

Flashlight: Property Assessment Visualization for the City of Boston

Tyler Brown, Sicheng Hao, Nischal Mahaveer Chand, Sumedh Sankhe

Summary

As a new home-buyer, it's easy to find out about your home but hard to get an understanding of your neighborhood. Flashlight makes it easier for you to see your potential neighborhood in Boston. This discrepancy is because current real estate websites emphasize individual properties rather than individual neighborhoods. Our group communicates the differences between Boston neighborhoods using an interactive data visualization called Flashlight.

Our dataset includes Property Assessment history from 2014-2017 [1] using Boston's Open Data Hub. We enriched the property assessment data with coordinates from Open Addresses [2], and neighborhood boundaries from Zillow [3]. These combined datasets provide unique insights to new home-buyers in Boston. As open data becomes more prevalent in cities across the United States, we can scale our insights and models.

Methods

We used methods for collecting, preparing, modeling, and presenting our data. Each step of the process is detailed here.

0.1 Data Collection

We started with Boston's Property Assessment data from 2014-2017 [1]. This dataset “[g]ives property, or parcel, ownership together with value information, which ensures fair assessment of Boston taxable and non-taxable property of all types and classifications.”[1]. We wanted to use this information because it helps us capture changes in Boston properties over time. For example, a remodeled property would change it's property tax assessment value we have this variable available to us.

After starting with the Property Assessment dataset, we brought in additional datasets to increase the value of our data collection. Neighborhoods in Boston were not named or geographically demarcated in the Property Assessment dataset so we brought in Neighborhood

Boundaries from Zillow [3] to make this distinction. Additionally, geographic coordinates for each assessed property’s address were occasionally not coded correctly or included at all for 2017 so we had to bring in those values using Open Addresses [2]. Once neighborhood names, boundaries, and missing coordinates were available, we were able to proceed to data preparation.

0.2 Data Preparation

There were a number of steps involved in data preparation.

0.2.1 Data Audits

The purpose of a data audit is to answer questions related to data quantity and quality. We started with our Property Assessment dataset by checking for the quantity of populated items within each variable. About 73% of variables were less than 70% populated. We were able to disregard these variables for our modeling purposes. Data quality checks are not as easily automated but value added.

Analyzing data quality helped us understand data problems up front such as having a non-unique primary key about 0.3% of the time, latitude and longitude were missing or corrupted about 35% of the time, and the ratio of unique addresses concatenated with coordinates to unique coordinates were about 3 : 1. We also found that the Property Assessment data did not map to defined neighborhoods in Boston. Understanding these shortcomings with a data audit allowed our group to plan remediation steps early in our analysis.

0.2.2 Geocoding Missing Addresses

During the data audit we found about 35% of addresses within the Property Assessment data were not matched to any coordinate pair. Our team work to mitigate this deficiency by leveraging the a “free and open global address collection” called the OpenAddresses [2] project. There are several strategies one can take when Geocoding (mapping address strings to coordinate pairs) addresses. The simplest would be to use Google Maps Geocoding API [4].

The Google Maps Geocoding API has a free usage tier which maxes out at 2,500 requests per day [4]. Our data required geocoding an order of magnitude more coordinates so this approach was out of the question. We resolved the issue by creating ‘address hashes’ for each address in the OpenAddress and Property Assessment dataset with missing coordinates. An ‘address hash’ was a string concatenated of concatenated values for street number, street name, city, and zip code. We excluded unit number as a simplifying assumption because we expected all unit numbers to be at the same property. Once ‘address hash’ had been computed for both datasets we performed an inner join.

In some cases, this join generated a many-to-many relationship between addresses in both datasets due to our exclusion of unit numbers. We resolved this issue by grouping and taking the first coordinate pair. Subsequent iterations of our analysis can be improved by adding robustness to our geocoding procedure.

0.2.3 Working with GeoJSON and Python

stuff about leaflet here.

0.3 Data Modeling

0.4 Data Presentation

Results

We had some results.

Discussion

Let's discuss what we did.

Statement of Contributions

Together everyone achieves more.

- **Tyler Brown:**
- **Sicheng Hao:**
- **Nischal Mahaveer Chand:**
- **Sumedh Sankhe:**

References

- [1] C. of Boston, “Property assessment - datasets - analyze boston.” <https://data.boston.gov/dataset/property-assessment>. (Accessed on 10/26/2017).

- [2] OpenAddresses, “Openaddresses.” <https://openaddresses.io/>. (Accessed on 12/09/2017).
- [3] Zillow, “Zillow neighborhood boundaries.” <https://www.zillow.com/howto/api/neighborhood-boundaries.htm>. (Accessed on 12/09/2017).
- [4] Google, “Google maps geocoding api.” <https://developers.google.com/maps/documentation/geocoding/>. (Accessed on 12/09/2017).

Appendices

Appendices here.