

Neighborhood Analysis in the Washington, D.C. Area for the Future Amazon HQ2

Capstone Project - Battle of the Neighborhoods

March 2, 2020

Contents

1	Introduction	1
1.1	Goal	1
1.2	Objectives	1
2	Data	2
2.1	Data Requirements	2
2.2	Data Collection	2
	Bibliography	4

1 Introduction and Business Problem

On November 13, 2008 Amazon announced their plans to establish a second headquarter office (named HQ2) in the United States, located in New York city and Washington D.C. ([1, 2]). This represents a great opportunity for the greater D.C. area to experience growth and to attract a significant number of both skilled workers and other businesses, which are expected to populate the DMV (D.C., Maryland and Virginia) area.

The National Landing in Arlington County was chosen as the future site of HQ2. This county houses multiple national landmarks (Arlington National Cemetery, the Pentagon), a thriving business economy (with strong presence of both private companies and federal agencies), and along with the neighboring Alexandria routinely top national rankings of best places to live and retire. All these characteristics make it a very desirable area to live and work on.

However, the D.C. metropolitan area already suffers from severe issues that are likely to worsen with the expected heavy affluence of workers. Specifically, this report will look into the housing prices (rent and sales), and the commuting time. This last issue is a major concern in the area, as the D.C. area commute regularly ranks within the national top 5 work commute times (regardless of whether the commute is analyzed in terms of actual time spent commuting, or just the time spent sitting in traffic).

For all the reasons above, this document aims to analyze the different regions in the DMV area in order to help understand what is the current situation prior to the arrival of Amazon's HQ2, in order to be able to better plan for the changes that the arrival will trigger, and in this way, allow all the interested stakeholders, from city planners to realtors and small business owners, to identify potential issues and areas to improve.

1.1 Goal

The goal of this analysis will be **to improve the understanding of the DMV area in terms of amenities, commute time, and housing costs**. These are significant aspects that will drive the decision of many of the new workers that will arrive to the area when deciding where to settle down, which will, in turn, affect back to those same issues (amenities, commuting, and housing costs).

It shall be noted that the goal of this analysis is not to perform a prediction on trends or metrics, as the arrival of HQ2 is expected to be significantly disruptive, and therefore using current trends to forecast the future seems not appropriate. Therefore, we will focus this report on the analysis and understanding of the current state of the DMV areas.

1.2 Objectives

In order to achieve the main goal described previously, we will be performing descriptive analyses on the DMV area. Specifically, the following objectives will be targeted:

- *Amenities Similarity*: In order to identify neighborhoods and areas that offer similar amenities to National Landing, we will look for areas that offer similar types and quantities of amenities in the whole DMV area.
- *Commute Time Distribution*: As National Landing sits inside the I-495 Beltway, and is bordered by the Potomac river, the time to commute for new potential workers may be significant. To analyze this, we will study not only the average expected commute time from each neighborhood or area to National Landing, but we will also focus the analysis on finding patterns regarding the distribution of this commute time.
- *Housing Costs*: The study of the costs of renting and purchasing a home in the different areas will also be analyzed, and, as with the commute time analysis, patterns and distributions will be displayed and analyzed.
- *Desirability Distribution*: After individually analysing each of the features previously described, a further study will be performed attempting to combine all these individual aspects in a single metric that tells us how similar in terms of desirability the different areas are.

2 Data

In this chapter we will review the data that will be used for achieving the goals described in Chapter 1. We will look at what type of data is needed (including the formats, ranges, and types), the sources we can use, the understanding of how the data fits into the project, and any preparation that may be needed.

2.1 Data Requirements

In order to achieve the goals of this project, we need several pieces of data:

- **Geographical Description of the Areas to Analyze:** The first piece we will need is the geographical definition of the areas that we will be using for the analysis. These areas should be separate areas (i.e. with no overlapping) that provide full coverage of the areas of District of Columbia, Maryland, and Virginia. Furthermore, we will need to be able to replace each area with a point of reference that can be used to query FourSquare for amenities. Given the large amount of National Parks and water bodies in the area, we will need to be careful to ensure that some areas do not get misrepresented because of the presence of one of these areas (which, for example, would have no amenities nearby if querying from the center of a National Park).
- **Amenities in Each Area:** In order to simplify the analysis, we will be focusing on the top-10 most popular types of amenities in an area, as provided by FourSquare.
- **Commute Times:** The analysis of the commute times will require the availability of the required time to commute from one point to the expected location of HQ2, as it is usually reported by apps like Google Maps, or Waze. Given that this expected commute time may vary significantly from one day of the week to another (e.g. depending on school hours or telecommuting in the area), we will need to consider the average commute time across a number of days and hours, to get a better representation of this metric. In order to simplify the analysis, we will look only at car commute times only, leaving walking, biking, or public transportation for a future extended analysis.
- **Housing Costs:** Finally, we will need to acquire a dataset that provides us with information regarding the cost of renting or buying a home in each of the areas analyzed. This information should be normalized to account for differences in terms of size, number of rooms, etc.

2.2 Data Collection

The main issue in collecting the data for this project may lay in the definition of the areas for the analysis. While D.C. defines neighborhoods, which seem to be a good partition for this type of analysis, it turns out that the definition of those neighborhoods is not clear, and there are overlaps ([3]). For this reason, and in order to have a division that can apply not only to D.C. but also to Maryland and Virginia, we decided to proceed with divisions based on ZIP codes tabulation areas (ZCTA), as defined by the U.S. Census Bureau ([4]). These ZIP code tabulation areas are full-coverage, non-overlapping areas that cover all of the U.S., and we can use the datasets available from the same U.S. Census Bureau to get the coordinates from a representative point of each area. The datasets are provided as Shapefiles, which can be easily converted to CSV files using any GIS application.

FourSquare will provide us with access to the most popular amenities near that representative point for each ZCTA, but we can also prevent the case where a ZCTA falls in a National Park or a similarly deserted area, and query the API with the ZIP code instead, thus helping us overcome one of the main issues of analysing the DMV area.

Expected commute times between points can be obtained through the Bing Distance Matrix API ([5]). Using the provided examples as a base, and automating the calls with Python, we can build our dataset of expected commute times for a range of hours each day between Monday and Saturday (inclusive).

2 *Data*

Finally, the housing costs will be obtained from the Zillow Home Value datasets, available at [6]. These datasets can be sorted by ZIP code only if they are provided as time series. Given that we are not interested in forecasting, we will have to filter the time series to process only the last value. The Zillow datasets already provide a single normalized value that accounts for footage and room number, so we will not have to perform any normalization nor create our own metric based on the price.

Bibliography

- [1] “Amazon’s Grand Search For 2nd Headquarters Ends With Split: NYC And D.C. Suburb.” <https://www.npr.org/2018/11/13/665646050/amazons-grand-search-for-2nd-headquarters-ends-with-split-nyc-and-d-c-suburb>. Accessed: 2020-03-01.
- [2] “Amazon Picks New York City, Northern Virginia for Its HQ2 Locations.” <https://www.wsj.com/articles/amazon-chooses-new-york-city-and-northern-virginia-for-additional-headquarters-1542075336>. Accessed: 2020-03-01.
- [3] “Neighborhoods in Washington, D.C..” https://en.wikipedia.org/wiki/Neighborhoods_in_Washington,_D.C.. Accessed: 2020-03-01.
- [4] “2019 TIGER/Line Shapefiles (machinereadable data files) / prepared by the U.S. Census Bureau, 2019.” <https://www.census.gov/geographies/mapping-files/time-series/geo/tiger-line-file.html>. Accessed: 2020-03-01.
- [5] “Distance Matrix API.” <https://www.microsoft.com/en-us/maps/distance-matrix>. Accessed: 2020-03-01.
- [6] “Zillow Home Value Index.” <https://www.zillow.com/research/data/>. Accessed: 2020-03-01.