

RHEIN-WAAL UNIVERSITY OF APPLIED SCIENCES

Machine Learning Based Restaurant Location Analysis in Hamburg

by

Vincent Meyer zu Wickern, Sachin Kumar, Viet-Hung Vu

in the

Faculty of Communication and Environment
Geoinformatics, WS 2018/19

February 2019

Contribution of authors to the project

This paper was implemented in collaboration of all three authors. This implementation is available on Github under the link https://github.com/Vincentxyz/Hamburg_Food_Geo. The documentation of the paper was split among the three authors with the following chapters written by the respective authors:

Vincent Meyer zu Wickern:

- Abstract
- chapter 1: Introduction
- section 4.2: Regression of performance measures
- chapter 6: Data analysis
- section 7.2: Discussion of results
- chapter 8: Conclusion

Sachin Kumar

- chapter 2: Factors of restaurant success
- section 4.1: Random forest regression
- section 7.1: Exploratory data analysis

Viet-Hung Vu

- chapter 3: Data sources
- chapter 5: Data extraction

We,

- Vincent Meyer zu Wickern
- Sachin Kumar
- Viet-Hung Vu

, declare that this paper and the work presented in it are our own. We confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this paper has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where we have consulted the published work of others, this is always clearly attributed.
- Where we have quoted from the work of others, the source is always given. With the exception of such quotations, this paper is entirely our own work.
- We have acknowledged all main sources of help.
- Where the paper is based on work done by ourselves jointly with others, we have made clear exactly what was done by others and what we have contributed ourselves.

Signed:

Signed:

Signed:

Date:

Abstract

Faculty of Communication and Environment

Geoinformatics, WS 2018/19

by Vincent Meyer zu Wickern, Sachin Kumar, Viet-Hung Vu

The location of restaurants is an important factor for their success and their prosperity of business. Although many location dependent factors influencing the success of a business have been explained in literature, there is no automated method of using these factors together with available data to find a beneficial location for a new restaurant. Therefore, this paper suggests a method to use machine learning methods for the recommendation of success contributing locations to potential restaurant owners. As Hamburg transparently publishes administration owned data to the public and is a location to many restaurants, the location recommendation method is based on data from the city of Hamburg and provides recommendations to locations in Hamburg. The success of a business in form of profit or revenue is most frequently not open to the public, so that the ratings and reviews in restaurant review portals are used as an indicator for the success of a restaurant in this paper.

The paper finds that the location of restaurants in Hamburg contributes to their success and the prediction of success given a location returns more accurate values than taking the mean of values as a guess for the success of a restaurant. With this indication of success contribution using locations, recommendations for potential restaurant locations are given using recommendation maps that depict each of Hamburg's locations' contribution to success. The five most beneficial locations in Hamburg are presented and single features leading to the predictions of success are discussed. Weaknesses of the automated restaurant location search using machine learning are found, especially the complexity of the models and the associated difficulty of understanding are posed as challenges that would benefit from further research in the future. Nevertheless, the contribution of the locations to the success of a potential restaurant provide a proven value to the location analysis and may be included into location searches of potential owners.

Contents

Abstract	iii
List of Figures	vi
List of Tables	vii
Abbreviations	viii
1 Introduction	1
2 Factors of restaurant success	3
2.1 Restaurant viability factors	3
2.2 Elements of a restaurant location analysis	3
3 Data sources	7
3.1 Transparency Portal Hamburg	7
3.2 Online restaurant portals	8
3.2.1 Tripadvisor	8
3.2.2 Google Places	8
3.2.3 Facebook	8
3.2.4 Yelp	9
4 Analysis methods	11
4.1 Random Forest Regression	11
4.1.1 General description	11
4.1.2 Random Forest Regression Algorithm	12
4.1.3 Estimation of the random forest error rate	12
4.1.4 Bagging and boosting	12
4.2 Regression performance measures	13
5 Data extraction	15
5.1 Hamburg district map	15
5.2 Yelp restaurant data extraction	15
5.2.1 Overview	15

5.2.2	Raw data extraction	16
5.2.2.1	Retrieving data within a certain area from the Yelp server	16
5.2.2.2	Determining latitudes and longitudes for restaurant searches	17
5.2.2.3	Extracting data of all restaurants in Hamburg	18
5.2.3	Data wrangling and storing	19
5.2.4	Restaurant success calculation	20
5.3	City district profiles Hamburg	22
5.4	Proximity to water	22
5.5	Restaurant density	24
5.6	Count of criminal cases	24
6	Data analysis	25
6.1	Feature selection and preparation	25
6.2	Regressor training	26
6.3	Creation of a recommendation map	29
6.4	Locations for restaurant recommendations for restaurants in Hamburg	32
7	Discussion	34
7.1	Exploratory data analysis	34
7.2	Discussion of results	38
8	Conclusion	40
	Bibliography	42

List of Figures

2.1	Impact of various factors on restaurant viability [1]	4
3.1	Statistics on Yelp review according to Yelp Inc. [2]	10
5.1	Administrative boundaries of Hamburg in QGIS	16
5.2	Yelp data extraction search grid	19
5.3	Un-wrangled data from JSON file	21
5.4	Wrangled data stored in namedtuple	21
5.5	Distribution of Restaurant Success Values	22
5.6	WMS Layer “Geologische Karte Hamburg”	23
6.1	Boxplot diagram of the R^2 values for training data predictions	27
6.2	Boxplot diagram of the Adjusted R^2 values for training data predictions	27
6.3	Boxplot diagram of the R^2 values for test data predictions of all regressors	28
6.4	Boxplot diagram of the R^2 values for test data predictions excluding linear regression	28
6.5	Boxplot diagram of the Adjusted R^2 values for test data predictions excluding linear regression	29
6.6	Recommendation grid centroids	30
6.7	Recommendation grid centroids filtered with the Hamburg districts	30
6.8	Greek restaurant success contribution map	32
6.9	Top 5 locations for restaurant success according to random forest predictions	33
7.1	Feature importance for random forest predictions	34
7.2	Boxplot diagram of restaurant price levels with success	35
7.3	Restaurant distance to water and success	35
7.4	Condominium Price and success	36
7.5	Income and success	36
7.6	Restaurant density and success	37
7.7	Population density and success	37
7.8	One-person households and success	38

List of Tables

6.1	Restaurant category frequencies in Hamburg	31
6.2	Top 5 locations for restaurant success according to random forest predictions (ETRS89)	32

Abbreviations

ALKIS© Amtliches Ligenschaftskatasterinformationssystem (Authoritative Real Estate Cadastre Information System)

ASCII American Standard Code for Information Interchange

API Application Programming Interface

CPT Central Place Theory

CRS Coordinate Reference system

CSV comma-separated values

ETRS89 European Terrestrial Reference System 1989

JSON JavaScript Object Notation

LR Linear Regression

MetaVer© MetadatenVerbund

OOB out-of-bag

RBF Radial basis function

RFR Random Forest Regression

SSE sum of squared errors of the regression model

SST sum of squared errors of the baseline model

SVR Support Vector Regression

WMS Web Map Service

Chapter 1

Introduction

The success of restaurants is determined by a variety of factors, both internal features, e.g. the management of the restaurant and the quality of food, as well as external features that are majorly determined by the location of a restaurant [1]. The extent, to which the location of a restaurant adds into its success has not been numbered, yet, and while many location success factors have been determined in literature, no model uses these factors for an automated location search for new restaurants. This paper strives to find a method to measure the contribution of the location in a city to the success of restaurants therein and implement an automated location search model.

A method to infer one variable from multiple other given attributes, is machine learning [3]. The prediction of continuous variables using a supervised learning is called regression, which is applied in this paper.

The administration of important location dependent data is regularly performed on a regional level and by the respective authorities of a region. As the German metropolitan city of Hamburg is home to many restaurants and because the city of Hamburg has a transparency law in place that obliges it to publish information of the city [4], the location success analysis model is based on the city of Hamburg.

Success is typically measured by financial figures, such as profit or revenue, however, these figures are not available to the public for small businesses, however, restaurants are small businesses in the majority of cases. Therefore, another success metric is used in this paper, which is visible to the public for all restaurants: reviews of online restaurant review portals. In these restaurant review portals, users can publish their opinion about restaurants and quantify their evaluation comparably to the opinions of other users [5]. While these reviews influence users in their decision to visit and find restaurants, the rating and number of reviews furthermore are viewed as a metric of the success of restaurants to both the restaurant owner and the consumers. Hence, the success metric of restaurants in this paper combines the number and rating of online reviews in online review portals.

To form a model for the prediction of restaurant success depending on a restaurant location, the factors of restaurant success are researched and described chapter 2 on page 3. In the next step, to gain information about the numbers of these factors for restaurants and locations in Hamburg, chapter 3 concentrates on the relevant data sources, which are available to the public. In chapter 4 on page 11, the analysis methods for the prediction of restaurant success are elaborated, mainly focusing on Random Forest Regression as a machine learning model to evaluate restaurant success based on feature sets and how the results thereof can be evaluated. Consequently, in chapter 5 on page 15, it is explained how the the data from multiple sources was extracted to be afterwards used for a prediction model. In the data analysis in chapter 6 on page 25, different regressors are trained, evaluated and a recommendation map for Hamburg is created. Chapter 7 on 34 discusses the predictions and results of this paper and provides an exploratory data analysis with a discussion of variable influences to restaurant success. In the conclusion in chapter 8 on page 40, the paper is summed up and statements on the results of the paper are presented.

Chapter 2

Factors of restaurant success

2.1 Restaurant viability factors

Success or failure of a restaurant depends on various factors. According to Parsa et. al [1], factors which impact the viability of restaurant are broadly classified into four categories. These categories are internal environment, external environment, family life cycle and organizational life cycle. The different factors under these categories are clearly classified by the model developed by Parsa et. al [1], which can be seen in figure 2.1.

For a successful restaurant, it is essential to secure a suitable location, maintain service levels and food quality, control the cost of labor, food and beverages, well defined goals, finance and marketing management, and strong leadership [1]. The restaurant may fail or succeed due to any of these factors. A restaurant location has significant impact on the success or failure of a restaurant , for example opening a nightclub restaurant near a police station thinking of safety may be a wrong selection of location as many customers will not feel comfortable due to police scrutiny. This project aims at selecting a beneficial location for a new restaurant while considering location dependent features which may impact restaurant success.

2.2 Elements of a restaurant location analysis

Location plays a critical role in the success or failure of a restaurant business [1] [6] [7]. Not every location may be a suitable location for every kind of restaurant. There are various factors such as demographic values of the restaurant neighborhood, accessibility, visibility and others, which need to be taken into account when selecting the location for opening a new restaurant. To find the best location for a new restaurant, location analysis technique can be used. Some of the elements of restaurant analysis technique are [8] [9]:

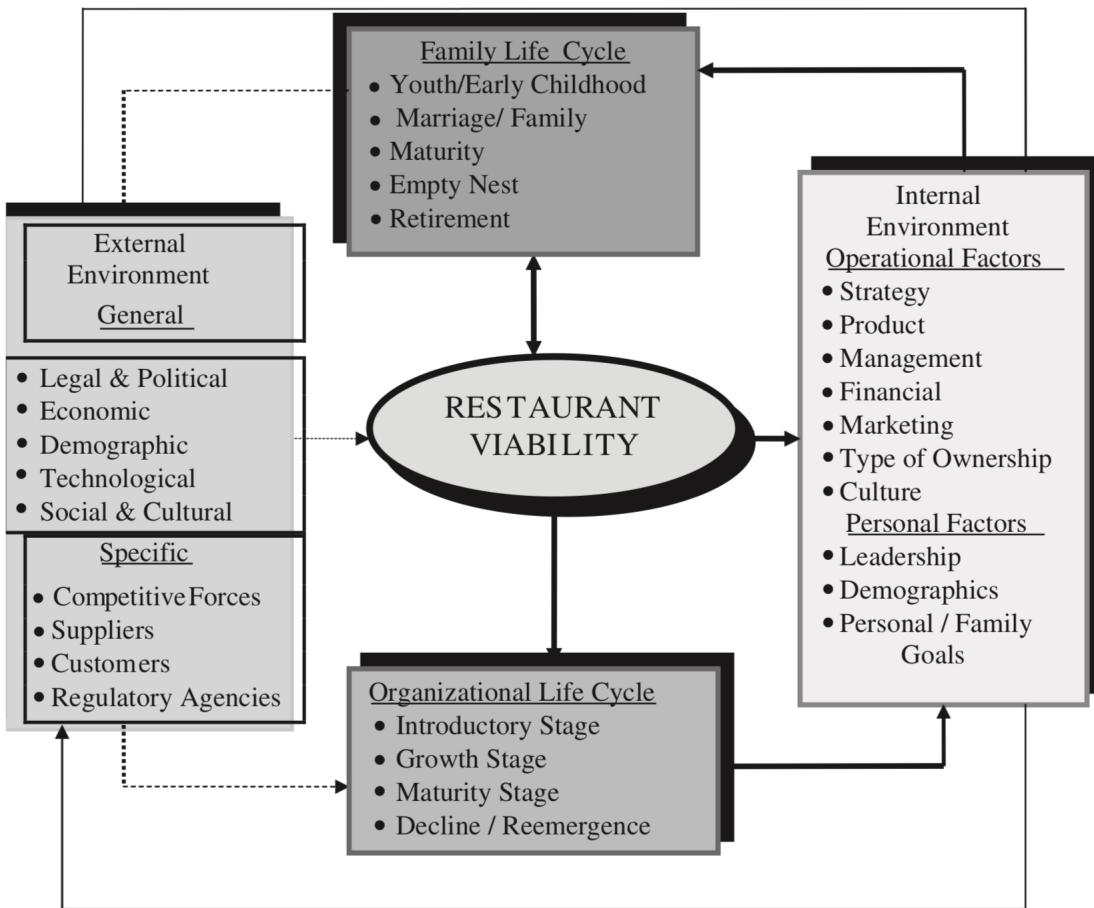


FIGURE 2.1: Impact of various factors on restaurant viability [1]

1. Demographics

Demographics information of the neighboring population such as age, gender, income, religion, relationship status, environment, and ethnicity have an important effect for potential restaurant owners on choosing a restaurant category like “fast food”, “casual dining”, “fine dining” or “bar & bistro” [8]. For example, fast food restaurants are most often favored by people with demographic criteria such as an age range between 15 and 35, with a low income, pedestrians and a fast food consumption is usually unplanned and connected to other nearby events like shopping. Fine dining restaurants, in contrast, are favored by people older than 35, couples, high income people, as pre-planned activities and usually mostly by vehicle traffic.

2. Psychographics

Psychographics information such as personality types and personal preferences of consumers support the decision of cuisine types [8]. For example, family-oriented and traditions-oriented consumers usually prefer Italian restaurants, whereas health-conscious consumer regularly choose organic or vegan grill restaurant.

3. Population

Information about the population and the population density in a specific area are necessary to launch any new business or product as defined by the Central Place Theory (CPT) [10]. A location must fulfill the criteria of CPT like range and threshold to be viable location for a new business launch. Range is referred to as the maximum distance that consumers travel for a desired meal and threshold is the minimum population required around the location to start a business or product (restaurant).

4. Customer Activity: Foot & Vehicle Traffic

For a location to be successful, it should have high levels of activities of potential customers around it, like in neighborhood of downtown and tourist places [8]. Foot traffic increases the flow of potential customer walking by a restaurant, whereas vehicle traffic enhances the information to customers who would potentially drive to the restaurant.

5. Competitor Analysis

Analyzing the competitors that are already active in an area can be beneficial in multiple ways [8]. If a location has too many restaurants in the neighborhood, on the one side this means that the regional market is healthy and supportive in that area and the chances of success are high in that area. However, on the other side, it may be difficult to enter the market and attract customers due to the high competition, so sometimes it is better to find a location with less competition, but in which demand still exists.

6. Labor Cost & Minimum Wage

Restaurant environment effects such as the cost of labor, minimum wage and the availability of potential employees vary with different locations [9]. These factors influence the profit and success of businesses to a large extent. Therefore, these pieces of information must be collected and taken into consideration before choosing the location for restaurant.

7. Accessibility and Visibility

For a location to be a potential spot for the opening of a restaurant, it should be easily accessible by foot traffic and car traffic with accessibility to nearby parking area [9]. Moreover, it should not be in a place where the traffic is too high and the restaurant would not be visible by drivers or pedestrians.

8. Proximity to suppliers

Proximity to suppliers is critical, as shipping cost and delivery cost of supplies might be higher than the original net procurement costs of goods, if it not chosen wisely [9]. The optimal choice of a location should minimize these expenses for the business to be successful.

9. Crime Rates

A potential restaurant location should only be finalized after researching in detail about the crime rate and the type of crime in the potential area [9]. Consumers generally do not prefer eating out in areas with high crime rates and it may be a major factor leading to failure of business.

10. Future Growth

Future plans and aspirations should be taken into consideration in the choice of restaurant locations [9]. In the case that the restaurant prospers and gains popularity, it may require an increase in the customer and employee capacity. Moreover, there should be sufficient parking spaces and storage options.

11. Health regulation and zoning

Different parts of cities have different regulations related to health and zones only allowed to certain types of businesses [9]. These regulations should be reviewed and researched properly before finalizing the location.

Chapter 3

Data sources

3.1 Transparency Portal Hamburg

As the first German federal state, Hamburg enacted a transparency law on October 6, 2012 [4]. Opposed to a right to request information, which all citizens had until this date, a new duty to inform the public was laid upon the state's administration offices. All information that would fall under this law, now had to be published in a freely available standard format on a centralized storage of information. The single pieces of information, which would fall under the law, varied highly in precision and the comprehensive term of “geodata” was requested opposed to precise datasets of geodata. A legal interpretation was worked out for all requested points and a plan for the release of geodata was designed consisting of the basic data for measurement administration and the technical geodata for special administration offices. The transparency law granted a period of two years for the technical implementation.

In October 2014, the “Transparency Portal” (<http://transparenz.hamburg.de/>) as the major component of the implementation of the transparency law was released [4]. With this portal, the Hamburg citizens have a multitude of data and documents available that was prior only available to Hamburg's administration. One important focus was the release of geodata that was even before the law in preparation for an “Open GeoData” model. In this “Open GeoData” model, geodata was split into two groups of data sets, one group able to be extracted with little effort, but free to the public and expected with a high use, and another group with expected high demand and high revenue on the sale of this data. For this second group of datasets, more effort with new measurements had to be arranged. With the transparency law in place, all datasets were merged into the Transparency Portal and yielded a much higher download count than the count of dataset sales before the portal was active. The Transparency Portal uses a standardized metadata repository called the MetadatenVerbund (MetaVer©) in collaboration with other German federal states.

3.2 Online restaurant portals

3.2.1 Tripadvisor

tripadvisor is one of the biggest rating portals for travel and travel related businesses, such as restaurants, with cumulated 600 million reviews and opinions until 2017 [11]. This made tripadvisor a potential portal for the analysis of restaurant reviews to extract data from. However, it was found that any scraping, download or copy of the data with automated or manual methods is legally prohibited from tripadvisor, which excluded tripadvisor as a data basis for the analysis of this paper.

3.2.2 Google Places

As Google is the most frequently used search engine on the world and since the research of consumers in the internet before any buying decision increases, the data processing company Google may have a high influence on purchase decisions [12] [13]. It further provides a service called “Places” for local businesses, among them restaurants, to present themselves with an opportunity of customers to leave reviews, ratings and answers on common questions [14]. With 100 million places, 25 million updates per day and one billion active users per month, the platform may provide insightful data to the restaurant landscape in Hamburg. The Places Application Programming Interface (API) provides an interface for developers to receive data about places, thereof restaurants, that may be used for an analysis.

However, when the API was explored more thoroughly, only the first five reviews of any restaurant could be retrieved and the review count was not extractable. Without the review count, there was information, how well a restaurant has been rated, but not how many customers have, in fact, rated it, which is an essential information. Therefore, Google Places was discarded as a potential source of restaurant information for this paper.

3.2.3 Facebook

Social media is important in the influence to customer choices of products, especially when customers prepare travel or plan the stay at a venue like a hotel or a restaurant [15]. In these cases, online reviews are often taken into consideration. Companies in the hospitality sector therefore have to carefully aim at the benefits of social media tools. Facebook, as one social media platform, has a large customer base with 2.32 billion monthly active users in the fourth quarter of 2018[16]. Out of these reasons, Facebook was considered as one portal to analyze restaurant data from. The Facebook Graph API provides access to content on Facebook for developers [17]. This Graph API is designed for apps that can read and write to Facebook and thereby connects a third-party service with Facebook.

Despite Facebook being a potential source of information of publicly visible restaurant data, the access to this information was impeded by an app review process. In May 2018, Facebook launched an enhanced developer app review, which made it necessary for developers to verify one's business, sign a supplemental terms contract and provide means of proving the business, such as utility bills or taxi ID numbers [18]. Since this was not available and since there was no intention of creating a business app on Facebook, Facebook as a source of restaurant information was discarded for this paper.

3.2.4 Yelp

Yelp is a popular and widely used social networking site for reviewing and sharing information about local businesses like restaurants, dentists and mechanics [2]. Yelp was founded in 2004 in San Francisco, California, and currently operates in 32 countries. It had an average of 164 million unique visitors every month (via the Yelp app, mobile web and desktop) and had 177 million reviews in 2018. Moreover, it presents information about events, special offers and provides a platform to connect and discuss among yelpers (registered users on yelp). The company Yelp Inc. sells advertisements to local businesses, but they claim that reviews do not get affected by advertisements, e.g. in form of added, manipulated, deleted comments. Figure 3.1 shows the statistics for the distribution of user ratings, recommendations and businesses reviewed by category for review across all categories as of 31st December 2018 [2].

From Figure 3.1, it can be seen that the restaurants have maximum reviews among all categories and the rating distribution and recommended distribution shows that yelpers use Yelp to share positive as well as negative experiences. The general, different types of information Yelp collects about restaurant businesses are ratings, reviews, price categories (inexpensive, medium, high, ultra-high), neighborhood, parking facilities, type of meal served (breakfast, brunch, lunch, dinner), smoking areas, reservation possibilities, delivery services and adequacy for children or groups, etc. [2].

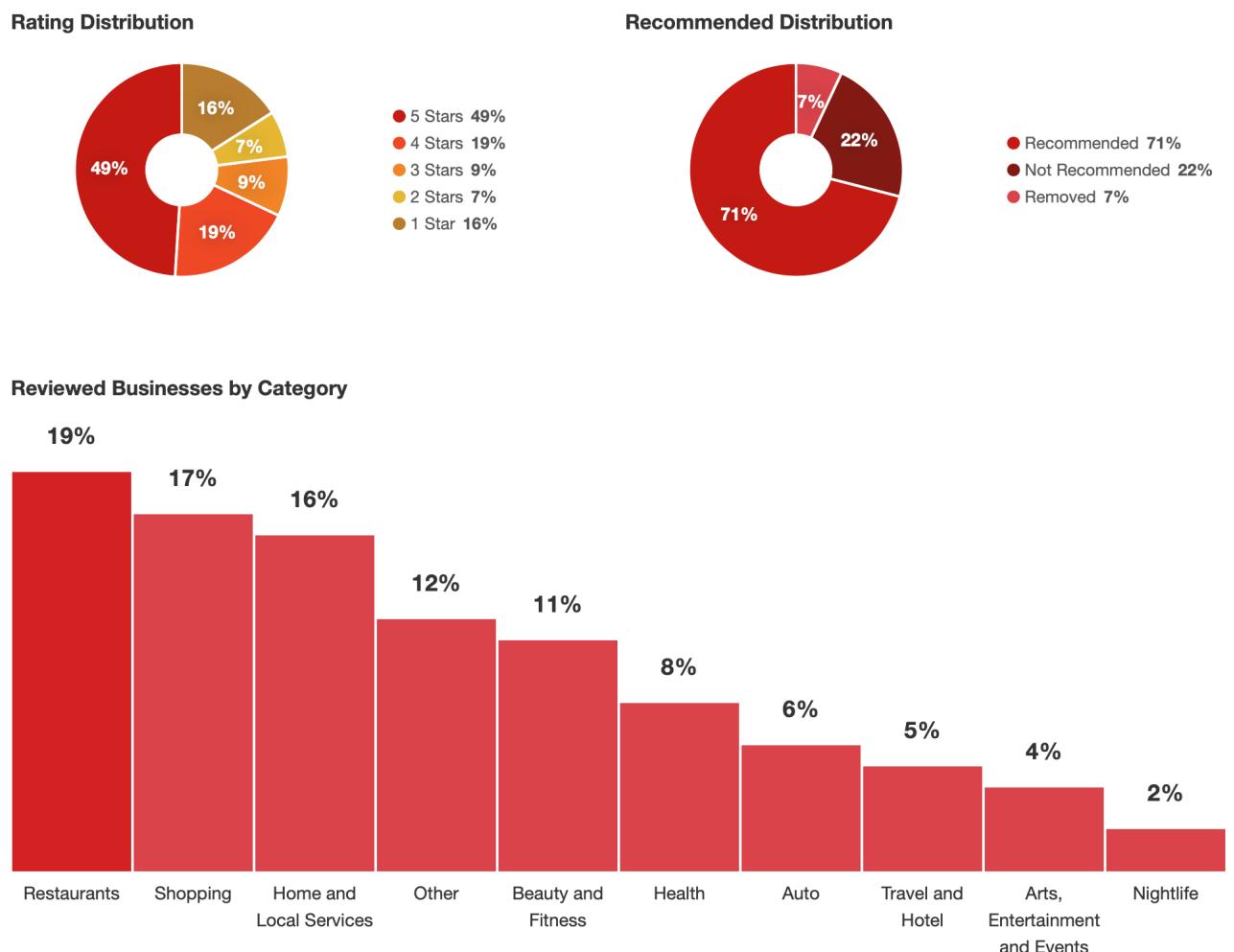


FIGURE 3.1: Statistics on Yelp review according to Yelp Inc. [2]

Chapter 4

Analysis methods

4.1 Random Forest Regression

4.1.1 General description

Random Forest Regression (RFR) is a supervised regression machine learning technique that applies an ensemble learning method [19]. The RFR builds so called “random forests”, which are groups of multiple decision trees per random forest. Predicted regression values of RFRs are an average of prediction values from each of the trees that they comprise. The term “random” stems from the fact that individual trees only use a random subset of the observations and the tree splits within a tree as well is based on a random subset of variables. Building these random subsets of observations and attributes, each creating another tree, is conducted in a defined number of iterations. The use of random feature subsets in the RFR technique makes RFR predictions more robust against overfitting than in other techniques, such as support vector machine and neural networks [20]. Moreover, RFR has been found to deliver more accurate predictions in cases, in which the number of variables is higher than the number of observations. In addition, RFR provides two important pieces of information:

1. Variable importance: RFR is able to measure the contribution of variables to predictions [20]. The measure of variable importance can be used to prune less important variables, since the likelihood of overfitting increases with the number of variables.
2. Proximity measure: Proximity measures the closeness of data points to each other and a proximity matrix can be generated [20]. The proximity matrix may then be used to study structure in the data.

4.1.2 Random Forest Regression Algorithm

The steps in the algorithm of RFR can be described as follows [20]:

1. From a given dataset, a defined number of bootstrap samples for trees is drawn.
2. Unpruned regression decision trees for each of the bootstrap samples are developed. Individual trees are based on random subsets of the observations and the splits within the trees are based on random subsets of variables after the defined number of iterations.
3. The dependent variable is predicted for new feature sets by calculating the average of the tree prediction.

4.1.3 Estimation of the random forest error rate

An estimation of the error rate in RFR is obtained with the help of the training data by conducting the following steps [20]:

1. Available data outside of the bootstrap sample (out-of-bag (OOB) data) is predicted at each iteration of bootstrapping.
2. The error rate is calculated by aggregating the OOB predictions. The result of this is known as the OOB estimate of error rate.

4.1.4 Bagging and boosting

Bagging and Boosting are two techniques to construct an ensemble of individual classifiers (trees) [21]. Individual trees of an ensemble are diverse and a higher accuracy of predictions can be obtained by the aggregation and voting of the results of single trees. In the bagging technique, each training set is constructed by replicating the bootstrap sample of the original training set. The single bootstrap samples in this process are called “examples”. The example creation in the bagging technique is conducted in parallel. In this process, a defined number of examples are taken from a given training set from the entire population of available data. These examples form a new training set. Multiple of these subsets of the training set are created to build multiple trees in the ensemble. Boosting techniques, in contrast, are based on assigning weights to original training set and adjusting these weights after each regressor prediction. Boosting is a sequential process, in which weights are increased for the examples which are predicted incorrectly and decreased for correct predictions. Moreover, the predictions are less affected by noise in bagging than in boosting. Additionally, bagging constructs diverse trees, only if small variations in the training set cause a large variation in the prediction outputs. RFR uses the

bagging technique in combination with techniques to randomize the internal decisions of the learning algorithm to create diverse trees regression trees.

4.2 Regression performance measures

A common performance measure for the quality of fit of regression models is the “coefficient of determination”, which is also called “ R^2 ” [22]. This coefficient of determination uses a baseline model, which is a model that consistently predicts the mean of all observations of the dependent variable. This baseline model is a model that a created regression model can be compared to, to see how much more accurate the own predictions were to a poor model.

Following are the calculation of the sum of squared errors of the regression model (SSE), sum of squared errors of the baseline model (SST) and R^2 [22]:

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$R^2 = 1 - \frac{SSE}{SST}$$

According to this calculation, R^2 always lies between 0 and 1, if the created regression model yields more accurate predictions than the baseline model [22]. The closer R^2 is to 1, the closer are the predictions by the regression model to the actual values. Additionally, the higher R^2 is, the more variations of the dependent variable are explained by the model.

A weak point of R^2 is that additions of more independent variables never lower the value of R^2 , even for variables with little or no information gain [22]. To cope with this problem, another performance measure called “adjusted R^2 ” was introduced, which introduces a negative effect on the measure for the inclusion of ineffective variables. It is calculated as follows [22]:

$$\bar{R}^2 = 1 - (1 - R^2) \frac{n - 1}{n - p - 1}$$

n = no of data points

p = no of independent variables in the model

R^2 therefore only improves (moves closer to 1), if significant variables are added to the model and deteriorates (moves closer to 0), if variables are added that are not valuable to the prediction of the dependent variable [22].

Chapter 5

Data extraction

5.1 Hamburg district map

Since usually cities raise important figures in aggregation per administrative area, in Hamburg being the single city districts, these administrative areas should be imported into QGIS to be able to link single restaurants to an administrative area and hence figures that could be important as dependent variables to predict restaurant success. The borders of the administrative areas are taken from a dataset of the Transparency Portal called “ALKIS Verwaltungsgrenzen Hamburg” [23]. This dataset is available in multiple data formats, is reported to have a 0% data deficit and a precision of 0.1 meters. It is part of the Amtliches Ligenschaftskatasterinformationssystem (Authoritative Real Estate Cadastre Information System) (ALKIS©), a digital combination of the real estate book information and a real estate map [24].

For this analysis, the GML version of the administrative boundaries were downloaded and imported to QGIS. Loading the administrative boundaries into QGIS, the European Terrestrial Reference System 1989 (ETRS89) Coordinate Reference system (CRS) defined by EPSG:25832 was selected, as was defined in the metadata of the dataset. All following elements that are loaded into QGIS are as well projected in this CRS. The imported boundaries can be seen in figure 5.1.

5.2 Yelp restaurant data extraction

5.2.1 Overview

Yelp provides developers with an API called Yelp Fusion API for accessing the Yelp database. The authentication of API calls is done via a private API key which can be obtained by:



FIGURE 5.1: Administrative boundaries of Hamburg in QGIS

- Registering for a Yelp developer account, and
- Creating an application

An extensive guide on the Yelp API authentication is provided on Yelp’s documentation page [25]. For this project, API calls were conducted using Python 3.7.

The extraction of data from the Yelp database consists of:

1. Loading raw data by performing search queries using Yelp’s Python API
2. Wrangling data into intermediate formats and storing wrangled data for further analyses

5.2.2 Raw data extraction

5.2.2.1 Retrieving data within a certain area from the Yelp server

The `yelpapi` Python module [26] must be installed in Python to perform API calls. The obtained private API key must be stored as a Python string:

```
from yelpapi import YelpAPI
api_key = 'pastetheapikeyhere'
yelp_api = YelpAPI(key)
```

Each search query is an https request to the Yelp server. The Python Yelp API exposes this as the following function:

```
yelp_api.search_query()
```

The list of parameters for this function is the same as the https query which is provided by Yelp [27]. For this project, this function was used as:

```
response = yelp_api.search_query(term='restaurants',
    latitude=latitude, longitude=longitude,
    limit=50, offset=0, radius = radius)
```

where latitude, longitude and radius defined the centroid of the search area. Note that for each query, only a maximum of 50 restaurants could be returned (default=20). This means that:

- The total number of restaurants had to be extracted from the returned object of the first function call
- The function had to be called multiple times to get data for all the restaurants within the desired search area

This process was built into the `get_restaurants(latitude, longitude, radius)` function in the Python file `yelp_hamburg_load.py`.

5.2.2.2 Determining latitudes and longitudes for restaurant searches

To receive data of all restaurants in Hamburg, the city had to be divided into tiles of size 2,000 meters by 2,000 meters. As Yelp only responded with less or equal than 1,000 restaurants, the tiles were even quartered, when exactly these 1,000 restaurants were sent back, in the Python scripts “analysis_big_tiles.py” and “get_big_tiles.py”. With smaller tiles, all restaurants could be retrieved. The vertices of the tiles were the latitudes and longitudes used by the function `get_restaurants()` mentioned in section 5.2.2.1. This process consisted of the following steps:

- Pre-determine the maximum latitude, maximum longitude, minimum latitude and minimum longitude of Hamburg. These points were manually picked out from a map of Hamburg in QGIS.
- Create a new empty list of latitudes and a new empty list of longitudes.
- Inner loop: loop from the minimum longitude to the maximum longitude, incrementing by 2,000 meters (length of tile’s side). Append each new longitude into the list of longitudes.
- Outer loop: loop from the minimum latitude to the maximum latitude, incrementing by 2,000 meters (length of tile’s side). Append each latitude into the list of latitudes.

Note that the set of pre-determined maximum/minimum latitudes and longitudes were of ETRS89 (EPSG:25832) CRS, while Yelp API calls require WGS84 (EPSG:4326). Therefore, a helper function [transform_coord()] using the Python module pyproj [28] was written to transform the formats of the latitudes and longitudes so that they could properly be passed into the Yelp API calls.

5.2.2.3 Extracting data of all restaurants in Hamburg

As mentioned in section 5.2.2.2, each Yelp search query requires a latitude, longitude and radius to determine the search location.

The radius of each search was calculated by the radius of all overlapping circles that could form a complete grid as shown in figure 5.2. To receive all restaurants in a grid of tiles having a square of 2,000 meters by 2,000 meters, the radius of the searches needed to at least contain these squares. The longest distance of a square center to any point in the square is given by the distance of the square center to the edges. This distance can be calculated with the law of Pythagoras. If ‘c’ is the distance of a triangle with ‘a’ being the distance of the center of the square to the center of the left side of the square and ‘b’ being the distance from the center of the left side of the square to the top left corner, then

$$a^2 + b^2 = c^2 \quad (5.1)$$

$$c = \sqrt{a^2 + b^2} \quad (5.2)$$

$$\text{radius} = \sqrt{\text{TileSize}/2 + \text{TileSize}/2} \quad (5.3)$$

The corresponding code for the calculation of the radius was given by the following formula:

```
radius = int(math.sqrt(math.pow(TILE_SIZE,2)*2)+1)
```

which was constant during all Yelp API calls.

Each {latitude, longitude} pair to be passed into each Yelp API calls was taken from the lists of latitudes and longitudes created in section 5.2.2.2.

The function get_restaurants(latitude, longitude, radius) was called for each {latitude, longitude} pair until all pair combinations were used.

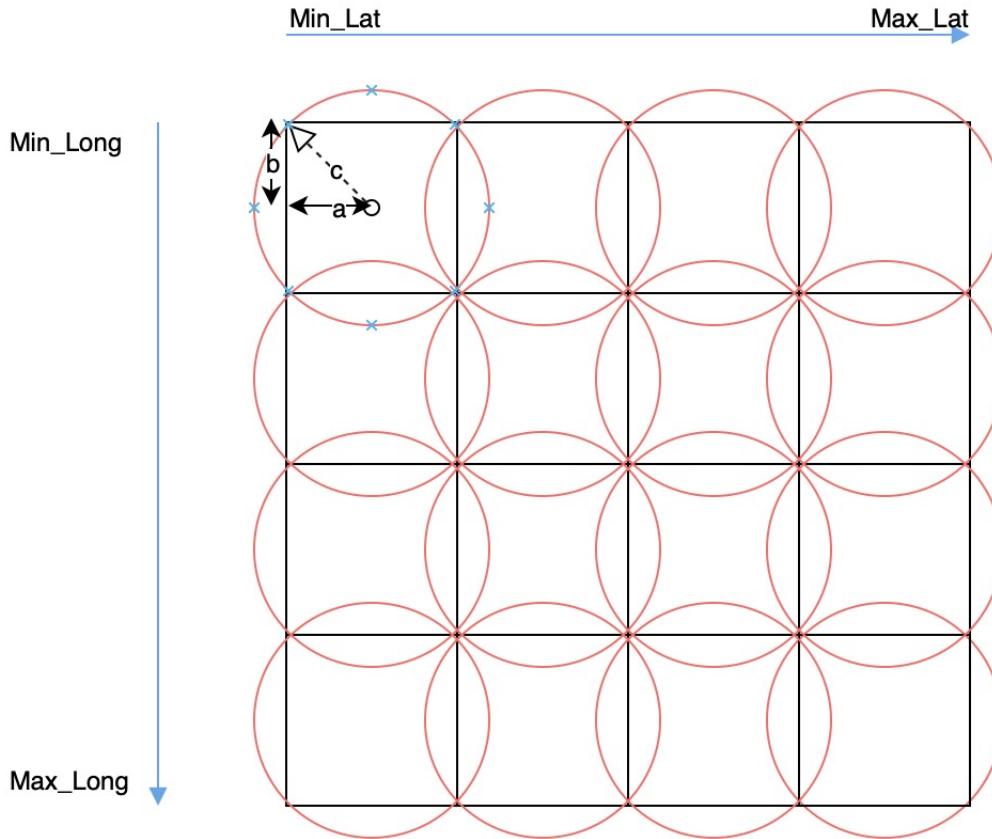


FIGURE 5.2: Yelp data extraction search grid

The result of each call was stored into a JavaScript Object Notation (JSON) file respective to the latitude and longitude used with the help of the Python JSON [29] module. All generated JSON files were subsequently merged into one JSON file for ease of data wrangling.

5.2.3 Data wrangling and storing

As data in the JSON file mentioned in section 5.2.2.3 was split, it had to be loaded back into Python for processing by:

```
with open('merged_extract.json') as json_data:  
    businesses = json.load(json_data)
```

A namedtuple “Restaurant” was created as an object model into which data from the JSON file would be passed into. The data fields for this namedtuple reflects the data retrieved for each restaurant from the search queries. The namedtuple’s data fields include:

- epsg25832_latitude

- epsg25832_longitude
- epsg4326_latitude
- epsg4326_longitude
- price_rating
- review_count
- is_closed
- zip_code
- category_restaurant_id
- category_alias
- category_title
- category_index

Originally the data loaded from the JSON file was a Python list of (nested) dictionaries. Therefore, data was un-wrangled using the dictionaries' keys. For example, the 'epsg4326_latitude' field for the each instance of the namedtuple Restaurant would be:

```
epsg4326_latitude = biz['coordinates']['latitude']
```

where 'biz' is a dictionary object representing data of one restaurant, which is part of the loaded list of dictionaries from the JSON file.

The list of namedtuples representing restaurants could then be passed into a pandas dataframes and subsequently stored into comma-separated values (CSV) files for further analysis.

5.2.4 Restaurant success calculation

The success of a restaurant is usually measured by financial figures, such as revenue, profit or business growth. However, this information is not publicly available for all small businesses, so a restaurant success has to be assumed from other metrics. Other metrics are important to businesses today are online reviews and ratings. They may influence the success of a restaurant and as well may mirror a restaurant's success. In fact, the rating and reviews a restaurant displays on a portal such as Yelp are their display of success to their customers on the Internet, for which reason the success measures in this paper should combine the two figures of review count and average rating.

```

✓  businesses = ([list] <Too big to print. Len: 37543>
  □  00000 = (dict) {'id': 'Apz9UhdlwKP4RGtjw7H6g', 'alias': 'de-krauler-kroog-hamburg-2', 'name': 'De Krauler Kroog', 'image_url': 'https://s3-media4.fl.yelpcdn.com/bphoto/4i5xiZ6T41NoeGwBNDEqw/o.jpg', 'is_closed': False, ...}
    □  'id' (139797216) = [str] 'Apz9UhdlwKP4RGtjw7H6g'
    □  'alias' (155670560) = [str] 'de-krauler-kroog-hamburg-2'
    □  'name' (155670720) = [str] 'De Krauler Kroog'
    □  'image_url' (155666456) = [str] 'https://s3-media4.fl.yelpcdn.com/bphoto/4i5xiZ6T41NoeGwBNDEqw/o.jpg'
    □  'is_closed' (155666496) = [bool] False
    □  'url' (155670912) = [str] 'https://www.yelp.com/biz/de-krauler-kroog-hamburg-2?adjust_creative=hBJOknQJ58S5goQyjmwRw&utm_campaign=yelp_api_v3&utm_medium=api_v3_business_search&utm_source=hBJOkn'...
    □  'review_count' (155666536) = [int] 6
  □  'categories' (155666576) = [list] [[alias: 'german', title: 'German']]
  □  'rating' (155671040) = [float] 5.0
  □  'coordinates' (155666616) = (dict) {'latitude': 53.3996, 'longitude': 10.2142}
    □  'latitude' (155666656) = [float] 53.3996
    □  'longitude' (155666696) = [float] 10.2142
    □  '_len_' = [int] 2
  □  'transactions' (155666736) = [list] []
    □  'price' (155671072) = [str] '€€€'
  □  'location' (155666776) = (dict) {'address1': 'Kraueler Hauptdeich 65', 'address2': None, 'address3': None, 'city': 'Hamburg', 'zip_code': '21037', 'country': 'DE', 'state': 'HH', 'display_address': 'Kraueler Hauptdeich 65, 21037 H...
    □  'phone' (155671424) = [str] '+49 40 7230368'
    □  'display_phone' (155667096) = [str] '+49 40 7230368'
    □  'distance' (155667176) = [float] 1646.4291237763646
    □  '_len_' = [int] 16
  □  00001 = (dict) {'id': 'NH8UdkeMrq2YFSQLMSerQ', 'alias': 'hof-eggers-hamburg-3', 'name': 'Hof Eggers', 'image_url': 'https://s3-media4.fl.yelpcdn.com/bphoto/jfzfR1d-xTAeF5dn9hAg/o.jpg', 'is_closed': False, 'url': 'https://w...
  □  00002 = (dict) {'id': 'OjgVmVhRkm62fICP646uQ', 'alias': 'doner-imbiß-beim-tv-nord-winsen', 'name': 'Doner Imbiß beim TÜV-Nord', 'image_url': 'https://s3-media3.fl.yelpcdn.com/bphoto/qkNvpocZhvRdwTchluSqv/o.jpg'...
  □  00003 = (dict) {'id': 'AwL7KooU6B09c4AQ7eN5w', 'alias': 'soetebier-winsen-3', 'name': 'Soetebier', 'image_url': 'https://s3-media1.fl.yelpcdn.com/bphoto/k2FZAEEfyPKbpAdqaRi-Bsg/o.jpg', 'is_closed': False, 'url': 'https://www.y...
  □  00004 = (dict) {'id': 'dzT3zUy7O8x0yL09p9_jgA', 'alias': 'bäckerei-u-café-zimmer-winsen', 'name': 'Bäckerei u. Café Zimmer', 'image_url': '', 'is_closed': False, 'url': 'https://www.yelp.com/biz/b%C3%A4ckerei-u-caf%C3%94-zimm...
  □  00005 = (dict) {'id': 'GorMSUXMnwrfUMzb1dVgCQ', 'alias': 'burger-king-winsen-aller', 'name': 'Burger King', 'image_url': 'https://s3-media3.fl.yelpcdn.com/bphoto/Vmgv2RbNbLf9eSUUHXIVQ/o.jpg', 'is_closed': False, 'url': 'https://w...
  □  00006 = (dict) {'id': 'CP15eMmPVJUHL363eNV7Ag', 'alias': 'croque-knut-winsen', 'name': 'Croque Knut', 'image_url': 'https://s3-media3.fl.yelpcdn.com/bphoto/hgxA8LaPmwhnicJkk&GSew/o.jpg', 'is_closed': False, 'url': 'https://w...

```

FIGURE 5.3: Un-wrangled data from JSON file

```

✓  businesses = ([list] <Too big to print. Len: 37543>
  □  00000 = (dict) {'id': 'Apz9UhdlwKP4RGtjw7H6g', 'alias': 'de-krauler-kroog-hamburg-2', 'name': 'De Krauler Kroog', 'image_url': 'https://s3-media4.fl.yelpcdn.com/bphoto/4i5xiZ6T41NoeGwBNDEqw/o.jpg', 'is_closed': False, ...}
    □  'id' (139797216) = [str] 'Apz9UhdlwKP4RGtjw7H6g'
    □  'alias' (155670560) = [str] 'de-krauler-kroog-hamburg-2'
    □  'name' (155670720) = [str] 'De Krauler Kroog'
    □  'image_url' (155666456) = [str] 'https://s3-media4.fl.yelpcdn.com/bphoto/4i5xiZ6T41NoeGwBNDEqw/o.jpg'
    □  'is_closed' (155666496) = [bool] False
    □  'url' (155670912) = [str] 'https://www.yelp.com/biz/de-krauler-kroog-hamburg-2?adjust_creative=hBJOknQJ58S5goQyjmwRw&utm_campaign=yelp_api_v3&utm_medium=api_v3_business_search&utm_source=hBJOkn'...
    □  'review_count' (155666536) = [int] 6
  □  'categories' (155666576) = [list] [[alias: 'german', title: 'German']]
  □  'rating' (155671040) = [float] 5.0
  □  'coordinates' (155666616) = (dict) {'latitude': 53.3996, 'longitude': 10.2142}
    □  'latitude' (155666656) = [float] 53.3996
    □  'longitude' (155666696) = [float] 10.2142
    □  '_len_' = [int] 2
  □  'transactions' (155666736) = [list] []
    □  'price' (155671072) = [str] '€€€'
  □  'location' (155666776) = (dict) {'address1': 'Kraueler Hauptdeich 65', 'address2': None, 'address3': None, 'city': 'Hamburg', 'zip_code': '21037', 'country': 'DE', 'state': 'HH', 'display_address': 'Kraueler Hauptdeich 65, 21037 H...
    □  'phone' (155671424) = [str] '+49 40 7230368'
    □  'display_phone' (155667096) = [str] '+49 40 7230368'
    □  'distance' (155667176) = [float] 1646.4291237763646
    □  '_len_' = [int] 16
  □  00001 = (dict) {'id': 'NH8UdkeMrq2YFSQLMSerQ', 'alias': 'hof-eggers-hamburg-3', 'name': 'Hof Eggers', 'image_url': 'https://s3-media4.fl.yelpcdn.com/bphoto/jfzfR1d-xTAeF5dn9hAg/o.jpg', 'is_closed': False, 'url': 'https://w...
  □  00002 = (dict) {'id': 'OjgVmVhRkm62fICP646uQ', 'alias': 'doner-imbiß-beim-tv-nord-winsen', 'name': 'Doner Imbiß beim TÜV-Nord', 'image_url': 'https://s3-media3.fl.yelpcdn.com/bphoto/qkNvpocZhvRdwTchluSqv/o.jpg'...
  □  00003 = (dict) {'id': 'AwL7KooU6B09c4AQ7eN5w', 'alias': 'soetebier-winsen-3', 'name': 'Soetebier', 'image_url': 'https://s3-media1.fl.yelpcdn.com/bphoto/k2FZAEEfyPKbpAdqaRi-Bsg/o.jpg', 'is_closed': False, 'url': 'https://www.y...
  □  00004 = (dict) {'id': 'dzT3zUy7O8x0yL09p9_jgA', 'alias': 'bäckerei-u-café-zimmer-winsen', 'name': 'Bäckerei u. Café Zimmer', 'image_url': '', 'is_closed': False, 'url': 'https://www.yelp.com/biz/b%C3%A4ckerei-u-caf%C3%94-zimm...
  □  00005 = (dict) {'id': 'GorMSUXMnwrfUMzb1dVgCQ', 'alias': 'burger-king-winsen-aller', 'name': 'Burger King', 'image_url': 'https://s3-media3.fl.yelpcdn.com/bphoto/Vmgv2RbNbLf9eSUUHXIVQ/o.jpg', 'is_closed': False, 'url': 'https://w...
  □  00006 = (dict) {'id': 'CP15eMmPVJUHL363eNV7Ag', 'alias': 'croque-knut-winsen', 'name': 'Croque Knut', 'image_url': 'https://s3-media3.fl.yelpcdn.com/bphoto/hgxA8LaPmwhnicJkk&GSew/o.jpg', 'is_closed': False, 'url': 'https://w...

```

FIGURE 5.4: Wrangled data stored in namedtuple

In a Python script called “add_restaurant_success.py”, the average rating and review count have been standardized to set them to a comparable range of values while still keeping the effects of outliers. The standardization has been conducted with the StandardScaler function of the scikitlearn-preprocessing package. Having the average rating and review count in comparable ranges the values were added and saved into a success variable to value them equally in their part of the success. The success, hence, is the sum of the standardized Yelp review count and average rating. In figure 5.5, the distribution of the success values can be seen, they are in the range between ca. -4.2 and ca. 13.3.

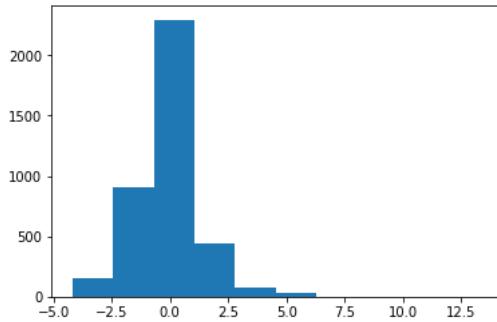


FIGURE 5.5: Distribution of Restaurant Success Values

5.3 City district profiles Hamburg

In this study, it was considered that the city district that a restaurant was located in, may have an influence on the success of this restaurant, e.g. through direct or indirect effects from the demographics in this district. Therefore, “city part profiles” published in the Transparency Portal of Hamburg were taken into account [30]. The dataset presents structure data for all of Hamburg’s city parts for the topic areas of population, living, council elections, social structure, infrastructure and traffic. The last dataset was published on the 19.03.2018, but dated back to data from 2016.

The data was available in XLSX-format, in which it was downloaded. Afterwards, the formatting was adjusted to be able to save the content as a CSV. With this CSV, all the necessary information, that would later be used in the analysis, was available as a spreadsheet that could be loaded into Python. In a Python script called “profile_translation.py”, all values of the city district profiles were translated from German to English.

5.4 Proximity to water

As Hamburg is an important port city and almost 10% of the city is covered by the harbor [31], water is an important consideration in the analysis of restaurant success. To know the water locations in hamburg is important, because restaurants are regularly not placed on water (restaurant boats are not considered as potential location candidates for simplicity). Additionally, water may influence the success of a restaurant with customers who may like to sit with a water view. Out of these reasons, the locations of water in Hamburg should be made out and the proximity of each restaurant to the next water location should be calculated.

For the extraction of water location, a geological ground map of Hamburg in the scale 1:5,000 [32] was added as a Web Map Service (WMS) layer into QGIS, as depicted in figure 5.6. The dark blue color on this WMS layer indicates that an area is on water ground, e.g. a river. Since the water spots

on WMS layer were not able to be processed as measurable points, yet, these had to be preprocessed before.

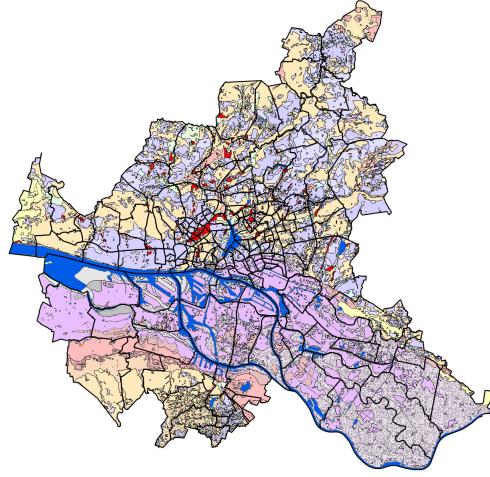


FIGURE 5.6: WMS Layer “Geologische Karte Hamburg”

In the first step, the WMS map was saved as a raster with the “Save as...” option of QGIS. The output was set to a “Rendered Image” in Format “GTiff”. The CRS stays ETRS89 and the resolution is first set to 100%, described by entering a 1 into both horizontal and vertical resolution. A new raster file is created.

The newly created raster file had a high resolution and needed to be reduced to be processed further. The resolution was reduced by conducting the QGIS raster conversion function “Translate (Convert format)”. In the upcoming dialog, the “outsize” resolution was set to 3% and the output location was set to a new GeoTIFF file. In the created execution code in the bottom of the dialog, the standard “of”-parameter “GMT” had to be changed to “GTIFF”. The resolution was after the resolution reduction still adequate to present the water areas.

In the third step, the new and reduced raster file was converted into a grid that would store a certain color value for each of the points on the raster together with their coordinates. This could be conducted by using the same QGIS conversion function “Translate (Convert format)” as in the previous step, however, in the corresponding dialog, the “outsize” resolution was not touched and the output format was set to an American Standard Code for Information Interchange (ASCII) gridded XYZ file. In the generated XYZ-grid, a delimited value storage of three attributes per record could be found. The first two attributes present the x and y values on the ETRS89 CRS and the third value a value for the colors with their luminosity is stored. The luminosity “0” is the darkest value in the grid and marks the prior dark blue water locations. This knowledge was used in a python script called “extract_water_locations.py” that was implemented to create a new dataset by looping through every record in the grid and to only keep the records, in which the luminosity would match “0”. Thereby only water locations would be part of the new dataset.

The fourth step consisted of the calculation of the distances of each restaurant to the next water spot by using a Python script called “calculate_distance_to_water.py”. In this script, all restaurants and water locations were loaded into the memory. Then, for each restaurant, the lowest distance to water was calculated by looping through all water points and calculating the Euclidean distance to them. At the end, a minimum water distance was found for each restaurant that was stored together with the restaurant ID in a CSV-file for the use in analysis.

5.5 Restaurant density

As described in chapter 2.2 on page 3, restaurant density is an important consideration in a selection of a location for a new restaurant. Therefore, this measure was included into the analysis of this paper. A pedestrian usually is willing to walk up to 400 meters to a nearby restaurant [33]. Hence, the restaurant density would present the number of all restaurants, which would be in a range of 400 meters to a location. To include the restaurant density as a feature into the restaurant analysis, a Python script called “get_restaurant_density.py” was created, in which the Euclidean distance of all restaurants to all other restaurants were calculated and the ones with a proximity of 400 meters or less were counted to the density of a restaurant.

5.6 Count of criminal cases

The criminal activity near a location of a restaurant can be an important factor in the success of the restaurant, as explained in chapter 2.2. The state office of criminal investigation in Hamburg uploaded a crime report in 2017, which presents different criminal activities per city district and additionally shows the total count of criminal activities per city district [34]. This total count of criminal activities was taken from the report and added as a record set for the analysis in the restaurant location analysis.

Chapter 6

Data analysis

6.1 Feature selection and preparation

To prevent overfitting, the number features in the training data should be kept to only the valuable features that support the prediction success in regression. The available data sources and features have been outlined in chapter 5 on page 15.

Apart from the dependent variable represented by the calculated success of a restaurant, included restaurant information from Yelp are the price level of a restaurant represented by one or multiple Euro-signs, and the restaurant category, e.g. a café, Vietnamese restaurant or a Korean grill. To use these two pieces of information, they were discretized and the restaurant category furthermore needed to be binarized. Additionally with the location of a restaurant, the proximity to the next water and restaurant density were calculated as explained in section 5.4 and section 5.5 and added as features for the regressor training.

As explained in section 5.3 and section 5.6, Hamburg provides city profiles with values representing each city district. Considering the factors for a restaurant location analysis explained in chapter 2.2, the following numbers were chosen to be included into the regressors training for a restaurant depending on the city district that the restaurant was located on:

- Count of criminal cases
- Population
- Share of under 18 year olds
- Share of foreigners
- Population with migration background

- Share of the population with migration background
- Number of households
- Share of one person households
- Share of households with children
- Share of households with single parents
- Population density
- Employment quote
- Share of unemployed people
- Sum of incomes per tax reliable person in EUR
- Prices for properties
- Prices for condominiums
- Share of students in Gymnasium
- Car density

To prevent the effects that numbers in different value ranges would have, the features were standardized.

6.2 Regressor training

In the training of a regressor, a regression technique is fitted to the training data to create a regressor that is able to predict values for unseen sets of features. In this paper, three different types of regression techniques were evaluated for their adequacy on the prediction of the restaurant success - RFR, Support Vector Regression (SVR) and linear regression.

The available data was split into a training set containing 70% of the data records and a test set containing 30% of the training records. Each of the regressors was fitted with the training set and then should predict both the values from the training set as well as the values from the test set. All three steps of the train-test-split, regressor fitting and prediction of test features were conducted 100 times to gain a stabler overview of the performance measures of the regressors. The R^2 and adjusted R^2 values were measured for both the training and test predictions for all regressors and are explained in the following.

Based on experiences made with the training and test data, the RFR regressor was set to consist of 1,000 trees and have a maximum depth of trees of five nodes. The SVR regressor was set to employ the Radial basis function (RBF) kernel.

In figure 6.1 the R^2 values for the prediction of training set values for all three regression techniques are shown. It can be seen that the random forest values vary between values of ca. 0.257 to 0.262, while SVR and Linear Regression (LR) show R^2 values between ca. 0.23 than 0.24.

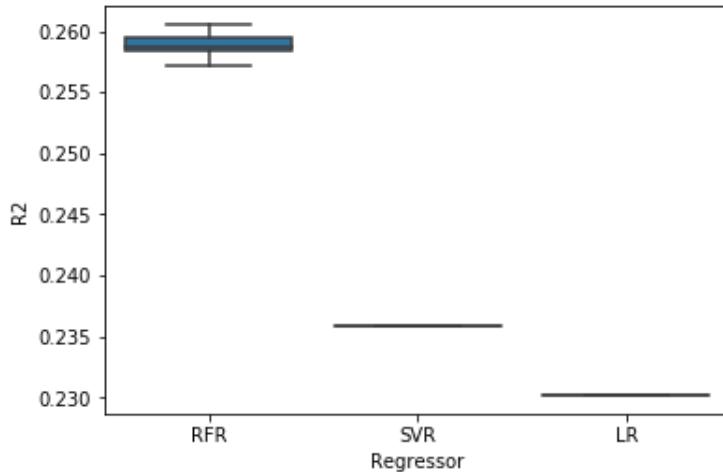


FIGURE 6.1: Boxplot diagram of the R^2 values for training data predictions

Figure 6.2 presents the Adjusted R^2 values of the training data predictions of the regressors. While the Adjusted R^2 values of all three regressors are a little lower than their R^2 values (deviation of ca. 0.05 for each regressor) for the predictions of the same data, the order of regressors in their performance is the same. Hence, RFR performs best on the training data among all three regressors.

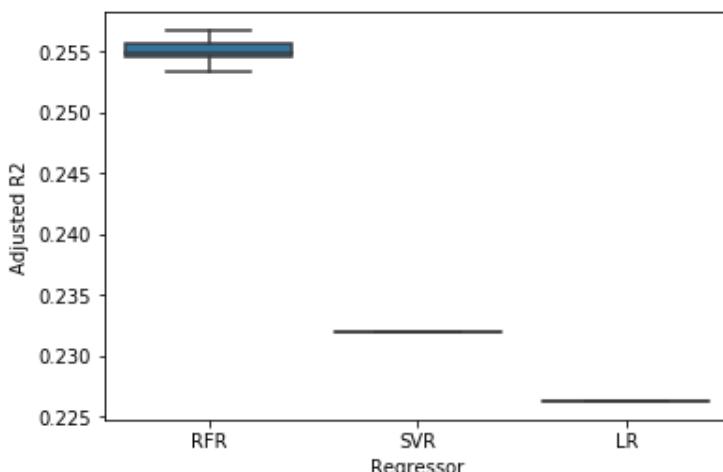


FIGURE 6.2: Boxplot diagram of the Adjusted R^2 values for training data predictions

The boxplot diagrams for the R^2 values of the test set predictions for all regressors are depicted in figure 6.3. However, the LR values are so low with values of $-4 * 10^{27}$ that the SVR and RFR values

are not displayed visibly. Therefore, the LR values have been excluded in figure 6.4. It is recognizable that the SVR and RFR predict the test values with a much higher accuracy than linear regression and that SVR performs even slightly better than RFR.

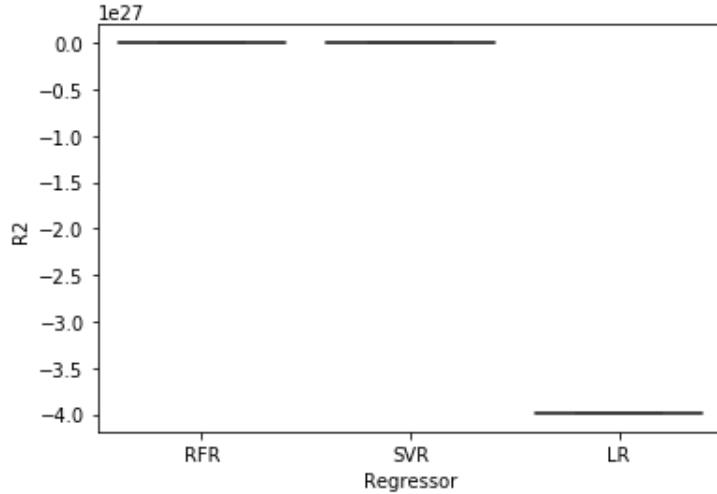


FIGURE 6.3: Boxplot diagram of the R^2 values for test data predictions of all regressors

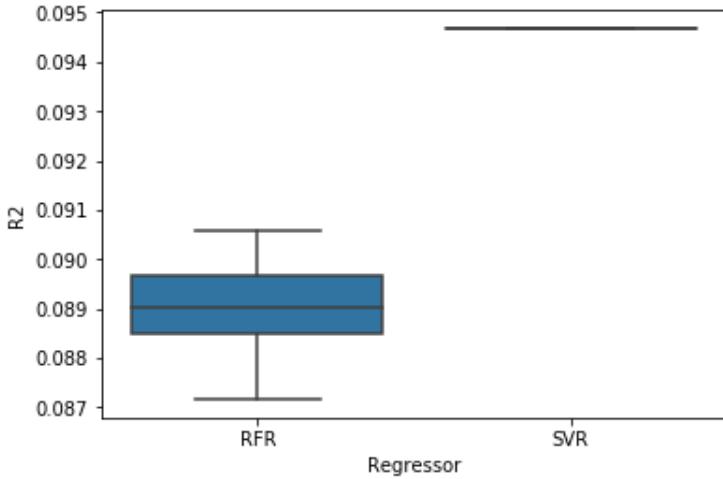


FIGURE 6.4: Boxplot diagram of the R^2 values for test data predictions excluding linear regression

Since a change of R^2 values to Adjusted R^2 values never increases the values, the the LR Adjusted R^2 values will still be very low for Adjusted R^2 and they are not depicted in the Adjusted R^2 values that are displayed in figure 6.5. In this boxplot diagram, it can be seen that SVR delivers slightly more accurate values than RFR for the test values, as prior was also shown by the R^2 values of the test set predictions.

These R^2 values, having values of more than 0, but still far away from 1, indicate that the feature analysis of location dependent features have a correlation with the restaurant success better than simple guessing of the mean values, but not as high as to be the one indicator for restaurant success. Values of ca. 0.2 for the training data and ca. 0.1 for the test data may give an indication of which

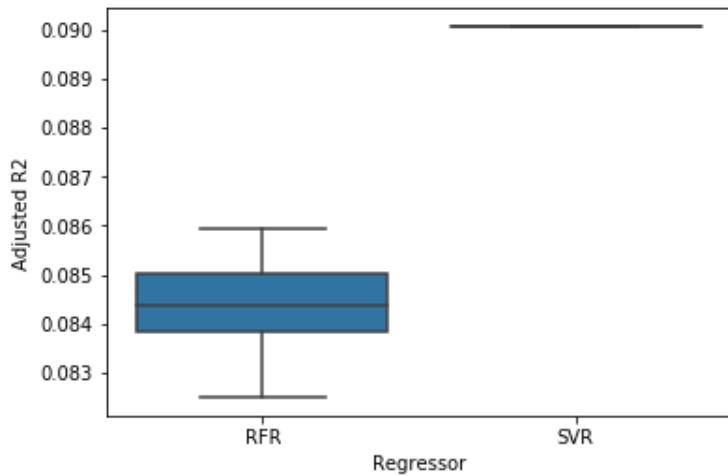


FIGURE 6.5: Boxplot diagram of the Adjusted R^2 values for test data predictions excluding linear regression

locations statistically lead to more success than others, but most of the restaurant success is still dependent on other factors, such as the quality of food, service, cleanliness and others, as explained in chapter 2.2.

The RFR regressor showed by far the most accurate predictions in the training set, while the SVR regressor delivered a slightly higher accuracy for the test set predictions. In consideration of the adequacy of RFR for restaurant location analyses given by the reasons explained in the section 4.1 and the overall acceptable performance in both the training data and the test data, further restaurant recommendations in the following sections are based on predictions by the RFR regressor.

6.3 Creation of a recommendation map

With the regressors that were trained in the section 6.2, potential beneficial points of opening up a restaurant in Hamburg should be deducted. To know, which locations represent a beneficial spot, the features given by points across all the map of Hamburg had to be analyzed and a chance for success had to be predicted. For this purpose, a recommendation grid was created that predicted the chance for success for centroids of squares (tiles) on the map of Hamburg. These centroids were determined with a Python script called “recommendation_grid_preparation”, which created centroids from the minimum x-value to the maximum y-value in the ETRS89 CRS. These tiles were created in a distance of 1,000 points on the x-axis and 1,000 points on the y-axis from each other. Therefore, each centroid was representative of a tile of 1,000 by 1,000 points/meters. The centroids were loaded into QGIS and are displayed in figure 6.6.

Since the grid centroids were created to a square map around Hamburg, not all grid centroids were, in fact, located within the boundaries of Hamburg. Therefore, the centroids were filtered to those on the

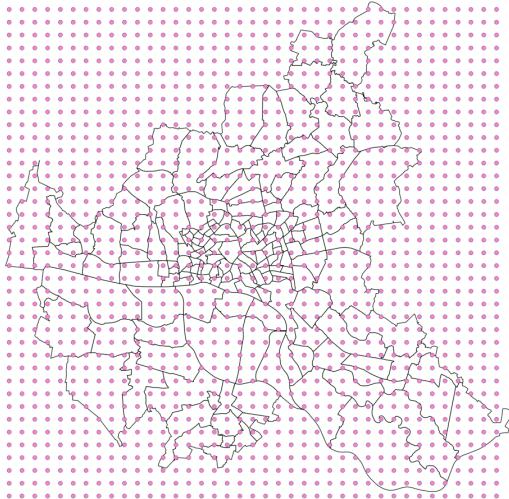


FIGURE 6.6: Recommendation grid centroids

map of Hamburg in QGIS and joined with the city districts, on which they were located, as displayed in figure 6.7.



FIGURE 6.7: Recommendation grid centroids filtered with the Hamburg districts

In the next step, all the features for the recommendation centroids were collected together and used to predict the success contributions of these locations. The location dependent features, which are described in section 6.1, were mainly location dependent and could be either extracted from the city district that a centroid was in, or calculated, like the distance to the next water bodies or the restaurant density in these centroids.

Two factors, however, were not location dependent and multiple values had to be fed into the feature sets to predict the success of each variation. One of these factors is the price level as a category from price level 1 to price level 5. Predictions for all of these price levels were created, however, most of

the times the price level 2 was selected as the price level with the highest success contribution, independent of the other features. As price has a relatively high contribution to restaurant success, the highest success contributing price was to be predicted for every point. The second factor, which was included and which was not location dependent, was the category of the restaurant. As there were many restaurants with different categories with the same success contributions, the “best” restaurant category was not predicted, but taken as an input parameter. The business question that the recommendation map therefore answers is, if a customer wants to open up a restaurant of a certain category, e.g. a Vietnamese restaurant, where the location with the highest success contribution would be. For this, the categories with the seven highest frequencies in Hamburg were taken into account and used for an exemplary analysis. The seven most frequent categories are presented in table 6.1. For each of these most frequent categories, a recommendation/prediction map was created.

restaurant category	number of restaurants
cafes	520
german	391
italian	387
hotdogs	245
greek	223
pizza	193
kebab	165

TABLE 6.1: Restaurant category frequencies in Hamburg

In the Python file “analysis.py”, predictions for all recommendation centroids were created in two for-loops that replicate the predictions for all seven top categories and all price levels. Afterwards, the prediction values were saved together with their locations and their feature values into CSV files that could later be loaded into QGIS. Additionally, the top five locations for every category map were extracted and saved into a separate CSV file to be able to highlight them on the maps.

For the creation of the recommendation maps, for each category, the recommendation centroids were loaded onto the QGIS Hamburg map together with their predicted success contribution and the additional features. These recommendation centroids were then saved in the same color and additionally converted to a raster that would show the success contribution as a color on the map. To save this raster, the attribute field was set to the predicted success and the raster resolution was set to 1,000 units per pixel. The color of the rasters were set to “Singleband gray” with a color gradient from black to white meaning the darker a location is, the less success is expected at this location and the brighter a location gets, the more success is expected. Furthermore, the transparency of the rasters were set to 40% to be able to see the Hamburg district map in the background. In addition, the top 5 centroids of each category were imported into QGIS and displayed on the map. In figure 6.8, the success prediction map for Greek restaurants is shown as an example. Enabling the “Identify

Features” option after clicking on the recommendation centroids layer of the respective category shows the features in a given location.

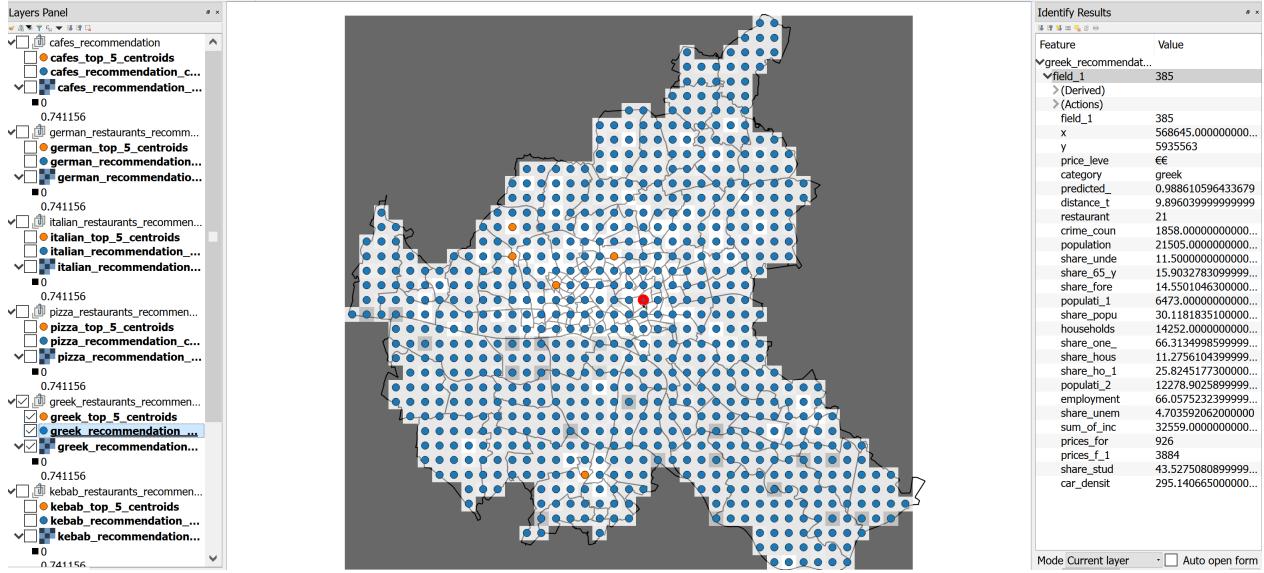


FIGURE 6.8: Greek restaurant success contribution map

6.4 Locations for restaurant recommendations for restaurants in Hamburg

In section 6.3, recommendation maps that indicate the success contribution of all points in Hamburg for a given restaurant category, were created. It was recognized that for the seven most frequent restaurant categories (cafes, german, italian, hotdogs, greek, pizza and kebab), always the five locations were found to be most success contributing. These five locations are displayed in figure 6.9. The x and y coordinates on the ETRS89 CRS are given in table 6.2. The predictions for these five locations are based on all the features, which were used in the RFR. Further, as price category, always price level 2 was predicted with the highest success contribution.

x	y
562645	5936563
566645	5938563
564645	5923563
559645	5938563
559645	5940563

TABLE 6.2: Top 5 locations for restaurant success according to random forest predictions (ETRS89)

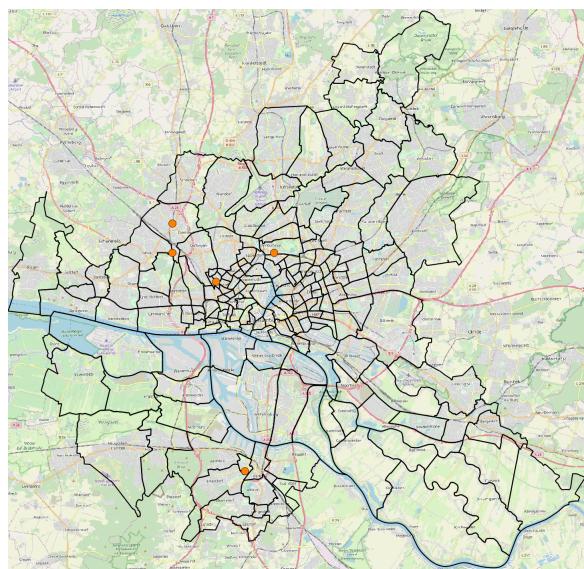


FIGURE 6.9: Top 5 locations for restaurant success according to random forest predictions

Chapter 7

Discussion

7.1 Exploratory data analysis

1. Variable Importance

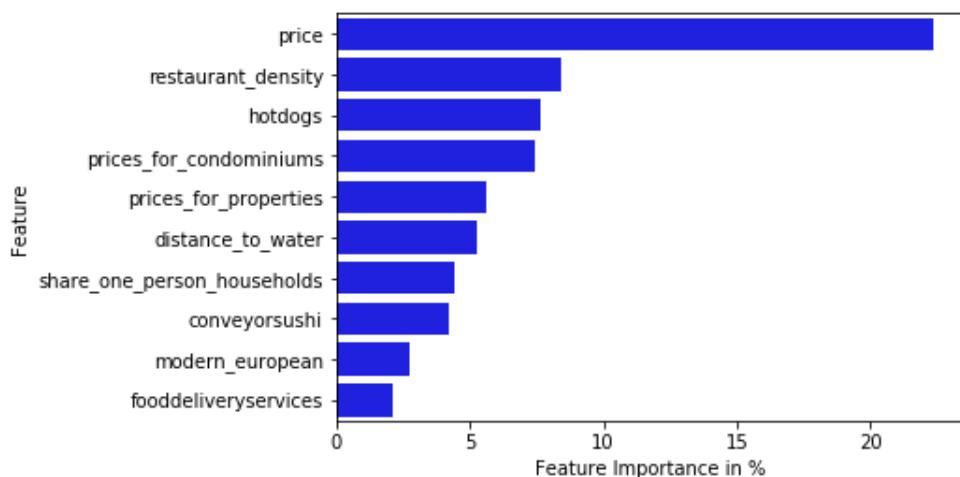


FIGURE 7.1: Feature importance for random forest predictions

Figure 7.1 shows the features which had high a importance in the random forest prediction. The overview indicates that the price is the most important feature for the success and failure of a restaurant. Further, the density of restaurants, type of food like hotdogs, price of condominiums, distance from water and households with one person play a significant role in the success of the restaurants.

2. Price level

Figure 7.2 shows the relationship between the price level and success of the restaurants in the dataset. It can be seen that the median of boxplots for different price levels increases together with the price level and it starts decreasing with a higher price level than 3. It indicates that the price level of

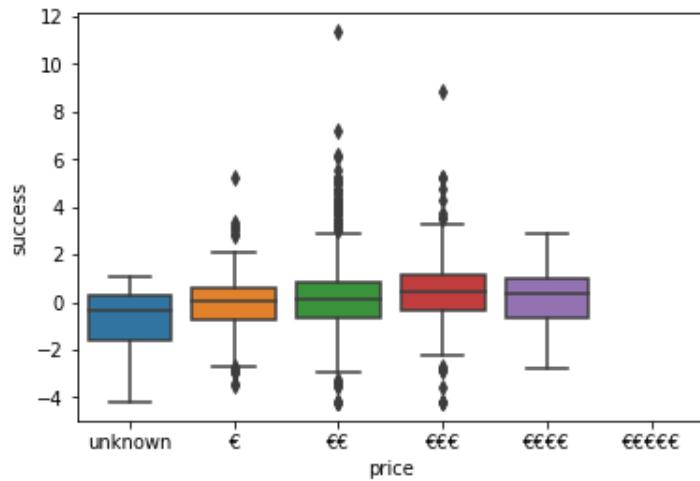


FIGURE 7.2: Boxplot diagram of restaurant price levels with success

restaurants in the extracted data had a correlation to their success. This could mean that there are consumers willing to pay high prices for restaurants and leave a good review, if they are satisfied and consider the prices justified. With price level 4, the overall restaurant success decreases, which could indicate either that too few consumers are willing to pay prices in this range or that the consumers were not satisfied with the restaurants' offerings. Price level 5 showed no records, which may either indicate a lack of demand in this price range or a lack of restaurants able to service a justified offering for price level 5.

3. Distance to Water

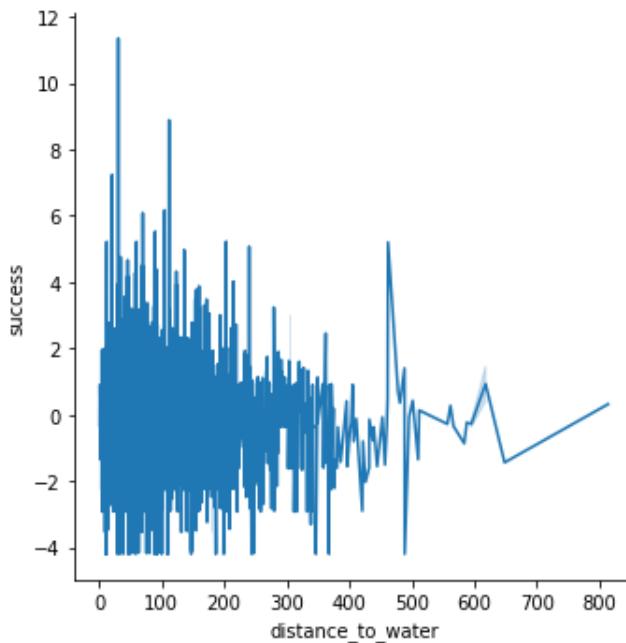


FIGURE 7.3: Restaurant distance to water and success

From figure 7.3, it can be seen that most of the restaurants are located within a range of 400 meters from water spots and that there is a huge variation in the success rate of these restaurant. However, the number of restaurants with the higher success rate is in the majority present in restaurants within a range of 200 meters near water spots.

4. Condominiums Price

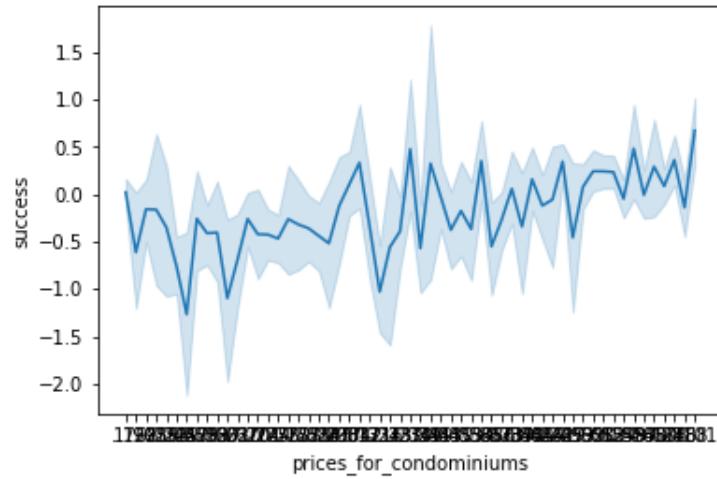


FIGURE 7.4: Condominium Price and success

Figure 7.4 shows that the price of condominiums and success of the restaurants in the dataset are slightly positively correlated, although this correlation varies strongly. This could mean that the prices of condominium are high in the area of high restaurant demand and supportive markets, such as downtown areas, which may lead to high restaurant success rates.

5. Income

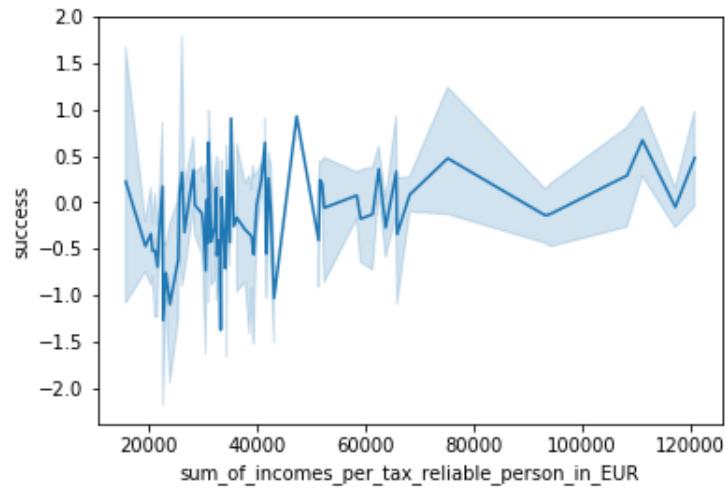


FIGURE 7.5: Income and success

From figure 7.5, it can be seen that the success of restaurants in the extracted data is proportional to the income of the inhabitants in the vicinity of a restaurant. Also, it could indicate that the number of restaurants is significantly less in extremely high-income areas of more than 70,000 euros per inhabitant and year, but still these restaurants have high success rate.

6. Restaurant density

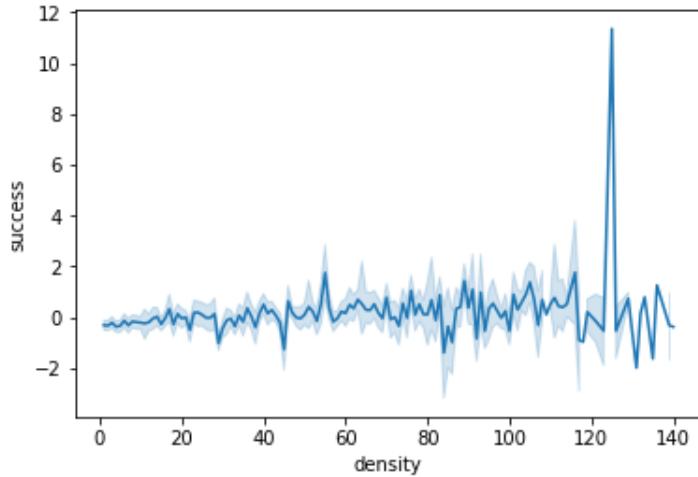


FIGURE 7.6: Restaurant density and success

Figure 7.6 shows that the overall success rate does not vary much with the increase in restaurant density. This could be due to the positive effect restaurant density through a high demand of the market providing enough customers to all restaurants and the negative effect of competition through other restaurants.

7. Population Density

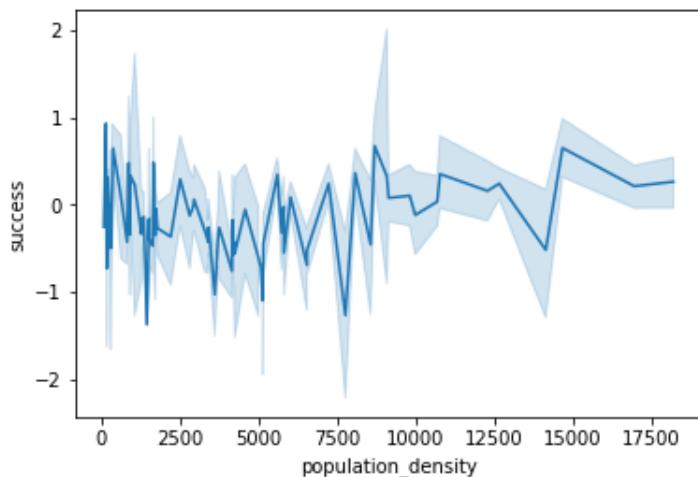


FIGURE 7.7: Population density and success

Figure 7.7 shows that the frequency and success rate of restaurants are higher in the low population density area with less than 2,500 inhabitants. This could be due to the reason that areas in the vicinity of city centers with more commercial activity are usually less residential area. Additionally, it can be seen that the frequency of restaurants in high population density areas with more than 9,000 inhabitants is lower, but the success rate tends to be higher which could be due to the residential area.

8. One-Person Households

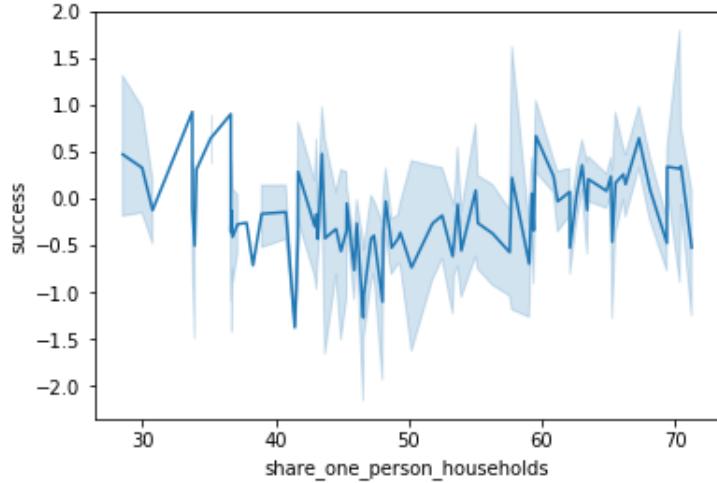


FIGURE 7.8: One-person households and success

It can be seen from the value frequency in figure 7.8 that the number of restaurants is higher in areas with a high share of one-person households of more than 45%. Additionally, the success rate tends to be higher in the areas with high shares of one-person households. This could be due to people in one-person households having their meal more frequently outside with others in company.

7.2 Discussion of results

In this paper, a recommendation map for restaurants together with beneficial locations for a new restaurant could be created. This recommendation map answers to any potential restaurant owner the question of where to place a restaurant. Since the R^2 and Adjusted R^2 values of the RFR were clearly positive, the predictions can bring value to the decision. Furthermore, the features that were used for the prediction of the success, can be delivered to potential restaurant owners, which enables them to evaluate these values themselves for their location decision.

However, the overview of 21 features and the restaurant category as input values for a location decision may be too complex for a human to handle. Since the final model of 1,0000 weighted decision trees within the random forest regressor is also too complex to understand for human beings, the grounds for the prediction of restaurant success are difficult to understand. This may leave a potential restaurant

owner with a decision on a location for a new restaurant without the understanding of why this location was chosen. The feature importance and exploratory data analysis presented in section ?? present the relationships between important features of the random forest regressor, but a decision model with more than 22 features is challenging to understand by two-dimensional models than can be visually understood by human beings.

Additional to the problem of model simplicity, the success of a restaurant is not only explainable by its location. As described in section 2.1, there are more factors to the viability of a restaurant than its location. This is recognizable in the R^2 and Adjusted R^2 values for the test sets of ca. 0.1. These values are high enough to add value to a decision, but are not high enough to rely the restaurant success on the predictions.

A challenge to the prediction of restaurant success, furthermore, was the success metric itself. Restaurant success would usually be measured by financial figures like the profit, revenue or customer count of this restaurant. In lack of this data, a success metric was created based on the Yelp reviews count and rating, which itself would be a valuable feature to predicting a restaurant's success. As the reviews of a restaurant both are a mirror of the restaurant success, but itself are a feature that can influence the restaurant's success, this metric is applicable for this paper, but has limitations. These limitations may also stem from the fact that a restaurant, which has high rated and many reviews, but which does not work profitably still would be seen as a successful restaurant.

Another limitation resulting from the given data is the restriction of the restaurant location analysis to one city. As one city's feature may be closer to each other than in the comparison to other cities or villages in the same or other countries, data from more cities would be beneficial to the prediction of success. However, as the provision and measurement of data is not performed equally in different countries and even within Germany, such an analysis would suffer from a questionable comparability of data and the availability of data. Nevertheless, the regression model found in this paper would benefit from future research into the verification of the model and review of the comparability of available data sources used in this paper for other areas.

Despite all these challenges, the predictions of restaurant success and creation of a recommendation may still be a beneficial technique for restaurant owners, e.g. to start or to verify their own analysis.

Chapter 8

Conclusion

This paper aims at a restaurant location analysis in Hamburg and strives to find a method of selecting a beneficial location for the success of restaurants. For this, the correlation of location dependent features to restaurant success has been analyzed and applied to a random forest regressor, which is able to predict the contribution to a restaurant success from the features given by a location. With this, a map of Hamburg could be created that presents recommended locations for the placement of a restaurant given a restaurant category. Five locations were named, which are predicted as beneficial to restaurant success on any of the seven most frequent restaurant categories in Hamburg.

Multiple challenges were found with the used method of analyzing data from a review portal and from city related sources to predict success for specific locations. Among these challenges was the complexity of RFR that impedes a clear understanding of the decision basis for choosing a location as beneficial for opening a restaurant. An exploratory data analysis in the review of two-dimensional graphs of features against each other supported the understanding, but suffered high variance and the little contribution that each of the features present in the overall complex prediction model. Moreover, as location dependent features are only one part of the factors for a restaurant success, only a contribution to restaurant success by locations could be recognized. Further, the restaurant success metric was a deducted metric from review portal reviews and, while presenting a mirror of a restaurant's success, review portal reviews are not regularly employed as the main indicator for restaurant success.

The regression model for location contribution to restaurant success was applied for the city of Hamburg. As differences of environment factors may be lower within a city than to other locations outside of a city, a future research could be conducted to cross-verify the model found for Hamburg in other areas. Another benefit from such an analysis would be the verification of the availability of data sources that could be found in Hamburg.

Notwithstanding the challenges of this research and the merger of multiple non-related data sources, a contribution of the location to the success of restaurants could be found and applied to new locations

within Hamburg. This may aid restaurant owners in their choice of a location, e.g. as the start of their location search analysis. Additionally, it may support future research in the exploration of location factors influences to restaurant success comprising which features can be used and how the data for these can be retrieved. In this manner, this paper adds knowledge about the composition of success factors to the restaurant location analyses and may represent one part of the way to a greater clarity in this area.

Bibliography

- [1] H. G. Parsa, John T. Self, David Njite, and Tiffany King. Why restaurants fail. *Cornell Hotel and Restaurant Administration Quarterly*, 46(3):304–322, 2005. ISSN 0010-8804.
- [2] Yelp Inc. An introduction to yelp metrics as of december 31, 2018, 2018. URL <https://www.yelp.com/factsheet>.
- [3] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*, volume 10. Springer series in statistics New York, 2001.
- [4] Roswitha Murjahn and Sascha Tegtmeyer. Open data/transparenzportal hamburg—grundlagen, umsetzung, erfahrungen, auswirkungen. *Zfv—Z Für Geodäsie Geoinformation Landmanagement*, 5(2016):330–335, 2016.
- [5] Kevin Mellet, Thomas Beauvisage, Jean-Samuel Beuscart, and Marie Trespeuch. A “democratization” of markets? online consumer reviews in the restaurant industry. *Valuation Studies*, 2(1):5–41, 2014. ISSN 2001-5992.
- [6] Gwo-Hshiung Tzeng, Mei-Hwa Teng, June-Jye Chen, and Serafim Opricovic. Multicriteria selection for a restaurant location in taipei. *International journal of hospitality management*, 21(2):171–187, 2002. ISSN 0278-4319.
- [7] Angelo A. Camillo, Daniel J. Connolly, and Woo Gon Kim. Success and failure in northern california: Critical success factors for independent restaurants. *Cornell Hospitality Quarterly*, 49(4):364–380, 2008. ISSN 1938-9655.
- [8] Evan Tarver. How to choose the best restaurant location for your business, 21.04.2017. URL <https://fitsmallbusiness.com/choose-a-restaurant-location/>.
- [9] Webstaurantstore.com. Restaurant location analysis, 25.07.2018. URL <https://www.webstaurantstore.com/article/81/restaurant-environmental-analysis.html>.
- [10] Li-Fei Chen and Chih-Tsung Tsai. Data mining framework based on rough set theory to improve location selection decisions: A case study of a restaurant chain. *Tourism Management*, 53:197–206, 2016. ISSN 0261-5177.

- [11] TripAdvisor Media Group. 2017 tripadvisor annual report and notice of 2018 annual meeting and proxy statement, 27.04.2018. URL <http://ir.tripadvisor.com/static-files/840c6d1c-9c17-46c9-b52f-3586ead2515f>.
- [12] Daying Zhao, Bin Fang, Huiying Li, and Qiang Ye. Google search effect on experience product sales and users'motivation to search: Empirical evidence from the hotel industry. *Journal of Electronic Commerce Research*, 19(4):357–369, 2018. ISSN 1938-9027.
- [13] Soyeon Shim, Mary Ann Eastlick, Sherry L. Lotz, and Patricia Warrington. An online prepurchase intentions model: the role of intention to search: best overall paper award—the sixth triennial ams/acra retailing conference. *Journal of retailing*, 77(3):397–416, 2001. ISSN 0022-4359.
- [14] Google LLC. Places. URL <https://cloud.google.com/maps-platform/places/>.
- [15] Linchi Kwok and Bei Yu. Spreading social media messages on facebook: An analysis of restaurant business-to-consumer communications. *Cornell Hospitality Quarterly*, 54(1):84–94, 2013. ISSN 1938-9655.
- [16] statista. Number of monthly active facebook users worldwide as of 4th quarter 2018 (in millions). URL <https://www.statista.com/statistics/264810/number-of-monthly-active-facebook-users-worldwide/>.
- [17] Facebook Inc. Graph api. URL <https://developers.facebook.com/docs/graph-api>.
- [18] Konstantinos Papamiltiadis. Enhanced developer app review and graph api 3.0 now live, 2018. URL <https://developers.facebook.com/blog/post/2018/05/01/enhanced-developer-app-review-and-graph-api-3.0-now-live/>.
- [19] Ulrike Grömping. Variable importance assessment in regression: Linear regression versus random forest. *The American Statistician*, 63(4):308–319, 2009. ISSN 0003-1305. doi: 10.1198/tast.2009.08199.
- [20] Andy Liaw and Matthew Wiener. Classification and regression by randomforest. *R news*, 2(3):18–22, 2002. ISSN 1609-3631.
- [21] Thomas G. Dietterich. An experimental comparison of three methods for constructing ensembles of decision trees: Bagging, boosting, and randomization. *Machine learning*, 40(2):139–157, 2000. ISSN 0885-6125.
- [22] R. Devasthal. Coefficient of determination (r-squared) explained, 2018. URL <https://towardsdatascience.com/coefficient-of-determination-r-squared-explained-db32700d924e>.
- [23] Alkis verwaltungsgrenzen hamburg, 28.02.2018. URL <http://suche.transparenz.hamburg.de/dataset/alkis-verwaltungsgrenzen-hamburg8?forceWeb=true>.

- [24] Authoritative real estate cadastre information system (alkis®). URL <http://www.adv-online.de/Products/Real-Estate-Cadastre/ALKIS/>.
- [25] Yelp Inc. Yelp fusion authentication, . URL <https://www.yelp.com/developers/documentation/v3/authentication>.
- [26] Geoffrey Fairchild. yelpapi github, 2018. URL <https://github.com/gfairchild/yelpapi>.
- [27] Yelp Inc. Yelp fusion /businesses/search, . URL https://www.yelp.com/developers/documentation/v3/business_search.
- [28] Jeff Whitaker. pyproj 1.9.6, 2018. URL <https://pypi.org/project/pyproj/>.
- [29] Python Software Foundation. json — json encoder and decoder¶. URL <https://docs.python.org/3/library/json.html>.
- [30] Stadtteil-profile hamburg, 2018. URL <http://suche.transparenz.hamburg.de/dataset/stadtteil-profile-hamburg2>.
- [31] Die stadt. URL <https://www.hamburgportal.de/die-stadt-hamburg/>.
- [32] Geologische karte 1:5 000, 11.12.2018. URL <http://suche.transparenz.hamburg.de/dataset/geologische-karte-1-5-00011?forceWeb=true>.
- [33] Yong Yang and Ana V. Diez-Roux. Walking distance by trip purpose and population subgroups. *American Journal of Preventive Medicine*, 43(1):11–19, 2012. ISSN 0749-3797.
- [34] Landeskriminalamt Hamburg. Polizeiliche kriminalstatistik 2017: Ausgewählte delikte nach bezirken / stadtteilen, 2017. URL <https://www.polizei.hamburg/contentblob/10538308/47b77ab8d33f4c8483d2efd29588ca5f/data/pks2017-stadtteilatlas-do.pdf>.