

# Using Crowdsourced Data to Study Crime and Place

David Buil-Gil and Reka Solymosi

Department of Criminology, University of Manchester, UK

25/05/2020

## Abstract

Crowdsourcing refers to the practise of enlisting the knowledge, experience or skills of a large number of people (the crowd) through some digital platform to collect data towards a collaborative project. Crowdsourcing can generate large volumes of data in relatively little time at a very small cost, and can be useful for research, strategic police management and many other purposes. To make effective use of crowdsourced data, it is important to understand its key strengths to emphasize, and limitations to mitigate. In this chapter we highlight the main strengths and weaknesses of crowdsourcing, and illustrate how to acquire, make sense of, and critically evaluate crowdsourced data to study crime and place. We present a step-by-step exemplar study using crowdsourced data from a platform called Place Pulse, where people rate their feelings of safety between different areas. Taking the case study of Atlanta, Georgia, we work through analyzing and interpreting these data while highlighting how to emphasize and evaluate the strengths and limitations of crowdsourcing. Exercises are presented using R software.

**Keywords:** Fear of crime, perceived safety, crime mapping, open data, GIS, Atlanta

**Full reference:** Buil-Gil, D., & Solymosi, R. (2020). Using crowdsourced data to study crime and place. In E. Groff & C. Haberman (Eds.), *The study of crime and place: A methods handbook*. Temple University Press.

**Contact details:** David Buil-Gil. G18 Humanities Bridgeford Street Building, Cathie Marsh Institute for Social Research, University of Manchester. E-mail address: [david.builgil@manchester.ac.uk](mailto:david.builgil@manchester.ac.uk)

**ORCID IDs:** David Buil-Gil: 0000-0002-7549-6317. Reka Solymosi: 0000-0001-8689-1526.

## Introduction

Crowdsourcing refers to the practise of enlisting the knowledge, experience or skills of a large number of citizens (*the ‘crowd’*) to achieve a common goal or cumulative result, usually via a platform powered by online technologies, mobile phones, social media or a website (Howe, 2006). Digital platforms allow recording large volumes of data in relatively little time at a small cost, and such data is often utilized for a variety of functions ranging from academic research to policy making and emergency management (Goodchild, 2007; Hecker et al., 2019).

For example, during the 2007-2009 wildfires in the Santa Barbara area, California, residents shared their real-time knowledge about the location of fires and emergency shelters via various online forums and websites, which proved to be an invaluable source of information for disaster response (Goodchild & Glennon, 2010). Specific to research in crime and place, crowdsourcing projects have been used to understand people’s experiences with crime and their perceptions about space and safety (e.g., Solymosi & Bowers, 2018; Williams et al., 2017).

In this chapter we present some examples of how crowdsourced data can be used creatively for criminological research, and specifically highlight the strenghts and limitations of data produced from crowdsourcing platforms. We present a step-by-step exemplar study in R software (R Core Team, 2020) using crowdsourced perceptions of safety in Atlanta, Georgia.

# Crowdsourcing to study crime and place

In criminological research, crowdsourcing has been primarily used to harness data about various forms of crime and antisocial behavior and to process information about citizens' perceptions and emotions about crime. Crowdsourced data provide new angles of insight into people's behaviors and perceptions, thus allowing researchers to devise new explanations of crime and perceived safety.

Public perceptions and emotions about crime have traditionally been analyzed by using surveys and interview-type qualitative approaches (Gabriel & Greve, 2003), but these methods are costly and may be limited in their ability to capture the time- and context-specific emotional reactions of fear (Castro-Toledo et al., 2017). They may also fail to record any behavioral responses to such emotions, such as avoiding certain places or situations, or acquiring alarm systems or weapons. As an alternative, some researchers have endorsed the use crowdsourcing to record data about the specific places and times in which episodes of fear of crime are more frequent (Solymosi et al., 2020; Solymosi & Bowers, 2018).

For example, Hamilton et al. (2011) developed a mobile phone app to record public perceptions of crime on public transportation in Melbourne, Australia. Similarly, Solymosi et al. (2015) designed an app and asked participants to report their worry about crime, which allowed authors to map the users' fear of crime across different areas of London, UK. Birenboim (2016) developed a mobile app to record data about the perceptions of security of attendees at a music festival in Jerusalem, Israel. And Gómez et al. (2016) designed a collaborative web-based tool that allowed the citizens of Bogotá, Colombia, to report those areas in which they feel less safe.

People can report not only their perceived emotions, but also about things they see in their environments. For example, Solymosi et al. (2017) analyzed secondary data recorded from FixMyStreet, an online problem-reporting website, where citizens can report graffiti, broken street lights, and other signals of neighborhood disorder in London, UK.

Concerned with the effect of the built environment on people's perceptions of crime, Salesses et al. (2013) designed a website which presented people with two images from Google Street View, and asked them to choose 'which place looks safer'. This platform is called '*Place Pulse*', and has received thousands of views, with people all over the world evaluating images of places based on their feelings of safety. Then, based on these evaluations, Salesses et al. (2013) produced a map of perceived safety in New York. It is this specific project from which we will be analyzing data later in this chapter.

These are only a few examples, but there are many other crowdsourcing platforms that have been designed and utilized to study emotions about crime (see a review in Solymosi et al., 2020). Moreover, open data recorded from social media and online forums [LINK TO TWITTER CHAPTER], which can be considered to be specific forms of crowdsourced data, enable detecting various forms of online crimes (e.g., hate speech towards minority groups; Miró-Llinares et al., 2018), and even associate patterns of online communication with offline disorder and serious offences (Bendler et al., 2014; Williams et al., 2017, 2020).

## Strengths and weaknesses of crowdsourced data

Crowdsourced data about public perceptions of space and crime have some key strengths over data recorded from traditional survey methods. Due to the data being provided by people in real-time, using technology which can record auxiliary information such as GPS or time-stamp, besides the information people report we also get precise spatial data, information about immediate environmental variables, and other relevant information, without any additional cost to researchers or participants. However, the mode of production of crowdsourcing is also associated with certain limitations or weaknesses that, if uncontrolled, may affect the validity of such measures and the reliability and generalizability of our results (Buil-Gil et al., 2020; Elliott & Valliant, 2017).

Solymosi et al. (2020) conducted a systematic review of 27 studies utilizing or discussing the use of crowdsourcing to study perceptions and emotions about crime. Here we will summarize the key strengths and weaknesses identified in their review.

## Strengths

The most frequent strength of crowdsourcing and app-based methods identified by researchers was that these techniques allow capturing the spatial-temporal specific nature of fear of crime. Unlike traditional survey instruments, crowdsourcing data collection can be designed to generate point-level location data, and accurate-to-the-second time-stamp data with each report (Solymosi & Bowers, 2018). This is very beneficial for anyone carrying out crime and place research.

Another strength relevant to crime research is the ability for people to record data about the architectural features and environmental characteristics of spaces where they report (Chataway et al., 2017; Traunmueller et al., 2015), which ultimately allows to “un-erroneously associate them [perceptions about crime] with elements of the environmental backcloth such as incivilities, crime, and disorder” (Solymosi et al., 2015, p. 198). Users can provide photos of their environments (e.g., with FixMyStreet; Solymosi et al., 2017), can be asked to evaluate photos (e.g., with Place Pulse; Salesses et al., 2013), or such data can be linked from other sources by a common information, like GPS location, to conduct on-site observation of places (e.g., with InseguridApp; Solymosi et al., 2020).

Crowdsourcing can also produce large sample sizes, often at a very low cost. Dubey et al. (2016), for example, analyzed more than 350,000 votes of perceived safety recorded from the Place Pulse platform; and Solymosi et al. (2017) analyzed more than 275,000 reports of disorder in London. These large samples are very costly to record by using traditional probability surveys. In this chapter we will illustrate how to download data about more than 1.5 million votes registered from the Place Pulse platform, and we will analyze more than 37,000 votes of perceived safety in Atlanta.

## Limitations and weaknesses

Perhaps the main weakness of data recorded from crowdsourcing is related to participants’ self-selection. Probability surveys are carefully designed to select participants randomly, which means that all units in the population have equal probabilities of being chosen; whereas crowdsourcing projects harness data from non-probability samples who decide when and where to share their perceptions and emotions, and whether they want to participate at all (Elliott & Valliant, 2017). The mode of production of crowdsourced data increases the risk of self-selection bias, and as a consequence males and young citizens tend to be overrepresented in these data (Chataway et al., 2017), and citizens from deprived areas are generally less represented than persons from wealthy neighborhoods (Solymosi & Bowers, 2018). For example, Salesses et al. (2013) observed that 78.3% of participants who informed about their gender when using the Place Pulse platform were males, and Solymosi et al. (2017) highlight that only 26% of those who informed about their gender when reporting instances of disorder via FixMyStreet were females.

Even within the sample of self-selected participants there may be unequal participation. Many crowdsourcing platforms allow users to submit data multiple times. This may lead to participation inequality (or unequal participation), where “few users are responsible for most crowdsourced information, while the majority participate only a few times” (Buil-Gil et al., 2020, p. 6). To illustrate this, Dubey et al. (2016) show that 6,118 of the 81,730 persons who used the Place Pulse platform participated only once, while 30 users participated more than 1,000 times and the most prolific user voted 7,168 times. Solymosi et al. (2017) also show that one fourth of all FixMyStreet reports are produced by one percent of participants, and 73% of participants contribute only once.

It is not only people who may be unequally represented in these sorts of data, there may also be a bias in the places and times that do or do not feature prominently. Since users of crowdsourcing projects can decide where and when to participate, certain types of areas and times can be underrepresented. For instance, it is possible that app-based platforms fail to capture data from high-crime-density areas, since participants may avoid those places where they feel more exposed to crime (Innes, 2015). The routine activities of participants are also reflected on an under-representation of data points at night, when people are less likely to be out and about (Blom et al., 2010).

Moreover, many of these projects rely on the enthusiasm of the crowd of participants, which may die down over time (Blom et al., 2010). It is important to consider when such projects are launched, as the beginning is much more likely to see more active participation than later on in the project.

Finally, there are ethical considerations that may arise from the use of crowdsourcing, which are related to the user's privacy (e.g., risk of participants' identification) but also concerns that these techniques may sensitize participant and increase their fear of crime by asking them to constantly think about crime-related risks (Jackson & Gouseti, 2015). These last points are not so much limitations as key concerns that all researchers working with crowdsourced data should keep at the forefront of their minds.

## Crowdsourcing perceptions of safety: A step-by-step example in R

In order to illustrate the use of crowdsourcing in criminological research, we present an exemplar study using data recorded by the Place Pulse 2.0 platform mentioned earlier (see Salesses et al., 2013). This section will introduce the Place Pulse project and provide annotated R scripts for downloading, cleaning and exploring this source of crowdsourced data. Then, we will analyze the spatial distribution of crowdsourced perceptions of space and safety in Atlanta, and illustrate with examples how to explore some of the known issues of crowdsourced data discussed above in a real-world dataset. Thus, in this exercise you will acquire a number of different skills, including:

- Downloading and exploring crowdsourced data from Place Pulse.
- Cleaning and wrangling Place Pulse data to analyze perceptions of safety.
- Obtaining spatial data about those places assessed by Place Pulse participants as more or less safe.
- Mapping perceptions of safety from crowdsourced Place Pulse data.
- Examining some of the known issues of crowdsourcing using the Place Pulse data, including:
  - quantifying the effect of participation by *'supercontributors'* in our dataset,
  - quantifying the overall participation inequality,
  - analyzing the under-representation of areas, and
  - checking the decrease of participation over time.

### The Place Pulse project

Place Pulse 2.0 was an online crowdsourcing platform designed to record data about citizens' perceptions of a variety of topics including safety, beauty, wealth, liveability, boredom and depression in urban areas. Each topic had a related question. To assess safety, two images were shown to participants, who then were asked to answer *'Which place looks safer?'* (see Figure 1). Participants could also be asked which of the two images looked wealthier, more beautiful, more boring, livelier or more depressing, but we will focus on perceptions of space and safety in this chapter.

The images were selected randomly from Google Street View across 56 cities from 28 countries. All images were originally taken between 2007 and 2012. All the data collected were stored on the Place Pulse open website (<http://pulse.media.mit.edu/>), and available to everyone there, but the platform closed in late 2019. We have been granted access to all the data recorded between May 28th 2013 and August 22nd 2019 to write this chapter. All data have also been uploaded onto an open repository with consent of the data producers.

### Download and explore Place Pulse data

We have saved all Place Pulse data (more than 1.5 million votes) in a data repository on Figshare: [https://figshare.com/articles/Place\\_Pulse/11859993](https://figshare.com/articles/Place_Pulse/11859993). You can download this directly into R by using the `read.csv()` function. It is a large file so it may take some minutes to read in.

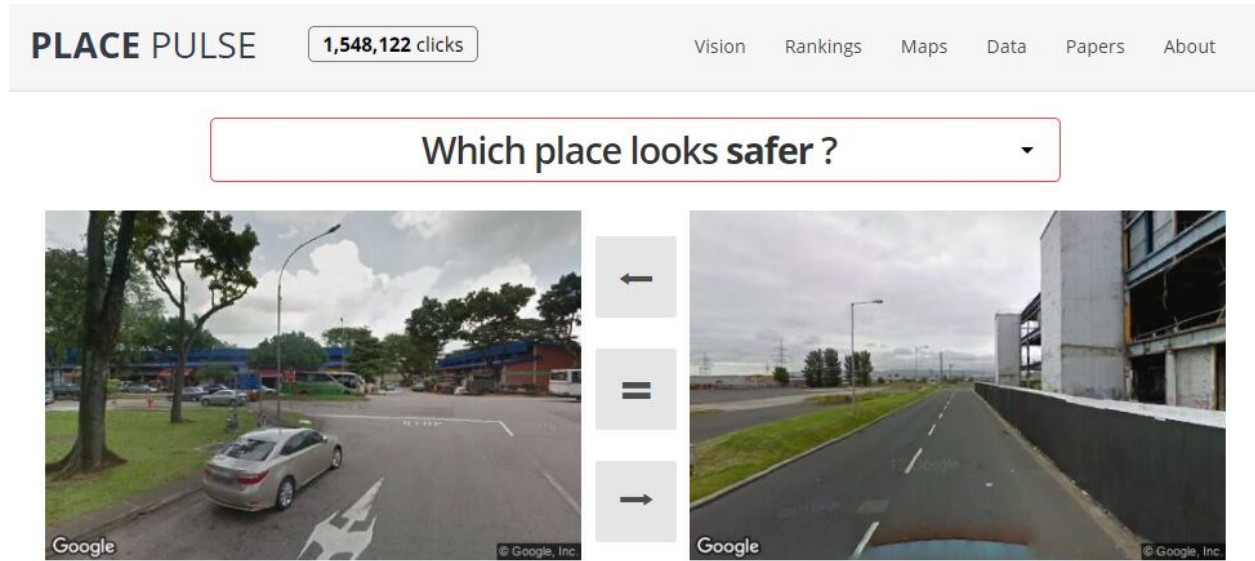


Figure 1: Figure 1: Place Pulse website

```
pp_data <- read.csv('https://ndownloader.figshare.com/files/21739137')
```

This dataset includes 17 variables, but we will only use the following variables:

- *'X'*: A unique identification code for each vote.
- *'right'*: A unique identification code for the image in the right side of the pairwise comparison.
- *'left'*: A unique identification code for the image in the left side of the pairwise comparison.
- *'voter\_uniqueid'*: A unique identification code given to each participant.
- *'place\_name\_left'*: The name of the city in the left side of the comparison.
- *'place\_name\_right'*: The name of the city in the right side of the comparison.
- *'choice'*: Which image was selected as 'safer' (left, right, or equal).
- *'study\_question'*: The variable of interest (e.g., safety, wealth, beauty).
- *'day'*: Day of the data point (vote).
- *'time'*: Time of the data point (vote).
- *'long\_right'*: Longitude for image on the right side of the comparison.
- *'lat\_right'*: Latitude for image on the right side of the comparison.
- *'long\_left'*: Longitude for image on the left side of the comparison.
- *'lat\_left'*: Latitude for image on the left side of the comparison.

The Place Pulse dataset has 1565723 observations. Each observation is one comparison between two images. Thus, our units of analysis here are comparisons between an image on the left and an image on the right.

We can start by answering some descriptive questions, for example: which cities are most frequently assessed within the Place Pulse platform? We know that for each row (each comparison) there are two images (**left** and **right**), and there are two columns that name the city for each image (**place\_name\_left**, **place\_name\_right**). To see how many times each city appears we can create two frequency tables (one for **place\_name\_left**, and one for **place\_name\_right**), and then join them and sum the two frequencies. We can use the **group\_by()** and **summarize()** functions from **dplyr** package (Wickham, François, et al., 2020) to achieve this:

```
library(dplyr) # load dplyr package

right_freq <- pp_data %>%
  group_by(place_name_right) %>% # categories based on cities on the right
  summarize(right_count = n()) # count number of units in each category

left_freq <- pp_data %>%
  group_by(place_name_left) %>% # categories based on cities on the left
  summarize(left_count = n()) # count number of units in each category

total_freq <- left_freq %>%
  left_join(., right_freq,
            by = c("place_name_left" = "place_name_right")) %>% # merge
  mutate(total_count = left_count + right_count) # create summatory column
```

Now we can see the top 3 most common cities using the `top_n()` function:

```
total_freq %>% top_n(3)

## # A tibble: 3 x 4
##   place_name_left left_count right_count total_count
##   <fct>           <int>      <int>      <int>
## 1 Atlanta         56992        57140       114132
## 2 Berlin          55265        55583       110848
## 3 Tokyo           53817        53514       107331
```

Atlanta appears to be the city with the largest number of votes. We can also check which variables (e.g., safety, beauty, wealth) were more frequently assessed by participants:

```
pp_data %>%
  group_by(study_question) %>% # categories based on study questions
  summarize(count = n()) %>% # count number of units in each category
  arrange(desc(count)) %>% # print in descending order
  filter(!is.na(study_question)) #remove NAs
```

```
## # A tibble: 6 x 2
##   study_question count
##   <fct>          <int>
## 1 safer          509961
## 2 livelier       366802
## 3 more beautiful 220604
## 4 wealthier      174758
## 5 more depressing 149355
## 6 more boring    144060
```

We see that ‘safety’ was the most commonly assessed variable, with 509,961 votes in total. In this chapter we will examine reports of safety in the city of Atlanta. Before analysing the data, however, we can also examine if participants were more inclined to vote for images in the left or right part of the platform. In other words, we analyze if responses were biased by the position in which images were shown on the website. This is an important step to confirm the validity of this study instrument. In survey design, much thought and research goes into elements like the ordering of the questions. It is important that we consider crowdsourced platforms with the same care and attention.

```
pp_data %>%
  group_by(choice) %>% # categories based on vote (right, left or equal)
  summarize(count = n()) %>% # count number of units in each category
  mutate(`%` = round(count/sum(count), 3)*100) %>% # compute percentage
  top_n(3) # print 3 most frequent categories
```

```
## # A tibble: 3 x 3
##   choice  count   `%`
##   <fct>   <int> <dbl>
## 1 equal  206147  13.2
## 2 left   668680  42.7
## 3 right  690792  44.1
```

The frequency of votes for left and right options is very similar; the left image wins about 43% of the time, the right one about 44% of the time, and the rest of users voted ‘equal’ (13%). Thus we can conclude that the position of the image on the website platform does not appear to have much of an affect on participants’ votes.

## Cleaning Place Pulse data

Our city of interest, Atlanta, is the capital city of the State of Georgia, United States. In 2018, its estimated population was close to 500,000 residents, and it is the 37th most populated city in the United States. It is also, as we have seen earlier, the city with the largest number of votes in the Place Pulse platform. There has been some considerable research looking into predictors of crime and fear of crime in this city, which can be used to interpret our findings later (see McNulty & Holloway, 2000; Tester et al., 2011).

In order to analyze perceptions of safety in Atlanta using our crowdsourced Place Pulse data, the first step is to clean the data to make it as complete and useful as possible to answer our research questions. For example, in this case, we want to map the perceived safety of areas in Atlanta. However, the Place Pulse dataset includes reports from all over the world, which means that we will need to select a subset of votes. Moreover, we would like our unit of analysis to be the Atlanta locations, rather than each comparison. There are a few steps that we need to take to make the data look like what we need to answer our questions. Specifically, we need to do the following:

- Select only votes about safety.
- Select only those votes that contain images of Atlanta.
- For each area in Atlanta, create a score from all the votes.

### Select only votes about safety

We can use the function `filter()` from `dplyr` to create a new dataframe that includes those pairwise comparisons which relate to safety:

```
pp_s <- pp_data %>%
  filter(study_question == "safer") # select votes of safety
```

This has created a new dataframe, `pp_s` which contains 509961 votes about perceived safety.

## Select votes that contain images of Atlanta

Now we want to select those votes where at least one of the images is from Atlanta. Remember that there are two images in each comparison (i.e., `left` and `right`). In order to select rows where Atlanta appears in either one or the other, or both, we can use the `'or'` operator (`|`) to ensure we get all comparisons which feature Atlanta on at least one side.

```
# select cases in which the image of the right or left is from Atlanta
pp_atl_s <- pp_s %>%
  filter(place_name_right == "Atlanta" | place_name_left == "Atlanta")
```

The new subset of cases has 37214 votes about the safety of places in Atlanta.

## Duplicate comparisons in which both images are from Atlanta and create new columns for analysis

In many cases, the dataset we need to answer our research question is slightly different to the dataset we have from the crowdsourced platform. This is often the case with much secondary data analysis, where data were created for another purpose initially. This often means that much data wrangling is to be done to make the data look like what we need.

We are interested in analyzing the proportion of 'safer' votes in each neighborhood of Atlanta. For this, we need the following information from each photo of Atlanta: 1) location (latitude and longitude) and 2) whether it won or lost (i.e., perceived as 'safer' or not) each vote.

First, we need every instance of a vote on each unique location in Atlanta. There are three possibilities for this:

- 1) Image on the right of the pairwise comparison is from Atlanta, but image on the left is somewhere else.
- 2) Image on the left of the pairwise comparison is from Atlanta, but image on the right is somewhere else.
- 3) Both images are from Atlanta.

For the first two cases, we can easily find the coordinates of the image in Atlanta, and whether it was perceived as 'safer' or not, using conditional statements in the `if_else()` function. We will create a dataset of only these votes and remove, for now, all rows in which both images are from Atlanta. We will merge all votes together once all data have been cleaned. Further, in order to assign photographs to their neighborhood, we need to create two new columns that specify the longitude and latitude of each image of Atlanta being assessed. We will name those columns `lat_Atl` and `long_Atl`. We will also create a new column that details whether each participant voted that the image of Atlanta was 'safer' than the other photograph or not (column `win`). In order to select those votes in which one of the images is from Atlanta, we can use the function `filter()` from `dplyr`, which we have also used above, but now, instead of the `'or'` operator (`|`) we use an `'and'` operator (`&`). We also compute the win score, and save a unique ID for each image:

```
atl_v_others <- pp_atl_s %>%
  # select where only one image is from Atlanta
  filter((place_name_right == "Atlanta" & place_name_left != "Atlanta") |
         (place_name_right != "Atlanta" & place_name_left == "Atlanta")) %>%
  # get coordinates from Atlanta side, label win, and ID of image
  mutate(lat_Atl = if_else(place_name_right == "Atlanta",
                           lat_right, lat_left),
         long_Atl = if_else(place_name_right == "Atlanta",
                           long_right, long_left),
         win = if_else((place_name_left == "Atlanta" &
```



```

        choice == "left") |
        (place_name_right == "Atlanta" &
         choice == "right"), 1, 0),
  unique_id = if_else(place_name_right == "Atlanta", right, left))

```

Some pairwise comparisons, however, assessed two different images from Atlanta, which means that we will need to duplicate those votes to account for both images. For this we can create a dataset for the right side and another dataset for the left side, also creating the new latitude (`lat_Atl`), longitude (`long_Atl`), and win (`win`) columns for each using `mutate()`:

```

#create dataset from Atlanta images on right side
right_side <- pp_atl_s %>%
  filter(place_name_right == "Atlanta" & place_name_left == "Atlanta") %>%
  mutate(lat_Atl = lat_right,
         long_Atl = long_right,
         win = if_else(choice == "right", 1, 0),
         unique_id = right)

#create dataset from Atlanta images on left side
left_side <- pp_atl_s %>%
  filter(place_name_right == "Atlanta" & place_name_left == "Atlanta") %>%
  mutate(lat_Atl = lat_left,
         long_Atl = long_left,
         win = if_else(choice == "left", 1, 0),
         unique_id = left)

```

Finally, we can merge all the data together with the `rbind()` function:

```
pp_atl_s <- rbind(atl_v_others, left_side, right_side)
```

We have a dataframe of 37892 votes about the safety in Atlanta that is ready to be analyzed, which now fits our criteria, in that:

- it contains only votes about safety,
- it contains only votes about Atlanta, and
- each row is a win (or no win) of a vote where an image from Atlanta was chosen as safer (or not safer).

We can now take this dataset and draw some conclusions.

## Environmental correlates of safety

Researchers may be interested in analyzing the environmental characteristics of those places assessed by Place Pulse users as safe or not safe. From such data we can get the coordinates of places rated, and use tools such as Google Street View to conduct a virtual environmental audit of these places. For example, we can look at the safest and least safe rated places using our dataset in Atlanta. For this, we need a dataset where each unique image is our unit of analysis, and we can compute a win score by considering the proportion of votes which that image has won:

```

unique_pp_wins <- pp_atl_s %>%
  group_by(unique_id) %>% # groups by images IDs
  summarize(winscore = mean(win, na.rm = TRUE), # proportion win

```

```

num_votes = n(), # count number of votes
latitude = first(lat_Atl), # save lat
longitude = first(long_Atl)) %>% # save long
arrange(desc(winscore)) # order by proportion 'safer' responses

```

This may be used by criminologists to observe the characteristics of those places assessed by Place Pulse participants as more or less safe. For example, we can use Google Street View [<https://www.google.com/maps>] to observe places rated as the least safe or safest amongst those rated at least 25 times.

```

unique_pp_wins %>%
  filter(num_votes >= 25) %>% # filter out pictures with less than 25 votes
  filter(row_number() == 1 | row_number() == n()) # highest and lowest score

```

```

## # A tibble: 2 x 5
##   unique_id          winscore num_votes latitude longitude
##   <fct>          <dbl>     <int>     <dbl>     <dbl>
## 1 513d7beffd9f03587006d1b 0.8529412         34 -84.46837  33.69525
## 2 513d9b54fdc9f035870079a2 0.16             25 -84.40463  33.73718

```

In this case, as shown in Figure 2, the least safe place is characterized by signs of physical disorder (i.e., graffiti, abandoned cars, rubbish lying around) which may increase negative emotions about crime (Toet & Schaik, 2012). Moreover, there are corners and hidden spaces between abandoned cars that can be perceived to offer concealment for possible criminals and obstruct the view onto certain spaces (see Fisher & Nasar, 1992); whereas the safest place is a wide street of a well-maintained residential area with green spaces, direct visual access to most places around it (large prospect), and natural surveillance from the house (Welsh & Farrington, 2004). We can do much more, but here we will focus on the specific issues to explore due to the crowdsourced nature of these data.

## Mapping Place Pulse data

We can use mapping techniques learned in other chapters (**LINK WITH MAPPING CHAPTERS**) to map crowdsourced data. We will be using the `sf` (Pebesma, 2020) and `ggplot2` (Wickham, Chang, et al., 2020) libraries in order to create a map of perceived safety of built environment across the areas of Atlanta.

First, we acquire a polygon which represents the census tracts of Atlanta. The Georgia Association of Regional Commission, for instance, publishes spatial data for Atlanta Region at the different spatial scales. We can go on their website to find out more about this boundary data: <https://opendata.atlantaregional.com/datasets/census-2000-tracts-atlanta-region>. We can download the shapefile directly using their Application Program Interface (or API) and the `st_read()` function from the `sf` package:

```

library(sf)

# download geojson from Georgia Association of Regional Commissions open data
atl <- st_read("https://opendata.arcgis.com/datasets/04b79404794f43959cda4f8c3f1817e6_49.geojson")

```

Note: if this link does not work, we saved a local copy of this file, in that case you can download from the following link: [https://raw.githubusercontent.com/maczkoni/crowdsourcing\\_pp\\_chapter/master/geojson/Census\\_2000\\_Tracts\\_Atlanta\\_Region.geojson](https://raw.githubusercontent.com/maczkoni/crowdsourcing_pp_chapter/master/geojson/Census_2000_Tracts_Atlanta_Region.geojson).

We can see what this file looks like by using the `plot()` and `st_geometry()` functions to plot the geometry of the 'atl' object we created:

**Least safe place among Atlanta places rated by at least 25 Place Pulse users**



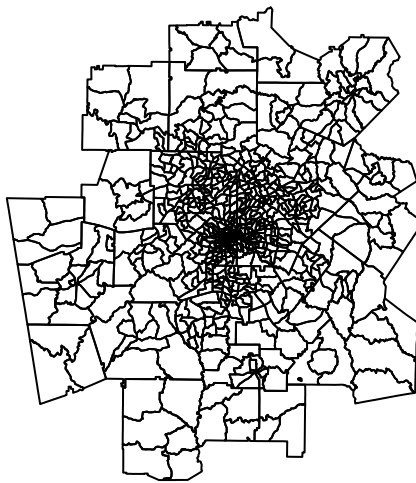
**Safest place among Atlanta places rated by at least 25 Place Pulse users**



Figure 2: Figure 2: Least safe and safest places rated by at least 25 Place Pulse users

```
plot(st_geometry(atl),
     main = "Atlanta Region census tracts")
```

## Atlanta Region census tracts



In order to plot the safety votes on this map, we first need to make our votes a spatial object, by specifying that the `'long_Atl'` and `'lat_Atl'` columns contain our longitude and latitude information. We use the `st_as_sf()` function for this:

```
#geocode votes
points_atl_s <- st_as_sf(pp_atl_s, coords = c("lat_Atl", "long_Atl"))
```

In order to plot both these spatial layers (i.e., votes recorded from Place Pulse and Atlanta census tracts) on the same map, their coordinate reference systems (CRS) need to match. We can assign the CRS from our polygon to our points layer:

```
st_crs(points_atl_s) <- st_crs(atl)
```

If we check, they should have the same CRS:

```
st_crs(points_atl_s) == st_crs(atl) #check if CRS is the same in both layers
```

```
## [1] TRUE
```

Now, to map 'safer' votes per census tract, we will compute a proportion of wins in each tract that will allow us to directly analyze the geographical distribution of perceived safety (see Buil-Gil et al., 2020).

```
p_atl_s_nhood <- st_intersection(atl, points_atl_s) %>% # points in areas
  group_by(TRACT) %>% # groups based on tracts
  summarise(winscore = mean(win, na.rm = TRUE), #proportion win
            num_votes = n()) %>% # count number of votes
  st_set_geometry(NULL) #remove geometry
```

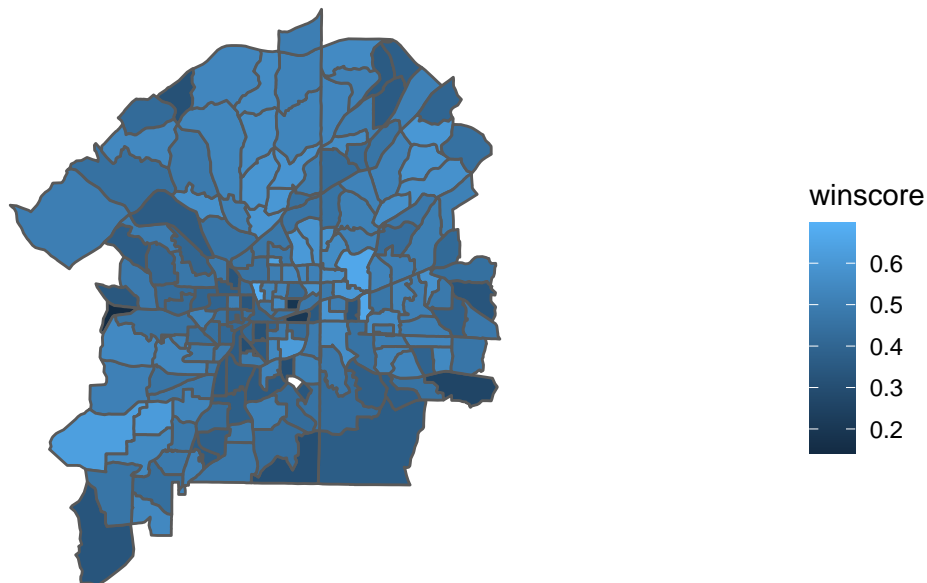
We have a dataframe which has the proportion of wins (`winscore`) for each tract. All that is left is to join this to our polygon of Atlanta (`atl`) and produce a map.

```
# merge census tracts and Place Pulse votes based on common 'TRACT' column
atl_pp_wins <- left_join(atl, p_atl_s_nhood, by = c("TRACT" = "TRACT")) %>%
  filter(!is.na(winscore)) # delete census tracts with 0 votes (NAs)
```

Finally, map the proportion of 'safer' votes in each census tract.

```
ggplot(data = atl_pp_wins) +
  ggtitle("Proportion of 'safer' votes per census tract") +
  geom_sf(aes(fill = winscore)) +
  coord_sf(xlim = c(-84.7, -84), ylim = c(33.5, 34), expand = FALSE) +
  theme_void()
```

Proportion of 'safer' votes per census tract



Areas in the south part of the city appear to have lower levels of perceived safety, whereas most areas in North Atlanta have higher values of perceived safety.



## Exploring known issues of crowdsourced data within Place Pulse

As we described earlier in the chapter, crowdsourced data comes with many possible issues which need to be properly understood. Here we illustrate how to explore some of these issues, working through examples from Place Pulse.

### How representative is the sample?

Place Pulse data contains votes from a self-selected sample, and unfortunately, the Place Pulse project did not record information about participants' demographic characteristics. Thus, we cannot directly examine the self-selection biases that may affect this dataset (Chataway et al., 2017; Elliott & Valliant, 2017). However, the sample's self-selection bias should be checked when possible. In this case, we do know that the first edition of Place Pulse did record some demographic variables from participants, and that in that iteration, 78.3% of participants were males, and only 21.7% were females, and the median self-reported age was 28 years (Saleses et al., 2013). We can expect that the version of Place Pulse we explore may have similar characteristics.

### Participation inequality (*'supercontributors'*)

Another issue is the participation inequality within the sample. Crowdsourced data tend to be affected by a few number of supercontributors that produce most votes (Dubey et al., 2016; Solymosi et al., 2017). In order to check if our dataset is affected by this, we can use the variable `voter_uniqueid` and produce a frequency table:

```
voter <- pp_data %>%  
  group_by(voter_uniqueid) %>% # create groups based on users unique id  
  summarise(num_votes = n()) # print the number of votes by user
```

We can have a look at this new dataframe using the `View()` function, and see that we have some *very active* participants. The top participant, for instance, has made 7168 votes on places. That is some very prolific participation. On the other hand, we can also see that 7494 of the participants made only one vote. We are definitely seeing signs of participation inequality in these data.

In fact, we can examine how many votes are produced by these 'supercontributors'. For example, we can assess the proportion of votes made by the top 1% of voters. We can do this using the `subset()` and `quantile()` functions:

```
# subset top 1% of most prolific participants  
top_1pc <- subset(voter, num_votes > quantile(num_votes, prob = 1 - 1/100))
```

We see that this new dataframe contains 954 people, who are our top 1% contributors to the Place Pulse dataset. We will now examine how much of the total number of votes are generated by the top 1% of users:

```
# proportion of votes by top 1% participants  
sum(top_1pc$num_votes) / sum(voter$num_votes) * 100
```

```
## [1] 17.87
```

That is a lot: 17.87% of the votes are made by the top 1% of contributors. We can also compute the proportion of votes made by the top 10% and 25% of participants in the same way. The top 10% contributors are responsible for 46.06% of votes, and the top 25% users contribute the 66.20% of all votes.

## Quantifying participation inequality

One way to quantify the extent to which participation inequality exists in our data is by using a Gini index, and visualizing it using a Lorenz curve. The Gini index (or Gini ratio) is a measure of statistical dispersion intended to measure inequality (Gastwirth, 1972). Although it is generally used to examine income inequality, it has also been frequently used to assess participation inequality in crowdsourcing platforms (see Solymosi et al., 2017; Solymosi & Bowers, 2018). Similarly, the Lorenz curve is a visual representation of inequality. For this we will need the `ineq` library (Zeileis & Kleiber, 2015). We can load this library and calculate the index using the `Gini()` function applied to our 'number of votes' column in our frequency table dataframe:

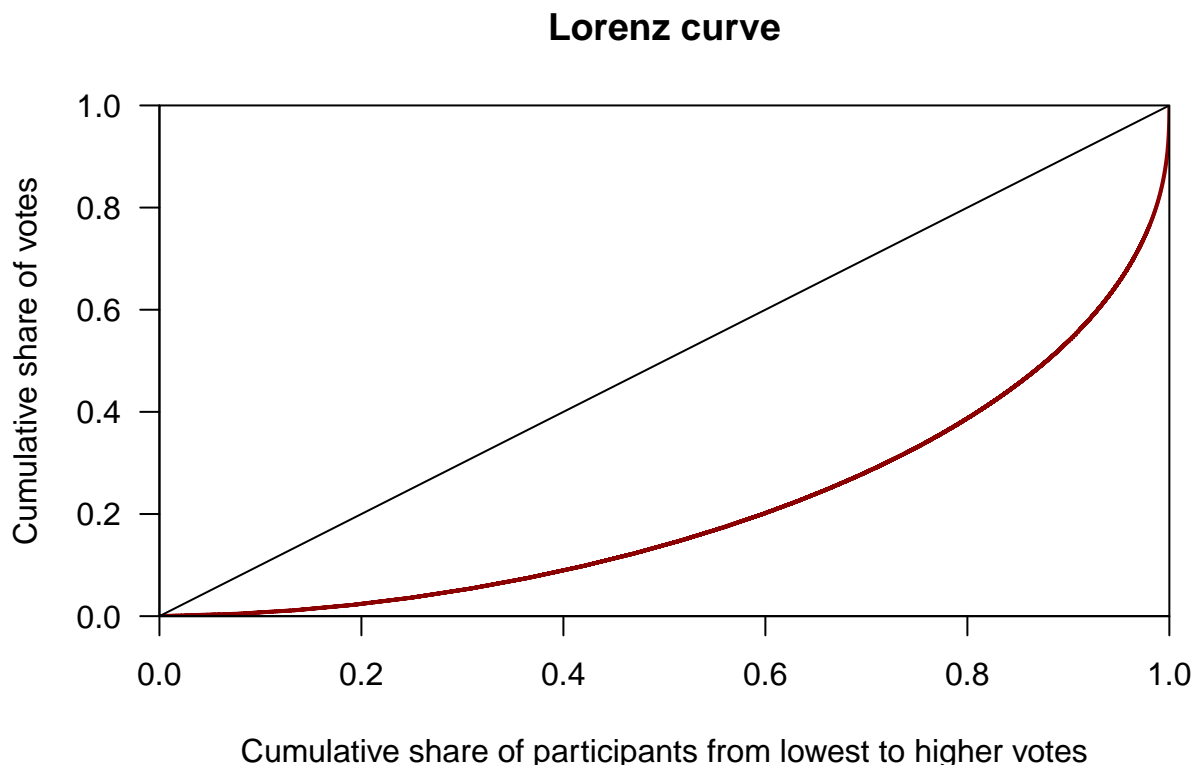
```
library(ineq)

Gini(voter$num_votes) # print Gini index
```

```
## [1] 0.58
```

To interpret this number, we can consider the following. A Gini index score of 0 represents perfect equality (everyone makes equal number of votes), while 1 shows perfect inequality (only one person making every single vote). Our answer of 0.58 shows some serious inequality. To put this into context, in 2017, according to the OECD, income inequality in the United States showed a Gini coefficient of 0.39. To visualize this we can use a Lorenz curve using the `plot()` and `Lc()` functions:

```
plot(Lc(voter$num_votes), # plot Lorenz curve
     xlab = "Cumulative share of participants from lowest to higher votes",
     ylab = "Cumulative share of votes", col = "darkred", lwd = 2)
```



The Lorenz curve (red line) shows how the top few percent of users contribute the majority of the reports. If we had perfect equality, we would expect to see the red line align perfectly with the black line with the slope of 1. With this information we can now quantify how severe the participation inequality is in our data, and compare with other crowdsourced data for context and understanding.

### Under-representation of certain areas

Another important consideration is the variation in the sample size across different areas. Some places might have many votes, while others not so much. While Place Pulse might not suffer from underreporting of high-crime-density areas due to people avoiding them (people are randomly presented with images from all over), it is still important to consider the number of votes in each census tract. We can analyze if certain areas are under-represented in our dataset by using the `summary()` function applied to the number of votes (`num_votes`) variable in our polygon object from earlier (`atl_pp_wins`).

```
summary(atl_pp_wins$num_votes)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##         2      85     151     193    248    1023
```

Whereas the average sample size per area is quite large (192.89), some tracts are clearly over-represented (the maximum number of votes is 1023) and others suffer from small representation (the minimum number of votes is only 2).

We may also want to know which areas suffer from severe under-representation in our dataset. We use the function `mutate()` to compute the number of votes divided by square miles in each census tract.

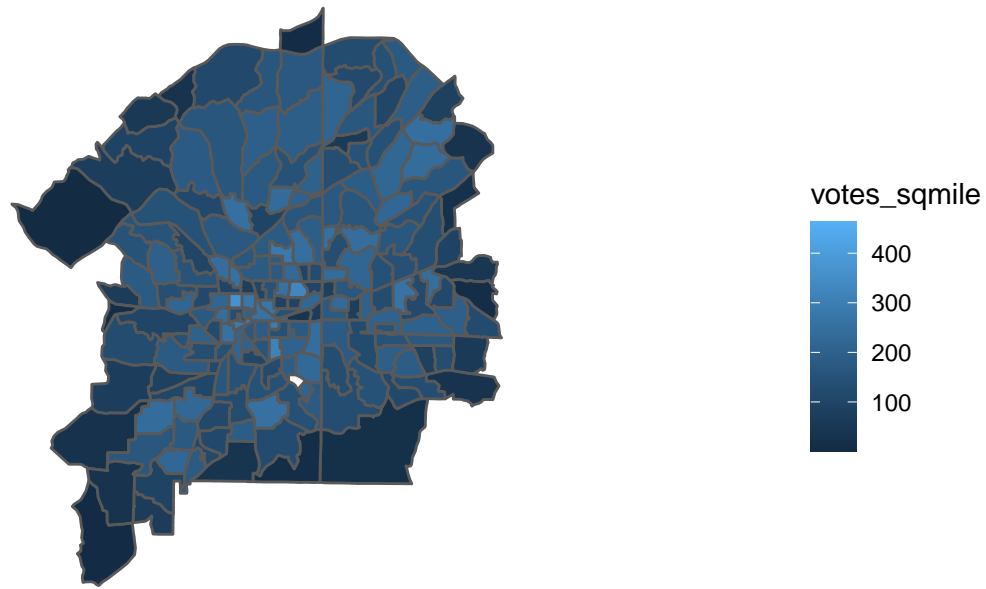
```
# compute new column of number of votes divided by square miles
atl_pp_wins <- atl_pp_wins %>%
  mutate(votes_sqmile = num_votes / SQ_MILES)
```

Then we can visualize the geographic distribution of the number of votes per census tract using the same code shown above to map perceptions of safety.

```
ggplot(data = atl_pp_wins) +
  ggtitle("Number of votes per square mile") +
  geom_sf(aes(fill = votes_sqmile)) +
  coord_sf(xlim = c(-84.7, -84), ylim = c(33.5, 34), expand = FALSE) +
  theme_void()
```



## Number of votes per square mile



We see that areas in the city center tend to have larger number of votes per square mile and are therefore well represented, whereas tracts in surrounding areas suffer from smaller sample sizes. Estimates of perceived safety in under-represented areas are likely to be affected by a small number of responses and may suffer from low precision. In order to increase the reliability of estimates produced from crowdsourced data for areas with small sample sizes, some researchers suggest using resampling and model-based techniques (see Arbia et al., 2018; Buil-Gil et al., 2020).

### Participation decrease

Finally, some researchers have identified that the number of users of crowdsourcing projects decreases over time: whereas the number of participants tends to be large during the first few days, users lose interest in the project if they do not obtain clear short-term benefits from using it (see Blom et al., 2010; Solymosi et al., 2020). We can also explore whether our Place Pulse dataset is affected by participation decrease.

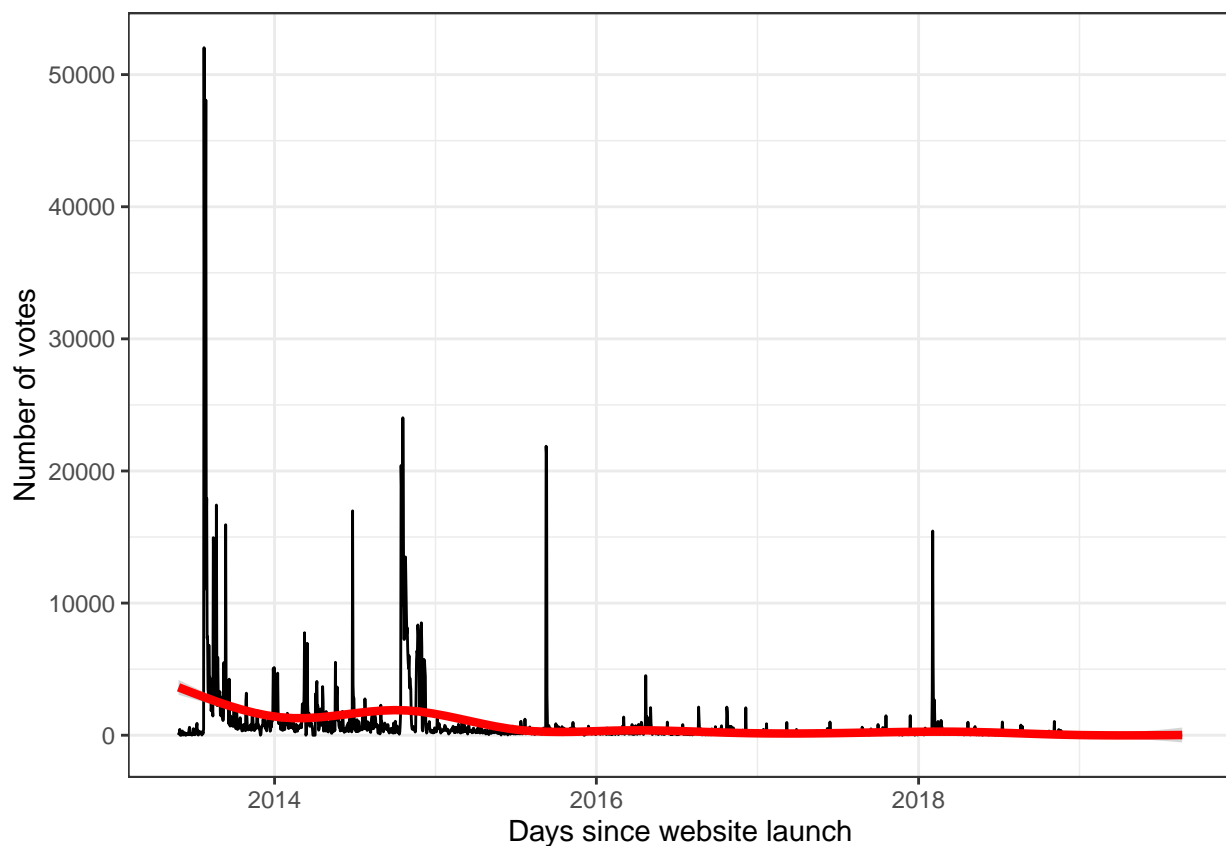
We will use various functions from `dplyr` and `ggplot2` packages (seen above) to visualize the number of votes between the Place Pulse website launch and its closure. But we also need to use other key functions: 1) the `ymd()` function from `lubricate` package is used to transform the dates in which votes took place into `Date` objects (Spinu et al., 2020), and 2) the `complete()` function from `tidyr` package is used to turn implicit missing dates into explicit missing dates and create a timeline of dates with and without votes (Wickham & Henry, 2020).

```
by_day <- pp_data %>%  
  mutate(day = ymd(day)) %>% # transform column into 'Date' object  
  group_by(day) %>% # create groups by days  
  summarise(num_votes = n()) %>% # count number of votes by day  
  complete(day = seq.Date(min(day),
```

```
max(day), by = "day")) %>% # complete days
mutate(num_votes = replace_na(num_votes, 0)) # replace NA by 0s
```

Once we have this dataset in the format we need to be able to look at reports by day, we can visualize the results using `ggplot`:

```
ggplot(by_day, aes(x = day, y = num_votes)) +
  geom_line() +
  geom_smooth(lwd = 1.5, col = "red") +
  theme_bw() +
  xlab("Days since website launch") +
  ylab("Number of votes")
```



The number of votes within the Place Pulse platform clearly decreased over time, but we also observe some peaks even years after the launching of the project. Some of these peaks match the dates of key publications using Place Pulse data and media releases, which shows that participation in crowdsourcing projects can be enhanced by periodic campaigns. For example, we observe a large peak beginning on July 24th 2013, date in which Salesses et al. (2013) published their paper and the Massachusetts Institute of Technology published a news article about the Place Pulse platform on their website: <http://news.mit.edu/2013/quantifying-urban-perceptions-0724>. We also observe another peak of participation beginning on October 15th 2014, just after the publication of Harvey (2014) Master's thesis about how to automate the study of the characteristics of streetscape skeletons and urban perceptions from Place Pulse data.

## Conclusions

The open data movement has provoked a revolution in social research methods, and will continue changing the way in which many social issues are researched, understood and managed. Digital technologies enable large volumes of data to become available for social researchers and data scientists, and crowdsourcing is becoming a key source of data to analyze and map social phenomena such as crime (Bendler et al., 2014) and perceptions of space and safety (Solymosi et al., 2015; Solymosi & Bowers, 2018). In this chapter we have described and explored the main strengths and weaknesses of using crowdsourced data for criminological research. Specifically, we have obtained access to a large dataset of more than 1.5 million votes about urban perceptions recorded from the Place Pulse project (Salesses et al., 2013), selected a sample of more than 37,000 votes of perceived safety for Atlanta, and studied the spatial distribution of perceptions of space and safety at a census tract level in this city. We have also shown how these data can be utilized to identify places assessed by participants as very safe or very unsafe; places in which researchers can then conduct observation to study those environmental features that make citizens feel fear of crime (Toet & Schaik, 2012).

Although crowdsourcing offers advantages over traditional survey methods to study perceptions and emotions about crime, data recorded from crowdsourcing is also affected by certain issues that, if uncontrolled, are likely to affect the validity of data and the reliability and generalizability of research outputs. For instance, we have observed how Place Pulse votes are largely produced by a few number of super-contributors (i.e., participation inequality), there is under-representation of certain areas outside the city center, and the number of votes decreases over time (i.e., participation decrease). These issues have also been observed in data produced from many other crowdsourcing and app-based projects (e.g., Chataway et al., 2017; Solymosi et al., 2015; Traunmueller et al., 2015). Other researchers have also highlighted that crowdsourced data tend to be affected by self-selection bias, which explains why males tend to participate more than females, and young persons more than adults (Salesses et al., 2013; Solymosi & Bowers, 2018); but the Place Pulse platform did not record demographic variables from participants and we have not directly assessed this issue here.

Due to the fact that crowdsourced datasets - and non-probability samples in general - may be affected by these potential sources of unrepresentativeness and bias, several researchers are exploring new techniques to enable obtaining reliable research outputs. Elliott & Valliant (2017), for example, present different methods to compute individual pseudo-sampling weights and adjust non-probability samples to target populations; Arbia et al. (2018) have developed a method to delete spatial outliers and calculate weights to adjust non-probability samples to optimal spatial samples; and Buil-Gil et al. (2020) investigate the use of resampling and model-based small area estimation techniques to allow producing reliable estimates at detailed spatial scales from crowdsourced data. Academics and practitioners will benefit from methods to mitigate the sources of bias in crowdsourced data, which may allow obtaining more precise and reliable - but also cheaper - findings and devise new explanations of crime, antisocial behavior and emotions about crime. In the context of crime analysis, bias-corrected crowdsourced data may become a key tool to understand crime patterns, anticipate crime trends and even provide assistance to police investigations (Bendler et al., 2014; Nhan et al., 2017).

## Authors bios

David Buil-Gil is a Research Fellow at the Department of Criminology of the University of Manchester, UK, and a member of the Cathie Marsh Institute for Social Research at this same university. His research interests cover small area estimation applications in criminology, environmental criminology, crime mapping, emotions about crime, crime reporting, new methods for data collection and open data.

Reka Solymosi is a Lecturer in Quantitative Methods at the Department of Criminology of the University of Manchester, UK, with interests in data analysis and visualization, crowdsourcing, rstats, fear of crime, transport, and collecting data about everyday life. As a former crime analyst, she is interested in practical applications to research and solving everyday problems with data.

## Acknowledgments

The authors would like to thank César A. Hidalgo for providing access to the data used in this chapter.

## References

- Arbia, G., Solano-Hermosilla, G., Micale, F., Nardelli, V., & Genovese, G. (2018). *Post-sampling crowd-sourced data to allow reliable statistical inference: The case of food price indices in nigeria*. Open Conference Systems.
- Bendler, J., Brandt, T., Wagner, S., & Neumann, D. (2014). *Investigating crime-to-twitter relationships in urban environments - facilitating a virtual neighborhood watch*. Association for Information Systems.
- Birenboim, A. (2016). New approaches to the study of tourist experiences in time and space. *Tourism Geographies*, 18(1).
- Blom, J., Viswanathan, D., Go, J., Spasojevic, M., Acharya, K., & Ahonius, R. (2010). *Fear and the city - role of mobile services in harnessing safety and security in urban contexts*. Association for Computing Machinery.
- Buil-Gil, D., Solymosi, R., & Moretti, A. (2020). Non-parametric bootstrap and small area estimation to mitigate bias in crowdsourced data: Simulation study and application to perceived safety. In C. Hill, P. Biemer, T. Buskirk, L. Japac, A. Kirchner, S. Kolenikov, & L. L. (Eds.), *Big data meets survey science*. John Wiley & Sons Ltd.
- Castro-Toledo, F. J., Perea-García, J. O., Bautista-Ortuño, R., & Mitkidis, P. (2017). Influence of environmental variables on fear of crime: Comparing self-report data with physiological measures in an experimental design. *Journal of Experimental Criminology*, 13.
- Chataway, M. L., Hart, T. C., Coomber, R., & Bond, C. (2017). The geography of crime fear: A pilot study exploring event-based perceptions of risk using mobile technology. *Applied Geography*, 86.
- Dubey, A., Naik, N., Parikh, D., Raskar, R., & Hidalgo, C. A. (2016). Deep learning the city: Quantifying urban perception at a global scale. In B. Leibe, J. Matas, N. Sebe, & M. Welling (Eds.), *Computer vision - eccv 2016* (pp. 196–212). Springer.
- Elliott, M. R., & Valliant, R. (2017). Inference for nonprobability samples. *Statistical Science*, 32(2).
- Fisher, B. S., & Nasar, J. L. (1992). Fear of crime in relation to three exterior site features. *Environment and Behavior*, 24(1).
- Gabriel, U., & Greve, W. (2003). The psychology of fear of crime. Conceptual and methodological perspectives. *British Journal of Criminology*, 43.
- Gastwirth, J. L. (1972). The estimation of the lorenz curve and gini index. *The Review of Economics and Statistics*, 54(3).
- Goodchild, M. F. (2007). Citizens as sensors: The world of volunteered geography. *GeoJournal*, 69.
- Goodchild, M. F., & Glennon, J. A. (2010). Crowdsourcing geographic information for disaster response: A research frontier. *International Journal of Digital Earth*, 3(3).
- Gómez, F., Torres, A., Galvis, J., Camargo, J., & Martínez, O. (2016). *Hotspot mapping for perception of security*. IEEE.
- Hamilton, M., Salim, F., Cheng, E., & Choy, S. L. (2011). *Transafe: A crowdsourced mobile platform for crime and safety perception management*. IEEE.
- Harvey, C. W. (2014). *Measuring streetscape design for livability using spatial data and methods* [Master's thesis]. The Faculty of the Graduate College, The University of Vermont.

- Hecker, S., Wicke, W., Haklay, M., & Bonn, A. (2019). How does policy conceptualise citizen science? A qualitative content analysis of international policy documents. *Citizen Science: Theory and Practice*, 4(1).
- Howe, J. (2006). The rise of crowdsourcing. *Wired Magazine*, 14(6).
- Innes, M. (2015). 'Place-ing' fear of crime. *Legal and Criminological Psychology*, 20(2).
- Jackson, J., & Gouseti, I. (2015). Psychological proximity and the construal of crime: A commentary on "mapping fear of crime as a context-dependent everyday experience that varies in space and time". *Legal and Criminological Psychology*, 20(2).
- McNulty, T. L., & Holloway, S. R. (2000). Race, crime, and public housing in atlanta: Testing a conditional effect hypothesis. *Social Forces*, 79(2).
- Miró-Llinares, F., Moneva, A., & Esteve, M. (2018). Hate is in the air! But where? Introducing an algorithm to detect hate speech in digital microenvironments. *Crime Science*, 7(15).
- Nhan, J., Huey, L., & Broll, R. (2017). Digilantism: An analysis of crowdsourcing and the boston marathon bombings. *British Journal of Criminology*, 57(2).
- Pebesma, E. (2020). *Sf: Simple features for r*. <https://CRAN.R-project.org/package=sf>
- R Core Team. (2020). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing. <http://www.R-project.org/>
- Salesses, P., Schechtner, K., & Hidalgo, C. A. (2013). The collaborative image of the city: Mapping the inequality of urban perception. *PloS One*, 8(7).
- Solymosi, R., & Bowers, K. (2018). The role of innovative data collection methods in advancing criminological understanding. In G. J. Bruinsma & S. D. Johnson (Eds.), *The oxford handbook of environmental criminology* (pp. 210–237). Oxford University Press.
- Solymosi, R., Bowers, K., & Fujiyama, T. (2015). Mapping fear of crime as a context-dependent everyday experience that varies in space and time. *Legal and Criminological Psychology*, 20(2).
- Solymosi, R., Bowers, K. J., & Fujiyama, T. (2017). Crowdsourcing subjective perceptions of neighbourhood disorder: Interpreting bias in open data. *British Journal of Criminology*, 58(4).
- Solymosi, R., Buil-Gil, D., Vozmediano, L., & Guedes, I. (2020). Towards a place-based measure of fear of crime: A systematic review of app-based and crowdsourcing approaches. *Environment and Behavior*.
- Spinu, V., Grolemond, G., & Wickham, H. (2020). *Lubridate: Make dealing with dates a little easier*. <https://cran.r-project.org/web/packages/lubridate/index.html>
- Tester, G., Ruel, E., Anderson, A., Reitzes, D. C., & Oakley, D. (2011). Sense of place among atlanta public housing residents. *Journal of Urban Health: Bulletin of the New York Academy of Medicine*, 88(3).
- Toet, A., & Schaik, M. G. van. (2012). Effects of signals of disorder on fear of crime in real and virtual environments. *Journal of Environmental Psychology*, 32(3).
- Traunmueller, M., Marshall, P., & Capra, L. (2015). Crowdsourcing safety perceptions of people: Opportunities and limitations. In T. Y. Liu, C. Scollon, & W. Zhu (Eds.), *SocInfo2015: Social informatics* (pp. 120–135). Springer.
- Welsh, B. C., & Farrington, D. P. (2004). Surveillance for crime prevention in public space: Results and policy choices in britain and america. *Criminology & Public Policy*, 3(3).
- Wickham, H., Chang, W., Henry, L., Pedersen, T. L., Takahashi, K., Wilke, C., Woo, K., Yutani, H., & Dunnington, D. (2020). *Ggplot2: Create elegant data visualisations using the grammar of graphics*. <https://CRAN.R-project.org/package=ggplot2>
- Wickham, H., François, R., Henry, L., & Müller, K. (2020). *Dplyr: A grammar of data manipulation*. <https://CRAN.R-project.org/package=dplyr>
- Wickham, H., & Henry, L. (2020). *Tidyr: Tidy messy data*. <https://CRAN.R-project.org/package=tidyr>

Williams, M. L., Burnap, P., Javed, A., Liu, H., & Ozalp, S. (2020). Hate in the machine: Anti-black and anti-muslim social media posts as predictors of offline racially and religiously aggravated crime. *British Journal of Criminology*, 60(1).

Williams, M. L., Burnap, P., & Sloan, L. (2017). Crime sensing with big data: The affordances and limitations of using open-source communications to estimate crime patterns. *British Journal of Criminology*, 57.

Zeileis, A., & Kleiber, C. (2015). *Ineq: Measuring inequality, concentration, and poverty*. <https://CRAN.R-project.org/package=ineq>