

# How many people are at intersections in Boston?

**Ben Lawson**

Department of Computer Science  
Boston University  
Boston, Massachusetts 02215  
balawson@bu.edu

*Social media is a constant factor in many people's lives in current times. We explore the usefulness of data collected from public resources to derive estimates of a subsection of people at different intersections in Boston. Using mainly Twitter data, we are able to apply sampling methods to determine an estimate of the monthly visitation rate of almost 60 intersections. This information can be incorporated into other systems to improve traffic congestion estimation.*

## 1 Introduction

In this project, we attempt to develop an understanding of human movement within the city of Boston. Using social media data from three companies, Brightkite and Gowalla, both of which are no longer active, and Twitter. To generate higher granular information, we used OpenStreetMap data to associate social media user's posts to specific intersections. Since not all people use social media, and thus are not represented in the datasets presented here, we must infer the actual amount of people that are present at these intersections in real life. Future work will attempt to cross-validate these methods with different types of observations, such as census population data and population counts at intersections derived from street cams and computer vision. Future work will also include using this data to solve classic problems, like max flow, in the pedestrian setting.

## 2 Data Resources

Many types of data were used in this project. Three sources of geosocial media data were used as well as geographical information from OpenStreetMap

### 2.1 Brightkite Dataset

This is a social media networking service that was acquired by a mobile social network, Limbo, in 2009. The dataset contains posts with a user id, geocoordinates, and a timestamp between 14 April 2008 and 18 October 2010.

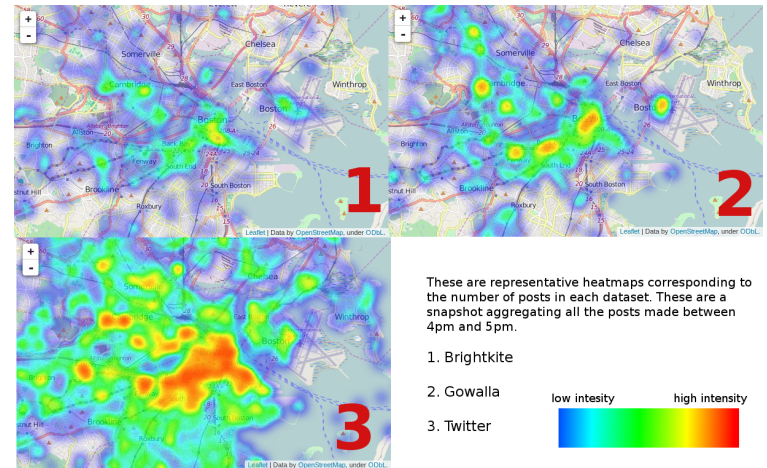


Fig. 1. Heatmap of tweets by direct count. Scaled proportionally total to each dataset.

### 2.1.1 Gowalla Dataset

This is a social media networking service that went out of business. Each post in the dataset has a user id, geocoordinates, and a timestamp, dating between 23 April 2009 and 22 October 2010.

### 2.1.2 Twitter Dataset

This is a micro blogging service that collects geological information about users's posts. This is still an active service and this data set was collected from 11 May 2015 until 2 April 2016. This data was collected via Twitter's streaming API, filtered by geolocation.

## 2.2 OpenStreetMap

This dataset was collected from OpenStreetMap, a collaborative, community based geographical open data resource. We use the labeled roads in the Boston area to create a chart of intersections.

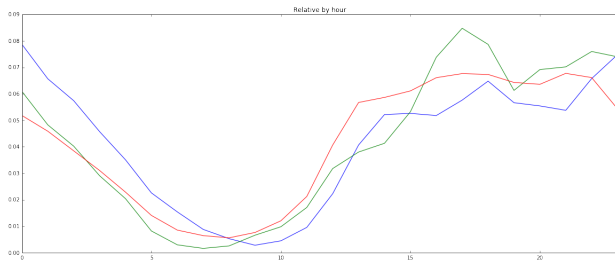


Fig. 2. Shows the percentage of tweets by hour. Demonstrates that tweeting behavior mimics the human sleep cycle. Red is Brightkite, Green is Gowalla, and Blue is Twitter.

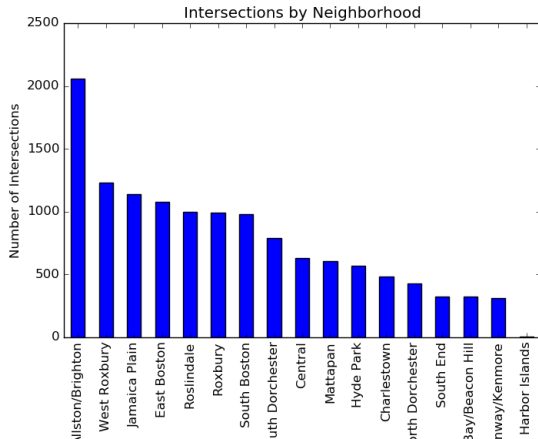


Fig. 3. The number of intersections per neighborhood

### 2.3 BostonMaps: Open Data

This dataset contained bounding boxes and geometric shapes for each of the neighborhoods in Boston. This was used to help give a general intuition for results.

### 3 Segmentation

Using the OpenStreetMap data, each Twitter post was associated with the closest intersection. This was done for each month in the Twitter dataset, so twelve months are represented from May 2015 to April 2016.

### 4 Sampling

**Capture & Recapture.** This type of sampling was first used when measuring animals in traps. Trappers would mark animals and then count how many of these animals returned to derive an estimate of the total population. We discovered only 94 intersections, of the almost 25,000, had five or more visitors each month. Only 57 of these intersections had visitors that returned during the capture/recapture period. The red markers show the intersections that estimates could be computed, scaled with the  $\log_2$  function. The blue markers show the intersections that did not have returning users. Fenway park had that max estimate with approximately 12,000 visitors per month.

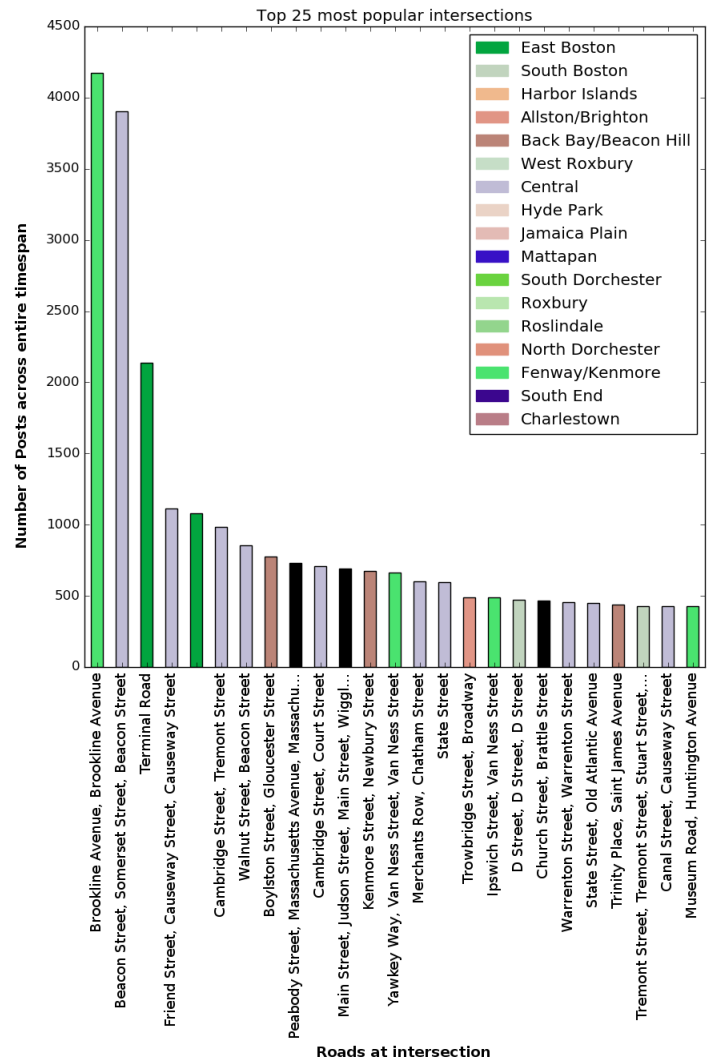


Fig. 4. The 25 most popular intersections overall (sheer volume) colored by neighborhood. Note missing road names and black color is due to missing data.

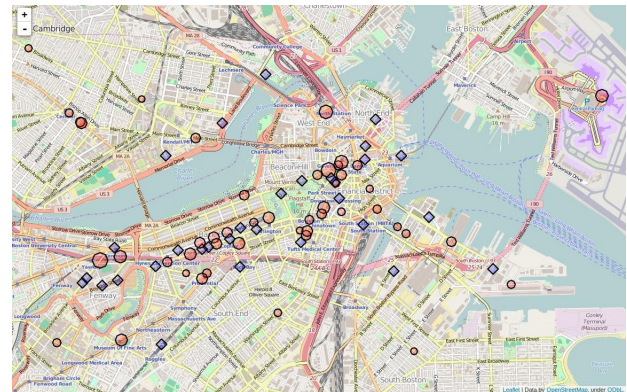


Fig. 5. Relative estimates for monthly social media users per intersection

$$\hat{N} = \frac{Kn}{k} \quad (1)$$

$\hat{N}$  is the estimation of the total population of social media users,  $n$  is the number of users during the first month,  $K$  is the number of users during the second month, and  $k$  is the number of users from the first month that are observed during the second month. With this formulation, we can derive estimations for the number of social media users that post near an intersection per month.

## 5 Conclusions

Although the social media data consisted of many posts, only a fraction of the intersections had data spanning the entire collection period. Of these intersections, only a fraction had entire data to compute the estimate via the capture & recapture method. Future work will be to include the Brightkite and Gowalla datasets in the estimation of intersection occupancy and analyzing flow problems associated with pedestrian traffic.

## Acknowledgements

Special thanks to Andrei Lapets who supported this project and helped flush out ideas about general project direction. Special thanks to Evimaria Terzi for help with the development of the Twitter dataset and Sofia Maria Nikolakaki for her work the segmentation of tweet segmentation by intersection.