

IBM Data Science Professional Certificate Capstone - Predicting Venue Like Count Using Natural Language Processing of Venue "Tip" Text

Andrew Barber

github.com/acbarber

Abstract

With such a large quantity of textual data on the internet today, there is an enormous opportunity for exploration and analysis. However, much of this data lacks the necessary context to conduct such research, making "labelled data", or data with an external measure of validation, extremely valuable. In this paper we attempt to leverage labelled textual data extracted from the Foursquare API (<https://developer.foursquare.com/>), making a model to predict venue like counts based on existing venue "tips". We found that our model successfully outperformed the baseline average like count metric, recording a mean squared error (MSE) of 10.99 compared to the baseline's MSE of 12.57.

1 Introduction

In many cases data is limited in scope and format. This is especially true now, with an abundance of unstructured data in the form of natural language - it can be hard to transform and extract meaning from the data in a valuable way. Labelled natural language processing (NLP) data presents a unique opportunity to extract meaning from natural language and hopefully extrapolate these findings to a wider category of unlabelled data. In the case of this paper, we have examined venue data from the Foursquare API including venue "tips" (reviews of the venue) and venue like count.

Our data was taken from venues scattered around the 40 major neighborhoods of Manhattan where we attempted to predict venue like counts based on the tips left for that venue. This NLP prediction task is used for user recommendation systems and may be useful for predicting the popularity of a venue based on unlabelled comments where another measure of popularity does not exist. It may also be used as a supporting method for

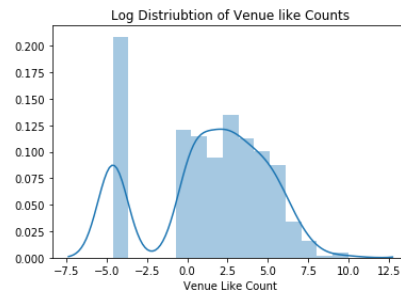


Figure 1: Log Distribution of Venue Like Counts

venue popularity evaluation when there are other methods of popularity such as like count or rating but the quantity of NLP data is far more extensive. If my model can successfully predict like counts in Manhattan venues, it could then be applied to popularity evaluation applications such as listed above or represent a viable base for future research.

2 Problem Definition and Data

As mentioned above, we used the Foursquare API for this paper. Specifically we looked at "tips" for various venues scattered across the 40 neighborhoods of Manhattan. Our problem addressed how to leverage this "tips" text in a way that could accurately predict a venue's like count. In order to retrieve a good mix of venues, we pulled the twelve nearest venues to each of the Manhattan neighborhoods. We then extracted all "tips" from the venues and concatenated them together to form the tip text for each venue. Rows with "N/A" values in either the "Venue Like Count" or "Venue Tip Text" fields were dropped. The final number of records was 459.

The challenge with this dataset came with the large number of zeros for like counts. In order to produce a model that wouldn't automatically predict 0 for all venues, putting the target field on

Model	MSE
Baseline - Mean Venue Log Like Count	12.53
Linear Regression (BOW)	11.00

Table 1: Evaluation of models

the narrow scope of the data. Looking at Manhattan venues exclusively limits the researcher’s ability to make generalizations regarding the results of their research. Expanding the scope of the data will then further increase the impact of the results.

Future models created for this task could vary in numerous ways. For one, different text-representation techniques could be used such as term-frequency independent-document-frequency (TF-IDF), word embeddings or pre-trained word embeddings. These NLP methods may produce better results than our BOW model. Different regression models may also be used to increase performance. Specifically, neural networks could be leveraged in the form of recurrent neural networks (RNNs) such as long-short-term models (LSTMs) or convolutional neural networks (CNNs).

6 Conclusion

Overall, we found that our research represents a promising jumping-off point for future discoveries in the context of our problem. After utilizing the ”tip” text of venues scattered throughout the 40 neighborhoods of Manhattan we found that we could more accurately predict the like count these venues received than simply assuming the like count for each venue was the average of the like counts for all venues. This discovery can now be used as a proof-of-concept for popularity-prediction in a variety of other datasets.