

Sydney Cyclability Analysis

Zijin Huang
Mengchen Ye
Wendi Tanaka

2019 May

Overview

This project aims to produce and analyze the cyclability of suburb areas around Sydney based on their respective statistics.

Dataset Description

Initially, there are five datasets provided. They contain some information on neighbourhoods of Sydney area, the specifics include but not limited to the population, number of dwellings, and the totals of the specific types of businesses. Also, the school provided a synthetic dataset of some bike pods and their respective attributes, longitude, and latitude. Since the primary goal of this project is to calculate the cyclability score of each area, all NULL values are treated as zero when performing score calculation.

In supplement to the datasets, the team has obtained the area shapefile¹ from Australian Bureau of Statistics (ABS). The dataset contains Statistical Area 2(SA-2), with polygon shape information for PostGis to process.

Also, the team collected data using web-scraping and scrape readable text from PDF file. The team scraped all the text data from 8 PDF files and 16 websites which are related to cycling in Sydney, then use NLP skills to grab and classify nouns from these data. After, a program counted the frequency of occurrence of area_name in StatisticalAreas.csv appears in these nouns.

¹<https://www.ausstats.abs.gov.au/ausstats>

Database Description

Schema

See Appendix.

Indexes

There are two normal indexes on this database:

1. neighbourhoods_land_area_idx
2. businessstats_num_businesses_idx

These indexes are created to enable fast search and sort on neighbourhoods based on its land area, and its business count. Since these two are used

Also, there are two spatial indexes:

1. bikepods_geo_index
2. neighbourhoods_geo_idx

These indexes are created for fast spatial join between bikepods and neighbourhoods.

Cyclability Analysis

Before we calculate the cyclability score of any given region, we first obtained the following statistics of each neighbourhood:

$$population_density(pd) = \frac{population}{land\ area} \quad (1)$$

$$dwelling_density(dd) = \frac{number\ of\ dwellings}{land\ area} \quad (2)$$

$$service_balance(sb) = \frac{education * 5 + food * 4 + retail * 3 + recreation * 2 + health}{number\ of\ service\ businesses} \quad (3)$$

$$bikepod_density(bd) = \frac{number\ of\ bikepods}{land\ area} \quad (4)$$

$$NLP_score(nlp) = The\ number\ of\ occurrences\ of\ the\ area_name\ in\ the\ fetched\ nouns. \quad (5)$$

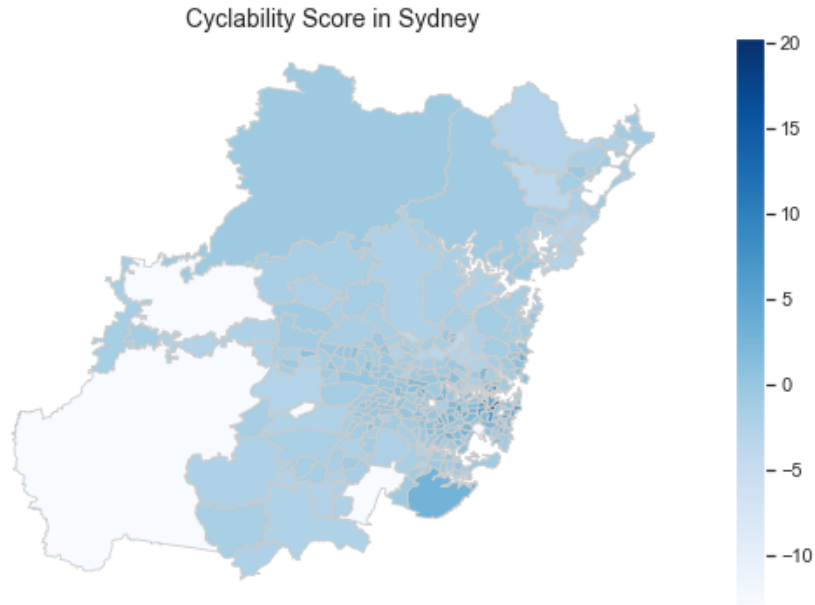
Details about service balance

Service_balance calculation derived from 5 businesses consists of education, food, retail, recreation and health. Those five businesses number scaled up by its usefulness.

After obtaining the above statistics of every neighbourhood, the cyclability score is calculated based on the given formula:

$$cyclability = z(pd) + z(dd) + z(sb) + z(bd) + z(nlp) \quad (6)$$

By implementing the above equations, we obtained the cyclability score of the Sydney suburbs:

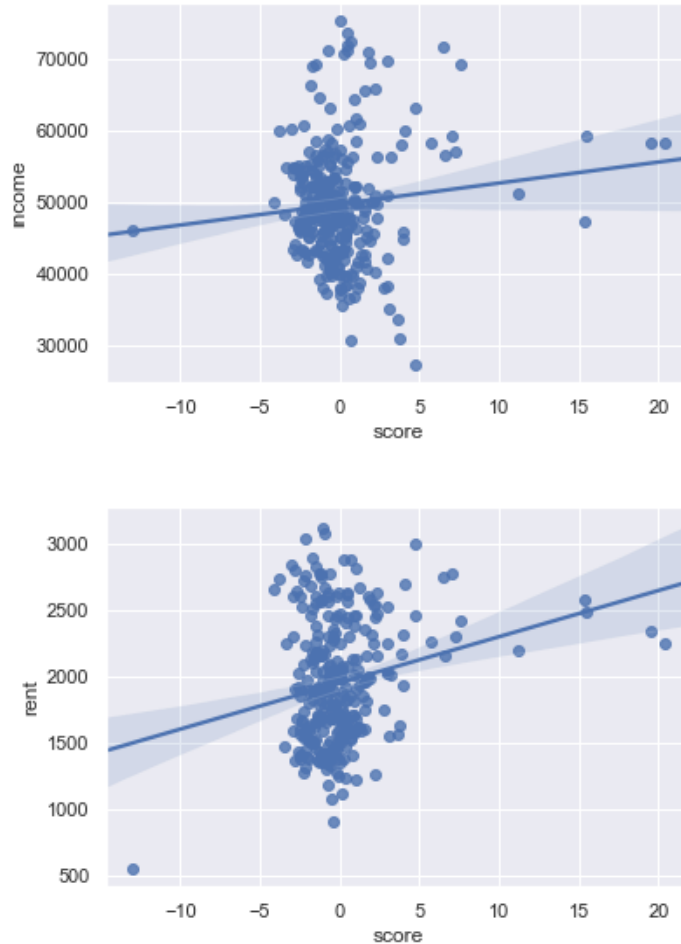


Web scraping and NLP

To enrich the datasets, we scrape readable text from government website and PDF which are related to the Sydney cycling infrastructures with BeautifulSoup and PyPDF2 library. Then the program processed the said text with NLP to classify the nouns. After removing the irrelevant entities from the existing nouns, the program counts the occurrence frequency of nouns related to each neighbourhood. Finally, the program computed NLP z score and added the result to the current cyclability score.

Correlation Analysis

Based on the statistics, we have obtained the following scatter plots:



We can see that the cyclability score is positively correlated with both average income and average weekly rent. The score and average annual income has a correlation coefficient of 0.109, and it has a correlation coefficient of 0.221 with neighbourhood average weekly rent.

Appendix

