



A circular word cloud centered on Craigslist items. The words are arranged in a circle, with larger words in the center and smaller words towards the perimeter. The words include: file, firewood, two foot, sofa, alert, car, metal, table, curb, cut, red, soil, size, nice, go, pallets, desk, fill, old, chair, X, iso, non, couch, dirt, chair, boxes, moving, tv, stand, new patio.

SPRINGBOARD CAPSTONE PROJECT TWO

WHAT'S FOR FREE ON CRAIGSLIST?

Author: Josh Mayer

Date: November 15, 2017

OUTLINE

- Introduction & Problem Statement
- The Dataset
- Analysis & Findings
- Statistical Inference
- Machine Learning
- Conclusions

★ La Z Boy love seat (Parker) ☺



★ FREE SCRAP METAL (Hudson) ☺



★ Free Wine Cabinet - holds 260+ bottles of wine (Stapleton) ☺



INTRODUCTION & PROBLEM STATEMENT

Craigslist was founded in 1995 and has become one of the most powerful sources for classified advertisements found on the web today.

THE PROBLEM

Oftentimes products are misclassified or the category itself is too broad to be useful.

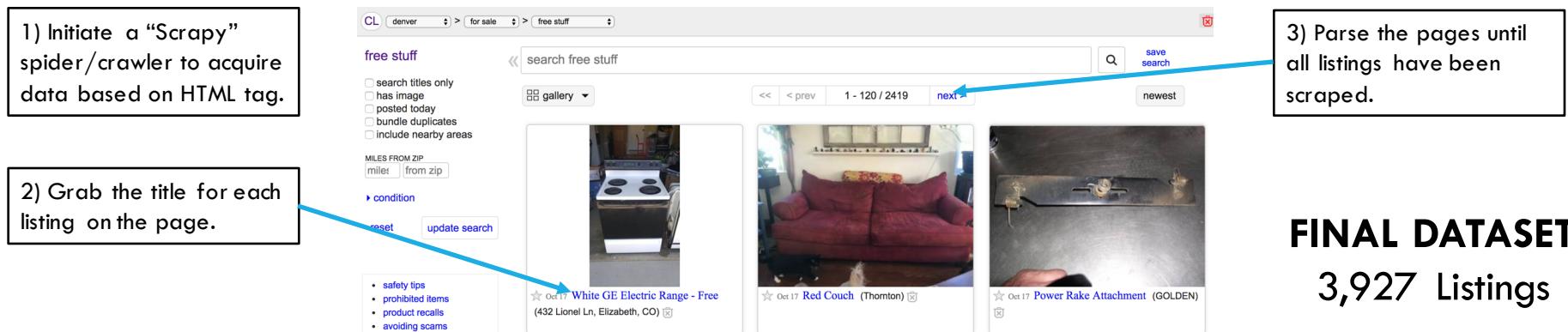
| denver, CO | | |
|--------------------------|-----------------|-----------------------|
| | | |
| community | housing | jobs |
| activities | local news | accounting+finance |
| artists | lost+found | admin / office |
| childcare | musicians | arch / engineering |
| classes | pets | art / media / design |
| events | politics | biotech / science |
| general | rideshare | business / mgmt |
| groups | volunteers | customer service |
| personals | for sale | education |
| strictly platonic | free | food / bev / hosp |
| women seek women | furniture | general labor |
| women seeking men | garage sale | government |
| men seeking women | general | human resources |
| men seeking men | heavy equip | internet engineers |
| misc romance | household | legal / paralegal |
| casual encounters | jewelry | manufacturing |
| missed connections | materials | marketing / pr / ad |
| rants and raves | motorcycles | medical / health |
| discussion forums | music instr | nonprofit sector |
| apple | photo | real estate |
| arts | help | retail / wholesale |
| atheist | history | sales / biz dev |
| autos | p.o.c. | salon / spa / fitness |
| beauty | housing | security |
| bikes | jobs | skilled trade / craft |
| celebs | psych | software / qa / dba |
| comp | jokes | systems / network |
| crafts | queer | technical support |
| diet | kink | transport |
| divorce | legal | tv / film / video |
| dying | linux | web / info design |
| eco | romance | writing / editing |
| | cell phones | [ETC] |
| | clothes+acc | [part-time] |
| | collectibles | |
| | computers | |
| | electronics | |
| | farm+garden | |

Source: <https://denver.craigslist.org/>

PRIMARY OBJECTIVE
 The goal is to explore the “For Sale/Free” category, leverage unsupervised learning to cluster and categorize products listed.

DATA ACQUISITION & WRANGLING

Data was scraped from the Denver, Boulder, Colorado Springs, and Fort Collins (Greater Denver Metro) Craigslist sites via the Python package “Scrapy”.



FINAL DATASET
3,927 Listings

DATA WRANGLING

- Duplicate Records – 133 records (3% of the acquired records) were dropped since they contained duplicate title AND timestamp. New timestamps are created if listings are “reposted”.
- Non-letter Removal – remove numbers, punctuation, emojis, etc. in order to have only letters.
- Lowercase Conversion – all letters were converted to be lowercase for comparison.

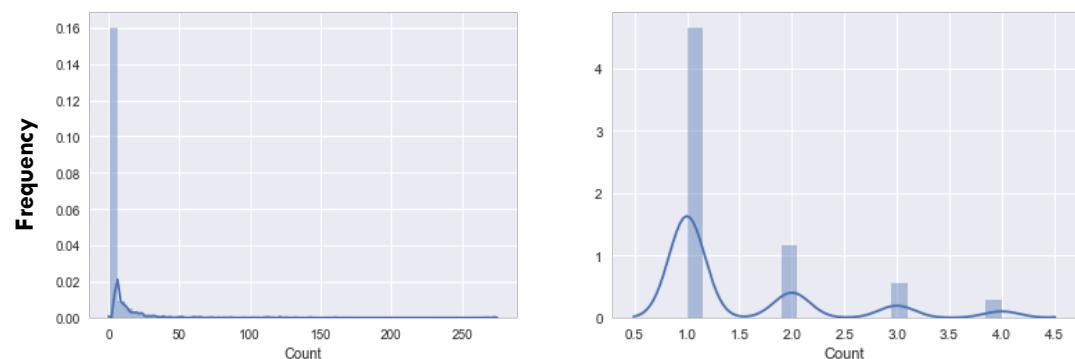
EXPLORATORY DATA ANALYSIS

The dataset acquired was highly skewed towards infrequently found words.

DATASET CHALLENGES

- 2,585 unique words were found in 3,927 unique listings.
 - The average word is found 4.51 times with a min count of 1 and max count of 274 counts.
- 57% of the words are found only one time throughout the many listings.
- 82% of the words are found fewer than 5 times.

FREQUENCY VS UNIQUE WORD COUNT



Charts above: frequency of all words (left) along side the frequency of words found fewer than 5 times (right).

FURTHER ANALYSIS

Analyzing the most (and least) frequent words begins to show some of the potential categories.

We can already start to see some of the potential categories based on word frequency alone (e.g. "wood" and "firewood" are amongst the top ten words found).



MOST FREQUENT VS INFREQUENT WORDS

| COUNT | TOP 10 WORDS |
|-------|--------------|
| 274 | wood |
| 267 | couch |
| 160 | tv |
| 143 | boxes |
| 132 | chair |
| 125 | dirt |
| 121 | pallets |
| 121 | curb |
| 115 | moving |
| 68 | firewood |

| COUNT | BOTTOM 10 WORDS |
|-------|-----------------|
| 1 | ability |
| 1 | newspaper |
| 1 | newer |
| 1 | neutral |
| 1 | network |
| 1 | net |
| 1 | nepro |
| 1 | needing |
| 1 | necklaces |
| 1 | mtb |

Charts above: the most frequent word found in the listings was "wood" at 274 (left chart). There are many words with a count of 1 (right chart).

MACHINE LEARNING

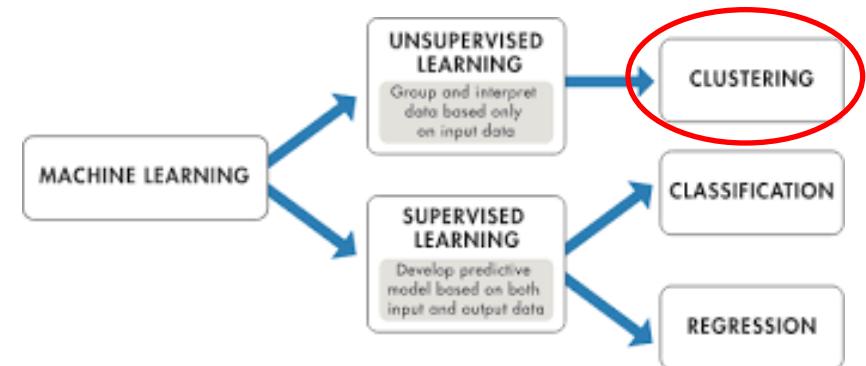
The primary goal of the machine learning (ML) section is to identify the best algorithm to properly classify the “free stuff” listings on Craigslist.

MACHINE LEARNING APPROACH

- 1) Text analysis was performed with the term frequency-inverse document frequency (tf-idf) algorithm and the vectorizers were then transformed to the tf-idf matrix.
- 2) Parameter tuning included setting min/max term frequency, “N grams”, stop words, and logarithmic data transformations.
- 3) Perform clustering with the “K-Means” algorithm.
- 4) Understand model results and identify the best model to name the clusters.

FOCUS ON UNSUPERVISED LEARNING

This is an unsupervised ML exercise and focused on testing and understanding key parameters as relates to the exploratory data analysis findings.



MACHINE LEARNING – TEXT ANALYSIS

Text analysis was performed with the term frequency-inverse document frequency (tf-idf) algorithm and the vectorizers were then transformed to the tf-idf matrix.

STEPS TO CREATE THE TF-IDF MATRIX

1. Count word (term) occurrences by document.
2. Transform into a document-term matrix.
3. Apply the term frequency-inverse document frequency weighting.
4. Perform sensitivity analysis on frequency weighting as well as some of the other key parameters.

VISUAL EXAMPLE TF-IDF MATRIX

| | Document 1 | Document 2 | Document 3 | Document 4 | Document 5 | Document 6 | Document 7 | Document 8 |
|-----------|------------|------------|------------|------------|------------|------------|------------|------------|
| Term(s) 1 | 10 | 0 | 1 | 0 | 0 | 0 | 0 | 2 |
| Term(s) 2 | 0 | 2 | 0 | 0 | 0 | 18 | 0 | 2 |
| Term(s) 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| Term(s) 4 | 6 | 0 | 0 | 4 | 6 | 0 | 0 | 0 |
| Term(s) 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| Term(s) 6 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| Term(s) 7 | 0 | 1 | 8 | 0 | 0 | 0 | 0 | 0 |
| Term(s) 8 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 |

Word Vector (Passage Vector)

Document Vector

MACHINE LEARNING – PARAMETER TUNING

Sensitivity analysis was performed on key model parameters to determine the impact.

KEY MODEL PARAMETERS

- Min/Max Frequency – this is the frequency within the listings a given term can have to be used in the tf-idf matrix. E.g. If the term is in greater than X% of the listings it probably carries little meaning.
- N grams: accounts for term relativity. As an example, setting this value equal to “2” would mean that I consider unigrams and bigrams, e.g. “leather” and “leather couch” which are both important to me.
- Stop Words: some words are too common (e.g. the, and) to be considered (stop words) for analysis and need to be removed. The only customization made for this model was to include the word “free” in the stop words list as it was found in nearly all listings.
- Logarithmic Transformation: setting this value to “true” provided a log10 transformation of the data, which in this case was sorely needed given the dataset was skewed towards infrequent words.

MACHINE LEARNING – K-MEANS CLUSTERING

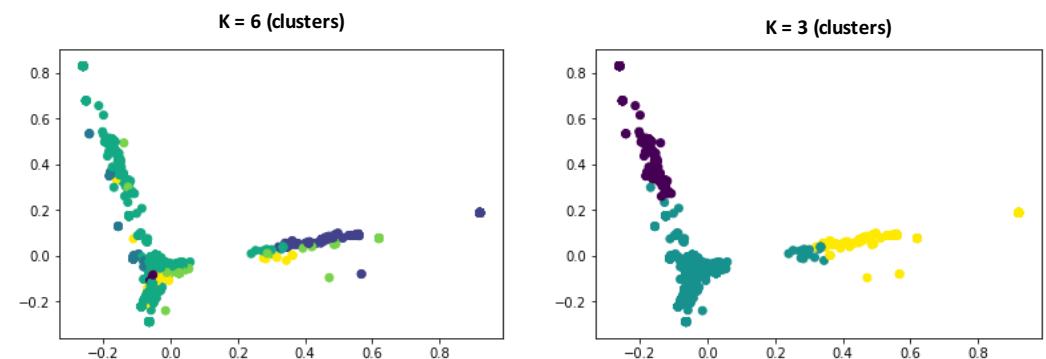
The clustering algorithm of choice for this model was “K-means” which required a pre-determined number of clusters for initialization.

ITERATING OVER POTENTIAL K VALUES

| K Value | Sum-of-Squares Score | Silhouette Score |
|---------|----------------------|------------------|
| 2 | -2,273 | 0.134 |
| 3 | -2,201 | 0.145 |
| 4 | -2,128 | 0.163 |
| 5 | -2,081 | 0.181 |
| 6 | -1,981 | 0.202 |
| 7 | -1,994 | 0.191 |
| 8 | -1,893 | 0.223 |
| 9 | -1,861 | 0.233 |

The silhouette score (SS) ranges from -1 (a poor clustering) to +1 (a very dense clustering) with 0 denoting the situation where clusters overlap. SS scores ranged from 0.134 to 0.233 depending on the K value.

COMPARING VISUAL RESULTS OF K CLUSTERS



Charts above: at k=6 (left) Interesting to see the supposed cluster on the right-hand side of the chart with multiple colors of dots next to each other. At k=3 (right) the clusters visually appear more relevant, however, yellow and green dots are still very close together on the right-hand side of the page. Looking back at the scores for k=3, the SS is 0.145 which means a weak structure.

MACHINE LEARNING – MODEL RESULTS

Performing sensitivity analysis on key model parameters produced a model with k=3 clusters and the following results.

A few key observations:

- As expected, some of the words are present across clusters. This was to be expected given the low silhouette scores and the visual of the many different colored dots close together.
- Only one bigram is present, meaning that the non-unigram words were not as important as initially thought.
- An attempt at naming the clusters might be:
 - Cluster 1 would be named “Moving Items”.
 - Cluster 2 would be named “Scrap Materials”.
 - Cluster 3 would be named “Furniture”.

NAMING THE CLUSTERS

| Cluster 1 “Moving Items” | Cluster 2 “Scrap Materials” | Cluster 3 “Furniture” |
|-----------------------------|--------------------------------|--------------------------|
| Wood | Firewood | Couch |
| Tv | Wood | Leather |
| Pallets | Pallets | Leather couch |
| Chair | Today | Loveseat |
| Boxes | Come | Recliner |
| Desk | Scrap | Chair |
| Dirt | Yard | Reclining |
| Stuff | Small | Sofa |
| Table | Mulch | Brown |
| Moving | Pallet | Bed |

Chart above: top 10 terms in each cluster of the model results of k=3.

RESULTS SUMMARY

KEY FINDINGS

- The data set acquired for the Greater Denver Metro “Free Stuff” analysis was heavily skewed by infrequent terms. Even a log 10 transformation was not able to properly redistribute or normalize the data.
- This does not mean it is “bad data”, only that the majority of the postings in the “free stuff” section represent very few terms.
- With a Silhouette Score of 0.14 on the training set and 0.15 on the test set, the K-means clustering model provided weak (at best) structure for the designated clusters.
- In this analysis of the “Greater Denver Metro” Craigslist sites, there was an abundance of items in the Craigslist “free stuff” that was classified into three categories:
 - Category 1 is “Moving Items”.
 - Category 2 is “Scrap Materials”.
 - Category 3 is “Furniture”.

RECOMMENDATIONS & FUTURE IMPROVEMENTS

Key results show that the Craigslist team could implement a clustering model to help users filter and find the products.

RECOMMENDATIONS & FUTURE IMPROVEMENTS

- Craigslist site administrators should create and enable “tags” or “categories” to assist in filtering the “free stuff” section.
- Perform unsupervised learning on the broader, national data set of all Craigslist site for better understanding of what is offered for free across the country.
- There could be a similar analysis done regarding miscategorized or mislabeled listings where the goal would be to find the outlying data points to understand their value.