*Article*

# The Geography of Taste: Using Yelp to Study Urban Culture

## Sohrab Rahimi [1,*], Sam Mottahedi [2] and Xi Liu [3]

[1]  Department of Architecture and Landscape Architecture, Pennsylvania State University, University Park, PA 16801, USA

[2]  Department of Architectural Engineering, Pennsylvania State University, University Park, PA 16801, USA; s.mottahedi@psu.edu

[3]  Department of Geography, Pennsylvania State University, University Park, PA 16801, USA; xiliu@psu.edu

*  Correspondence: sur216@psu.edu; Tel.: +1-781-2965152

check for updates

**Abstract:** This study aims to put forth a new method to study the sociospatial boundaries by using georeferenced community-authored reviews for restaurants. In this study, we show that food choice, drink choice, and restaurant ambience can be good indicators of socioeconomic status of the ambient population in different neighborhoods. To this end, we use Yelp user reviews to distinguish different neighborhoods in terms of their food purchases and identify resultant boundaries in 10 North American metropolitan areas. This dataset includes restaurant reviews as well as a limited number of user check-ins and rating in those cities. We use Natural Language Processing (NLP) techniques to select a set of potential features pertaining to food, drink and ambience from Yelp user comments for each geolocated restaurant. We then select those features which determine one's choice of restaurant and the rating that he/she provides for that restaurant. After identifying these features, we identify neighborhoods where similar taste is practiced. We show that neighborhoods identified through our method show statistically significant differences based on demographic factors such as income, racial composition, and education. We suggest that this method helps urban planners to understand the social dynamics of contemporary cities in absence of information on service-oriented cultural characteristics of urban communities.

## 1. Introduction

Socioeconomic polarization is a defining characteristic of cities in the global economy [1,2]. In global markets where economic regulations are minimized, social polarization is an inevitable consequence given the relatively small proportion of the population involved in this growing affluence [3]. In case of the U.S., this social polarization is also ethnic/racial as the prosperous economy in the U.S. was accompanied by massive immigration waves from other countries adding more dimensions to the long-lasting Black and white dichotomy. Not surprisingly, immigrants targeted large cities where most industries were located at and this, in part, led to more diversity in urban population. The multitude of cultural/ethnic groups led to cultural polarization and fragmentation of these global cities where every ethnic group occupied a piece of land [4]. Therefore, the American metropolis is plagued by both cultural and economic polarization [3].

During the past four decades, the debate over the definition and qualities of urban communities in developed countries grew significantly. Overall, scholars have different opinions regarding the strength of communities. Some believe that the notion of community is lost, some believe it has not

changed significantly and other say that it's been liberated from their constraints [5]. However, many of the recent studies have shown that the liberated hypothesis is more representative of the state of modern communities [5–7]. These studies assert that telecommunication and mobility has encouraged dispersed networks of friendship, kinship or communities of interest. Under this condition, the individual's network is a personal choice that she is free to choose from. Even though telecommunication has facilitated broad networks over space, the spatial segregation instigates sharp borders between communities in American cities. Emphasis on diversity and seeing the city as a melting pot, which is championed by postmodern thinking, has not addressed the gaps between ethnic and economic groups [8].

Many studies have attempted to fathom the sociospatial complexities that emerged in post-war American cities. Most of classic studies of this kind were based on the Census data [9,10]. Although the U.S. Census data provide valuable information about cities, these data hardly inform us about lifestyles, consumption behavior, cultural factors, and space-use patterns. The past two decades have seen a rapid advancement in the field of urban and social studies partly due to emergence of new crowd-sourced data sources and computation techniques [11]. The new data sources have enabled the researchers to go beyond basic demographics such as race and income and delve into a multitude of sociospatial phenomena in modern cities.

This study aims to contribute to this line of studies by proposing taste as an indicator of social status which integrates different facets of culture, economy and social networks of urban inhabitants. In sociology, taste refers to an individual's personal and cultural choice and preference patterns. Taste enables individuals to draw distinctions between a variety of cultural products such as styles, manners, consumption patterns, and art [12]. Taste is closely knit with social relations and human dynamics and many studies have studied taste in relation with aesthetics [13], consumption patterns [14], and social classes [15].

In this study, we also argue that using businesses as sensors can provide new insight into practiced taste in a region and consequently, the intricate social structure of the American metropolis. More specifically, this research aims to answer the following questions:

- To what extent is taste a good indicator of socioeconomic status of communities in American cities?
- By utilizing the concept of taste, can we use restaurant-as-sensor instead of citizen-as-sensor to examine the socioeconomic dynamics of neighborhoods? This issue is especially important to us since business data is far more accessible and plentiful than individual-level data.
- Are American cities comprised of regions with different dominant taste cultures? Are different regions in every city similar to regions from other cities?

## 2. Materials and Methods

### 2.1. Literature Review

#### 2.1.1. Previous Attempts to Define Sociospatial Boundaries

Recently, many studies have addressed these problems by using heterogeneous data sources that are updated frequently and exist at the scale of buildings or individuals [11]. Some investigate the communities on a large scale. For example, one study used vehicle GPS traces in Pisa, Italy to build a network and used community detection algorithms (i.e., Infomap) to identify non-overlapping communities of people at the county and municipality scale [16]. A similar study was conducted on a larger scale in Great Britain using telecommunications data [17]. Recently, detecting communities on urban scale has been more popular. For example, one study uses human mobility between different regions and the Points of Interest (POI) data to find the dominant functions of each urban region using topic-based inference model [18]. Using this model, this research identifies nine functional regions using clustering techniques.

Most often, urban studies that use crowd-sourced data to study the sociospatial structure of cities incorporate Location-Based Social Network (LBSN) techniques, that is, a network consisting of people in a social structure who share location-embedded information [11]. Much of research in this area uses social media data which includes the geographical location as well as their tagged images, videos, and texts. Common examples of data used for LBSNs include GPS trajectories of taxis, Twitter, Call Data Records (CDRs), Flickr geo-tagged photos, and Foursquare check-in data. Georeferenced crowd-sourced data such as tweets, photos, and check-ins can help understand people's lifestyles (e.g., likes and dislikes [19–21], cities' sociospatial structure [22], neighborhood functions and characteristics [23] and behavioral patterns [24] in cities.

One of the common techniques for studying urban structure is identifying similarities between users in terms of their use of urban spaces [22,25–28]. For example, among the most well-known studies of this kind is the Livehoods project, which uses check-in data to identify the zones where their establishments (e.g., restaurants and bars) share similar clientele [22]. This study uses 18 million check-ins collected from Foursquare, a location-sharing service where users share their location by checking in via their smart phones. By using clustering techniques, this study identifies clients with similar points of interest (POI). In another study, the authors studied the semantics of different locations by analyzing different categories of POIs in many neighborhoods [28].

Although the state of the art techniques used in these studies have dramatically improved our understanding of cities, they still have some limitations. First, accessing data that include individuals' behavior is often hard and these data are not freely available to the public. For example, companies which maintain a great inventory of georeferenced social networks do not share such information due to privacy issues. Second, the data is not often representative of the entire population. For example, not everyone has a Foursquare account and not all those account owners use Foursquare every time they visit a place [22]. Third, these studies only address one aspect of an individual's life, for example, Foursquare only covers check-in data and points of interest (POIs), and taxi data cover some travel patterns. While these datasets have proven helpful, multiple data sources need to be fused to provide an understanding of urban lifestyle.
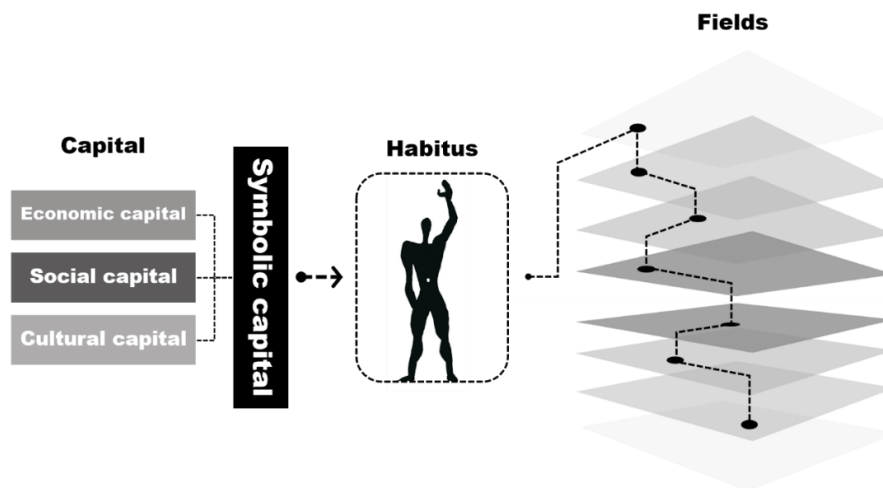
### 2.1.2. Taste as an Indicator of Urban Culture

The social construct of taste is a well-studied topic especially in the age of Internet where individual preferences are available to information-based companies (e.g., Amazon, Facebook, Spotify). In fact, many recommender systems (i.e., algorithms made for recommending products to users) are designed under the same assumption that people of same social groups are likely to consume similar products [29,30]. The underlying mechanism of the relationship between social groups and taste was discussed by Pierre Bourdieu in his well-recognized book Distinction: A Social Critique of the Judgment of Taste [15].

In Bourdieu's view, both cultural and economic capital are the most important forms of capital. Economic capital has to do with individuals' access to economic resources while cultural capital is a collection of non-material traits in a person, such as knowledge and skills, attitudes, philosophical views, use of vocabulary, and language skills. Bourdieu believes that taste is the means of identifying class distinction. He argues that these differences are most obvious in the routine everyday choices in taste of food, furniture, and clothing as they are representative of the pure taste. For example, he argues, children of a lower social status like plentiful and good meals while those of higher status go for original and exotic. These choices, according to Bourdieu, become intrinsic to one's personality and thereon he/she rejects the tastes of other groups. Bourdieu argues that high-taste is characterized by how far it is from pure necessities. The upper classes in this regard use taste as the ideal weapon in strategies of distinction [15].

Many studies followed Bourdieu's theory of distinction to determine how demographic factors were correlated with taste. For example, some studies showed that generally, people of higher economic status read more literature and quality papers [31] and have different taste in art [32].

More recent studies on Facebook and MySpace datasets, argue that people with similar social networks share similar tastes of music, movies, TV shows and books [20,33]. In all these studies, taste is seen as a means of distinction between different groups of people, which further supports Bourdieu's argument. According to Bourdieu, individuals may play different roles in different fields of a society (i.e., sub-spaces of society such as friend groups and institutions). The quality of these roles relies heavily on an individual's symbolic capital, which Bourdieu defines as a combination of social, economic and cultural capital. As discussed earlier, Bourdieu believes that taste best reflects the symbolic capital, which is the main reason of distinction in societies (Figure 1).



**Figure 1.** Bourdieu's theory of distinction. Fields refer to different sub-spaces of society such as family groups and work groups. Individuals' role in these fields is influenced by her symbolic capital.

2.1.3. How Can Information about Restaurants Help Us Understand the Socioeconomic and Cultural Structure of Cities?

Businesses are an effective type of sensors that can reflect what is accessible and offered to a neighborhood. Theoretically, it is not surprising to expect geographically concentrated clusters of similar tastes between individuals in American metropolitan areas: first, as discussed earlier, these cities are characterized by highly fragmented social fabric with segregated communities of different taste, culture, ethnicity and economic status. Second, their economies are global and products of all types belonging to all different cultures and nationalities are offered in the marketplace and therefore the consumer is offered a variety of goods from which she can choose [34]. Third, in case of the U.S., the rise of individualism and diversity along with the economic growth of the post-war period has generated a dominant landscape in cities known as consumption spaces. These spaces gradually took the place of production spaces such as factories after the era of industrialization [35–38]. The emergence of these spaces is a result of the increasing impact of consumerism, pushing the individuals towards consuming goods and certain types of services [4,39–42].

Restaurants are one of the most common and frequently used consumption spaces [43]. In the U.S., restaurant expenditures exceed spending in higher education, computers, books, magazines, newspapers, movies and recorded music [44]. Data on consumption behaviors in restaurants is available in different social media venues such as Yelp. Yelp is a web-based application which maintains crowd-sourced reviews of local businesses (i.e., mostly restaurants, coffee shops and bars). In this platform, reviewers sign in with pseudonyms and provide reviews and evaluate the performance of local businesses. Yelp users have generated nearly 127 million reviews for different businesses across the world [45]. As in any crowd-sourced dataset, Yelp suffers from some biases the most important of which is the risk of fake review. Although Yelp has certain strategies to classify fake reviews from the real ones, some commentators claim that approximately 20% of the reviews are fake [46] and Yelp is still dealing with constant complaints from time to time leading to lawsuits [47]. Despite these

potential biases, due to its extensive geographic coverage as well as large review corpus, Yelp has been used extensively in different studies [48–51] in a variety of research topics including dietary and food health [51,52], hospitality and tourism [53], document modeling [54] and sentiment analysis [55].

Georeferenced crowd-sourced data sources such as Yelp, TripAdvisor, and Airbnb can play supplemental roles to common datasets such as U.S. Census data, to inform us about the city dynamics. Although the U.S. Census data provides valuable information from cities, this data hardly informs us about lifestyles, consumption patterns, cultures, and space use patterns. In addition to providing a different understanding from the sociospatial dynamics of urban areas, new data sources are current and often collected with high spatial resolution (i.e., geographic coordinates) whereas census data is aggregated on different sets of spatial units that often lack desired accuracy in many cases. Moreover, Census conducts a national level data collection only every ten years which is a long period under the extremely dynamic social life of postmodern cities. Although new datasets cannot substitute the valuable demographic data provided by the Census, they have proved to be great assets to further our understanding of the complexities of global cities [56].

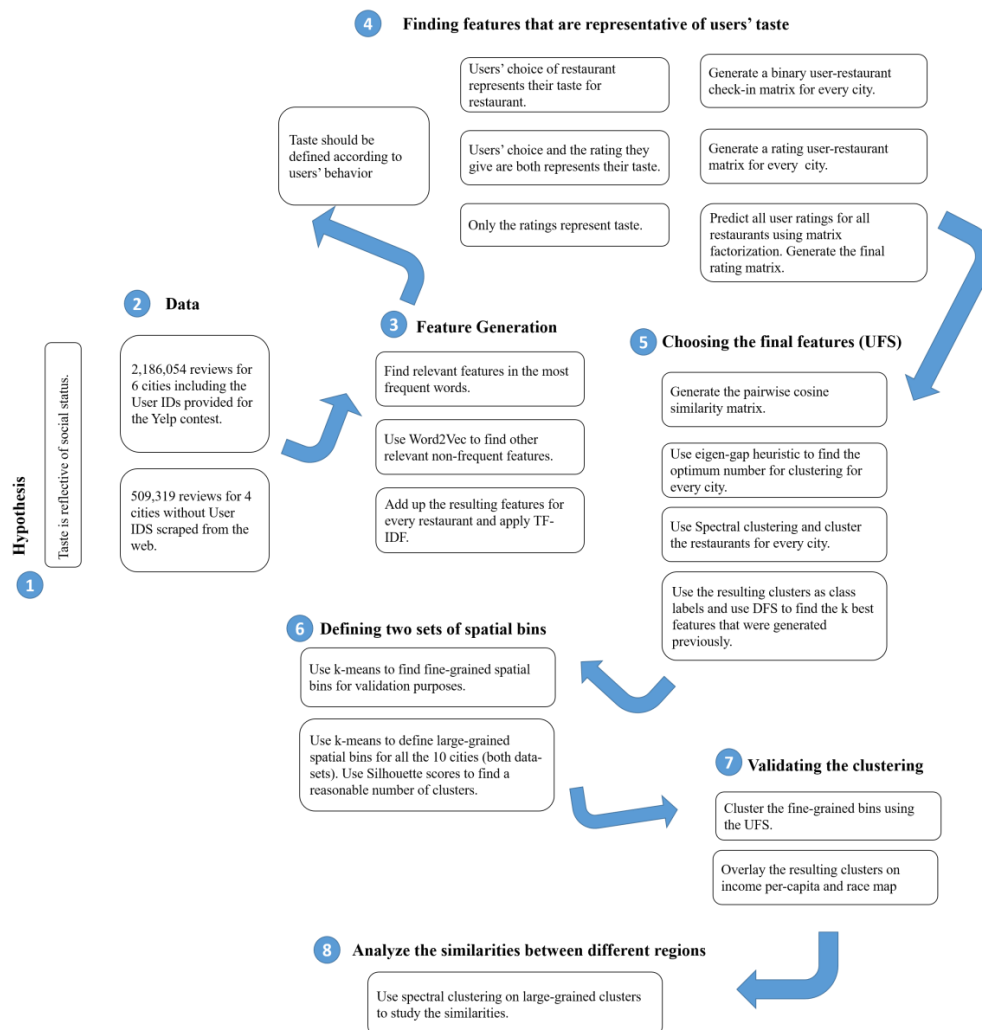*2.2. Data*

Two sets of data were used in this research:

- Data provided by Yelp [57] which includes 11 cities, 8 of which are in North America (i.e., Cleveland, Pittsburgh, Charlotte, Urbana-Champaign, Phoenix, Las Vegas, Toronto, and Montreal). This data includes 4.1 M reviews by 1 M users for 144 K businesses as well as 1.1 M business attributes (e.g., hours, parking availability and ambience). For the case of Montreal, of 86,054 reviews 11,284 were in French, as identified through langdetect 1.0.7 package in Python [58]. Since English reviews may not equally represent all Montreal neighborhoods, demographics, and resident population, we considered Montreal as an outlier and removed it from our analysis. For this study we were only interested in restaurants in English-speaking North American metropolitan areas therefore, we filtered out Montreal and Urbana-Champaign (a small city) as well as points that fell out of the metropolitan boundaries. Also, we only used businesses tagged as restaurants. This process resulted in 2,186,054 reviews for 34,231 restaurants. This data includes the following fields: Business ID, User ID, Reviews, Business Name, Star Rating, Address, City, State, Zip code, Business Category, Review Count, Longitude, Latitude. The geographic coordinates represent the location of businesses.
- As we discussed in the introduction section, we intended to see if we can characterize the socioeconomic status of urban communities without having information about users. This is very important, because although it is possible to scrape data from different websites such as Yelp, the user IDs are often not provided in the interface and cannot be scraped easily. In other words, extracting information from businesses from the web is often easier than finding individual-level data. To investigate the extent to which business-level data scraped from the web and stripped from user IDs can inform us about neighborhoods, we scraped restaurant reviews and attributes for Boston, Washington D.C., Detroit and Philadelphia metropolitan areas. The data collection process took a few months in 2017, the same year as the Yelp contest data. All these cities are characterized by high segregation as well as ethnic and cultural diversity. This data includes 509,319 reviews for 120,801 restaurants. Using the earlier dataset, we expect to be able to study the communities in this dataset where the user IDs are absent. In addition, the four cities are important metropolitan areas and studying the sociospatial dynamics of these cities can be useful per se. Table 1 provides a summary of these the data.

**Table 1.** Number of reviews and restaurants for the 10 cities.

| City | MSA Population (2017) | Number of Restaurants | Number of Reviews | Reviewers' User-ID | Number of Reviewers |
|---|---|---|---|---|---|
| Boston | 4,836,531 | 44,597 | 172,401 | Not available | Not available |
| Charlotte | 2,525,305 | 2780 | 139,188 | Available | 39,813 |
| Cleveland | 2,058,844 | 3996 | 139,824 | Available | 21,939 |
| DC | 6,216,589 | 8206 | 40,420 | Not Available | Not available |
| Detroit | 4,313,002 | 35,823 | 81,301 | Not available | Not available |
| Las Vegas | 2,204,079 | 6312 | 826,358 | Available | 275,012 |
| Philadelphia | 6,096,120 | 29,045 | 91,660 | Not available | Not available |
| Phoenix | 4,737,270 | 9692 | 731,744 | Available | 97,476 |
| Pittsburgh | 2,333,367 | 3130 | 124,170 | Available | 33,268 |
| Toronto | 6,417,516 | 11,451 | 357,940 | Available | 58,355 |

## 2.3. Methodology

In using the Yelp dataset, our assumption is that when a person talks about a food or drink in her comment, she has purchased or at least considered that food or drink and therefore, it can be used as an indicator of one's choice of food or drink. In the following sections we explain our methods for this research. Figure 2 summarizes our workflow.



**Figure 2.** Research workflow.

2.3.1. Feature Generation

In this study we use the text provided by Yelp reviewers when they post restaurant reviews on Yelp.com. We use a bag-of-words model to define features for every restaurant. In this model the existence of a word, regardless of the way its embedded in the comment, is considered. A bag-of-words model is suitable for our case, as we are only interested in the frequency of these words and not the way they're used in the sentence. According to Bourdieu's theory of distinction, food, drink, and interior decoration are among the best indicators of taste reflecting one's everyday choice [15]. We are, therefore, interested in three categories of features: foods and drinks (e.g., pizza, martini), adjectives used to describe foods (e.g., fried, steamed), and adjectives described for ambience (e.g., rustic, minimalist). We assumed that ambience is an equivalent of decoration. Ambience are among those concepts that are frequently discussed in Yelp reviews along with food, price, and service [59] and provide an overview of the restaurants atmosphere and decorative features such as classy, intimate, romantic, hipster and so forth. In choosing features we avoided selecting words that have multiple connotations or are too general (e.g., nice, green).

In order to select relevant features from reviews, we used a four-step process:

1. First, we used English stop-words to remove commonly-used words [60] and then, chose features among the top 1000 frequent words. Forty-five features of the three categories (i.e., foods and drinks, food adjectives, and ambience adjectives) were selected at this step (Appendix A).

2. Although frequent features can provide much information for restaurants, we expect to get more specific words from the comments. For example, different types of fish (e.g., haddock, tilapia) or different adjectives used to describe an ambience (e.g., divey, hipster) are not among frequent words. To address this problem, we used the Word2Vec model. This open-source model was developed by Google in 2013 which transforms words in a document to high-dimensional spatial vectors by using a Neural Network Language Model (NNLM) [61,62]. Given $N$ user comments and the $n$-th word in the comment $w_n$ and the window size of the context centered on the $n$-th word as $C$, the maximum likelihood function of the NNLM model will be as follows:

$$I(\theta) = \frac{1}{N} \sum_{i=1}^{N} log \ p(w_n | w_{n-c}^{n+c}) \tag{1}$$

where $w_{n-c}^{n+c}$ represents a set of words at the center of which is $w_n$ with context sampling window size of $c$. Word2Vec suggests two mathematical frameworks for solving Equation (1) i.e., Continuous Bag-of-Words (CBOW) and Skip-Gram. In summary, Skip-Gram uses stochastic processes to sample from the words whereas CBOW offers a continuous input and training mechanism. In this study, we use CBOW to train the model as some studies suggest it has a better performance at characterizing the words [63]. We trained our Yelp corpus with this model and every word was turned into a 100-dimensional vector. As an example, Table 2 shows the closest words to the word classy. It is noteworthy that the model does not necessarily return synonyms of classy but rather, it considers the way word classy is used in a sentence and therefore, it returns all adjectives that are used to describe an ambience. The 45 words chosen in the last step were given as input to this model to find the 20 closest words in cosine distance. However, not all these 20 words were relevant to food, drink, or ambience. Accordingly, we went through all the 900 words (i.e., 45 × 20) and selected related words subjectively. It is important to note that Word2Vec model significantly simplified the filtering process and instead of going through all the words in the corpus, we just went through the Word2Vec outputs that is 900 words total. At the end of this step, a total of 454 features were selected.

3. We binarized the number of words selected from the last step in each comment (1 word exist 0 otherwise) and aggregated them for every restaurant. Given that these words are not equally

common we use Term Frequency-Inverse Document Frequency model (TF-IDF) to weight these features:

$$idf(t, D) = log \frac{N}{1 + |\{d \in D : t \in d\}|} \tag{2}$$

where $N$ is the total number of restaurants in the corpus and $|\{d \in D : t \in d\}|$ is the number of times that term $t$ appears in the restaurant $d$. We can then multiply IDF by the Term Frequency (TF) that we previously generated. After this step, for every restaurant, we will have 454 features that are properly weighted.

4. The features generated in the previous steps can sometimes fall into categories which can be even more important than the individual features themselves. For example, specific fish types (e.g., salmon) might be important but less informative than the combination of all types of fish. This information tells us that seafood is popular in a certain area. Appendix B indicates the groups of features that we combined in order to generate new features. By including these new features, a total of 477 potentially-unnecessary features remain (e.g., does the word "water" really explain anything about a community's taste?). In the next step, we explain our methodology for reducing the dimensionality and choosing the most important features.

**Table 2.** Top 10 most similar words to classy.

| Word2Vec Output | Similarity to Classy |
| --- | --- |
| swank | 0.87688 |
| trendy | 0.86152 |
| chic | 0.85917 |
| posh | 0.84972 |
| elegant | 0.84592 |
| stylish | 0.84019 |
| cozy | 0.83344 |
| modern | 0.80526 |
| contemporary | 0.78569 |
| homey | 0.77934 |

### 2.3.2. User's Taste and the Curse of Dimensionality

In the feature generation process, we took an inclusive approach and considered all features that could possibly represent user taste. Considering all these features for clustering is problematic due to high dimensionality. It is also unclear whether these features represent people's taste. In other words, we are interested in a subset of features that distinguishes between different groups of users in terms of their practiced taste. For example, the word water may be used equally in all restaurants. In this case, considering water not only does not add any additional information about different neighborhoods but also increases the dimensionality. Therefore, it is important to only select those features that have to do with people's taste.

Recall that the data-set provided by Yelp includes User IDs as well as user-generated ratings for rated restaurants. This data can assist us to select a subset of the 477 features that actually has to do with users' taste of food, drink, and decoration. Therefore, we examined three scenarios to select the best features related to taste:

**1. Users' choice of restaurant represents their taste for food, drink and decoration:** Under this assumption, a person's taste is only reflected in the type of restaurants she chooses to visit. Therefore, if we find clusters of restaurants that have been visited by similar users, we should be able to find distinguishing features between these clusters. To this end, we first create a matrix for every city showing whether a user has visited a restaurant (1) or not (0). We generate this matrix for each city separately to reflect how a user living in one city is more likely to go to restaurants in the same city. By separating the cities, the effect of geography is minimized and we can draw our focus on the effects of restaurant attributes on users' choice of restaurant. Of all the 525,863 Yelp reviewers,

311,866 reviewers have provided only one review. We removed users with only 1 review since first, these reviews are more likely to be biased and have extremely high or low ratings and we will use the ratings in the next steps. Second, excluding these would reduce the computational costs and also increase the accuracy of our clustering, which we will explain in the next steps. From these matrices, we generated a pairwise similarity matrix using cosine distance:

$$\cos(A, B) = \frac{AB}{|A||B|} = \frac{\sum_{i=1}^{n} A_i B_i}{\sqrt{\sum_{i=1}^{n} A_i^2} \sqrt{\sum_{i=1}^{n} B_i^2}} \tag{3}$$

where restaurant A and restaurant B are n dimensional vectors with n being the number of Yelp users in each city. Every element of A and B is 1 if a given user has reviewed that restaurant and 0 otherwise.

In the next step, we used spectral clustering [64] to first find the restaurants with similar clientele. This method constructs a graph from the similarity matrix, where the data points (i.e., restaurants) are the nodes and the similarity between them are presented as weighted edges. The algorithm finds partitions of the similarity matrix by detecting low-weight edges. More specifically, this algorithm first performs a dimensionality reduction and then applies a k-means clustering [65] on the low-dimensional embedding. To reduce the dimension, the algorithm first generates a Graph Laplacian $L$ [66]:

$$L = 1 - D^{-1}W \tag{4}$$

where $D$ is the degree matrix with diagonal terms, $d_i = \sum_{j=1}^{n} W_{ij}$, and $W$ is the adjacency weight matrix of an undirected graph. The Laplacian matrix $L$, in fact, is used to calculate the eigenvalues for the matrix. The k-means clustering will then be applied to these eigenvalues, which represent an image of the similarity matrix in a lower-dimension space. Since the k-means is applied to a reasonably lower dimension, the resulting clusters are expected to be more distinguishable and informative. To ensure an optimal number of clusters, we use eigen-gap heuristic method [64] to find the largest difference between two consecutive eigenvalues of the Laplacian matrix and set the number of clusters equal to the rank of the eigenvalues (Figure 3). The check-in row in Figure 4 shows the resulting eigen-gaps for different number of clusters. As we can see, for Pittsburgh for example, 2 is the best number of clusters for the check-in matrix.

We then select the k best features (from those 477 features) that affect the membership status of a restaurant in one of those previously defined clusters. In other words, we discover which subset of the 477 features actually distinguishes between the clusters using a Deep Feature Selection (DFS) model [67] to select features at the input level of the deep network. The DFS model used in this study has the following network structure $477 \rightarrow 477 \rightarrow 256 \rightarrow 64 \rightarrow 16$ with a softmax output layer. The first one-to-one linear layer $w$ between the input layer and the first hidden layer with linear activation function is regularized using an elastic-net [68]. The resulting sparse one-to-one layer weights $w$ only selects those features corresponding to none-zero terms in $w$. The model parameters are learned by minimizing this Equation (5).

$$\begin{aligned} min_\theta f(\theta) &= l(\theta) + \lambda_1 \left( \frac{1-\lambda_2}{2} \|w\|_2^2 + \lambda_2 \|w\|_1 \right) \\ &+ \alpha_1 \left( \frac{1-\alpha_2}{2} \sum_{k=1}^{K+1} \|W^{(k)}\|_F^2 + \alpha_2 \sum_{k=1}^{k+1} \|W^{(k)}\|_1 \right) \end{aligned} \tag{5}$$

where $l(\theta)$ is the log-likelihood of the data, the matrix $W^{(k)}$ is the $k$th hidden layer weights and $\lambda_{1,2} \in [0, 1]$ is the parameter that controls the sparsity of $w$ and the term $\alpha_{1,2}$ is another elastic-net like term that reduces the model complexity and increases the speed of optimization.

To find the best subset of features, we tuned hyper-parameters $\alpha_{1,2}$ and $\lambda_{1,2}$ corresponding to the sparsest model with the highest prediction accuracy measured using $F_1$ score which is a weighted harmonic mean of the precision and recall metrics described below:
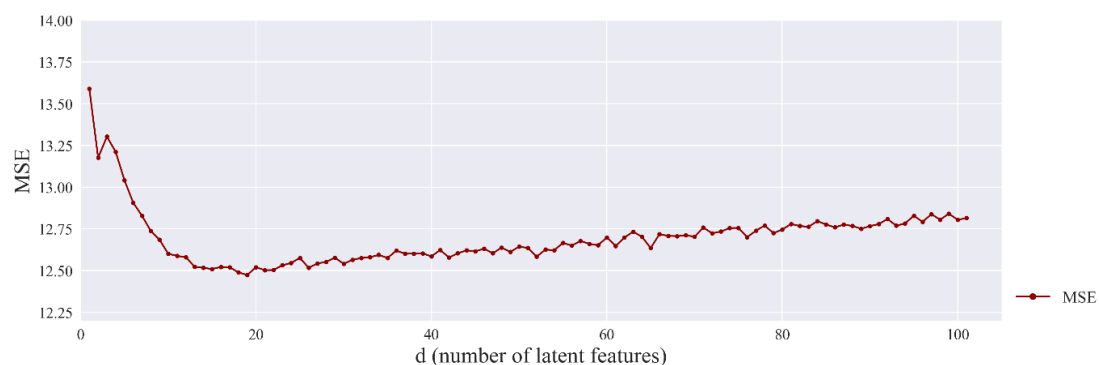
$$\text{recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \text{ and precision} = \frac{\text{TP}}{\text{TP} + \text{Fp}} \qquad (6)$$

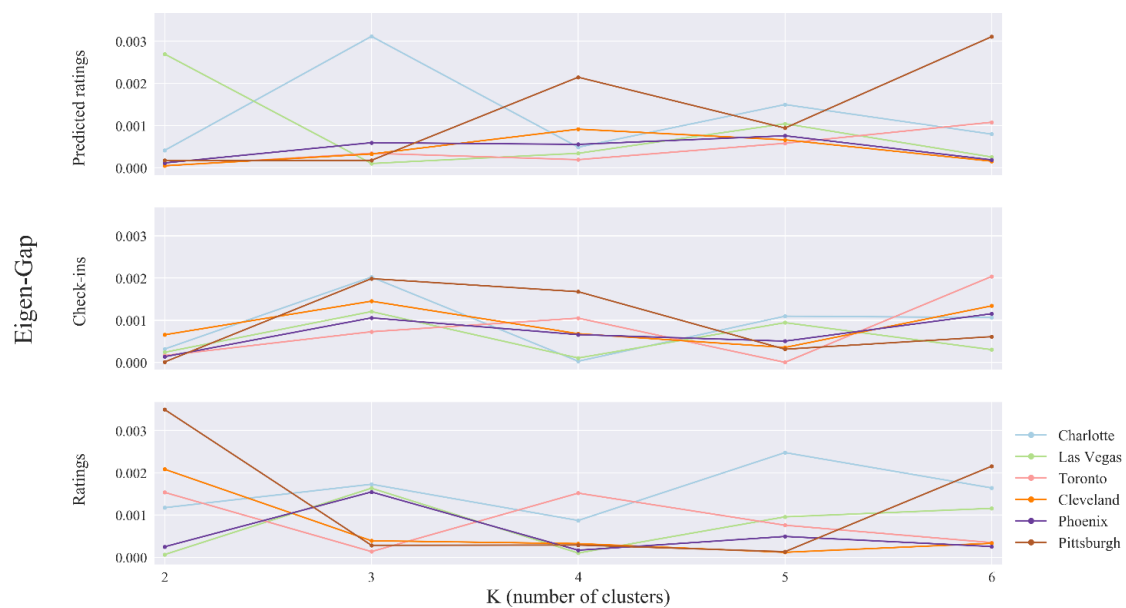$$F_1 = 2\,\frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} \qquad (7)$$

where TP, FP and FN stand for true positive, false positive and false negative respectively [69]. Since the data for each city is moderately small, 10-fold cross-validation was performed to prevent over-fitting to the training data set.

**2. Users' choice and the rating they provide both affect their taste for restaurant:** The only difference between this hypothesis and the first one is that the rating that one provides for a restaurant acts as a weight to the check-in matrix from the last hypothesis. Accordingly, in this hypothesis, not all restaurants visited by the user are equally important, but rather, we assume those that the user rates higher are more important in determining one's taste.

**3. Only the users' ratings determine their taste:** In the second assumption we assumed that taste is reflected in the way people rate a restaurant. The only difference here from the last assumption is that we try to see what would happen if every user rated every restaurant. Under this assumption, however, a problem arises: The rating matrix is sparse and many ratings for many restaurants are missing. Using the original rating matrix cannot help us identify how would every user like every restaurant. Therefore, we will need to predict the ratings by using matrix factorization method [70]. The fundamental assumption of this method is that there are d latent features in restaurants that affect the users' ratings. The advantage of this method is that without having to know what those d features are; we can predict how users might rate restaurants which they have not yet reviewed. We use Singular Value Decomposition (SVD) method to factorize the rating matrix [71]. To find the best number for *d*, we used 10-fold cross validation. The results indicate that there are approximately 20 latent features (*d* = 20) that affect one's rating for a restaurant. The Mean Square Error (MSE) decreases significantly up to *d* = 20 and gradually increases afterwards due to being over-fit (Figure 3). After predicting the rating matrix with 20 latent features for every city, we repeat the steps described in the last two hypotheses. In all three hypotheses above, we selected the number of clusters with the largest eigen-gaps (Figure 4) for every city. Table 3 shows the final number of clusters selected for different matrices and different cities.



**Figure 3.** 10-fold cross-validation results for rating predictions.

**Figure 4.** Eigen-gaps for different number of clusters and different matrices.

**Table 3.** Selected number of clusters for different matrices and cities.

| City | Predicted Matrix | Check-in Matrix | Rating Matrix |
|------|------------------|-----------------|---------------|
| Charlotte | 3 | 3 | 5 |
| Cleveland | 4 | 3 | 2 |
| Las Vegas | 2 | 3 | 3 |
| Phoenix | 5 | 3 | 3 |
| Pittsburgh | 6 | 3 | 2 |
| Toronto | 6 | 6 | 2 |

### 2.3.3. Defining the Spatial Bins

The features generated from the previous steps reflect Yelp reviewers' preferences in different urban areas. We next aggregate restaurant features on some spatial units fabric to ensure that nearby restaurants will fall in the same spatial cluster. Aggregating restaurants on geographic units will enable us to minimize the impact of outliers and noise. It also enables us to get an overall sense of taste preference given all different types of restaurants in a region. This practice, however, raises a new challenge which is, choosing the best spatial unit for this purpose.
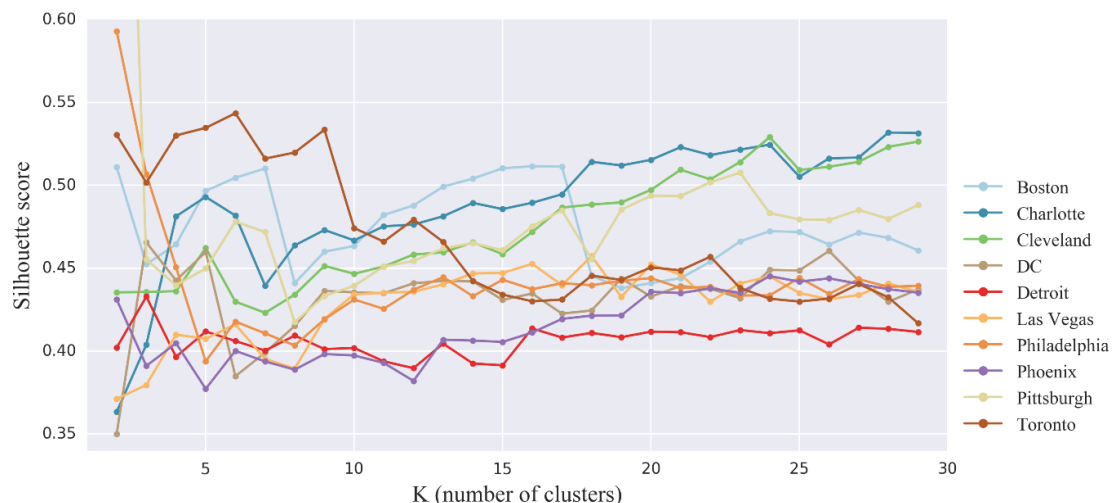
Accounting for sociospatial boundaries have long been an issue in the literature. Harvey proposes that appropriate scales for assessment of spatiotemporal patterns can be found by seeking the scale or spatial region across which the pattern ceases to be significant, i.e., description of the system is no longer accurate [72]. More recently, Kwan and Weber discuss the scale effect in spatial analysis as it relates to design of accessible and sustainable land use and transportation schemes [73]. Through this work, an offshoot of the scale problem known as the Modifiable Areal Unit Problem (MAUP), is used to describe the implications of how features in the built environment as well as the people who inhabit that system are aggregated for analysis. Specifically, MAUP concerns itself with how areas or regions are defined such that the aggregation of data about the contents of that region (e.g., people, localized places, things) is relevant to understanding a research question [73–75]. Traditionally, however, a large body of urban literature on crime in cities, for example, use Census tracts as the main spatial unit of analysis [76–79] which does not address the MAUP.

Since our sensors are restaurants, we define these geographical units based on their density and configuration and avoid using administrative boundaries e.g., block groups. Two sets of spatial bins are required to answer our research questions:

**1. Large-grained spatial bins:** These spatial bins enable us to compare different parts of cities together as to see how different cities interact in terms of food, drink, and decoration related attributes. The existing administrative boundaries are too small for this purpose. For example, we are looking at dividing up Washington DC to 3–6 parts and conventional administrative boundaries are too fine-grained for this purpose. Also, we intend to have reasonable spatial bins that are actually representative of the city form. The number of these bins is actually a matter of preference, however, for visualization and simplification purposes we choose large-grained clusters. Accordingly, we use k-means clustering on the restaurants' geographic coordinates to find reasonable spatial clusters. To find the best number of clusters for each city, we use the silhouette scores [80] for different number of clusters for every city. Silhouette score measures the extent of tightness and separation for each cluster. In other words, it specifies which objects are within their clusters and which ones are somewhere in between:

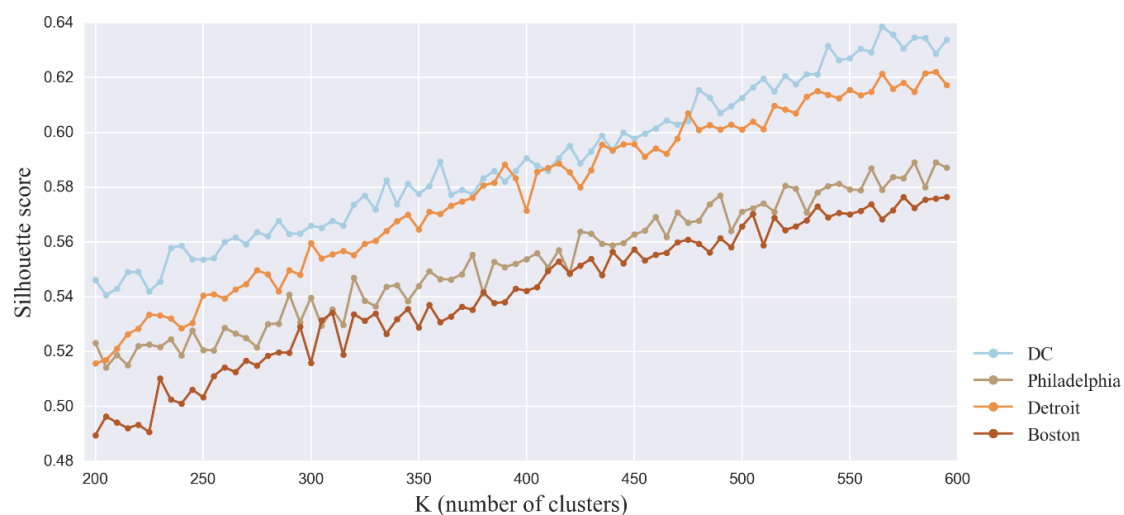$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \tag{8}$$

where $a(i)$ is the average dissimilarity of datum $i$ with all other data points and $b(i)$ is the lowest average dissimilarity of $i$ to any other cluster. We then average $s(i)$ over all data points, a measure that we used for goodness of clustering. Silhouette score ranges from $-1$ to 1, where 1 means that the clustering configuration is appropriate. Figure 5 shows the Silhouette scores when we divide each city to less than 30 clusters. At this point, we make a compromise between the number of restaurants in every city, area of the city as well as the Silhouette score.



**Figure 5.** Silhouette scores for different Ks for different cities (large-grained spatial bins).

**2. Small-grained spatial bins:** To validate the results and compare it with other demographic datasets, more fine-grained spatial bins are needed. Administrative boundaries are not helpful in this case either since these boundaries do not consider the formality of the built environment. For example, restaurants located on the Woodward Ave and East 9 Mile Rd cross section in Detroit, MI have been divided between four Census tracts, whereas they are all located near the same cross section and are very close to one another. Another problem with the administrative boundaries is that their sizes are not consistent with the distribution of the restaurants. For example, as we move to the suburbs of Detroit we can see tracts which contain one or two restaurants in them. Accordingly, same as the last step, we use k-means clustering and Silhouette scores to define these spatial bins. This method enables us to consider for the distribution of restaurants while defining the spatial bins. Figure 6 shows the Silhouette scores for the four cities. As we can see, for all these cities the Silhouette score improves as we increase the number of clusters. At this point, Silhouette scores are not useful for our purposes as they do not suggest any optimum number of clusters. Therefore, we base our decision on the number of restaurants and city area. Given the number of restaurants we have for every city (Table 1),

we expect about 200 clusters for Washington D.C., 500 for Detroit and Philadelphia and 600 for Boston. It is important to note that there are more census tracts in these areas than the number of clusters that we determined. For example, Detroit metropolitan area has 909 census tracts however, as discussed earlier, due to the uneven spatial distribution of restaurants, our spatial bins are larger than census tracts in the suburban areas with low number of restaurants, but smaller than block-groups in the city centers. It is important to note that the size and number of these spatial bins can change depending on one's research question as well as spatial resolution of the original dataset (i.e., Yelp in this case).



**Figure 6.** Silhouette scores for different Ks for different cities (small-grained spatial bins).

We use the small-grained and large-grained spatial bins defined in the last step in two different ways. The small-grained clusters are for validation purposes. Our purpose is to see if we can find any clear spatial pattern by clustering these fine-grained clusters. Using small bins enables us to assess the accuracy of this method and compare it with other high-resolution data sources. We will first average the selected set of features from the last step on these spatial bins, scale the features using min-max scaling for every bin, and then calculate the pairwise cosine similarity between the fine-grained bins separately for every city which we did not have information about user IDs (i.e., Philadelphia, DC, Detroit, Boston), using Formula (3). To calculate the similarities, we will use principal components instead of the actual features, to further reduce the dimension and improve the clustering results. For every resulting matrix, we will use spectral clustering method [64] as described in Section 3.2. We will then overlay the resulting clusters on the block-group level map of 2017 income per capita provided by Tableau 10.0 software for those four cities. At this point, we expect to see a geographic pattern in our clustering as well as a reasonable alignment between the clustering results and the block-group income per capita layer.

After validating, we can use the selected set of features from the last steps to study the interactions between different regions in cities. It is important to note that this capacity is the advantage of this set of features over using user ID data since, at this point, this feature set only relies on the aggregated comments for every restaurant and not the users' check-ins and ratings. To this end, we will average these features on the large grained clusters, calculate the pairwise similarities and cluster, same as the last step. Due to the extreme cultural, economic, and racial divisions in the American metropolis [2,4] we expect to see different clusters in every city and due to the global nature of these cities [2] we expect some regions from some cities to be similar to other regions in another cities.

## 3. Results

### 3.1. Selected Features

We took the steps described in Section 2.3.2, to reduce the dimension of the data set and only focus on those features that are actually representative of users' choice of food, drink and ambience. Figure 7 shows the resulting $F_1$ scores for the three hypotheses (i.e., check-ins, ratings, and predicted ratings) and the six cities where the user IDs were available (Table 1). For every city we selected a taste scenario that returned the highest $F_1$ score. The resulting features along with the scenarios that returned the highest $F_1$ scores, as well as the $F_1$ scores are presented for every city in Table 4. By considering all these features, we will have a total of 105 features which we call the Universal Feature Set (UFS). We can now use the UFS to study those cities where the user data is not available. The underlying assumption here is that the six cities that we have based the UFS on, are diverse enough that cover the types of food, drink and ambience that one expects to find in the four other cities where the user IDs are not available.
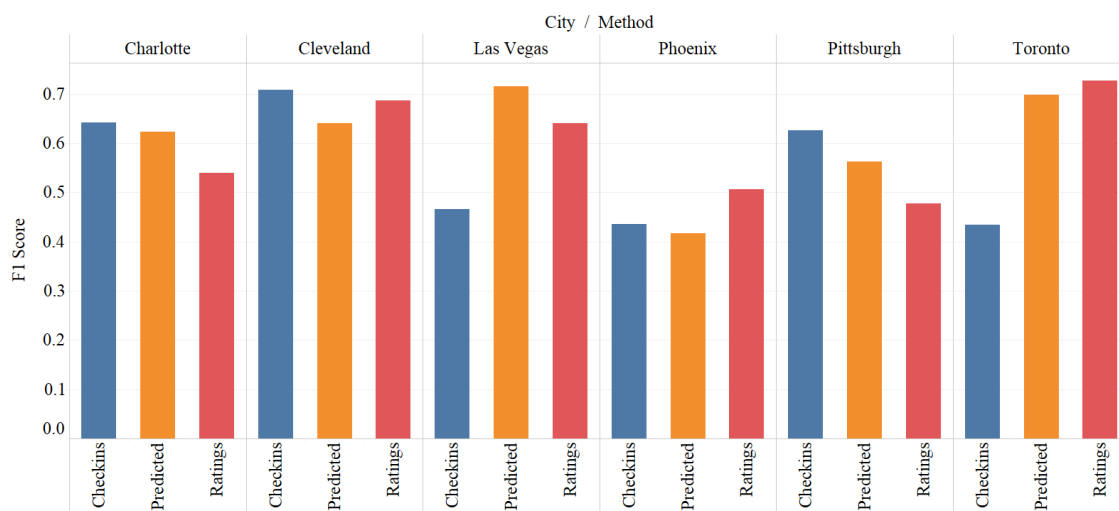


**Figure 7.** $F_1$ scores resulting from classification for different cities.

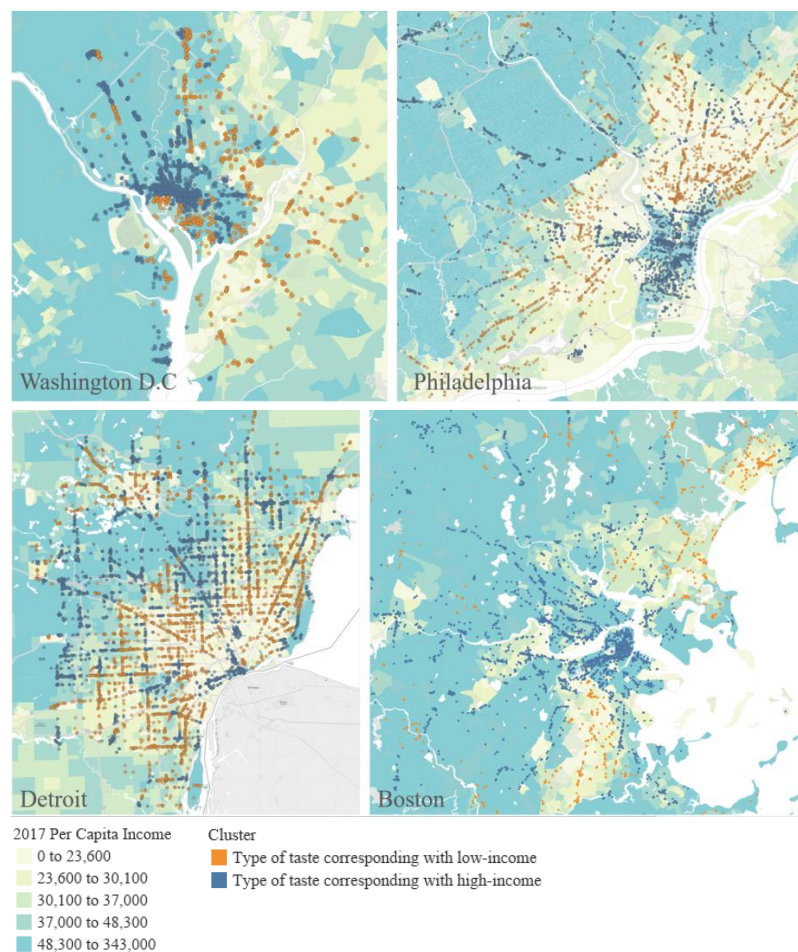**Table 4.** Selected features for different cities.

| City | Best Method | $F_1$ Score | Selected Features |
|---|---|---|---|
| Charlotte | Check-ins | 0.64244 | salty, vegetarian, creamy, hipster, divey, dessert, calamari, asparagus, vodka |
| Cleveland | Check-ins | 0.70865 | sweet, spicy, hipster, tomato, lime, meat_types, vegie_types, herb_types |
| Las Vegas | Predicted Ratings | 0.71563 | braised, seared, salty, creamy, intimate, classy, modern, casual, upscale, elegant, rice, soup, wine, crab, salmon, lobster, lamb, dessert, duck, cocktail, calamari, martini, ranch, steak_types, vegie_types, herb_types, hardliq_types, sofliq_types, sweet_types, asian_types, seafood_types, pos_ambience, neg_ambience, style_types |
| Phoenix | Ratings | 0.50608 | spicy, upscale, wine, pos_ambience |
| Pittsburgh | Check-ins | 0.62651 | crispy, vegetarian, hipster, romantic, rice, noodle, curry, sausage, cocktail, tofu, coleslaw, wing, cheesesteak, lettuce, provolone, ranch, fast_food, dressing_types, pos_ambience, style_types |
| Toronto | Ratings | 0.72686 | fried, Chinese, salty, Asian, Japanese, steamed, oily, hipster, rice, beer, soup, pork, shrimp, wine, tea, noodle, seafood, cocktail, sashimi, soy, squid, milk, sesame, Fanta, meat_types, softliq_types, Asian_types, soda_types, seafood_types, ethnic_food |

### 3.2. Clustering Results

Results derived from clustering the small-grained spatial bins with the selected set of features reveal clear geographic patterns which correspond with block-level per capita income for the four cities where the user IDs were absent (Figure 8). We set the number of clusters on two (k = 2) for the ease of comparison (Table 5).
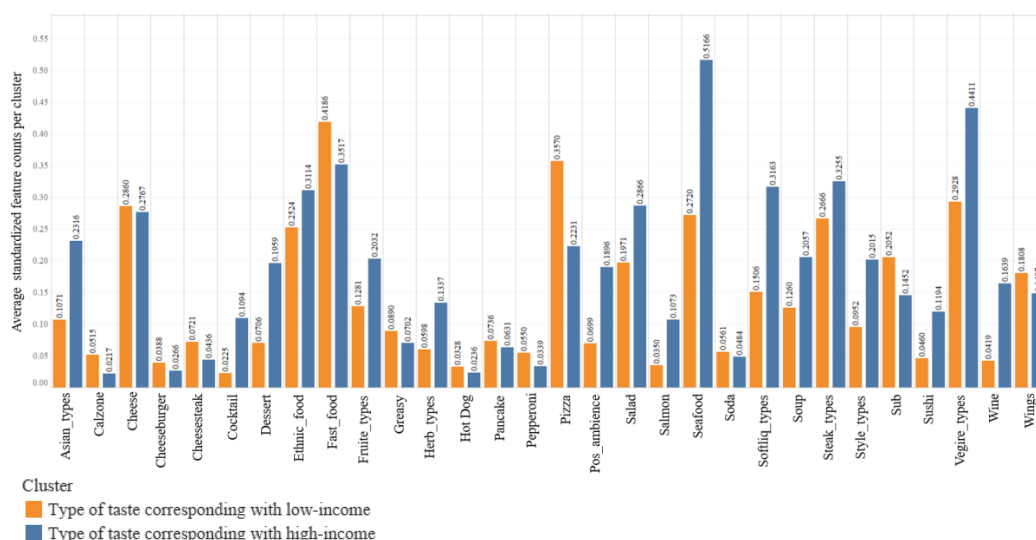
The difference between the type of tastes practiced between the two clusters is shown in Figure 9. This figure shows the top 30 features with highest average difference between the two clusters. As we can see, features such as seafood, salad, ethnic foods, vegetables, fruits, and Asian food types show higher values in high-income communities whereas the low-income cluster shows higher consumption of fast food.



**Figure 8.** Clustering results overlaid on per capita income map for four cities. As we can see the two clusters clearly correspond with block-group level income per capita map from Census.

**Table 5.** Number of restaurants in the two clusters for different cities.

| Cluster | Boston, MA | Detroit, MI | Philadelphia, PA | Washington, D.C. |
|---------|-----------|-------------|------------------|------------------|
| Cluster 1 | 16,827 | 17,226 | 13,849 | 2780 |
| Cluster 2 | 27,770 | 18,597 | 15,180 | 5419 |

**Figure 9.** A comparison between a set of features for the two clusters.

The fact that this spatial distribution has been derived from small spatial bins indicates the high accuracy of taste as indicator. These maps show that income can be an important factor in determining a communities' taste. To see empirically how our clusters, correspond with demographic factors, we considered racial composition, educational status, and annual household income at the block-group level for the four cities. Block-group level data is the highest spatial resolution available on Census for these demographics. The data was collected from the American Community Service (ACS) website [81]. We defined educational ratio as the ratio of population that have a bachelor degree or higher, in each block-group. The racial composition was defined as the population ratio of Black/A.A., White, and Asian for different block-groups. The income variable is the annual household income in U.S. Dollars. All these demographic factors were estimates provided by the ACS for 2016. We spatially joined the restaurants to the block-groups and conducted *t*-tests to evaluate the extent to which our clustering results compare with these demographic factors.
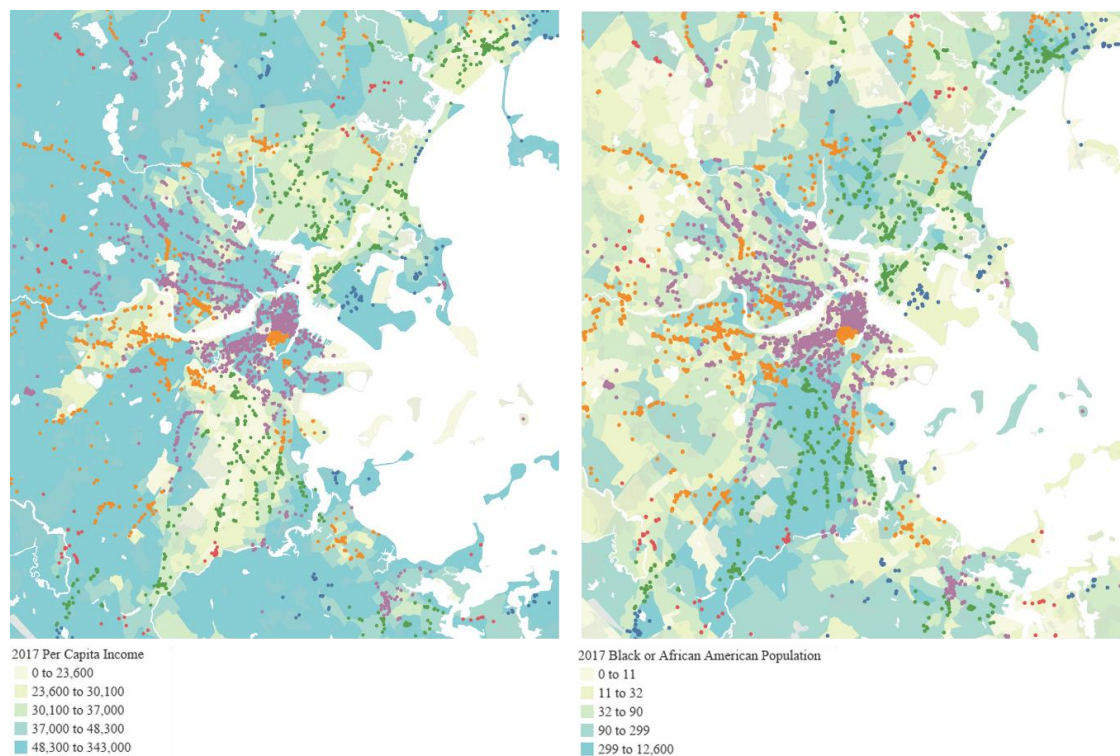
Table 6 provides a summary of the results. As we can see, the two clusters show significantly different demographic features in all four cities. Looking at all four cities together, we can see that education is the most different demographic factor between the two clusters. Considering the restaurants in all four cities, we can see that education and the Asian population ratio are the most distinctive factors with the highest T-statistics. As we consider each city individually, we can see that the order of importance for different demographic factors differs among different regions. For example, in Boston, the top distinctive factors are education and Asian population ratio whereas in Washington D.C. the Black population ratio and annual household income have the highest T-statistics. It is important to note that all the four cities show clear spatial boundaries separating the two clusters. In other words, this method proves to be capable of identifying spatial segregation patterns that may have different demographic reasons in different regions (e.g., education level and Asian population in Boston, MA versus income and Black/A.A. population ratio in Washington D.C.).

Figure 10 illustrates the clustering result with 5 clusters for Boston, MA. In this case, as well, we can see clear geographic patterns. For example, we can see orange and green points are both clustered together around the low-income areas. By overlaying these clusters on the African American population, we can see that most of the green points are located in areas with high concentration of African American population. On the other hand, many purple points are located at areas with high income and high concentration of African Americans. This issue gets to the heart of Bourdieu's argument [15], that taste as an indicator of social status, is not merely a construct of economic capital, but rather it's derived from symbolic capital, which is in turn, a combination of social, cultural and economic capitals. Accordingly, using taste as an indicator of symbolic capital can shed light on

different aspects of communities' lifestyles which may not be explained similarly with conventional demographic indicator (e.g., income, race) for different geographic and cultural contexts.

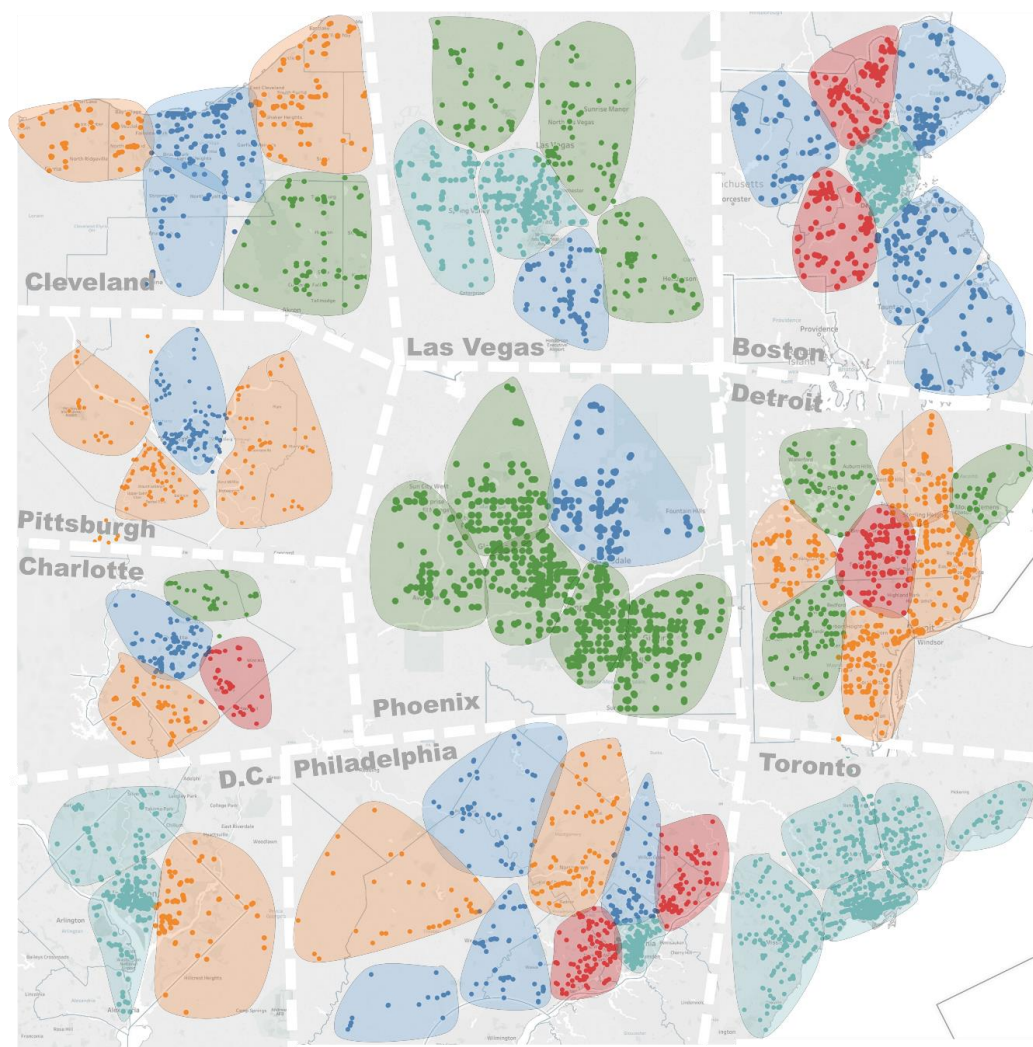**Table 6.** *t*-Test results between the two clusters for demographic variables.

| City | Factor | Mean Value in Cluster 1 | Mean Value in Cluster 2 | T Statistic (Absolute Value) | *p* Value |
|---|---|---|---|---|---|
| **Boston, MA** | | | | | |
| | Educated population ratio | 0.06 | 0.10 | 97.46 | 0.000 |
| | Annual household income (USD) | 66,985.93 | 68,655.57 | 5.47 | 0.000 |
| | Black/A.A. population ratio | 0.41 | 0.40 | 13.49 | 0.000 |
| | White population ratio | 0.53 | 0.56 | 32.97 | 0.000 |
| | Asian population ratio | 0.02 | 0.07 | 73.04 | 0.000 |
| **Detroit, MI** | | | | | |
| | Educated population ratio | 0.04 | 0.07 | 59.39 | 0.000 |
| | Annual household income (USD) | 50,359.52 | 61,600.40 | 40.69 | 0.000 |
| | Black/A.A. population ratio | 0.41 | 0.38 | 21.84 | 0.000 |
| | White population ratio | 0.55 | 0.60 | 35.75 | 0.000 |
| | Asian population ratio | 0.01 | 0.03 | 43.08 | 0.000 |
| **Philadelphia, PA** | | | | | |
| | Educated population ratio | 0.05 | 0.09 | 72.51 | 0.000 |
| | Annual household income (USD) | 55,067.55 | 64,436.73 | 25.42 | 0.000 |
| | Black/A.A. population ratio | 0.39 | 0.35 | 24.14 | 0.000 |
| | White population ratio | 0.55 | 0.62 | 41.79 | 0.000 |
| | Asian population ratio | 0.03 | 0.05 | 33.30 | 0.000 |
| **Washington, D.C.** | | | | | |
| | Educated population ratio | 0.11 | 0.15 | 25.48 | 0.000 |
| | Annual household income (USD) | 53,222.42 | 80,220.32 | 28.74 | 0.000 |
| | Black/A.A. population ratio | 0.36 | 0.22 | 41.42 | 0.000 |
| | White population ratio | 0.55 | 0.68 | 23.94 | 0.000 |
| | Asian population ratio | 0.02 | 0.04 | 15.19 | 0.000 |
| **All four cities combined** | | | | | |
| | Educated population ratio | 0.06 | 0.09 | 134.74 | 0.000 |
| | Annual household income (USD) | 57,322.86 | 66,673.53 | 51.06 | 0.000 |
| | Black/A.A. population ratio | 57,322.86 | 0.37 | 29.46 | 0.000 |
| | White population ratio | 0.40 | 0.60 | 69.42 | 0.000 |
| | Asian population ratio | 0.54 | 0.05 | 91.96 | 0.000 |



2017 Per Capita Income
- 0 to 23,600
- 23,600 to 30,100
- 30,100 to 37,000
- 37,000 to 48,300
- 48,300 to 343,000

2017 Black or African American Population
- 0 to 11
- 11 to 32
- 32 to 90
- 90 to 299
- 299 to 12,600

**Figure 10.** Clustering results with 5 clusters for Boston.

Having set our new indicator, we can now use this indicator to study the socioeconomic interactions between different regions in different cities. We use the large-grained spatial bins that we previously defined for all cities and choose five clusters for simplification purposes (Figure 11). The results are consistent with our understanding of global cities. American cities are comprised of spatially separated cultural groups [4]. We can also see that the distribution of these cultural clusters is consistent with our knowledge of some cities. For example, we know that the racial and economic segregation pattern for Phoenix, Pittsburgh, and Washington D.C. approximately corresponds with our results. In some cases, the clusters do not necessarily match with racial and economic measures of those regions. For example, the north-eastern side of Phoenix is in the same cluster as downtown Cleveland while the two regions are demographically different. The earlier is dominantly white and high-income whereas the latter is a low-income mixed-race region. Another anomaly is Toronto which seems to have all its regions in the same cluster colored in cyan. Clusters shown in cyan signify high-income multicultural areas with a variety of restaurant types and cultural groups. This issue might be due to the fact that Toronto does not suffer from extreme racial and economic segregation as is the case for American metropolitan areas [82].



**Figure 11.** Cultural interactions between different cities. Similar colors across cities indicate similar tastes.

## 4. Conclusions

In this study, by using Natural Language Processing techniques, we extracted several features from the Yelp reviews corpus which characterize the practiced taste in a given region in 10 major American cities. By clustering these features and comparing them with block-group level demographic and economic variables derived from the U.S. Census in four of these cities, we showed that our definition of taste can be used as an indicator for studying the socioeconomic structure for the four cities where we did not have the user IDs. We found a clear alignment between areas of low-income and high-income and our clusters for all the four cities (Figure 8). We also showed statistically that the two clusters are significantly different based on different demographic factors representing income, education and racial composition. We showed that education is the most distinctive factor between the two clusters once we consider all four cities combined. We also showed that the two clusters in different cities, while forming clear spatial boundaries, are different in terms of demographic differences between the two clusters. For example, we found that Education and Asian population ratio are the most distinctive factors in Boston while in case of Washington D.C., Black/A.A population ratio and annual household income are the main distinctive factors.

Once we increased the number of clusters we still observed a geographic pattern (Figure 10) which results from a combination of demographic factors such as race and income. This issue reflects the multifaceted nature of taste as argued by Bourdieu [15]. We showed that this method also works well for more than two clusters, although the performance of this method depends highly on the quality of data and number of reviews. Lastly, we used the selected set of features to study the inter-regional similarities for 10 North American cities. Our results showed that all the nine American cities were comprised of regions that are less similar to one another and more similar to some regions in other cities. This observation is close to our understanding of the global cities as described in the literature [4]. In case of Toronto, all the spatial bins were in the same cluster which might be due to the fact that extremely disadvantaged neighborhoods for different racial groups do not exist compared to the U.S. metropolitan areas [82].

As discussed earlier, we do not expect to see a direct relationship between clusters derived from taste and racial and economic segregation patterns in all cultural and geographic contexts: First, commonly used foods and drink in a White community in one city might be quite popular in the African American communities in another. From a theoretical point of view, the taste index assists us to see cities regardless of their mere economic and racial composition, but rather the symbolic capital of the inhabitants which results from social, economic, and cultural capital, combined. Second, reviews provided by Yelp users in a region might not have necessarily been authored by the residents of that region. It is not surprising to see that a considerable number of reviews in downtown Cleveland, for example, have been authored by visitors who do not reside in that region. This issue can be seen as both a limitation or potential [83]. It is a limitation in a sense that restaurants-as-sensors, may fail to capture the cultural characteristics of the resident population in a neighborhood as these restaurants may target the visitors and not the resident population. On the other hand, it could be a potential since most of the information collected by different agencies such as Census are collected from residents while ignoring the ambient population. This issue has also been discussed by other studies [83] that argue about the mismatch between density of tweets and residents' population. The taste index, therefore, enables us to see the cultural preferences practiced by the ambient population who actually are the clientele of these restaurants. Using ambient population can help urban planners to gain a better understanding of the people who actually use urban spaces and design spaces accordingly [84].

Working with socially sensed data comes with many limitations. First, Yelp reviewers may be a biased sample of the population and therefore, the comments that they provide might not be reflective of the entire population's judgment for a restaurant. Second, our definition of taste was limited to the types of food, drinks and restaurants' ambience. Although this definition may reflect the characteristics of neighborhoods to some extent, additional data on people's lifestyle such as the interior decorations, grocery purchases, and types of movies they watch will provide a more accurate

understanding of different neighborhoods. The extent of these limitations for different geographic contexts may affect the final results, significantly. In case of Phoenix for example, we can see that the final $F_1$ score, according to Table 4 was low (i.e., 0.50608) compared to other cities, which may be due to data bias or similar food tastes between different user groups.

Despite all these limitations, our method uses community-authored comments scraped from the web at no cost with a reasonable spatial and temporal resolution. Given the variety and accessibility of business data [56], the information derived from this method can complement the conventional demographic data of the cities and provide a multifaceted understanding of cities which integrate economic, social and cultural components at once. Using datasets with high temporal and spatial resolution such as Yelp, to better understand the transitional nature of global cities in an ever-changing economical and societal setting at no significant cost.

Future work can focus on improving the methods by using a multitude of crowd-sourced datasets other than Yelp. Concatenating several datasets would minimize the potential biases that may be specific to Yelp. Amazon reviews, Spotify, Instagram, and Flickr are some examples of potential datasets that can be used in combination with Yelp and enrich our understanding of taste in different neighborhoods. The validity of the methods that we used are entirely reliant on the quality and quantity of the reviews used to characterize different regions. Accordingly, using different other geographies to enrich the study sample can also significantly improve the specificity and sensitivity of our methods. Future studies can also examine applying our methods to a variety of different fields of research such as hospitality and tourism, community planning, health and nutrition, marketing, cultural studies, and other fields where characterizing the taste of a region can play a pivotal role.

**Author Contributions:** S.R., S.M. conceived and designed the methodology; S.R., S.M. and X.L. collected the data; S.R. and S.M. cleaned and analyzed the data; S.R. provided the theoretical discussions and literature review; S.R., S.M. and X.L. contributed materials and analysis tools; S.R. and S.M. wrote the paper.

**Conflicts of Interest:** The authors declare no conflict of interest.

**Appendix A.**

**Table A1.** List of 45 Features Used as Seed to Word2Vec Model.

| Category | List of Features |
|---|---|
| Food | chicken, pizza, ketchup, cheese, salad, hot dog, burger, bacon, burrito, mushroom, fish, wings, strawberry |
| Drink | coffee, tea, beer, soda, water, wine, cocktail, alcohol, smoothie |
| Food adjectives | Mexican, Italian, Chinese, sweet, fried, spicy, vegetarian, greasy, homemade, juicy, organic, stuffed, crispy |
| Ambiance | cozy, hipster, trendy, classy, modern, homey, intimate, romantic, upscale, divey |

**Appendix B.**

**Table A2.** List of Features Generated by Aggregation.

| New Feature | List of Combined Features |
|---|---|
| steak_types | meatloaf, Barclay, flank, wagyu, kalbi, tenderloin, striploin, bavette, rib, brisket, mignon, steak, ribeye |
| meat_types | chicken, meat, beef, pork, lamb, veal, duck, turkey, steak |
| sweets_types | yogurt, gelato, pudding, cupcake, biscuit, pie, tiramisu, crepe, custard, tart, sorbet, Nutella, cheesecake, cream, cannoli, muffin, donut, cookie, cake, shake |

**Table A2.** *Cont.*

| New Feature | List of Combined Features |
| --- | --- |
| **fast_food** | pizza, hot fog, sandwich, burger, chips, pepperoni, max, finger, cheeseburger, cheesesteak, calzone, meatball, hoagie, poutine, blt, Rueben, wing |
| **vegie_types** | turnip, lettuce, celery, seaweed, parsley, scallion, eggplant, broccoli, zucchini, kale, cilantro, veggie, ceasar, cabbage, cucumber, basil, vegetable, mushroom, sprout, carrot, asparagus, bean, onion, tomato, coleslaw, avocado, spinach, artichoke |
| **breakfast_types** | bacon, sausage, egg, benedict, scramble, omelet, bagel, pancake, croissant, pretzel, syrup, waffle, roast |
| **fruite_types** | pineapple, peach, strawberry, raspberry, blueberry, coconut, apple, mango, banana, orange |
| **nut_types** | walnut, pecan, peanut, almond |
| **herb_types** | oregano, thyme, fennel, sumac, paprika, garnish, herb, radish, chive, dill, arugula, mint |
| **dressing_types** | ranch, ketchup, mayo, gravy, marinara, sriracha |
| **coffee_types** | espresso, cappuccino, decaf, americano, mocha, latte |
| **soda_types** | Pepsi, Fanta, spirit, coke, soda |
| **softliq_types** | champagne, beer, wine, margarita, sangria, mimosa, cider |
| **hardliq_types** | tequila, whiskey, vodka, martini, bourbon, shot |
| **ethnic_food** | Thai, Chinese, Mexican, Italian, Asian, Indian, Japanese, Vietnamese, Hawaiian, Sicilian, Arabic, Middle Eastern, Korean, Taiwanese, Persian, Greek, Lebanese, Portuguese, Ethiopian, Spanish |
| **latin_types** | salsa, burrito, quesadilla, taco, carnitas, tamale, guacamole, tapa, enchilada, tortilla, fajita, carne, jalapeno, nacho, ceviche, empanada |
| **Italian_types** | pastrami, panini, lasagna, bruschetta, pasta, prosciutto, stromboli, vermicelli, risotto, spaghetti, pesto, chorizo, gnocchi |
| **Asian_types** | fusion, sesame, wonton, spring roll, omakas, sushi, aman, tofu, kimchi, nigiri, sashimi, mushi, noodle, teriyaki |
| **Mideast_types** | shawarma, flatbread, pita, naan, hummus, falafel |
| **pos_ambience** | cozy, homey, classy, trendy, artsy, urbane, posh, swanky, upscale, festive, romantic, eclectic, elegant, chic, stylish |
| **neg_ambience** | casual, divey, kitschy, masculine |
| **style_stypes** | hipster, hippie, bohemian, rustic, modern, minimalistic, contemporary, retro, deco, quaint |
| **material_types** | wooden, hardwood, marble, concrete, mosaic, metal, steel, brick |

## References

1. Musterd, S.; Ostendorf, W. *Urban Segregation and the Welfare State: Inequality and Exclusion in Western Cities*; Routledge: Abingdon, UK, 2013.
2. Sassen, S. *The Global City*; Princeton University Press: Princeton, NJ, USA, 1991.
3. Badcock, B. Restructuring and spatial polarization in cities. *Prog. Hum. Geogr.* **1997**, *21*, 251–262. [CrossRef]
4. Dear, M.; Flusty, S. Postmodern Urbanism. *Ann. Assoc. Am. Geogr.* **1998**, *88*, 50–72. [CrossRef]
5. Wellman, B. The Community Question: The Intimate Networks of East Yorkers. *Am. J. Sociol.* **1979**, *84*, 1201–1231. [CrossRef]
6. Wenger, G.C. A Comparison of Urban with Rural Support Networks: Liverpool and North Wales. *Ageing Soc.* **1995**, *15*, 59–81. [CrossRef]
7. Hampton, K.; Wellman, B. Neighboring in Netville: How the Internet supports community and social capital in a wired suburb. *City Community* **2003**, *2*, 277–311. [CrossRef]
8. Madanipour, A.; Cars, G.; Allen, J. *Social Exclusion in European Cities: Processes, Experiences, and Responses*; Psychology Press: Abingdon, UK, 2000; Volume 23.
9. Lyons, W.E.; Lowery, D. Citizen Responses to Dissatisfaction in Urban Communities: A Partial Test of a General Model. *J. Politics* **1989**, *51*, 841–868. [CrossRef]
10. Bracken, I.; Martin, D. The generation of spatial population distributions from census centroid data. *Environ. Plan. A* **1989**, *21*, 537–543. [CrossRef] [PubMed]

11. Zheng, Y.; Capra, L.; Wolfson, O.; Yang, H. Urban Computing: Concepts, Methodologies, and Applications. *ACM Trans. Intell. Syst. Technol.* **2014**, *5*, 1–55. [CrossRef]

12. Hennion, A. Those things that hold us together: Taste and sociology. *Cult. Sociol.* **2007**. [CrossRef]

13. Kant, I. *Critique of Judgment*; Hackett Publishing Company: Indianapolis, IN, USA, 1787.

14. Cummings, J. The Theory of the Leisure Class. *J. Political Econ.* **1899**. [CrossRef]

15. Bourdieu, P. *Distinction: A Social Critique of the Judgment of Taste*; Routledge: Abingdon, UK, 1984; Volume 1.

16. Rinzivillo, S.; Mainardi, S.; Pezzoni, F.; Coscia, M.; Pedreschi, D.; Giannotti, F. Discovering the Geographical Borders of Human Mobility. *Künstl. Intell.* **2012**, *26*, 253–260. [CrossRef]

17. Ratti, C.; Sobolevsky, S.; Calabrese, F.; Andris, C.; Reades, J.; Martino, M.; Claxton, R.; Strogatz, S.H. Redrawing the map of Great Britain from a network of human interactions. *PLoS ONE* **2010**, *5*. [CrossRef] [PubMed]

18. Yuan, J.; Zheng, Y.; Xie, X. Discovering regions of different functions in a city using human mobility and POIs. In Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Beijing, China, 12–16 August 2012; pp. 186–194. [CrossRef]

19. Yuan, N.; Zhang, F.; Lian, D.; Zheng, K. We know how you live: Exploring the spectrum of urban lifestyles. In Proceedings of the First ACM Conference on Online Social Networks, Boston, MA, USA, 7–8 October 2013; pp. 3–14. [CrossRef]

20. Liu, H. Social network profiles as taste performances. *J. Comput. Commun.* **2007**, *13*, 252–275. [CrossRef]

21. Rahimi, S.; Liu, X.; Andris, C. Hidden style in the city: An analysis of Geolocated Airbnb rental images in Ten Major Cities. In Proceedings of the 2nd ACM SIGSPATIAL Workshop on Smart Cities and Urban Analytics, Burlingame, CA, USA, 31 October–3 November 2016.

22. Cranshaw, J.; Hong, J.I.; Sadeh, N. The Livehoods Project: Utilizing Social Media to Understand the Dynamics of a City. In Proceedings of the 6th International AAAI Conference on Weblogs and Social Media, Dublin, Ireland, 4–7 June 2012; pp. 58–65.

23. Yin, Z.; Cao, L.; Han, J.; Zhai, C.; Huang, T. Geographical topic discovery and comparison. In Proceedings of the 20th International Conference on World Wide Web, Hyderabad, India, 28 March–1 April 2011; pp. 247–256.

24. Wakamiya, S.; Lee, R. Crowd-sourced Urban Life Monitoring: Urban Area Characterization based Crowd Behavioral Patterns from Twitter Categories and Subject Descriptors. In Proceedings of the 6th International Conference on Ubiquitous Information Management and Communication, Kuala Lumpur, Malaysia, 20–22 February 2012. [CrossRef]

25. Li, Q.; Zheng, Y.; Xie, X.; Chen, Y.; Liu, W.; Ma, W.Y. Mining user similarity based on location history. In Proceedings of the 16th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems, Irvine, CA, USA, 5–7 November 2008.

26. Hung, C.-C.; Hung, C.; Chang, C.-W.; Chang, C.; Peng, W.; Peng, W.-C. Mining Trajectory Profiles for Discovering User Communities. In Proceedings of the 2009 International Workshop on Location Based Social Networks, Seattle, WA, USA, 3 November 2009. [CrossRef]

27. Zheng, Y.; Xie, X.; Ma, W. GeoLife: A Collaborative Social Networking Service among User, Location and Trajectory. *IEEE Data Eng. Bull.* **2010**, *33*, 32–40.

28. Xiao, X.; Zheng, Y.; Luo, Q.; Xie, X. Finding similar users using category-based location history. In Proceedings of the 18th SIGSPATIAL International Conference on Advances in Geographic Information Systems, San Jose, CA, USA, 2–5 November2010. [CrossRef]

29. He, J.; Chu, W.W. *A Social Network-Based Recommender System (SNRS)*; Springer: Boston, MA, USA, 2010; Volume 12, ISBN 9781441962867.

30. Bonhard, P.; Harries, C.; McCarthy, J.; Sasse, M. Accounting for taste: Using profile similarity to improve recommender systems. In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, Montréal, QC, Canada, 22–27 April 2006; pp. 1057–1066. [CrossRef]

31. Knulst, W.; Kraaykamp, G. Trends in leisure reading: Forty years of research on reading in The Netherlands. *Poetics* **1998**, *26*, 21–41. [CrossRef]

32. Van Eijck, K. Social Differentiation in Musical Taste Patterns. *Soc. Forces* **2001**, *79*, 1163–1185. [CrossRef]

33. Lewis, K.; Kaufman, J.; Gonzalez, M.; Wimmer, A.; Christakis, N. Tastes, ties, and time: A new social network dataset using Facebook.com. *Soc. Netw.* **2008**, *30*, 330–342. [CrossRef]

34. Zukin, S. Urban Lifestyles: Diversity and Standardisation in Spaces of Consumption. *Urban Stud.* **1998**, *35*, 825–839. [CrossRef]

35. Harvey, D. *The Condition of Postmodernity*; Blackwell: Oxford, UK, 1991; Volume 67, ISBN 0-631-16292-5.

36. Lash, S.; Urry, J. *Economies of Signs and Space*; Sage: Newcastle upon Tyne, UK, 1994; ISBN 0803984723.

37. Zukin, S. *The Cultures of Cities*; Wiley-Blackwell: Hoboken, NJ, USA, 1996; Volume 25, ISBN 1557864373.

38. Mullins, P.; Natalier, K.; Smith, P.; Smeaton, B. Cities and Consumption Spaces. *Urban Aff. Rev.* **1999**, *35*, 44–71. [CrossRef]

39. Bocock, R. *Consumption*; Routledge: London, UK; New York, NY, USA, 1993.

40. Featherstone, M. *Consumer Culture and Postmodernism*; Sage: Newcastle upon Tyne, UK, 1991; ISBN 9781412910132.

41. Clarke, S. *The World of Consumption*; Routledge: London, UK; New York, NY, USA, 1994.

42. Miller, D. Consumption Studies as the Transformation of Anthropology. In *Acknowledging Consumption*; Routledge: New York, NY, USA, 1995; pp. 272–301. ISBN 0415106893.

43. Neal, Z.P. Culinary deserts, gastronomic oases: A classification of US cities. *Urban Stud.* **2006**. [CrossRef]

44. Schlosser, E. *Fast Food Nation: The Dark Side of the All-American Meal*; Houghton Mifflin Harcourt: Boston, MA, USA, 2012.

45. Yelp. *Yelp Information*; Yelp: San Francisco, CA, USA, 2017.

46. Luca, M.; Zervas, G. Fake It Till You Make It: Reputation, Competition, and Yelp Review Fraud. *Manag. Sci.* **2016**. [CrossRef]

47. Zukin, S.; Lindeman, S.; Hurson, L. The omnivore's neighborhood? Online restaurant reviews, race, and gentrification. *J. Consum. Cult.* **2017**. [CrossRef]

48. Rahimi, S.; Andris, C.; Liu, X. Using Yelp to Find Romance in the City: A Case of Restaurants in Four Cities. In Proceedings of the 3rd ACM SIGSPATIAL Workshop on Smart Cities and Urban Analytics, Redondo Beach, CA, USA, 7–10 November 2017.

49. Hu, B.; Ester, M. Spatial topic modeling in online social media for location recommendation. In Proceedings of the 7th ACM Conference on Recommender Systems, Hong Kong, China, 12–16 October 2013.

50. Harrison, C.; Jorder, M.; Stern, H.; Stavinsky, F.; Reddy, V.; Hanson, H.; Waechter, H.; Lowe, L.; Gravano, L.; Balter, S. Using Online Reviews by Restaurant Patrons to Identify Unreported Cases of Foodborne Illness—New York City, 2012–2013. *Morb. Mortal. Wkly. Rep.* **2014**, *63*, 441–445.

51. Griffis, H.M.; Kilaru, A.S.; Werner, R.M.; Asch, D.A.; Hershey, J.C.; Hill, S.; Ha, Y.P.; Sellers, A.; Mahoney, K.; Merchant, R.M. Use of social media across US hospitals: Descriptive analysis of adoption and utilization. *J. Med. Internet Res.* **2014**. [CrossRef] [PubMed]

52. Nsoesie, E.O.; Kluberg, S.A.; Brownstein, J.S. Online reports of foodborne illness capture foods implicated in official foodborne outbreak reports. *Prev. Med.* **2014**. [CrossRef] [PubMed]

53. Xiang, Z.; Du, Q.; Ma, Y.; Fan, W. A comparative analysis of major online review platforms: Implications for social media analytics in hospitality and tourism. *Tour. Manag.* **2017**. [CrossRef]

54. Tang, D.; Qin, B.; Liu, T. Document Modeling with Gated Recurrent Neural Network for Sentiment Classification. In Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, Lisbon, Portugal, 17–21 September 2015.

55. Salinca, A. Business Reviews Classification Using Sentiment Analysis. In Proceedings of the 17th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing, Timisoara, Romania, 24–27 September 2016.

56. Arribas-Bel, D. Accidental, open and everywhere: Emerging data sources for the understanding of cities. *Appl. Geogr.* **2014**, *49*, 45–53. [CrossRef]

57. Yelp Dataset Challenge Yelp Dataset Challenge. 2017. Available online: https://www.yelp.com/dataset_challenge (accessed on 15 February 2017).

58. Danilak, M. Langdetect 1.0.7: Language Detection Library Ported from Google's Language-Detection. Available online: https://github.com/Mimino666/langdetect (accessed on 15 January 2018).

59. Sajnani, H.; Saini, V. Classifying Yelp Reviews into Relevant Categories. Available online: http://www.ics.uci.edu/vpsaini/ (accessed on 15 January 2018).

60. Bird, S.; Klein, E.; Loper, E. *Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit*; O'Reilly Media, Inc.: Newton, MA, USA, 2009.

61. Mikolov, T.; Sutskever, I.; Chen, K.; Corrado, G.S.; Dean, J. Advances in neural information processing systems. In Proceedings of the Neural Information Processing Systems Conference, Lake Tahoe, NV, USA, 5–10 December 2013; pp. 111–3119.

62. Mikolov, T.; Chen, K.; Corrado, G.; Dean, J. Efficient estimation of word representations in vector space. *arXiv*, 2013; arXiv:1301.3781.

63. Yu, M.; Dredze, M. Improving Lexical Embeddings with Semantic Knowledge. In Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, Baltimore, MD, USA, 22–27 June 2014. [CrossRef]

64. Von Luxburg, U. A tutorial on spectral clustering. *Stat. Comput.* **2007**, *17*, 395–416. [CrossRef]

65. Hartigan, J.A.; Wong, M.A. Algorithm AS 136: A K-Means Clustering Algorithm. *J. R. Stat. Soc. C* **1979**, *28*, 100–108. [CrossRef]

66. Belkin, M.; Niyogi, P. Towards a theoretical foundation for Laplacian-based manifold methods. *J. Comput. Syst. Sci.* **2008**, *74*, 1289–1308. [CrossRef]

67. Li, Y.; Chen, C.-Y.; Wasserman, W.W. Deep feature selection: Theory and application to identify enhancers and promoters. *J. Comput. Biol.* **2016**, *23*, 322–336. [CrossRef] [PubMed]

68. Zou, H.; Hastie, T. Regularization and variable selection via the elastic net. *J. R. Stat. Soc. Ser. B* **2005**, *67*, 301–320. [CrossRef]

69. Powers, D.M.W. Evaluation: From Precision, Recall and F-Measure To Roc, Informedness, Markedness & Correlation. *J. Mach. Learn. Technol.* **2011**, *2*, 37–63.

70. Koren, Y.; Bell, R.; Volinsky, C. Matrix Factorization Techniques for Recommender Systems. *Computer* **2009**, *42*, 42–49. [CrossRef]

71. Sarwar, B.; Karypis, G.; Konstan, J.; Riedl, J. Incremental Singular Value Decomposition Algorithms for Highly Scalable Recommender Systems. In Proceedings of the Fifth International Conference on Computer and Information Science, Seoul, Korea, 28–29 November 2002; pp. 27–28.

72. Harvey, D. Pattern, Process, and the Scale Problem in Geographical Research. *Trans. Inst. Br. Geogr.* **1968**, *45*, 71–78. [CrossRef]

73. Kwan, M.-P.; Weber, J. Scale and accessibility: Implications for the analysis of land use-travel interaction. *Appl. Geogr.* **2008**, *28*, 110–123. [CrossRef]

74. Jelinski, D.E.; Wu, J. The modifiable areal unit problem and implications for landscape ecology. *Landsc. Ecol.* **1996**, *11*, 129–140. [CrossRef]

75. Kwan, M.-P. How GIS can help address the uncertain geographic context problem in social science research. *Ann. GIS* **2012**, *18*, 1–11. [CrossRef]

76. Peterson, R.; Krivo, L.; Harris, M. Disadvantage and neighborhood violent crime: Do local institutions matter? *J. Res. Crime Delinq.* **2000**, *37*, 31–63. [CrossRef]

77. Bellair, P.E. Informal surveillance and street crime: A complex relationship. *Criminology* **2000**, *38*, 137–170. [CrossRef]

78. Rountree, P.W.; Warner, B.D. Social ties and crime: Is the relationship gendered? *Criminology* **1999**, *37*, 789–814. [CrossRef]

79. Scribner, R.; Cohen, D.A.; Farley, T.A. A Geographic Relation Between Alcohol Availability. *Sex. Transm. Dis.* **1998**, *25*, 544–548. [CrossRef] [PubMed]

80. Rousseeuw, P.J. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *J. Comput. Appl. Math.* **1987**, *20*, 53–65. [CrossRef]

81. U.S. Census Bureau. American Community Survey 5-Year Estimates. 2016. Available online: https://factfinder.census.gov/faces/nav/jsf/pages/index.xhtml (accessed on 15 January 2018).

82. Fong, E.; Gulia, M. Differences in neighborhood qualities among racial and ethnic groups in Canada. *Sociol. Inq.* **1999**, *69*, 575–598. [CrossRef]

83.  Jiang, B.; Ma, D.; Yin, J.; Sandberg, M. Spatial distribution of city tweets and their densities. *Geogr. Anal.* **2016**, *48*, 337–351. [CrossRef]
84.  Jacobs, J. *The Death and Life of Great American Cities*; Alexander, C., Ishikawa, S., Silverstein, M., Eds.; Macat Library: New York, NY, USA, 1961; Volume 71.