. We report the share of respondents that took each action using victimization weights.

Table A.15: 2SLS - Alternative SXSWS Controls

Control variable(s)	None	Pooled	Pooled	Individual	Individual	Individual	Individual	Other festivals
Control period	—	2006	2006-Feb. 2007	Feb. 2007	2006	2006-Feb. 2007	2007	2007
Control counties	—	67	121	59	67	121	25	25
Corr(March 2007, Control), average	—	0.77	0.83	0.72	0.49	0.54	0.33	0.33
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(7)
<b>Panel A: Reduced form</b>								
SXSWS measure, March 2007	0.096*** (0.022)	0.070** (0.032)	0.066* (0.036)	0.083** (0.034)	0.071** (0.030)	0.061* (0.035)	0.094*** (0.023)	
SXSWS measure, control (linear combination)	—	0.059 (0.055)	0.047 (0.040)	0.036 (0.058)	-0.120 (0.233)	-0.109 (0.230)	0.272 (0.225)	
<b>Panel B: 2SLS</b>								
Log(Twitter users)	0.141*** (0.031)	0.121** (0.054)	0.165** (0.082)	0.137*** (0.048)	0.117** (0.048)	0.125* (0.068)	0.151*** (0.035)	
Weak IV 95% AR confidence set	[0.085; 0.198]	[0.021; 0.221]	[-0.003; 0.318]	[0.040; 0.216]	[0.027; 0.206]	[-0.013; 0.251]	[0.085; 0.216]	
SXSWS measure, control (linear combination)	—	0.032 (0.067)	-0.026 (0.071)	0.008 (0.060)	0.060 (0.222)	0.065 (0.227)	0.199 (0.206)	
Observations	3,107	3,107	3,107	3,107	3,107	3,107	3,107	
Mean of DV	0.018	0.018	0.018	0.018	0.018	0.018	0.018	
Robust F-stat.	241.22	86.97	36.73	68.04	133.64	66.58	161.03	

Notes: This table presents county-level OLS and IV regressions where the dependent variable is the log change in hate crimes against Muslims between 2010 and 2017. *Log(Twitter usage)* is instrumented using the measure described in the top rows; column 2 plots the baseline specification. “Pooled” controls refer to one variable for the entire control period and “individual” to a vector of individual variables for each control period (e.g. one variable for March 2006, one variable for April 2006, etc.). *SXSWS measure, control (linear combination)* is the estimate for the SXSWS control variable(s). In case of multiple controls, we plot linear combinations of coefficients and standard errors. In those cases, we also plot the *average* of the correlation of the individual controls with the March 2007 measure. All regressions control for population deciles, state fixed effects and demographic controls that include population growth between 2000 and 2016 as well as age cohort controls for the share of people aged 20-24, 25-29, 30-34, 35-39, 40-44, 45-49, and those over 50. Weak IV 95% Anderson-Rubin (AR) confidence sets are calculated using the two-step approach of Andrews (2018) using the Stata package from Sun (2018). For the just-identified case we study here, the “robust” *F*-stat. is equivalent to the “Kleibergen-Paap” or the “effective” *F*-statistic of Olea & Pflueger (2013). Robust standard errors in parentheses are clustered by state. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .



**Table A.16: Robustness - Alternative Measures of Twitter Usage**

	Survey # households using Twitter (1)	Survey % households using Twitter (2)	GESIS Tweets (Pre-Trump) (3)	GESIS Twitter users (4)
<b>Panel A: First stage - Twitter usage</b>				
Log(SXSW followers, March 2007)	0.440*** (0.041)	0.080*** (0.018)	0.443*** (0.061)	0.461*** (0.061)
<b>Panel B: OLS - Hate crimes against Muslims</b>				
Twitter measure	0.061*** (0.016)	0.024** (0.010)	0.017*** (0.005)	0.018*** (0.005)
<b>Panel C: 2SLS - Hate crimes against Muslims</b>				
Twitter measure	0.160** (0.068)	0.877** (0.387)	0.159** (0.073)	0.152** (0.070)
Weak IV 95% AR confidence set	[0.033; 0.286]	[0.159; 10.745]	[0.023; 0.308]	[0.023; 0.295]
Log(SXSW followers, Pre)	0.040 (0.062)	0.007 (0.083)	0.034 (0.069)	0.034 (0.069)
Observations	3,106	3,106	3,107	3,107
Mean of DV	0.018	0.018	0.018	0.018
SD of Twitter measure	1.474	0.549	1.925	1.908
Robust F-stat.	114.10	20.59	53.15	58.04

*Notes:* This table presents county-level OLS, reduced form, and IV regressions where the dependent variable is the log change in hate crimes against Muslims between 2010 and 2017. *Twitter usage measure* is the measure listed in the top row, instrumented using the number of users who started following SXSW in March 2007 (in log with 1 added inside). *SXSW followers, Pre* is the number of SXSW followers who registered at some point in 2006 (in log with 1 added inside). All regressions control for population deciles and state fixed effects, as well as demographic controls including population growth between 2000 and 2016 as well as age cohort controls for the share of people aged 20-24, 25-29, 30-34, 35-39, 40-44, 45-49, and those over 50. Weak IV 95% Anderson-Rubin (AR) confidence sets are calculated using the two-step approach of Andrews (2018) using the Stata package from Sun (2018). For the just-identified case we study here, the “robust”  $F$ -stat. is equivalent to the “Kleibergen-Paap” or the “effective”  $F$ -statistic of Olea & Pflueger (2013). Robust standard errors in parentheses are clustered by state. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

**Table A.17: Further Robustness - Social Media and the Rise in Hate Crimes Against Muslims**

	Pop. weights (1)	Change since 1990 (2)	Log hate crimes (3)	Drop zero change counties (4)	Drop potentially nonreporting counties (5)	Drop counties with few Muslims (6)	Only neighbouring counties (7)	Drop zero follower counties (8)	Diff. Trump campaign start (9)
<b>Panel A: OLS</b>									
Log(Twitter users)	0.095*** (0.023)	0.056*** (0.012)	0.146*** (0.034)	0.059** (0.024)	0.043*** (0.009)	0.068*** (0.018)	0.053*** (0.012)	0.066** (0.030)	0.038*** (0.010)
<b>Panel B: Reduced form</b>									
Log(SXSW followers, March 2007)	0.104** (0.043)	0.138*** (0.035)	0.302*** (0.066)	0.097** (0.047)	0.072** (0.033)	0.073** (0.035)	0.077** (0.036)	0.098** (0.045)	0.091*** (0.030)
<b>Panel C: 2SLS</b>									
Log(Twitter users)	0.163** (0.065)	0.189*** (0.051)	0.519*** (0.108)	0.176** (0.084)	0.123** (0.056)	0.132** (0.062)	0.138** (0.064)	0.180* (0.091)	0.156*** (0.049)
Weak IV 95% AR confidence set	[0.043; 0.284]	[0.132; 0.343]	[0.318; 0.720]	[0.024; 0.329]	[0.020; 0.227]	[0.018; 0.246]	[0.019; 0.269]	[0.035; 0.376]	[0.065; 0.247]
Log(SXSW followers, Pre)	0.000 (0.065)	0.007 (0.065)	0.061 (0.162)	0.030 (0.087)	0.039 (0.067)	0.049 (0.073)	0.021 (0.075)	0.027 (0.074)	0.039 (0.071)
Observations	3,107	3,107	3,107	394	2,185	586	1,167	172	3,108
Mean of DV	0.158	0.022	0.088	0.137	0.025	0.080	0.038	0.153	0.027
Robust F-stat.	63.74	86.97	86.97	53.29	93.14	90.16	74.52	47.70	109.24

*Notes:* This table presents county-level OLS and IV regressions where the dependent variable is the log change in hate crimes against Muslims between 2010 and 2017 in all columns except columns 2 and 3. In column 2, the dependent variable is the log change between 1990 and 2017; in column 3, it is the log number of hate crimes against Muslims in a county after the start of Donald Trump's presidential run on June 16, 2015. *Log(Twitter usage)* is instrumented using the number of users who started following SXSW in March 2007. *SXSW followers*, *Pre* is the number of SXSW followers who registered at some point in 2006. All regressions control for population deciles, state fixed effects (except in column 1), and demographic controls that include population growth between 2000 and 2016 as well as age cohort controls for the share of people aged 20-24, 25-29, 30-34, 35-39, 40-44, 45-49, and those over 50. Column 4 drops all counties for which the change in hate crimes between 2010 and 2017 was zero. Column 5 drops all counties which never report a hate crime between 1990 and 2017. Column 6 drops all counties for which the (rounded) share of Muslims in the county population is zero according to Census data. Column 7 only keeps neighbouring counties that differ in whether they have SXSW followers in March 2007 or not. Column 8 only keeps counties with any SXSW follower in March 2007 or the pre-period. In column 9, the dependent variable is the log-difference in anti-Muslim hate crimes after Trump's campaign start (June 16, 2015) until 2017, compared to the period from 2010 until his campaign. Weak IV 95% Anderson-Rubin (AR) confidence sets are calculated using the two-step approach of Andrews (2018) using the Stata package from Sun (2018). For the just-identified case we study here, the "robust" *F*-stat. is equivalent to the "Kleibergen-Paap" or the "effective" *F*-statistic of Olea & Pflueger (2013). Robust standard errors in parentheses are clustered by state. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

**Table A.18: Robustness - Alternative Estimators**

	IV Probit (1)	IV Poisson (2)	Inverse Hyperbolic Sine (3)	Index Dependent Variable (4)
<b>Panel A: OLS</b>				
Log(Twitter users)	0.049*** (0.008)	0.242*** (0.027)	0.030*** (0.007)	0.042** (0.018)
<b>Panel B: Reduced form</b>				
Log(SXSW followers, March 2007)	0.036** (0.015)	0.138*** (0.033)	0.070** (0.032)	0.191*** (0.074)
<b>Panel C: 2SLS</b>				
Log(Twitter users)	0.081*** (0.031)	0.288*** (0.094)	0.183** (0.083)	0.396*** (0.152)
Weak IV 95% AR confidence set	[0.348; 10.182]		[0.030; 0.336]	[0.115; 0.678]
Log(SXSW followers, Pre)	-0.014 (0.030)	-0.016 (0.078)	0.035 (0.061)	-0.082 (0.144)
Observations	2,648	2,648	3,106	3,106
Mean of DV	0.094	0.264	0.023	0.032

*Notes:* This table presents county-level OLS, reduced form, and IV regressions where the dependent variable is measure of hate crimes against Muslims. Column 1 reports the results from an IV probit regression estimated using maximum likelihood, where the dependent variable is a dummy for counties with an increase in hate crimes against Muslims (and 0 otherwise). Column 2 estimates a Poisson regression, where the dependent variable is the total number of hate crimes after Trump’s presidential campaign start. Column 3 replaces the dependent variable with the change in the inverse hyperbolic sine of hate crimes, and the Twitter variables with their inverse hyperbolic sine (instead of  $\log(1+)$ ). Column 4 recodes the dependent variable into an index equal to 1 for increases, -1 for decreases, and 0 for no changes in hate crimes. All regressions control for population deciles and state fixed effects, as well as demographic controls, geographical controls, and race and religion controls, and socioeconomic controls. Weak IV 95% Anderson-Rubin (AR) confidence sets are calculated using the two-step approach of Andrews (2018) using the Stata package from Sun (2018). For the just-identified case we study here, the “robust”  $F$ -stat. is equivalent to the “Kleibergen-Paap” or the “effective”  $F$ -statistic of Olea & Pflueger (2013). Robust standard errors in parentheses are clustered by state. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

**Table A.19: Social Media and Types of Hate Crimes**

	Any (1)	Vandalism (2)	Theft (3)	Burglary (4)	Robbery (5)	Assault (6)
<b>Panel A: OLS</b>						
Log(Twitter users)	0.030*** (0.007)	0.011** (0.005)	0.001 (0.001)	0.001 (0.002)	0.002 (0.002)	0.030*** (0.008)
<b>Panel B: Reduced form</b>						
Log(SXSW followers, March 2007)	0.070** (0.032)	0.020 (0.020)	0.001 (0.004)	0.006 (0.009)	-0.000 (0.003)	0.068** (0.032)
<b>Panel C: 2SLS</b>						
Log(Twitter users)	0.121** (0.054)	0.034 (0.034)	0.002 (0.007)	0.010 (0.015)	-0.001 (0.005)	0.117** (0.055)
Weak IV 95% AR confidence set	[0.021; 0.221]	[-0.035; 0.097]	[-0.012; 0.015]	[-0.017; 0.038]	[-0.010; 0.009]	[0.014; 0.219]
Log(SXSW followers, Pre)	0.032 (0.067)	0.058 (0.046)	-0.002 (0.007)	-0.014 (0.015)	0.022 (0.015)	0.020 (0.066)
Observations	3,107	3,107	3,107	3,107	3,107	3,107
Mean of DV	0.018	0.007	0.000	0.000	0.001	0.013
Robust F-stat.	86.97	86.97	86.97	86.97	86.97	86.97

*Notes:* This table presents county-level OLS and IV regressions where the dependent variable is the log change in hate crimes against Muslims of the type in the top row between 2010 and 2017. *Log(Twitter usage)* is instrumented using the number of users who started following SXSW in March 2007. *SXSW followers, Pre* is the number of SXSW followers who registered at some point in 2006. All regressions control for population deciles and state fixed effects (not shown). Demographic controls include population growth between 2000 and 2016 as well as age cohort controls for the share of people aged 20-24, 25-29, 30-34, 35-39, 40-44, 45-49, and those over 50. Race and religion controls contains the share of people identifying as white, African American, Native American or Pacific Islander, Asian, Hispanic, or Muslim. Socioeconomic controls include the poverty rate, unemployment rate, local GINI index, the share of uninsured individuals, log median household income, the share of highschool graduates, the share of people with a graduate degree, as well as the employment shares in agriculture, information technology, manufacturing, nontradables, construction and real estate, utilities, business services, or other sectors. Media controls include the viewership share of Fox News, the cable TV spending to population ratio, and the prime time TV viewership to population ratio. Election control is the county-level vote share of the Republican party in 2012. Crime controls are the rates of violent or property crime from the FBI. Geographical controls include the linear distance from the SXSW festival location (Austin, Texas), population density, and the natural logarithm of county size. Weak IV 95% Anderson-Rubin (AR) confidence sets are calculated using the two-step approach of Andrews (2018) using the Stata package from Sun (2018). For the just-identified case we study here, the “robust” *F*-stat. is equivalent to the “Kleibergen-Paap” or the “effective” *F*-statistic of Olea & Pflueger (2013). Robust standard errors in parentheses are clustered by state. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

**Table A.20: Social Media and Hate Crimes - Alternative Standard Errors**

	Robust SE (1)	Bootstrap robust SE (2)	Bootstrap state cluster SE (3)	Spatial SE (4)
<b>Panel A: OLS</b>				
Log(Twitter users)	0.030*** (0.007)	0.030*** (0.007)	0.030*** (0.007)	0.030*** (0.007)
<b>Panel B: Reduced form</b>				
Log(SXSW followers, March 2007)	0.070** (0.028)	0.070*** (0.025)	0.070** (0.032)	0.070** (0.031)
<b>Panel C: 2SLS</b>				
Log(Twitter users)	0.121** (0.050)	0.121** (0.049)	0.121** (0.057)	0.121** (0.055)
Log(SXSW followers, Pre)	0.032 (0.057)	0.032 (0.060)	0.032 (0.069)	0.032 (0.067)
Observations	3,107	3,107	3,107	3,107
Mean of DV	0.018	0.018	0.018	0.018
Robust F-stat.	68.22	68.22	85.64	71.45

*Notes:* This table presents county-level OLS and IV regressions where the dependent variable is the log change in hate crimes against Muslims between 2010 and 2017. *Log(Twitter usage)* is instrumented using the number of users who started following SXSW in March 2007. *SXSW followers, Pre* is the number of SXSW followers who registered at some point in 2006. All regressions control for population deciles and state fixed effects (not shown). Demographic controls include population growth between 2000 and 2016 as well as age cohort controls for the share of people aged 20-24, 25-29, 30-34, 35-39, 40-44, 45-49, and those over 50. Spatial standard errors are based on the method proposed in Colella et al. (2019), implemented in Stata as *acreg*, using a 200 miles cutoff. For the just-identified case we study here, the “robust” *F*-stat. is equivalent to the “Kleibergen-Paap” or the “effective” *F*-statistic of Olea & Pflueger (2013). Standard errors are computed as indicated in the top row. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

**Table A.21: Social Media and Hate Crimes - Split by Number of Perpetrators**

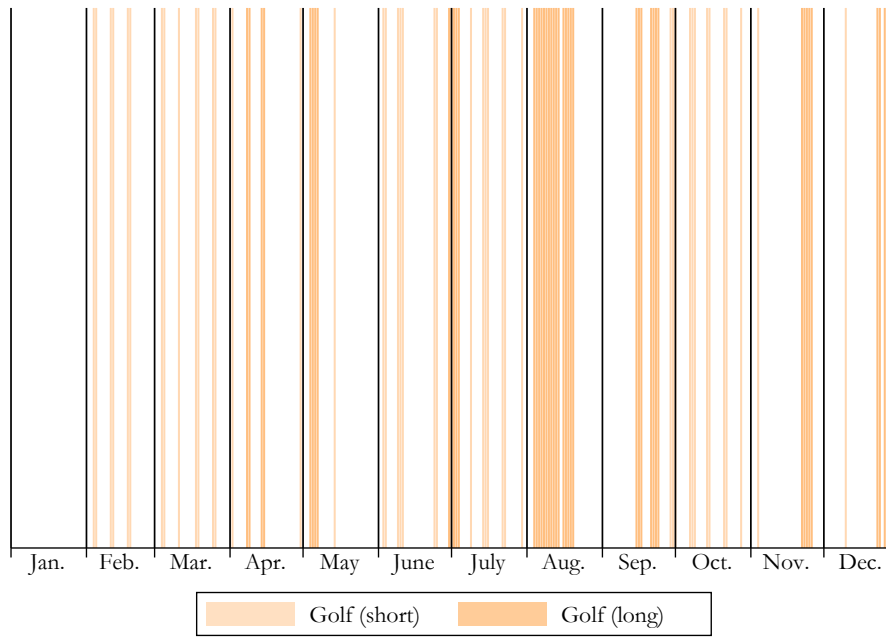
	Muslim bias		Hispanic bias	
	One offender (1)	Multiple offenders (2)	One offender (3)	Multiple offenders (4)
<b>Panel A: OLS</b>				
Log(Twitter users)	0.024*** (0.006)	0.004 (0.005)	-0.000 (0.008)	-0.008 (0.006)
<b>Panel B: Reduced form</b>				
Log(SXSW followers, March 2007)	0.052* (0.031)	0.014 (0.014)	0.073*** (0.026)	-0.003 (0.019)
<b>Panel C: 2SLS</b>				
Log(Twitter users)	0.108* (0.064)	0.028 (0.030)	0.151*** (0.051)	-0.006 (0.039)
Weak IV 95% AR confidence set	[-0.010; 0.226]	[-0.027; 0.084]	[0.056; 0.246]	[-0.086; 0.058]
Log(SXSW followers, Pre)	0.040 (0.064)	0.019 (0.022)	-0.058 (0.059)	0.000 (0.036)
Observations	3,106	3,106	3,106	3,106
Mean of DV	0.012	0.003	-0.005	-0.004
Robust F-stat.	76.10	76.10	76.10	76.10
Share of hate crimes	81%	19%	78%	22%

*Notes:* This table presents county-level OLS and IV regressions where the dependent variable is the log change in hate crimes with the indicated number of offenders between 2010 and 2017. We have information on the number of perpetrators for 62% of hate crimes in our sample. The bottom row reports the percentage of hate crimes falling into the one and multiple offender categories for incidents for which we have information. *Log(Twitter usage)* is instrumented using the number of users who started following SXSW in March 2007. *SXSW followers, Pre* is the number of SXSW followers who registered at some point in 2006. All regressions control for population deciles and state fixed effects (not shown). We also control the full set of controls. For the just-identified case we study here, the “robust” *F*-stat. is equivalent to the “Kleibergen-Paap” or the “effective” *F*-statistic of Olea & Pflueger (2013). Standard errors are clustered by state. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

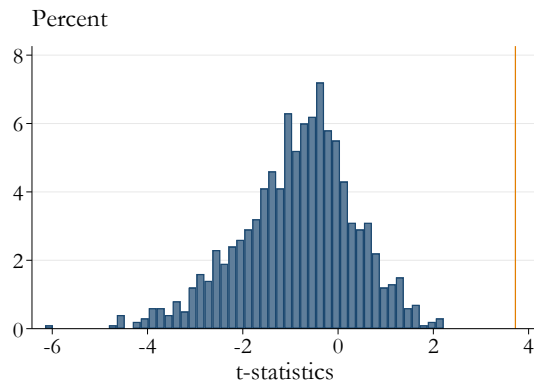
### A.3. Appendix 4: Additional Time Series Evidence

Figure A.7: Trump's Golf Days

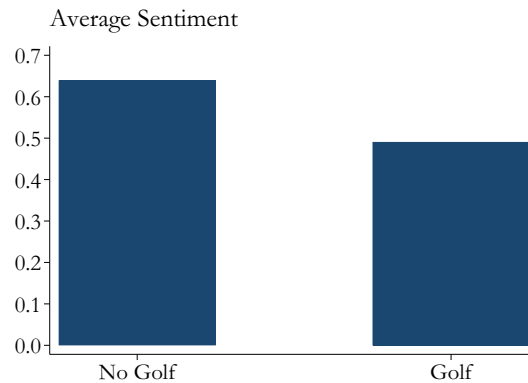
(a) Trump's Golf Days in 2017



(b) Randomization Test for Golf Days



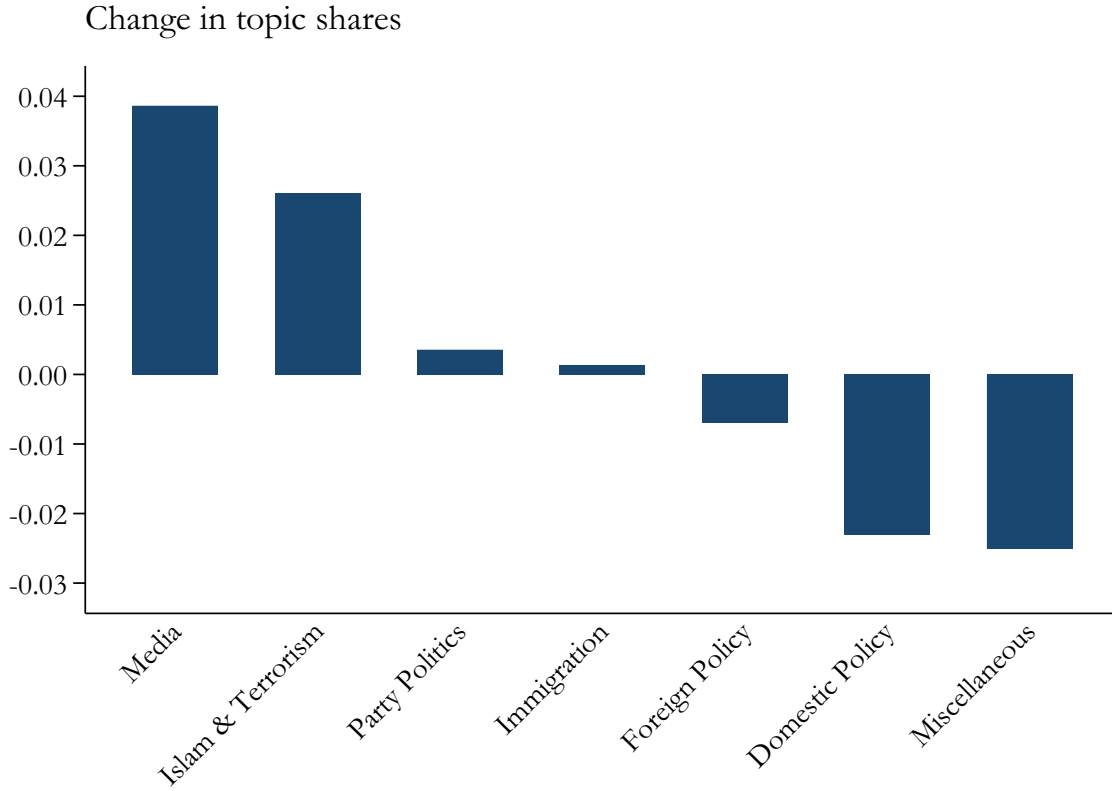
(c) Golf Days and Tweet Sentiment



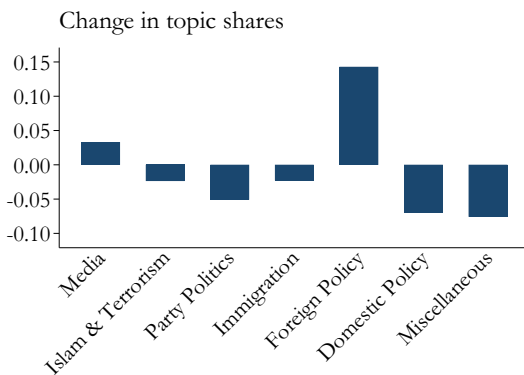
*Notes:* Panel (a) plots the days in 2017 when Donald Trump played golf. *Golf (long)* indicates three or more consecutive days of golfing. Panel (b) visualizes the distribution of  $t$ -statistics from a randomization test of the first stage regression of Trump's tweets about Muslims on placebo golf days. In particular, we create 1,000 placebo sets of 92 golf days, which is the number of times Trump golfed in 2017. We then regress the log number of Trump's tweets about Muslims on these dummies using the baseline specification in Equation (4) and report the resulting  $t$ -statistics. The orange line marks our baseline point estimate. Panel (c) plots the average sentiment of Donald Trump's tweets on golf and non-golf days. Lower values mean more negative sentiment. The sentiment was independently hand-coded using a scale from -2 (very negative) to 2 (very positive).

**Figure A.8: Shift in Topics of Trump’s Tweets During Events**

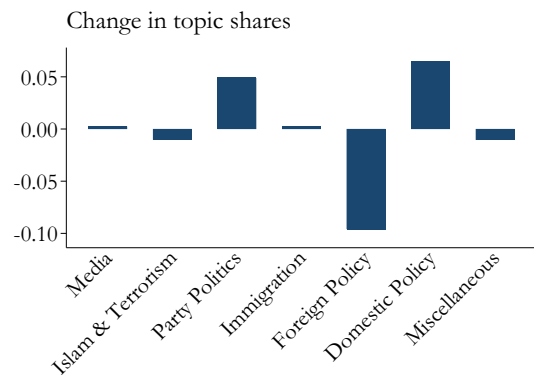
**(a) Golf Days**



**(b) Travel Abroad**



**(c) Policy Briefing**



*Notes:* This figure shows how the content of Donald Trump’s tweets changes on days when he plays golfs (Panel a), he is traveling abroad (Panel b) or receives a policy briefing (Panel c), based on the official presidential schedule. Topics are based on the independent hand-coding of three research assistants.



**Table A.22: Summary Statistics for Time Series**

Variable	Mean	SD	p50	Min	Max	N
<b>Trump tweets</b>						
Log(Muslim Trump tweets)	0.08	0.25	0.00	0.00	1.79	365
Log(Trump tweets)	1.95	0.58	0.00	1.95	3.30	365
Muslim Trump tweets (dummy)	0.09	0.29	0.00	0.00	1.00	365
<b>Hate crimes against Muslims (1 + natural logarithm)</b>						
All types	0.43	0.45	0.00	0.69	1.61	365
Assault	0.29	0.40	0.00	0.00	1.61	365
Vandalism	0.14	0.29	0.00	0.00	1.39	365
Theft	0.01	0.09	0.00	0.00	1.10	365
Burglary	0.01	0.07	0.00	0.00	0.69	365
Robbery	0.01	0.09	0.00	0.00	0.69	365
<b>Other hate crimes (1 + natural logarithm)</b>						
All hate crimes	2.91	0.27	2.08	2.94	3.58	365
Other ethnicity	0.38	0.45	0.00	0.00	1.79	365
Race	2.22	0.37	0.69	2.30	3.00	365
Sexual orientation	1.23	0.48	0.00	1.39	2.40	365
Religion (excl. Muslims)	1.28	0.50	0.00	1.39	2.83	365
<b>TV news coverage (1 + natural logarithm)</b>						
Muslim mentions (total)	3.71	0.64	0.69	3.69	5.26	365
Muslim mentions (Fox News)	2.75	0.66	0.00	2.77	4.29	365
Muslim mentions (CNN)	2.24	0.94	0.00	2.30	4.29	365
Muslim mentions (MSNBC)	2.75	0.66	0.00	2.77	4.26	365
<b>Trump's golfing</b>						
Trump golfs	0.25	0.43	0.00	0.00	1.00	365
Trump golfs (NYT only)	0.24	0.43	0.00	0.00	1.00	365
Trump golfs (alternative coding)	0.25	0.44	0.00	0.00	1.00	365
Golf holiday	0.16	0.37	0.00	0.00	1.00	365
Golf in previous week	0.75	0.43	0.00	1.00	1.00	365
<b>Other control variables</b>						
Google searches about Muslims (PC)	-0.27	1.98	-2.11	-0.59	21.51	365
Terror attack in the West	0.03	0.17	0.00	0.00	1.00	365

*Notes:* This table presents descriptive statistics for the IV sample. The sample year is 2017.  $1+\log$  or  $1+\text{natural logarithm}$  means that the logarithm of any variable is calculated with 1 added inside. The data on hate crimes come from the FBI hate crime statistics. Data on Trump's golfing come from the New York Times, the official White House presidential schedule, and trumpgolffcount.com. *Google searches about Muslims (PC)* is the first principal component of Google trends for the key words "islam", "mosque", "muslim", "refugee", "sharia", and "terror". We use these same keywords as measures of TV news attention based on data from the internet archive. The sources for the number of terror attacks is the Global Terrorism Database. See the online appendix for more details on data and variable construction.

**Table A.23: Summary Statistics by Day of Week (2017 only)**

Day of week		Hate crimes against Muslims	Tweets about Muslims	Trump golfs
Monday	Sum	43	3	4
	Mean	0.83	0.06	0.08
Tuesday	Sum	33	6	3
	Mean	0.63	0.12	0.06
Wednesday	Sum	43	10	4
	Mean	0.83	0.19	0.08
Thursday	Sum	43	6	6
	Mean	0.83	0.12	0.12
Friday	Sum	36	12	13
	Mean	0.69	0.23	0.25
Saturday	Sum	36	4	30
	Mean	0.69	0.08	0.58
Sunday	Sum	42	6	32
	Mean	0.79	0.11	0.60
<b>Total</b>	Sum	276	47	92
	Mean	0.76	0.13	0.25

*Notes:* This table presents descriptive statistics by day of week for the number of anti-Muslim hate crimes, the number of Trump’s tweets about Muslims and the number of Trump’s golf outing for the sample used in the instrumental variable regressions (2017 only).

Table A.24: Robustness Time Series 2SLS Regressions

	Baseline (1)	Add 7 lagged dependent variables (2)	Add golf holiday control (3)	Add previous week golf control (4)	Use Trump Tweet dummy (5)	Use only NYT golf count (6)	Use alternative golf count (7)
<b>Panel A: Log(Hate crimes against Muslims) in <math>t+2</math></b>							
Log(Muslim Trump tweets)	1.609** (0.791)	1.677* (0.948)	1.607** (0.820)	1.616** (0.791)	1.391* (0.727)	1.736** (0.831)	1.566* (0.821)
<b>Panel B: Log(News reports about Muslims) in <math>t</math></b>							
Log(Muslim Trump tweets)	2.701** (1.114)	2.689** (1.208)	2.614** (1.075)	2.673** (1.110)	2.335** (1.082)	2.672** (1.184)	2.869** (1.150)
<b>Panel C: Log(New content about Muslims by Trump Twitter followers) in <math>t</math></b>							
Log(Muslim Trump tweets)	1.151** (0.469)	1.008* (0.561)	1.233** (0.451)	1.155** (0.467)	0.996** (0.437)	1.032* (0.541)	1.206** (0.503)
Weak IV 95% AR confidence set	[0.177; 20.219]	[-0.048; 20.397]	[0.296; 20.348]	[0.184; 20.219]	[0.088; 20.076]	[-0.415; 20.157]	[0.160; 20.352]
Fixed effects (month, day of week)	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Time trends	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Observations	364	359	364	364	364	364	364
Robust F-stat.	13.02	12.08	12.80	13.40	11.81	10.45	12.56

Notes: This table presents IV regressions where the dependent variable is listed in the panel header. We use a dummy for days on which President Donald Trump golfs used as an instrument for his tweets about Muslims. Column 2 controls for seven lags of the dependent variable. Column 3 controls for the temperature on the golf day in Washington, D.C.. Column 4 controls for whether Trump golfed in the previous week. Column 5 replaces the number of Muslim Trump tweets with a dummy for whether Trump sends any tweet about Muslims. Column 6 replaces the main measure *Trump golfs* with one that only uses information from the New York Times (ignoring that contained in his presidential schedule). Column 7 uses an alternative golf count that incorporates information from *trumpgolfcount.com*. The sample year is 2017, for which we have information on Trump's golfing. All regressions include day-of-week and year-month dummies, linear and quadratic time trends as well as a dummy for whether Trump's golfing is the first of a series of golf days. See online appendix for more details on data and variable construction. Newey-West standard errors are reported in parentheses except in column 8. Weak IV 95% Anderson-Rubin (AR) confidence sets are calculated using the two-step approach of Andrews (2018) with the Stata package from Sun (2018). \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

**Table A.25: Time Series - Split By Pre-Existing Sentiment**

	No terror attacks (1)	Fox News Muslim Coverage	
		Low (2)	High (3)
<b>Panel A: First stage - Log(Trump tweets about Muslims)</b>			
Trump golfs	0.078*** (0.025)	0.085** (0.042)	0.121*** (0.047)
<b>Panel B: OLS - Log(Hate crimes against Muslims) in t+2</b>			
Log(Muslim Trump tweets)	0.194** (0.095)	0.120 (0.102)	0.074 (0.127)
<b>Panel C: Reduced form - Log(Hate crimes against Muslims) in t+2</b>			
Trump golfs	0.162** (0.077)	0.154* (0.082)	0.165 (0.118)
<b>Panel D: 2SLS - Log(Hate crimes against Muslims) in t+2</b>			
Log(Muslim Trump tweets)	2.094* (1.183)	1.815* (1.115)	1.370 (1.268)
Fixed effects (month, day of week)	Yes	Yes	Yes
Time trend	Yes	Yes	Yes
Observations	322	192	171
Robust F-stat.	8.78	3.59	5.84

*Notes:* This table presents OLS and IV regressions where the dependent variable is the number of hate crimes against Muslims on any given day based on FBI data. We use a dummy for days on which President Donald Trump golfs used as an instrument for his tweets about Muslims. Column 1 drops days with terror attacks from the sample. Columns 2 and 3 divide the sample based on whether the coverage of Muslim-related topics on Fox News on the day before the Trump tweet/golfing is above or below its median value. The sample year is 2017. All regressions include day-of-week and year-month dummies, linear and quadratic time trends as well as a dummy for whether Trump's golfing is the first of a series of golf days. See online appendix for more details on data and variable construction. Newey-West standard errors are reported in parentheses. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

Table A.26: Time Series - Split by Type of Hate Crime

	Any (1)	Vandalism (2)	Theft (3)	Burglary (4)	Robbery (5)	Assault (6)
<b>Panel A: OLS</b>						
Log(Muslim Trump tweets)	0.109 (0.071)	0.027 (0.053)	0.023 (0.033)	0.093** (0.042)	0.011 (0.014)	0.009 (0.063)
<b>Panel B: Reduced form</b>						
Trump golfs	0.164** (0.069)	0.136** (0.055)	-0.003 (0.014)	0.022 (0.015)	-0.007 (0.013)	0.071 (0.069)
<b>Panel C: 2SLS</b>						
Log(Muslim Trump tweets)	1.609** (0.791)	1.338** (0.621)	-0.033 (0.132)	0.216 (0.148)	-0.065 (0.131)	0.693 (0.712)
Weak IV 95% AR confidence set	[0.278; 40.036]	[0.293; 30.245]	[-0.308; 0.268]	[-0.091; 0.581]	[-0.441; 0.156]	[-0.646; 20.595]
Fixed effects (month, day of week)	Yes	Yes	Yes	Yes	Yes	Yes
Time trends	Yes	Yes	Yes	Yes	Yes	Yes
Observations	363	363	363	363	363	363
Robust F-stat.	13.15	13.15	13.15	13.15	13.15	13.15

*Notes:* This table presents OLS and IV regressions where the dependent variable is the number of hate crimes against Muslims on any given day based on FBI data. We use a dummy for days on which President Donald Trump golfs used as an instrument for his tweets about Muslims. The sample year is 2017, for which we have information on Trump's golfing. All regressions include day-of-week and year-month dummies, linear and quadratic time trends as well as a dummy for whether Trump's golfing is the first of a series of golf days. See online appendix for more details on data and variable construction. Newey-West standard errors are reported in parentheses. Weak IV 95% Anderson-Rubin (AR) confidence sets are calculated using the two-step approach of Andrews (2018) with the Stata package from Sun (2018). \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

**Table A.27: Time Series - Split by Motivating Bias**

	All (1)	Hispanic (2)	Other Ethnicity (3)	Race (4)	Sexual Orientation (5)	Religion (excl. Muslims) (6)
<b>Panel A: OLS</b>						
Log(Muslim Trump tweets)	0.108** (0.049)	0.033 (0.076)	0.236** (0.100)	0.015 (0.072)	0.051 (0.071)	0.136* (0.073)
<b>Panel B: Reduced form</b>						
Trump golfs	0.035 (0.049)	-0.149** (0.064)	0.046 (0.078)	0.054 (0.060)	0.007 (0.067)	0.056 (0.063)
<b>Panel C: 2SLS</b>						
Log(Muslim Trump tweets)	0.343 (0.465)	-1.465* (0.769)	0.450 (0.749)	0.529 (0.586)	0.065 (0.666)	0.547 (0.580)
Weak IV 95% AR confidence set	[-0.717; 10.310]	[-30.975; -0.323]	[-10.255; 20.007]	[-0.689; 10.864]	[-10.188; 10.714]	[-0.887; 10.752]
Fixed effects (month, day of week)	Yes	Yes	Yes	Yes	Yes	Yes
Time trends	Yes	Yes	Yes	Yes	Yes	Yes
Observations	363	363	363	363	363	363
Robust F-stat.	13.15	13.15	13.15	13.15	13.15	13.15

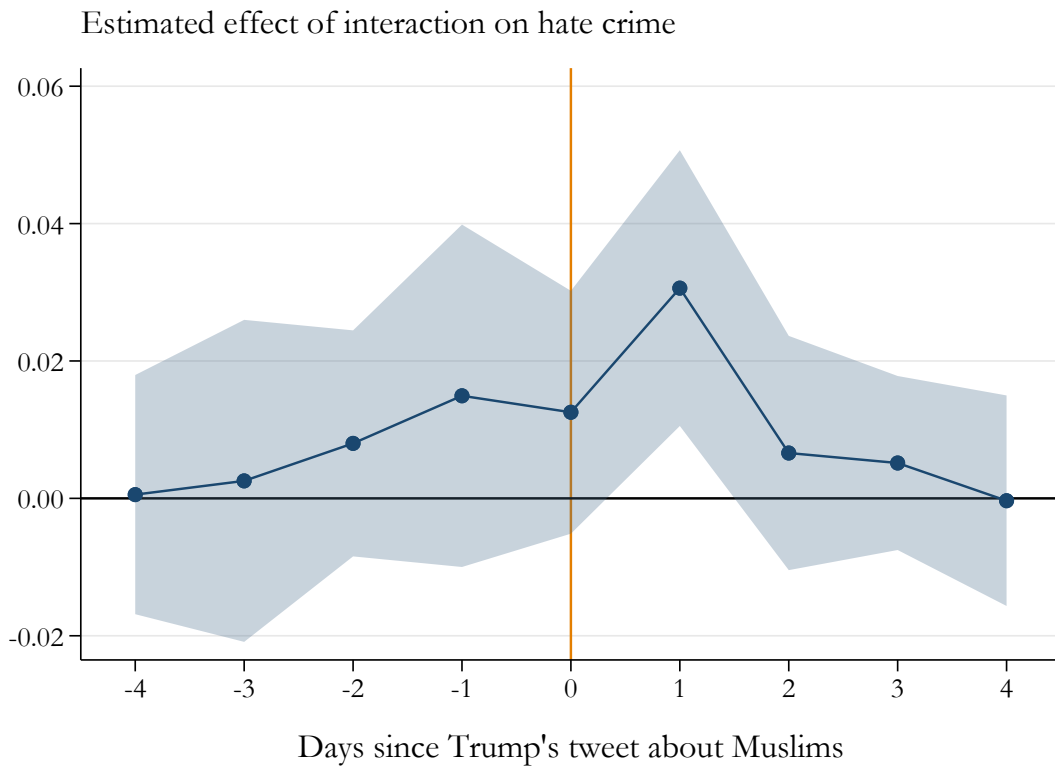
*Notes:* This table presents OLS and IV regressions where the dependent variable is the number of hate crimes against the group in the top row on any given day based on FBI data. We use a dummy for days on which Trump golfs used as an instrument for his tweets about Muslims. The sample year is 2017, for which we have information on Trump's golfing. All regressions include day-of-week and year-month dummies, linear and quadratic time trends, as well as a dummy for whether Trump's golfing is the first of a series of golf days. See online appendix for more details on data and variable construction. Newey-West standard errors are reported in parentheses. Weak IV 95% Anderson-Rubin (AR) confidence sets are calculated using the two-step approach of Andrews (2018) with the Stata package from Sun (2018). \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

**Table A.28: Time Series Regression Full Period**

	Baseline (1)	Add lagged dependent variable (2)	Add total tweets control (3)	Use Trump Tweet dummy (4)
<b>Panel A: Before campaign announcement</b>				
Log(Muslim Trump tweets)	0.009 (0.007)	0.009 (0.007)	0.007 (0.007)	0.028 (0.035)
Fixed effects (year, month of year, day of week)	Yes	Yes	Yes	Yes
Time trends	Yes	Yes	Yes	Yes
Observations	2,234	2,233	2,234	2,234
$R^2$ (partial)	0.00	0.00	0.00	0.00
<b>Panel B: After campaign announcement</b>				
Log(Muslim Trump tweets)	0.039** (0.016)	0.037** (0.016)	0.035** (0.016)	0.121** (0.057)
Fixed effects (year, month of year, day of week)	Yes	Yes	Yes	Yes
Time trends	Yes	Yes	Yes	Yes
Observations	1,295	1,294	1,295	1,295
$R^2$ (partial)	0.01	0.01	0.01	0.01

*Notes:* This table presents OLS regressions where the dependent variable is the number of hate crimes against the group in the top row on any given day based on FBI data. The sample is split into the period before and after June 16, 2015 when Trump announced his presidential campaign. All regressions include day-of-week and year-month dummies as well as linear and quadratic time trends. Partial  $R^2$  excludes these controls. See online appendix for more details on data and variable construction. Newey-West standard errors are reported in parentheses. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

Figure A.9: Panel Event Study - Trump Tweets, Twitter Usage, and Hate Crimes

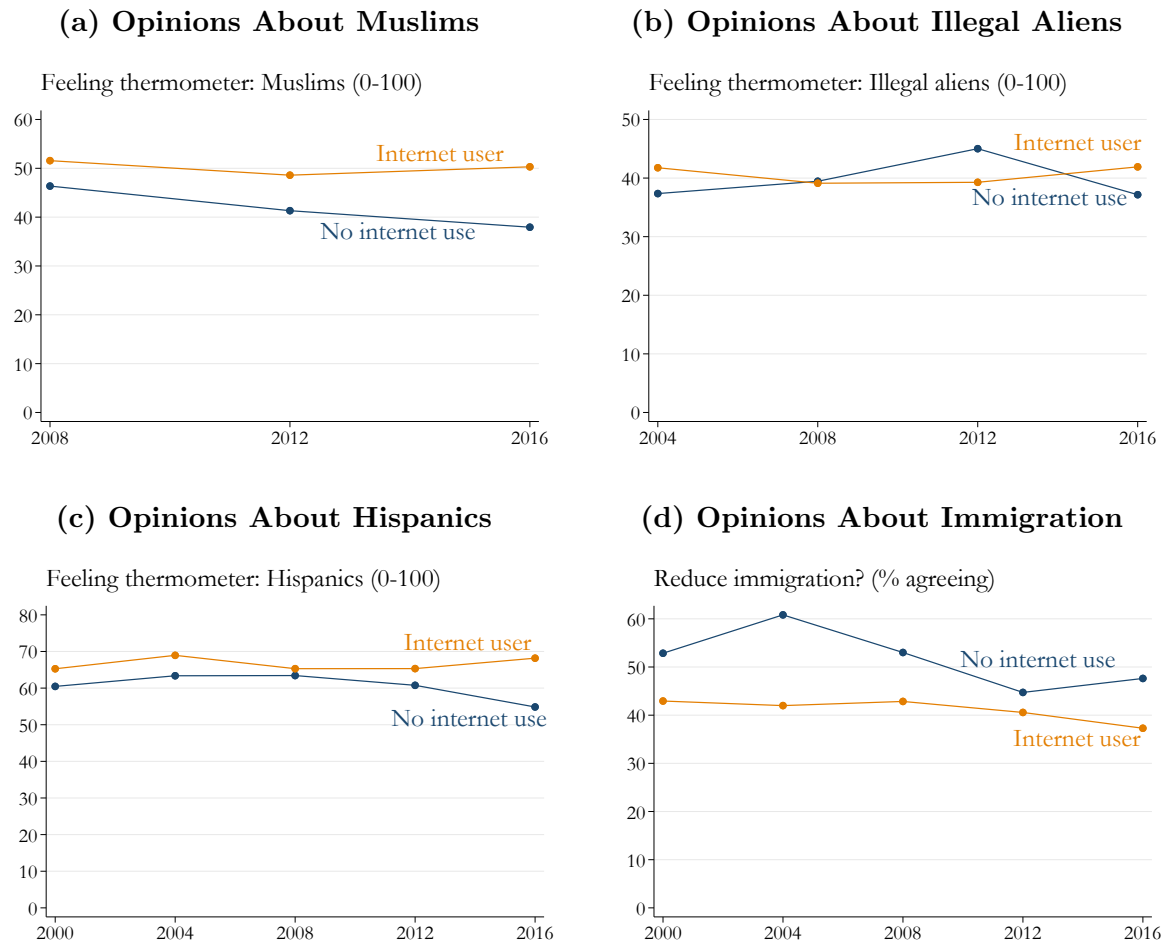


*Notes:* This figure plots the coefficients  $\beta_t$  from a dynamic version of Equation (6), where we allow values of  $t$  between  $-4$  and  $4$  days around Donald Trump's tweets about Muslims. The dependent variable is an indicator for anti-Muslim hate crimes in county  $i$  on day  $t$ . The coefficients are multiplied by 100 for readability. The regression also includes population controls, interacted with day dummies, state  $\times$  day fixed effects, and county  $\times$  day-of-week fixed effects, and county  $\times$  day-of-month fixed effects. The shaded areas are 95% confidence intervals based on standard errors clustered by state.



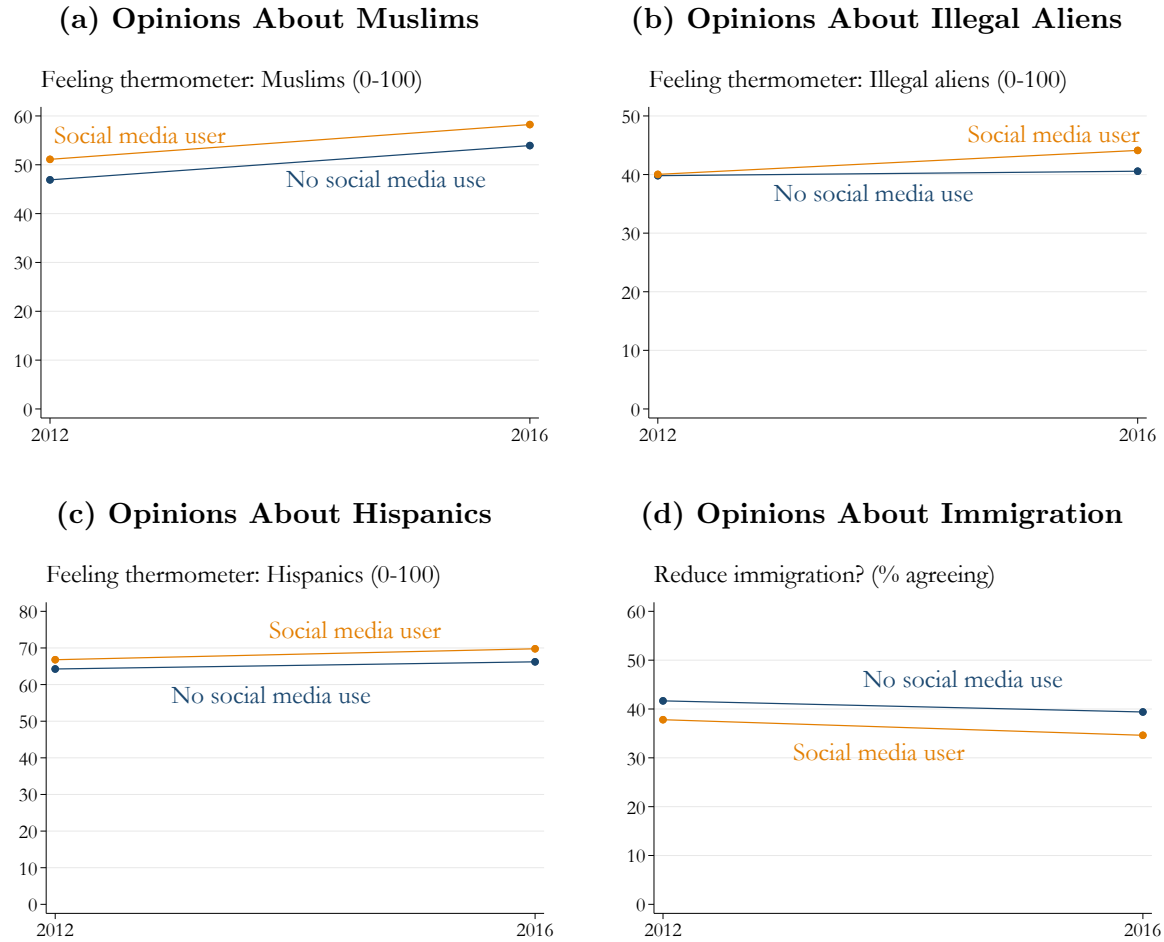
## A.4. Appendix 5: Additional Evidence for Mechanism

Figure A.10: Americans' Attitudes Towards Minorities, By Internet Use



*Notes:* These figures show Americans' attitudes towards Muslims, illegal aliens, Hispanics, and immigration more generally over time based on data from the American National Election Studies (ANES). The feeling thermometers in Panels (a) through (c) proxy for people's general attitude towards minorities. Panel (d) is coded from the question "Should the number of immigrants permitted to come to the U.S. be increased, decreased, or should number be the same as now?". We report the share of respondents agreeing with reducing immigration based on survey weights.

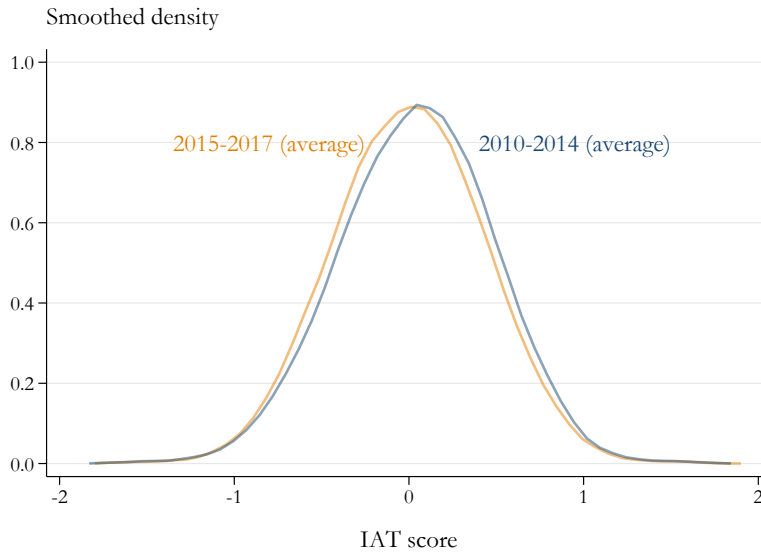
**Figure A.11: Americans' Attitudes Towards Minorities, By Social Media Use**



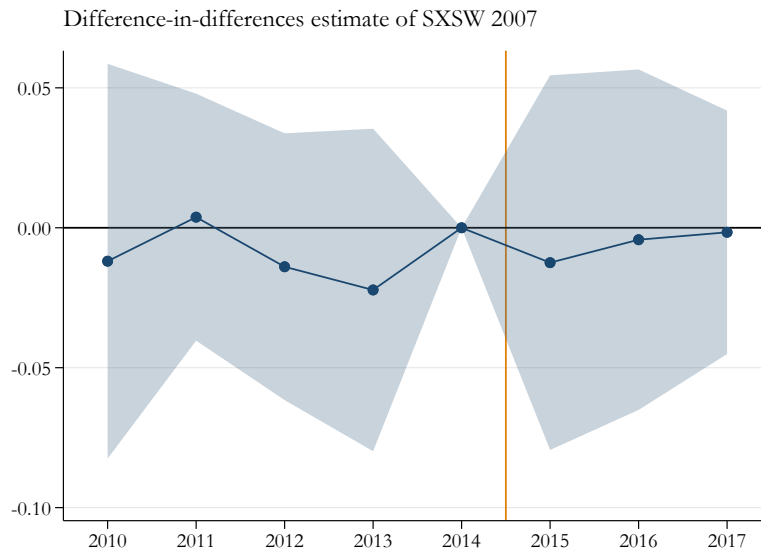
*Notes:* These figures show Americans' attitudes towards Muslims, illegal aliens, Hispanics, and immigration more generally over time based on data from the American National Election Studies (ANES). The feeling thermometers in Panels (a) through (c) proxy for people's general attitude towards minorities. Panel (d) is coded from the question "Should the number of immigrants permitted to come to the U.S. be increased, decreased, or should number be the same as now?". We report the share of respondents agreeing with reducing immigration based on survey weights.

## Figure A.12: Evidence from IAT Scores

### (a) Implicit Bias Against Muslims, 2010-2017



### (b) Reduced-Form Effect of Social Media on Implicit Bias



*Notes:* Panel (a) shows two density plots for the distribution of IAT scores measuring implicit bias against Muslims for the periods 2010-2014 (before Donald Trump’s political rise) and 2015-2017. The source of the data is Project Implicit’s Arab-Muslim module. Panel (b) plots the coefficients from a panel event study regression as in Equation (1). The dependent variable is the mean county-level IAT score from Project Implicit, which measures implicit bias against Muslims, residualized with regard to education, gender, ethnicity, race, age, and age squared (as in column 2 of Table A.29). We plot the coefficient for the log number of SXSW followers who joined in March 2007 and control for county and year fixed effects as well as SXSW followers who joined before the 2007 festival. We standardize the variables to have a mean of zero and standard deviation of one. The vertical line indicates the approximate start of the 2016 presidential primaries. The shaded areas are 95% confidence intervals based on standard errors clustered by state.

Table A.29: Social Media and Changes in Implicit Bias Against Muslims

	Raw IAT scores (1)	Residual IAT scores (2)	Only conservatives (3)	Only whites (4)	Only Christians (5)	Only non-Muslims (6)	Only obligatory tests (7)	At least 10 tests (8)
<b>Panel A: OLS</b>								
Log(Twitter users)	0.002 (0.016)	-0.001 (0.015)	0.001 (0.030)	-0.014 (0.016)	-0.008 (0.014)	-0.005 (0.017)	-0.014 (0.020)	0.008 (0.008)
<b>Panel B: Reduced form</b>								
Log(SXSW followers, March 2007)	0.003 (0.012)	-0.006 (0.012)	-0.000 (0.024)	0.001 (0.012)	-0.003 (0.016)	-0.003 (0.011)	0.005 (0.016)	0.009 (0.007)
<b>Panel C: 2SLS</b>								
Log(Twitter users)	0.006 (0.026)	-0.014 (0.028)	-0.001 (0.057)	0.003 (0.026)	-0.006 (0.036)	-0.007 (0.024)	0.012 (0.036)	0.025 (0.020)
Weak IV 95% AR confidence set	[-0.047; 0.049]	[-0.070; 0.032]	[-0.105; 0.103]	[-0.050; 0.046]	[-0.071; 0.060]	[-0.055; 0.032]	[-0.054; 0.078]	[-0.010; 0.065]
Log(SXSW followers, Pre)	-0.025 (0.019)	-0.025 (0.021)	-0.035 (0.043)	-0.048** (0.022)	-0.004 (0.020)	-0.025 (0.019)	-0.030* (0.017)	-0.022** (0.009)
Observations	2,113	2,075	1,181	2,002	1,832	2,086	1,570	833
Mean of DV	-0.031	-0.026	-0.022	-0.035	-0.024	-0.030	-0.042	-0.042
Robust F-stat.	70.02	72.71	84.18	74.94	82.83	71.06	87.28	62.38

Notes: This table presents county-level OLS and IV regressions where the dependent variable is the change in average Implicit Association Test (IAT) scores that measures implicit bias against Muslims between 2010-2014 and 2015-2017. Higher scores reflect more bias. *Log(Twitter usage)* is instrumented using the number of users who started following SXSW in March 2007. *SXSW followers*, *Pre* is the number of SXSW followers who registered at some point in 2006. All regressions control for population deciles and state fixed effects (not shown). In column 2, IAT scores are residualized with respect to age and its squared term, as well as a full set of fixed effects for educational attainment, race, sex, and ethnicity. In columns 3 through 6, the sample is restricted to respondents as indicated in the top row. Column 7 only includes tests that are obligatory, e.g. as part of a work program. Column 8 restricts the sample to counties with at least 10 IAT tests before and after Trump's presidential run. Weak IV 95% Anderson-Rubin (AR) confidence sets are calculated using the two-step approach of Andrews (2018) using the Stata package from Sun (2018). For the just-identified case we study here, the "robust" *F*-stat. is equivalent to the "Kleibergen-Paap" or the "effective" *F*-statistic of Olea & Pflueger (2013). Robust standard errors in parentheses are clustered by state. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .

**Table A.30: Heterogeneous Effects - Hate Groups and Hate Crimes**

Dependent variable:	(1)	(2)	(3)	(4)
Log(Anti-Muslim hate crimes)	No hate groups	Any hate group	Few hate crimes	Many hate crimes
Log(Twitter Usage) x Year=2010	-0.01* (0.01)	0.01 (0.03)	-0.00 (0.00)	0.00 (0.03)
Log(Twitter Usage) x Year=2011	-0.00 (0.01)	0.00 (0.03)	0.00 (0.00)	0.00 (0.03)
Log(Twitter Usage) x Year=2012	0.00 (0.01)	-0.01 (0.04)	0.00 (0.00)	-0.01 (0.04)
Log(Twitter Usage) x Year=2013	-0.00 (0.00)	-0.00 (0.03)	-0.00 (0.00)	0.00 (0.03)
Log(Twitter Usage) x Year=2015	0.01 (0.01)	0.09*** (0.03)	0.00 (0.00)	0.11*** (0.03)
Log(Twitter Usage) x Year=2016	0.01 (0.01)	0.14*** (0.03)	0.01* (0.00)	0.15*** (0.03)
Log(Twitter Usage) x Year=2017	-0.00 (0.01)	0.06* (0.03)	0.00 (0.00)	0.04 (0.04)
County FE	Yes	Yes	Yes	Yes
Year FE	Yes	Yes	Yes	Yes
Pop. deciles x Year FE	Yes	Yes	Yes	Yes
Observations	22,024	2,832	22,344	2,504

*Notes:* This table presents panel event study regressions where the dependent variable is the log number of hate crimes against Muslims (with one added inside). We standardized the variables to have a mean of zero and standard deviation of one. The sample period is 2010 to 2017. 2014 is the excluded period.  $\text{Log}(SXS\text{W followers})$  is the number of local SXS\text{W} followers that joined Twitter in March 2007. The existence of hate groups is based on data from the Southern Poverty Law Center (SPLC). The number of hate crimes in the pre-period is based on the total number of hate crimes per capita the FBI registered in a county from 2010 until 2015, split at the 90th percentile. All regressions control for the interaction of population deciles with year dummies. Standard errors in parentheses are clustered by state. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .