

TRANSFERRING KNOWLEDGE FROM MONITORED TO UNMONITORED AREAS FOR FORECASTING PARKING SPACES

ANDREI IONITA

*Computer Science, RWTH Aachen University
Aachen, Germany
andrei.ionita@rwth-aachen.de*

ANDRÉ POMP

*Institute of Information Management in Mechanical Engineering, RWTH Aachen University,
Aachen Germany
andre.pomp@ima.rwth-aachen.de*

MICHAEL COCHEZ

*Fraunhofer Institute for Applied Information Technology FIT, Aachen, Germany
Department of Computer Science, Vrije Universiteit Amsterdam, Netherlands
Faculty of Information Technology, University of Jyväskylä, Finland
michaelcochez@gmail.com*

TOBIAS MEISEN

*Chair of Technologies and Management of Digital Transformation, University of Wuppertal,
Wuppertal, Germany
meisen@uni-wuppertal.de*

STEFAN DECKER

*Computer Science 5, RWTH Aachen University, Germany
Fraunhofer Institute for Applied Information Technology FIT, Aachen, Germany
stefan.decker@dbis.rwth-aachen.de*

Preprint of an article to be published in Int J. on Artificial Intelligence Tools (IJAIT) ©2019 [copyright World Scientific Publishing Company] <https://www.worldscientific.com/worldscinet/ijait>.

Smart cities around the world have begun monitoring parking areas in order to estimate available parking spots and help drivers looking for parking. The current results are promising, indeed. However, existing approaches are limited by the high cost of sensors that need to be installed throughout the city in order to achieve an accurate estimation. This work investigates the extension of estimating parking information from areas equipped with sensors to areas where they are missing. To this end, the similarity between city neighborhoods is determined based on background data, i.e., from geographic information systems. Using the derived similarity values, we analyze the adaptation of occupancy rates from monitored- to unmonitored parking areas.

Keywords: smart parking, machine learning, semantic annotation, data mining

1. Introduction

Parking problems and overall traffic congestion are commonplace in cities nowadays. While the number of cars continues to increase¹, studies show that about 30% of the traffic in cities is caused by cars that are actively searching for parking². Often it is the lack of planning on behalf of cities that does not accommodate parking facilities proportional to building developments³, which leads to double-parking, more accidents due to distracted drivers, more busy traffic, and, ultimately, a waste of fuel⁴. As infrastructure solutions are not always optimal and take a long time to implement, there are other strategies to overcome these issues.

Parking can be managed more efficiently if parking spaces would work on an allocation basis, with drivers reserving a spot of their choosing. Accounting for each individual spot and managing its reservations is, however, unrealistic and unsustainable at the moment. Merely providing an overview of areas with free parking spaces would be a big help. Suppose a driver would have access to such a service that indicates free parking spaces at the time when she is arriving in the area. The system would take into account the usual parking levels in the respective area depending on the day of the week and the time of day. Additional information such as current traffic, event data and weather would improve its estimations. Being aware of this sort of parking information beforehand, the driver would pick a traveling path that is less busy through the city and spend significantly less time finding a parking spot that is limited to the area indicated.

This work introduces an approach inspired by the vision to produce parking occupancy estimations for city areas without any previous measurements by taking into account the city infrastructure and the parking data from other areas. This paper extends our previous work^{5,6}, by enriching the evaluation of, and exploring alternative assumptions as, the original approach. The approach distinguishes itself from related research by including open Geographical Information System (GIS) data into the prediction computation.

This paper is organized as follows. Following the introduction, an overview of the research landscape in city parking is presented, before outlining the main assumptions behind our approach. The approach itself is broken down and described in detail, after which its evaluation is set up and carried out. Finally, further possible extensions are outlined and conclusions are drawn.

2. Smart Parking Research

Improving the parking situation using sensor data has been subject of research, especially since around 2000. In a 2017 survey^{7,8}, Lin compiles an overview of the advances in smart parking by splitting the results into three categories: *information collection*, *system deployment*, and *service dissemination*. We follow their categorization below.

2.1. Information Collection

Under information collection, the survey lists techniques to acquire parking information. Static sensors are usually mounted around parking meters, in the ground or on nearby lamp posts. The dynamic sensors are usually managed through a wireless network infrastructure and are present inside cars, such as taxis, which travel through the city and collect roadside information on parking. In both cases, the transmitted information is the occupancy situation: either the car has left or arrived at a parking space. A number of types of sensors are usually being used in the process: example include infrared sensors, ultrasonic sensors, accelerators, optical sensors, inductive loops, piezoelectric sensors, cameras, and acoustic sensors. Most of the time, the captured data requires post-processing in the form of image or audio recognition before arriving at the target occupancy information. Some captured information also raises privacy issues as it contains sensitive data about the car and driver. Smartphones provide means to collect data and have a great impact through crowdsourcing if the users are given incentives to enable the respective smartphone functions that automatically collect their information. Conversely, parking applications may give drivers incentives to initiate the data transmission themselves and report the parking situation on site.

Data collection based on smartphone-integrated sensors has been studied in numerous research publications, as it spared the authors to produce and mount special-purpose sensors. Xu et al.⁹ makes real-time parking availability estimations based on a system that aggregates the data coming in from mobile phones. The system uses algorithms based on statistical weighted schemes and Kalman filters. Additionally, the authors create parking availability profiles based on historical data and using statistical algorithms.

Chen et al.¹⁰ developed an Android application that finds a parking location at park-and-ride facilities by calculating the probability of parking availability and taking in consideration the shortest travel time. The authors employ fuzzy logic to model the uncertainty of parking availability, with the fuzzy membership function being linear. The authors proposed multiple criteria in finding the best parking location, such as train frequency, service quality, and park-and-ride price.

PocketParker is a crowdsourcing system, proposed by Nandugudi et al.¹¹ that uses smartphone data to predict parking availability. The system is used for parking lots. It requires no input from the user since it notices automatically when a user starts to drive or stops, i.e., departure and arrival events. Based on these, the system builds a probability distribution model that is used to answer queries about parking availability. PocketParker has proved robust to hidden parkers, i.e., parking vehicles that are not using the application. In the authors' simulation, it has reached 94% rate for parking availability prediction with 105 users over 45 days.

Koster et al.¹² propose a smartphone-based solution that recognizes when drivers arrive or leave parking spaces. A Bayesian approach and Hidden Markov Models (HMM) are used to model the parking spaces and respond to user queries for the

next parking space. The HMMs are based on gathered historical data. The answer to the user query is a parking space nearby and the probability that it is free at the respective time. The authors emphasize the non-intrusive nature of their solution where drivers only have minimal interaction with their phones to get a recommended free parking space.

2.2. System Deployment

Regarding system deployment, the smart parking survey⁷ refers to the varieties of parking systems, looks into how well they scale, and touches upon the data analysis side. The parking system software is the interface between data sources and the users. Software systems are often in the form of reservation systems, typically run by municipalities or private car parks. These systems may also provide guiding assistance in arriving to the desired parking lot or at the individual parking space. The vacancy prediction component informs the drivers about availability of parking spaces at the destination, either in real-time or at a specific date and time.

There are several publications that investigate parking system worth mentioning. Rajabioun and Ioannou¹³ introduce an information system for parking guidance that enables communication between vehicles and the infrastructure. It proposes a prediction algorithm that forecasts the availability of parking locations based on real-time parking information. It takes into account parameters such as parking duration, arrival time, destination, pricing, walking distance, parking capacity, rates of vehicles occupying and leaving parking spots, time restrictions, parking rules, events that disrupt parking availability, etc. Their algorithm uses a probabilistic density distribution model. The parking data was collected both from on-street parking meters and off-street garages in Los Angeles and San Francisco, USA. In a follow-up, Rajabioun and Ioannou¹⁴ propose a multivariate autoregressive model that considers the temporal and spatial correlations of parking availability when making predictions.

Tiedemann et al.¹⁵ present the development of a prediction system that estimates occupancy of parking spaces. The occupancy data is collected online via roadside parking sensors and the predictions are realized using neural gas machine learning combined with data threads. The authors notice that some factors play a significant part in the predictions, such as holidays, weather and use the neural gas clustering to separate the data before the data thread method is applied.

Richter et al.¹⁶ address the parking prediction problem with the focus on model storage in vehicles. The authors train models of various granularity that would predict parking availability based on the information contained: A one-day model per road segment, a three-day model per road segment, and a seven-day model per road segment. Additionally, models based on regions and time intervals computed by clustering are tried out. Hierarchical clustering with complete linkage is employed. The models are evaluated on street data from the SF *park* project¹⁷. The application of clustering before building the models shows a 99% decrease of model storage

space. The prediction success rate is at about 70%.

With iParker, Kotb et al.¹⁸ propose a system that handles parking reservations. It achieves resource allocation so that drivers pay less for parking, while parking managers receive more resource utilization and hence reach higher revenue. The system is based on mixed-linear programming (MILP). The system uses dynamic resource allocation and pricing models to achieve its goal. In its evaluation, cost cuts for drivers of 28% were reported, while achieving a 21% increase in resource utilization and an increased total revenue for parking management by 16%.

Shin and Jun¹⁹ propose an algorithm for smart parking that assigns cars to parking facilities in the city. The criteria based on which the assignment is realized includes driving distance to the parking facility, walking distance from the parking facility to the destination, parking cost, and traffic congestion. The real-time data is collected from parking facilities and from sensors that are integrated in cruising cars. The data is transferred from the central server, where it is managed through a wired/wireless telecommunication network. The authors tested their approach in Luxembourg City. The results of the simulations show improved figures for average driving duration, average walking distance, parking failing rate, parking utilization rate, average standard deviation on the number of guided cars to each parking facility, average occupancy ratio of parking facility, and for the parking facility occupancy rate.

ParkNet, developed by Mathur et al.²⁰ is a system made up of vehicles that capture parking space information while driving. Every ParkNet vehicle is equipped with a GPS receiver and an ultrasonic sensor facing sideways. The latter determines whether it passes by parking spaces and whether they are occupied. The data is sent to a central server that aggregates it in order to build parking space occupancy maps in real-time. The information is queried by clients that search for a free parking space. The system was evaluated in Highland Park, New Jersey and San Francisco on 500 miles road-side parking data and yield 95% accurate parking maps and 90% parking occupancy accuracy. The authors show that the system can further be improved if the sensors are fitted into taxicabs or city buses.

2.3. Data Dissemination

Under data dissemination, the survey⁷ addresses the capability of sharing parking information. This scenario occurs in decentralized parking systems, where the cars find out about free parking spaces in an area where other cars merely drive by and report the real-time situation.

A selected group of publications is driven by the information exchange approach. Caliskan et al.²¹ model the prediction of available parking spaces as a vehicular ad-hoc network (VANET). The network disseminates parking data in order to help the estimation of future occupancy of parking lots. The pieces of disseminated data are timestamp, total capacity of parking lot, number of parking spaces that are currently occupied, the arrival rate, and the parking rate. The latter two are used

in the modeling of continuous-time homogeneous Markov chains. The approach is otherwise based on queuing theory.

Klappenecker et al.²² builds on the result of Caliskan and uses an improved version of continuous-time Markov chains for predicting availability of parking spaces. Predictions are communicated between cars in an ad-hoc network. The approach simplifies the computations of transitional probabilities inside a Markov chain model. The system applies to parking lots that are connected to the ad-hoc network. These communicate the number of occupied spaces, capacity, arrival and parking rate.

Also based on VANETs is Szczurek et al.²³ work, which propose a novel approach that combines machine learning with the information disseminated in ad-hoc vehicular networks. The building blocks of the system are parking reports, which are issued by vehicles leaving a parking space and comprise a report identifier, a location, and a timestamp. The parking reports are being learned by a model, which then indicates whether a parking is available for a specific vehicle. A conditional relevance is used to determine whether a particular report is useful for a specific vehicle. This is modeled using a Naive Bayes method. A parking availability report R is labeled relevant by vehicle V , if the parking space referenced in R is available when V reaches it. Upon evaluation of the methods, the authors reported an improvement in parking discovery times for vehicles.

3. The SFpark Project

To implement and evaluate our proposal, we use the data from the *SFpark* project. This project was realized by the San Francisco Municipal Transportation Agency (SFMTA), the city agency that manages the city's transportation, which includes on-street parking^{24,25}, with the goal to improve parking availability. The SFMTA had the possibility of changing parking rates for on-street parking meters on short notice. Before the project started, parking rates were the same all day, every day, independent of the parking demand. By implementing a demand responsive pricing scheme, parking availability improved dramatically.

Dynamic pricing is a way to control parking occupancy. Parking prices are raised in areas that are almost fully occupied, whereas areas with low parking rates get assigned a lower price. A more advanced version adjusts the prices when enough demands received by the parking system would point to a future parking overload in the respective area.

In conducting the project, nine pilot areas were chosen for monitoring. Out of these areas, seven were selected to have new pricing policies, while two were control areas. The number of metered spaces used was 6,000, which amounts for 25% of the city's total. The meters allow rates to be deployed remotely, and they transmitted data to a central server through a wireless connection.

The data was collected using parking sensors. These provided the central server with the information needed to calculate the demand-responsive parking rates and

provided real-time parking availability information. A parking sensor is a magnetometer that detects changes in the earth's electromagnetic field. A total of 11,700 sensors were deployed, resulting in 8,000 spaces that were equipped with one or two sensors. The sensors delivered valid data from April, 2011 to December, 2013. *SFpark* made available real-time information on parking rates and parking occupancy through a smart phone application. The scale and scope of the *SFpark* project and its freely available data sets played an important role when choosing to base our project on it.

4. Rationale Behind Our Approach

The approach presented offers a solution to estimating parking occupancy without the help of sensor data. It is based on the observation that parking is determined by the specificities of city areas. Two residential neighborhoods of similar sizes, perhaps far apart from each other, will have very similar parking occupancies: high during nighttime and low during daytime. This will likely differ significantly from office areas, which tend to have most parking spaces occupied during the day and free during the night. Restaurants or shopping centers may represent another distinct category, where customers park usually during the evenings and on weekends, while in the other times they are not very busy, therefore producing a low parking occupancy.

Looking at a city, we can identify a pattern: the types of buildings and the time people spend there determines parking behavior. The presented approach builds on this pattern in order to estimate the level of parking occupancy. Specifically, it uses amenity types and time-spent data to complement established machine learning algorithms in order to arrive at parking levels in places where such forecasts cannot be made only with straightforward models.

The approach does *not* infer parking levels solely based on building metadata and busy times. This information is currently not enough and other factors regarding the city would be needed to arrive at a direct result. Some cities have better parking infrastructures while fitting the same number of people in the offices as cities with scant parking facilities. In the former parking are likely concentrated around the offices, while in the latter the cars are probably distributed uniformly around a larger area around the offices. To circumvent these inconsistencies between cities, we focus on single cities, where parking infrastructure is likely the same given the type of amenities and their dimension. This could be extended to cities in a region or a whole country, depending on the specifics.

The approach therefore uses the (dis)similarities between city areas, with their respective amenity types and time-spent information to help infer parking occupancy. It is assumed that an estimation model can be transferred from a source city area A to a target city area B without any amendment, provided A and B are perfectly similar according to their parking profile. In contrast, the parking occupancy estimation would very much differ if A and B are dissimilar. Specifically, the

approach does not offer a precise result in this case, resorting to an interval that expresses the possible parking level.

4.1. Motivational Example

The dimension of the problem we are attempting to solve here is best illustrated with a concrete scenario.

Bob is excited about the interview with a big IT company in his city. He will be driving to the office building located in one of the several office sites in the city. Bob does not like being late, even more so on this occasion, and he wants to leave himself enough buffer time before he arrives at the company reception desk. He has no idea about the parking situation on site, however. In this city, he could spend up to half an hour to find a free parking space. Therefore, Bob uses a new parking app, that can estimate the parking levels at almost every location; the system does not employ sensors everywhere, instead it works by extending the parking behavior from one site to another depending on their specificity, be it offices, restaurants, shopping, or residential. Bob likes the idea and enters his estimated time of arrival at the site and sees that the parking occupancy there will be between 60% and 80%. This is good enough for him, he knows that at least 1 out of 5 spaces will be free on average and will likely find a spot in a few minutes. He is suddenly more confident about his punctuality and can now drive more assured to the interview.

5. Approach

The approach, in its generic form, is split into several steps. It begins by acquiring access to data that contains information about parking occupancy for a desired area. The parking data is mapped geographically using OpenStreetMap (OSM). Next, the points of interest (POIs) contained in the OSM data are spatially clustered so that the individual clusters are of about the same size. Afterwards, machine learning models are trained on parking data for the computed clusters. For each cluster, mathematical representations are constructed based on the OSM data, which are then used to compute similarity values using *cosine similarity* and *earth mover's distance*. Finally, estimations of parking occupancy are computed by applying the models on areas without parking data with the similarity values factored in. We detail the process step by step below.

5.1. Overview

(i) **Get access to appropriate parking data**

To have a solid analysis foundation, it is essential to find a well-defined spatial area for which parking measurements over a continuous period of time have been made. Regular status updates, usually by hour, are preferred, if not as soon as they happen. In case multiple distinct data sources for the spatial area and time period are available, limiting oneself to the richest data source is recommended,

as multiple sources tend to have different time and space references, and can be inconsistent with regard to sensor errors.

(ii) **Map the parking data to OpenStreetMap layers**

Another important part is geographically referencing the parking data. OpenStreetMap layers such as points, lines, and polygons that include city artifacts and geographical coordinates together with their metadata are suitable. These are downloaded and associated with the parking occupancy information. The information on amenities included in OSM layers together with the time people spent in amenities is collected as well.

(iii) **Cluster the spatially-referenced data into multiple city areas**

Splitting the data corresponding to an entire city into multiple groups is a prerequisite of the approach. Especially, having city areas without parking data completely separated from the city areas with parking data so that the latter can later serve as estimation basis for the former. Splitting is performed spatially. Including any other property, such as OSM metadata, in the clustering algorithm results in noncontinuous areas, which would defeat the purpose of a driver finding a parking space inside a certain radius. Furthermore, the resulting clusters should be of about the same size, as this helps to make inferences later in the process. Averaging the occupancy among the parking spaces inside a cluster, for instance, is less representative for another cluster that has a number of parking spaces of a different order of magnitude.

(iv) **Build machine learning models for the clustered city areas**

A freshly split city area that contains parking data will have a occupancy estimation model. In the training process, the predictor variables includes the measurement timestamp, parking lot capacity, and parking price, while the target variable is the parking occupancy. Methods used for building the models are *decision trees*, *support vector machines*, *multilayer perceptrons*, and *boosted trees*.

(v) **Build mathematical representations for the city areas**

Amenities and time spend information are organized in mathematical objects to reflect the parking demand in their respective city areas. *Cluster vectors*, which serve in *cosine similarity* computation and *cluster Gaussians* that contribute to calculating *earth mover's distance* between city areas are used.

(vi) **Compute similarity values between any two city areas**

The built mathematical representations make it possible to compute similarity measures between city areas. *Cosine similarity* and *earth mover's distance* are defined and computed for every pair of city areas.

(vii) **Apply models on city areas that do not have parking information**

When computing the occupancy in clustered city areas with no parking data, the elements built up to now come together. Basically, the machine learning models are applied to the clustered areas without parking data. In the result, the similarity measure between the originating model area and the target area is factored in. In practice, this means that input data for the models will need to be constructed. The result output of the model will be extended in form of an

interval upon applying the similarity value: the smaller the similarity, the more the interval will be stretched around the original occupancy result. Occupancy values are expressed between 0% and 100%.

5.2. Getting Access to Parking Data

We consider the following types of data as parking data:

- (i) *parking occupancy* contains information on the availability of parking spaces at a defined location
- (ii) *traffic data* contains information regarding the street traffic intensity
- (iii) *weather data* contains temperature and rain information for the geographical location considered
- (iv) *event data* contains information relating to events such as street closures which may have an impact on parking
- (v) *parking revenue data* contains economic information on parking pricing
- (vi) *fuel price data* contains prices of fuel in the region

A detailed overview of the *SFpark* dataset is shown in table 1.

Each piece of data is geographically referenced by a *location unit*, i.e., street block, street, district or entire city. In the prospect of using the parking data for training estimation models, the different location unit poses a problem. For the traffic and events data sets, it is entire streets; in the parking revenue dataset, it is city districts, while in case of the weather and fuel price datasets, the location reference is valid for the whole city of San Francisco. Hence, there is a need to align the datasets before they can be used together. In the cases when aggregating values associated to street blocks to the street level is performed, the aggregated values leads to a poorer training performance. Even more so when aggregating street blocks to the city level. Therefore, we were forced to continue without traffic-, events-, parking revenue-, fuel price-, and weather data and rely strictly on the occupancy data further in the process.

The occupancy data amasses 1.05 million entries with measurements between 04.2011 and 07.2013. The original file provided by *SFpark* is about 192M large. The *SFpark* data are visualized in fig. 1 using a Leaflet application built as part of this work.

5.3. Mapping Parking Data to OpenStreetMap Layers

We complement the parking data by downloading OpenStreetMap^a data corresponding to the location where the parking data belongs to. OSM data is generally available as shapefiles containing the geometry layers: points, polylines, and polygons. We extract the *points of interest* (POIs), which, among multiple attributes, contain

^a<https://www.openstreetmap.org> The maps used in this article are ©OpenStreetMap contributors.

Table 1. Overview of the properties available in the data used from the SF*park* project.

Parking Occupancy		Traffic	
timestamp	Recorded at full hours	timestamp	Recorded at full hours or in periodic time intervals
parking capacity	The total number of parking spaces at the given location	traffic value	Expressed as average traffic road occupancy, average vehicle count, median speed, or average speed of the traveling cars
parking price	The price of a ticket in dollars at the certain location and the given time	location unit	Given as entire street
parking occupancy	Expressed either as rate (sub-unitary fraction or percent) or in absolute numbers		
location unit	Given as street block		
Events		Weather	
date and time	Expressed as calendar date or time interval within a day	date	Expressed as calendar date
event name class	Given as the name of the event and its class: road closure or rise of parking demand	temperature	Expressed as the maximum value of the day
location unit	Given as entire street	precipitation	Expressed in the quantity of rain or snow for the corresponding time interval
		location unit	Given as entire city
Fuel Price		Parking Revenue	
type of fuel	Provided as gasoline, diesel, etc.	payment type	Expressing the way the driver opted to pay for parking: cash or credit card
price per unit	Provided as the price per gallon	payed amount	Expressed as the amount in US dollars
location unit	Given as entire city	location unit	Given as city district

the *amenity* attribute indicating the public service, facility, or type of building located at this position as it was annotated by the OSM users (cf. fig. 2). The types of the public amenities collected from the POIs are listed in table 2. The polylines layer contains artifacts mostly in linear form, such as streets or foot paths. Polylines are less interesting for our problem and therefore we ignore them. The polygons layer contains artifacts of polygon shapes such as buildings, parks, university campuses, etc. Polygon objects may contain an amenity attribute as well, in practice the authors have found it often empty, however. When the attribute is present, it enables us to compute the area of the amenity and make an inference towards the capacity of the building.

For the San Francisco area corresponding to the SF*park* project, the OSM file containing the above described layers amounts to about 173M. Inside there are

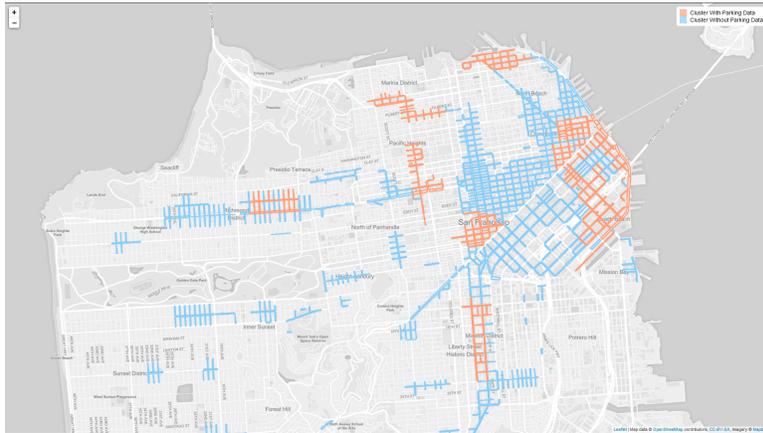


Fig. 1. The blocks accounted in SFpark. The light blue ones are blocks *without* parking data, the light red ones are *with* parking data ²⁶.

30,798 POI entries, out of which 5,462 have a non-empty amenity attribute. The number of polygon entries is 147,881.

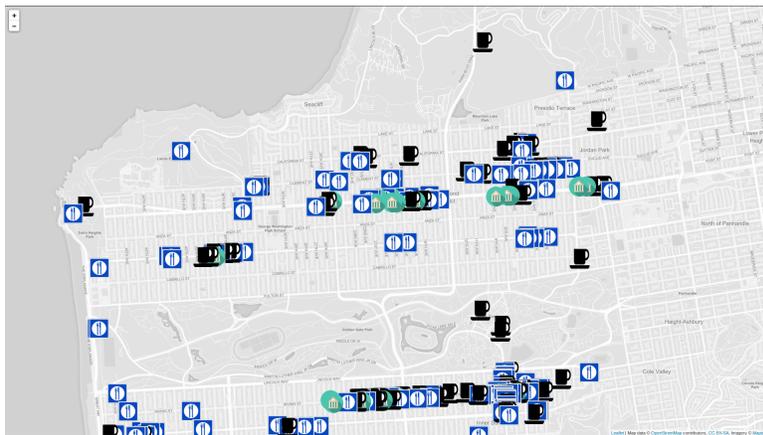


Fig. 2. A map indicating public amenities (cafes, restaurants, banks) found at points of interest in OSM (customly built by the authors).

Furthermore, we collect the visiting duration corresponding to the amenities. We have found that this information is offered by Google Places and FourSquare. The latter data is available via an API, however the access is not free of charge. We used Google Places instead and collected data corresponding to the parking data's location. An example of the service is found within Google Maps for smartphones. It displays typical visiting duration or *time spent* values and popularity of the place for specific time intervals, obtained by Google using a crowdsourcing approach

Table 2. List of all OSM amenities found in the SF *park* blocks.

arts_centre	dojo	marketplace	shelter	conference_centre
bank	embassy	music_rehearsal_place	shop	fire_station
bar	fast_food	music_school	spa	fuel
biergarten	grocery	nightclub	stripclub	parking
bureau_de_change	gym	pet_grooming_shop	studio	place_of_worship
cafe	hookah_lounge	pharmacy	training	social_centre
clinic	ice_cream	police	veterinary	swimming_pool
clothes_store	karaoke	post_office	vintage_and_modern_resale	theatre
community_centre	lan_gaming_centre	pub	bus_station	training
dentist	laundry	restaurant	car_rental	bicycle_parking
doctors	library	salon	childcare	car_wash
brokarage	community_centre	courthouse	fountain	nursing_home
recycling	social_facility	toilets		

that averages the values received from users’ smart phone location (cf. fig. 3). To obtain the *time spent* values, we manually extract information from 470 places in San Francisco, for which a maximum duration of stay was provided (the minimum duration is not always given)^b. The results are shown in table 3 and the *time spent* values are provided in minutes and have been rounded to the nearest integer. We have included only amenities for which at least two stay duration sources were found.



Fig. 3. An example of *time spent* information found on Google Places ²⁷.

In order to combine the parking and city data, both datasets require a common location unit. For parking occupancy it is street blocks that are provided in latitude and longitude for the coordinate reference system EPSG 4326. The POIs inside the OSM data are expressed in the same geometry reference system and therefore a *merge distance* that matches a parking space to a public amenity can be defined. The *merge distance* can be intuitively understood as the radius around a public

^bThis piece of information is not accessible yet via the Google Places API. Google Feature Request: <https://issuetracker.google.com/issues/35827350>

Table 3. All amenities listed with their corresponding mean *time spent* information as collected from Google Places. The *cat* column indicates whether the average visiting time is under half an hour (1), 31 to 90 minutes (2), or more than 1.5 hours (3).

amenity name	mean	stdev	cat	amenity name	mean	stdev	cat
arts_centre	110	37	3	laundry	78	16	2
bank	42	65	2	library	83	13	2
bar	121	38	3	music_school	120	30	3
cafe	76	39	2	nightclub	189	20	3
clinic	100	29	3	pharmacy	25	20	1
clothes_store	41	37	2	post_office	16	2	1
community_centre	119	40	3	pub	135	21	3
dentist	104	35	3	restaurant	135	32	3
doctors	60	42	2	salon	141	53	3
embassy	75	24	2	shelter	90	0	2
fast_food	31	15	2	shop	43	21	2
grocery	20	10	1	spa	161	54	3
gym	100	22	3	stripclub	140	46	3
hookah_lounge	130	17	3	studio	60	0	2
ice_cream	23	7	1	veterinary	67	29	2
karaoke	188	15	3	vintage_modern_resale	38	32	2

amenity. It is defined to represent the parking area that is relevant for a particular public amenity, or, more straightforward, the walking distance from the parked car to e.g., the restaurant, the office, the bank, etc. Based on the above rationale, concrete instances of the merging distance are set to 100m, 200m, and 400m. In section 6.1, we discover which one delivers the best results.

5.4. The Clustering Process

By splitting into city areas, we are making sure that smaller regions lead to more representative parking profiles and therefore parking estimations. As we want an exclusively location-based separation, we may employ K-Means, DBSCAN or OPTICS to cluster the city areas. The distance is calculated between (latitude, longitude) -pairs of location unit coordinates corresponding to one street block. Since having control over the number of clusters is the goal here, we choose to use K-Means, where we provide the number of expected clusters as input. In practice, sklearn’s *kmeans* module is used, specifically the K-Means++ algorithm, which initializes the centroids purposefully distant from each other and therefore achieves faster convergence.

There are two clustering processes executed, one for the city area *with* parking data, another one for the city area *without* parking data. The number of clusters chosen in each area is kept proportional to the number of total street blocks that each area contains. It turns out that for SF*park* data the proportion is approximately 2.6, following the division between the total number of blocks from each group. We have chosen two numbers of clusters to run the evaluation, namely 8 clusters and 16 clusters. The area without parking data will therefore have 20 and 41 clusters, respectively. In the evaluation, we will refer to the number of clusters *with* parking

data as *the number of clusters*.

After running the K-Means clustering process, the Leaflet application map reveals the individual clusters by highlighting them on mouse-over. The clusters *with* parking data will turn dark red, while the clusters *without* parking data will appear in dark blue (cf. fig. 4).

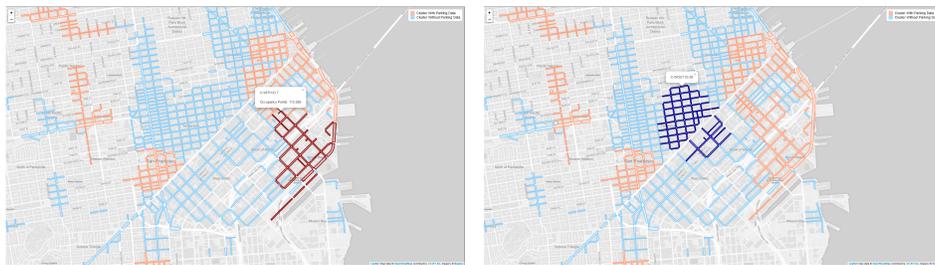


Fig. 4. Highlighted cluster *with* parking data on the left side and a cluster *without* parking data on the right side ²⁶.

5.5. Building Estimation Models

The estimation of parking occupancy is realized using machine learning. We choose to explore this methodology following the solid results machine learning models have delivered for the various smart parking investigations discussed in section 2. A machine learning model \mathcal{M} will be trained for every cluster *with* parking data.

The training data are composed of the parking occupancy data. Specifically, the timestamp, parking capacity and parking price act as predictor variables, while the occupancy rate is the target variable. Before training, data is aggregated across all blocks so that it becomes comparable to other clusters and can be used when training and testing models. The averaging is performed per timestamp, i.e., if multiple blocks have an occupancy record for the same time and block, the occupancy rate will be averaged across both of these. Features such as price and parking capacity per block are averaged as well. See table 4 for an example of this process. This means that the original collection of data records shrinks, which should decrease the training time. In table 5, the shrinking rate is shown for various number of clusters. The resulting parking capacity, parking price, and occupancy values are taken as cluster representatives.

The model training and evaluation is performed in Python via the *scikit-learn* library. During the training phase, we evaluate models such as *decision trees*, *support vector machines*, *multilayer perceptrons* and *boosted trees* using:

- (1) `sklearn.tree.DecisionTreeRegressor`,
- (2) `sklearn.svm.SVR`,
- (3) `sklearn.neural_network.multilayer_perceptron.MLPRegressor`, and

Table 4. Example of aggregating datapoints

block id	timestamp	price rate	total spots	occupied
902	2011-04-02 7:00:00	0	46	58
32800	2011-04-02 7:00:00	0	32	2
33005	2011-04-02 7:00:00	3	36	12
902	2011-04-02 8:00:00	2	46	54
32800	2011-04-02 8:00:00	4	32	5
33005	2011-04-02 8:00:00	3	36	22

timestamp	price rate	total spots	occupied
2011-04-02 7:00:00	1	38	24
2011-04-02 8:00:00	3	38	27

Note: In the first subtable, three distinct blocks belonging to a cluster are transformed into two entries by averaging *price rate*, *total spots* and *occupied* attributes for the two distinct timestamps (second subtable).

Table 5. Number of datapoints aggregated per timestamp vs. all datapoints alongside the shrinking rate for 8, 16 and 32 clusters.

cluster size	aggregated datapoints	all datapoints	shrinking rate
8	9741	128525	12.3
16	8409	73332	8.3
32	6257	29355	4.6

Note: Values have been averaged across clusters.

(4) *xgboost.XGBRegressor* respectively.

As error metric, we use *root mean square error* (RMSE) and perform a five-fold cross-validation. A model will be evaluated on other clusters with parking occupancy data.

Specific details on the actual training of the models can be found in the table table 6. Upon training the models, the clusters can be visualized with the web application as shown in fig. 5 as screenshot. Figure 6 displays the presented table in more detail.

5.6. Building Mathematical Representations for City Areas

Using amenity data from OSM and time-spent information, we build *cluster vectors* and density estimation kernels.

To form the *cluster vectors*, we first divide all amenities into categories $Cat_1, Cat_2, \dots, Cat_n$. The criteria for division will be their average *time spent* values. Each cluster gets represented by an n -dimensional vector, whose components correspond to the amenity categories. The magnitude of component i is equal to the number of amenities of category Cat_i that can be found in that particular cluster. For example, a short duration category of up to 30 minutes, a medium duration between 31 and 90 minutes and a large duration of above 90 minutes stay. Compare fig. 7 for a general representation.

In the case of San Francisco, the categories are based on the *time spent* mean and

Table 6. Overview of the machine learning methods used, their inputs and parameters.

Predictor Variables		Target Variable	
timestamp	split into <i>year</i> as integer, <i>calendar week</i> : [1 - 52] as integer, <i>weekday</i> : [1 - 7] as integer, and <i>hour</i> : [0 - 23] as integer	parking occupancy	: [0 - 100] as floating point number
parking capacity	as floating point number, when aggregated, otherwise as integer		
parking price	as floating point number		
Decision Tree Parameters		SVM Parameters	
Model Selection	randomized search on hyper-parameters for 10 iterations	Model Selection	epsilon-support vector regression model with fixed parameters
<i>min_samples_split</i>	: [2, 3, 4, 5] as the minimum number of samples required to split an internal node	<i>kernel</i>	radial basis function
<i>min_samples_leaf</i>	: [0.03 - 0.1] as the minimum number of samples required to be a leaf of a node	<i>C</i>	penalty parameters equal to 1
<i>max_features</i>	: [0.7, 0.8, 0.9, 1] as the number of features (as fraction from all available features) to consider at each split	<i>gamma</i>	kernel coefficient for the kernel equal to 0.01
<i>criterion</i>	: [mean squared error, mean absolute error] as the function to measure the quality of a split		
<i>min_weight_fraction_leaf</i>	: [0, 0.1, 0.2] as the minimum weighted fraction of the total sum of weights from all the input samples required to be a leaf node		
MLP Parameters		XGB Parameters	
Model Selection	multi-layer perceptron regressor with fixed parameters	Model Selection	exhaustive search over specified parameter values for an extreme gradient boosting model
<i>hidden_layer_sizes</i>	(7, 11) as tuple representing the number of neurons in each hidden layer	<i>max_depth</i>	: [2, 3] as maximum tree depth for base learners
<i>max_iterations</i>	500 as the number of iterations the solver iterates until convergence	<i>n_estimators</i>	: [50, 100] as the number of trees to fit
		learning_rate	: [0.1, 0.25] as boosting learning rate (η)

are split in three categories. The assigned partitions for every amenity are shown in table 3.

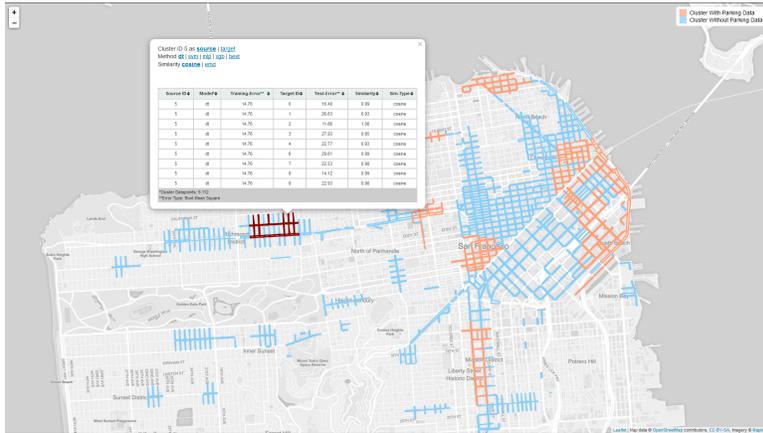


Fig. 5. Selected cluster *with* parking data and the pop-up table in the Leaflet application ²⁶.

Source ID	Model	Training-Error**	Target ID	Test-Error**	Similarity	Sim-Type
5	dt	14.76	2	11.66	1.00	cosine
5	dt	14.76	8	14.12	0.99	cosine
5	dt	14.76	0	16.40	0.99	cosine
5	dt	14.76	7	22.53	0.98	cosine
5	dt	14.76	4	22.77	0.93	cosine
5	dt	14.76	9	22.93	0.98	cosine
5	dt	14.76	1	26.63	0.93	cosine
5	dt	14.76	3	27.93	0.85	cosine
5	dt	14.76	6	29.61	0.99	cosine

*Cluster Datapoints: 9.112
**Error Type: Root Mean Square

Fig. 6. The pop-up table for the Leaflet application view of fig. 5 ²⁶.

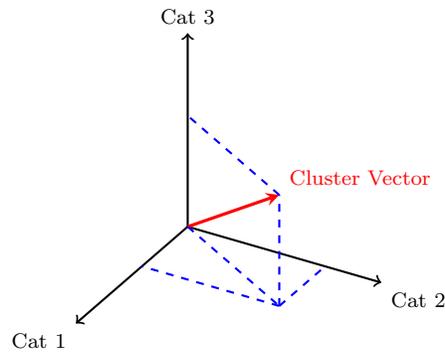


Fig. 7. An example of a *cluster vector* representing amenity *time spent* information composed from for three categories.

- (i) < 30 min,
- (ii) 30 to 90 minutes, and
- (iii) > 90 minutes.

A *cluster Gaussian* is a *kernel density estimation* among the probability distribution of an amenity’s *time spent* value. We use a Gaussian kernel to express the probability distribution, hence the name. To construct a *cluster Gaussian*, we first collect the mean and standard deviation of the individual amenities’ *time spent* values and then we construct its corresponding Gaussian curve. Multiple amenities, each appearing multiple times, will result in a curve that is the linear combination of the individual representations of the amenities as normal distribution curves. Compare fig. 8 for a visualization of the summing process.

$$emd(\mathcal{C}_i) = \sum_{j=1}^{|\text{amenities}|} K_{ij} \times A_j \tag{1}$$

$\forall i \in \{1, ..|\text{clusters}|\}$ and $\forall j \in \{1, ..|\text{amenities}|\}$
 where A_j is an amenity that appears K_{ij} times in the cluster \mathcal{C}_i .

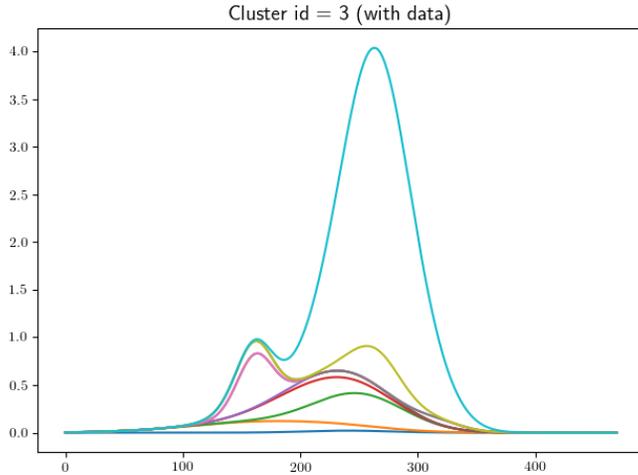


Fig. 8. The iterative summing of Gaussian curves representing the amenity *time spent* information resulting in a *cluster Gaussian* as the outer hull (graph obtained inside the Python application using the matplotlib package).

In practice, the computation is discretized using bins. A bin represents a unit on the X axis, the same on which the *time spent* value is expressed. We will take a number of bins equal to the maximum amenity mean and buffer them with $3 \times$ the largest standard deviation, as it is known that within $3 \times$ standard deviation

on both sides of the mean over 99% of the Gaussian sum is covered. Moreover, an offset on the X -axis equal to $3\times$ the maximum standard deviation is used. This way, we are sure the landscape of summed Gaussians will easily fit into the number of bins. The Gaussians are computed using Python’s Optimal Transport package (OT), specifically using the `ot.datasets.make_1D_gauss` module, by providing the number of bins, mean and standard deviation values.

5.7. Computing Similarities Between City Areas

The *cosine similarity* between two vectors is defined as the cosine of the angle between the two vectors. The *cosine similarity* implementation uses the direct mathematical formula by plugging in the magnitudes of the respective vector components.

$$\cos(\theta) = \frac{A \cdot B}{\|A\|_2 \|B\|_2} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}} \quad (2)$$

where A_i and B_i are the components of vector A and respectively.

The *earth mover’s distance* (emd) is a measure used in statistics that roughly expresses the difference between position and magnitude of two curves. It is best explained by regarding the curves as the hull of earth piles. For two separate earth piles, emd computes the minimum effort of rearranging a pile so that the shape of the other pile is obtained. Moving P particles over a distance D is equal to the effort $P \times D$. A prerequisite for this operation is that the two piles need to contain the same quantity of earth.

More rigorously, emd is better known in mathematics as Wasserstein Metric. Given two normal distributions $\mu_1 = \mathcal{N}(m_1, C_1)$ and $\mu_2 = \mathcal{N}(m_2, C_2)$, where m_1 and $m_2 \in \mathbb{R}^n$ are their respective expected values and C_1 and $C_2 \in \mathbb{R}^{n \times n}$, their 2-Wasserstein distance between μ_1 and μ_2 is:

$$W_2(\mu_1, \mu_2)^2 = \|m_1 - m_2\|_2^2 + \text{trace}(C_1 + C_2 - 2(C_2^{1/2} C_1 C_2^{1/2})^{1/2}) \quad (3)$$

Notice that emd is applicable only when the sum under both Gaussian curves is equal. Therefore, all *cluster Gaussians* will get normalized before emd is computed. In practice, after summing up the respective Gaussians and arriving at the *cluster Gaussians*, the `scipy.stats` module is used to compute emd by means of the `wasserstein_distance` method. Provided here are the bin ranges on the X axis and their Y axis values corresponding to the computed *cluster Gaussians*.

5.8. Computing Parking Occupancy Estimations

Once all models \mathcal{M} have been built for the clusters *with* parking data, making estimations on parking occupancy in these areas is straightforward. The input data fed to the model is manufactured from averaging the predicted variables from the

areas with parking data. However, we want to apply these models on the clusters that are missing parking data. We derive the *estimation interval* for cluster \mathcal{C}_{wout}^j based on the model of cluster \mathcal{C}_{with}^i as follows.

For *cosine similarity*:

$$E(\mathcal{C}_{with}^i, \mathcal{C}_{wout}^j) = [\mathcal{M}(\mathcal{C}_{with}^i) - (1 - sim_{ij}), \mathcal{M}(\mathcal{C}_{with}^i) + (1 - sim_{ij})] \quad (4)$$

$$\text{where } sim_{ij} = sim(\mathcal{C}_{with}^i, \mathcal{C}_{wout}^j) \in [0, 1]$$

For *emd*:

$$E(\mathcal{C}_{with}^i, \mathcal{C}_{wout}^j) = [\mathcal{M}(\mathcal{C}_{with}^i) - emd_{ij}, \mathcal{M}(\mathcal{C}_{with}^i) + emd_{ij}] \quad (5)$$

$$\text{where } emd_{ij} = emd(\mathcal{C}_{with}^i, \mathcal{C}_{wout}^j) \in [0, 1]$$

$$\forall i \in \{0, \dots, |\mathcal{C}_{with}| - 1\} \text{ and } \forall j \in \{0, \dots, |\mathcal{C}_{wout}| - 1\}$$

The result is an *estimation interval* that *stretches* the punctual estimation into an interval depending on the similarity value. The lower the similarity value is, the larger the length of the resulting estimation interval will be.

Furthermore, we define an *estimation intersection interval*, whose purpose is to narrow down the computed estimation interval. An estimation intersection interval for the clusters \mathcal{C}_{with}^i and \mathcal{C}_{wout}^j is computed by intersecting the *estimation intervals* that have a better similarity among the clusters with data $\mathcal{C}_{with}^0, \dots, \mathcal{C}_{with}^{i-1}$ and the same cluster without data \mathcal{C}_{wout}^j .

$$EII(\mathcal{C}_{with}^i, \mathcal{C}_{wout}^j) = \bigcap_{k=0}^{i-1} EI(\mathcal{C}_{with}^k, \mathcal{C}_{wout}^j) \quad (6)$$

where

$$sim(\mathcal{C}_{with}^k, \mathcal{C}_{wout}^j) < sim(\mathcal{C}_{with}^i, \mathcal{C}_{wout}^j), k \in \{0, \dots, i-1\} \text{ for emd} \quad (7)$$

$$sim(\mathcal{C}_{with}^k, \mathcal{C}_{wout}^j) > sim(\mathcal{C}_{with}^i, \mathcal{C}_{wout}^j), k \in \{0, \dots, i-1\} \text{ for cosine} \quad (8)$$

$$\forall i \in \{0, \dots, |\mathcal{C}_{with}| - 1\} \text{ and } \forall j \in \{0, \dots, |\mathcal{C}_{wout}| - 1\}$$

6. Evaluation

We evaluate various pieces of the system that has been presented. Firstly, we establish the machine learning method that achieves best results on average across clusters. Afterwards, we compare the cluster models’ test errors with the independently-computed similarity values between clusters. More specifically, a *source* cluster’s model will be tested on a *target* cluster and the error is correlated to the similarity between the *source* and the *target* clusters. Both *cosine* and *emd* functions will be used. The correlations will be expressed as Pearson- and Spearman’s rank coefficients. Afterwards, we take a look at the results of applying the models to clusters *without parking data* and visualize the results.

Furthermore, some alternative methods are investigated. Firstly, we look at the model test errors and correlation results by skipping the aggregating step, i.e., instead of averaging the datapoints over timestamp per cluster, we build the cluster models using the entire occupancy data directly. Secondly, we use the *amenity area* as the basis for the similarity functions in calculating the correlations between model test errors and similarity values. Finally, we question whether the similarity function approach is the most efficient and transfer its purpose to the machine learning phrase. The model will receive absolute *cosine vectors* and *cluster Gaussian* values as additional features and its model test error and correlation values will be compared to the ones from the original approach.

6.1. Best model method

Of the models were trained using the four methods (*decision trees*, *support vector machines*, *multilayer perceptrons*, and *gradient boosted trees*), table 7 shows the distribution of best machine learning methods in case of 8 and 16 clusters. The values were obtained by summing up the number of times a method produced the smallest estimation error, i.e., RMSE, among the four methods for all combinations of clusters with parking data ($\mathcal{C}_{source}, \mathcal{C}_{target}$). Extreme gradient boosting claims the first spot in both cases. In further experiments in the evaluation we shall only report on models trained using extreme gradient boosting.

Table 7. The fraction of best models among *decision trees*, *support vector machines*, *multilayer perceptrons* and *extreme gradient boosting* measured as RMSE when applied on all pairs of clusters.

	dt	svm	mlp	xgb
8 clusters	24.6%	17.5%	12.3%	45.6%
16 clusters	14.6%	13.8%	13.8%	57.9%

6.2. Similarity Values vs. Estimation Errors

The central goal of this work is estimate the occupancy of clusters where no parking data is available by using model predictions and pair-wise cluster similarity values. For the purpose of evaluating the models, we shall use clusters with parking data that were left out from the training dataset as the application target of the prediction models, while the computed similarity values will serve to confirm the prediction errors. The higher the absolute correlation between the cluster similarity values and the model test errors, the better the accuracy of the cluster similarity.

Concretely, for every pair of clusters $(\mathcal{C}_{source}, \mathcal{C}_{target})$, a model $\mathcal{M}_{\mathcal{C}_{source}}$ is trained on \mathcal{C}_{source} and its test error $(\mathcal{M}_{\mathcal{C}_{source}}(\mathcal{C}_{target}))$ shall be correlated with the cosine and emd similarity values between \mathcal{C}_{source} and \mathcal{C}_{target} . Here, we use two correlation coefficients: the Pearson correlation coefficient and Spearman’s rank correlation coefficient, which results in four correlation measures:

- (i) cosine (Pearson) correlation
- (ii) cosine (Spearman’s) rank correlation
- (iii) emd (Pearson) correlation
- (iv) emd (Spearman’s) rank correlation.

The evaluation was performed for configurations of 8 and 16 clusters respectively. Additionally, we varied the *merge distance* between 100m, 200m, and 400m to see how the correlation behaves. In table 8 the final results are shown. A correlation result is averaged across all clusters. Due to their mathematical meaning, the *cosine similarity* values below zero express a positive correlation, whereas emd values above zero express a positive correlation.

We notice that the *cosine similarity* achieves better results than emd for the same testing configuration, peaking for 8 clusters and 100m *merge distance*. Its average Pearson coefficient is -0.55 , while the mean Spearman rank coefficient is -0.49 . emd positively correlates the most for the same testing configuration (8 clusters, 100m), when the average Pearson coefficient is at 0.28 and Spearman’s rank coefficient equals 0.23 . There is a clear descending trend in correlations, as the *merge distance* increases. Also, the results for 8 clusters are superior to the ones when the city is split in 16 clusters. Further in the evaluation runs we fix the *merge distance* to 100m.

6.3. Estimations for clusters without parking data

We apply the models trained on SF *park* data on clusters *without* parking data. The testing data records are composed of values equal to the averages of the respective data types in the clusters *with* parking data. This is the case for *parking price* and *parking capacity*. One piece of data that still needs to be provided so that the estimation is computed is the timestamp. For convenience, we choose the next day relative to when we ran the experiment and 8 times spread throughout the day. A sample of the input values fed to the model is show in table 9, while the results of the

Table 8. Correlations between similarity values and model estimations errors for pairs of clusters *with* parking data (C_{source}, C_{target}).

	8 clusters			
merge distance	cosine	rank_cosine	emd	rank_emd
100m	-0.55	-0.49	0.28	0.23
200m	-0.34	-0.30	0.26	0.23
400m	-0.23	-0.08	0.25	0.27

	16 clusters			
merge distance	cosine	rank_cosine	emd	rank_emd
100m	-0.20	-0.17	0.10	0.11
200m	-0.13	-0.11	0.02	0.02
400m	-0.17	-0.17	0.08	0.11

Note: For *cosine similarity* the values show a negative correlation tendency, while for the correlation based on emd similarity expresses a positive correlation tendency. The correlations are measured using Pearson coefficient and Spearman's rank coefficient.

estimation which includes the similarity values is visualized in the web application as in fig. 9.

C-WOUT ID 19 Timepoint: 2017-11-04 09:00:00 ▾
 Similarity [cosine](#) | [emd](#)

C-WTIIH ID	Similarity	Sim-Type	Estimation	Interval	Intersection
6	1.00	cosine	45.27	[45.19 - 45.35]	[45.19 - 45.35]
7	1.00	cosine	45.15	[45.01 - 45.29]	[45.19 - 45.29]
8	0.99	cosine	57.73	[56.99 - 58.46]	empty
9	0.99	cosine	50.84	[49.84 - 51.83]	empty
5	0.99	cosine	66.71	[65.49 - 67.93]	empty
0	0.98	cosine	52.52	[50.58 - 54.46]	empty
1	0.98	cosine	36.32	[34.10 - 38.54]	empty
2	0.98	cosine	55.45	[53.03 - 57.87]	empty
4	0.97	cosine	58.11	[55.53 - 60.68]	empty
3	0.92	cosine	45.26	[37.51 - 53.01]	empty

Fig. 9. The pop-up table of a cluster without data. Notice the drop-down list from which the time can be selected. There are 8 times to select, equally spaced throughout the day that followed our experiment. The similarity values are rounded off. The results are expressed in intervals of occupancy. The successive intersection of intervals succeeds only for the first two values, afterwards it is empty ²⁶.

Table 9. Input data for models when predicting occupancy for city areas without parking data. The date has been chosen arbitrarily with several times equally spaced throughout the day. The price rate and total spots values are equal to the approximate averages of the entire parking dataset.

Date	Time	Price Rate	Total Spots
2017-11-04	00:00	1.0	20
	03:00	1.0	20
	06:00	1.0	20
	09:00	1.0	20
	12:00	1.0	20
	15:00	1.0	20
	18:00	1.0	20
	21:00	1.0	20

6.4. Models built on all occupancy datapoints

Up to now, the evaluation involved models trained and tested- on aggregated datapoints. We ask ourselves, however, whether a model trained on all datapoints performs better than when tested on an aggregated cluster. Or whether an aggregate model delivers better results on an all-datapoints cluster than a model trained on all datapoints? Here, we experiment these combinations by training models on both all- and aggregate datapoints and apply them on both types of aggregation forms.

See table 10 for an overview of test errors for 8 and 16 clusters. One observes that the errors from models applied on aggregated datapoints are about 5 unit points (5%) smaller than the errors from models applied on all datapoints. This is naturally accounted for by the smaller spread of occupancy values that the aggregation brings with itself. As far as source models are concerned, there is no significant difference between the aggregate and all datapoint models, i.e., the margin is under 1%. Regarding the number of clusters, the values for 8-cluster models are slightly better than 16-cluster models for aggregated datapoints, but lose when the testing bed is equal to all datapoints.

In table 11 the resulting correlations of testing errors with cosine- and emd similarity values are listed for models of 8 and 16 cluster configurations. It follows that the cosine- and rank cosine correlation values are on average closer to -1 in the 8-cluster case. The same applies for emd- and rank emd correlation values, which are closer to 1 in this case. The correlation values for 16 clusters are obviously weaker than for 8 clusters. In both sections, the aggregate datapoints target is superior to the all datapoints target.

6.5. Amenity area as similarity basis

We considered *time spent* information to complement the amenity information as the basis for creating the mathematical representation. However, there are other

Table 10. Test error for ML models build alternatively with all- and aggregated datapoints.

cluster size	datapoints source	datapoints target	test error
8	aggregate	aggregate	20.19
8	all	aggregate	21.37
8	aggregate	all	26.25
8	all	all	26.68
16	aggregate	aggregate	20.47
16	all	aggregate	21.32
16	aggregate	all	25.97
16	all	all	26.52

Note: Test errors for the same number of clusters can be accurately compared when the *datapoints target* is the same. All models above were build using extreme gradient boosting.

Table 11. Resulting correlation values for ML models built using all- and aggregated datapoints.

cluster size	datapoints source	datapoints target	cosine	rank_cosine	emd	rank_emd
8	aggregate	aggregate	-0.53	-0.52	0.30	0.17
8	all	aggregate	-0.53	-0.43	0.37	0.27
8	aggregate	all	-0.35	-0.36	0.20	0.14
8	all	all	-0.41	-0.43	0.34	0.25
16	aggregate	aggregate	-0.16	-0.11	0.10	0.05
16	all	aggregate	-0.18	-0.17	0.22	0.17
16	aggregate	all	-0.09	-0.06	0.08	0.00
16	all	all	-0.10	-0.11	0.17	0.08

factors that affect the parking demand towards an amenity. Obviously, one of them is the amount of people visiting the amenity. As we do not have data about this aspect, we use the area of the amenity as a proxy, assuming that larger places would have more visitors. OpenStreetMap provides a polygon layer for a certain geographic bounding box, which contains information across all the surfaces in that region. See fig. 10 for a visualization of the polygons and their areas in OSM.

To extract this information we investigated two options.

- (i) **Polygon containing POI.** Matching the amenities' POIs with the containing OSM polygon and then computing the polygon areas per amenity was the option tried first. This has several drawbacks. The relation POI : polygon is in practice by no means 1 : 1. Many cases arose where multiple POIs were contained by the same polygon, in which situation the area was split between them; a POI might also be on the edge of several polygons, in which case we have to either (arbitrarily) assign it to the first polygon or to all. The deciding factor against this approach was, however, the fact that the *coefficient of variation*, i.e., the ratio between the standard deviation and the mean of the sample, is larger than 1, i.e., 2.1 to be precise.
- (ii) **Amenity attribute in polygon layer.** The other option was to use the *amenity* attribute from the polygon layer of the region. We could avoid the



Fig. 10. OSM screenshot emphasizing polygons as buildings and the amenities that are housed by them ²⁸.

cumbersome matching by leveraging solely the polygon layer and calculating the amenity area mean and its standard deviation. The results are listed in table 12. On top of that, the coefficient of variation is 0.9 in this case, significantly lower than before.

Note that we have reduced the values in the table by a factor of 20, as it turned out that the actual mean and standard deviation were large enough to make the emd Gaussian computation extremely slow. As the standard deviation is linear with regard to the mean, both mean and standard deviation values were reduced conveniently.

By applying the procedure using the amenity area as similarity basis, we obtain the correlations values listed in table 13 for 8- and 16 cluster configurations. The superiority of the correlation values using *time spent* value is observed both for cosine- and rank cosine correlations, which are closer to -1, and for the emd- and rank emd correlations respectively, which are nearer to 1. All correlation values for 16 clusters are, however, relatively weak in absolute measures.

6.6. Extended prediction models

An alternative to building mathematical representation of city areas and applying similarity functions is to let the machine learning model find out the similarities by itself. One can choose to add the city data as further training information for clusters. The purpose is to produce better predictions by leveraging unknown patterns in the city data. Hence, features representing the *cosine* and *emd* functions are added to the model, as follows:

- (i) k features corresponding to the k categories the amenities are split in, i.e. the magnitudes of the vectors for each category;

Table 12. Amenity area values gathered and averaged from OMS polygon layer for the SF*park* region.

amenity name	mean	stdev	cat	amenity name	mean	stdev	cat
arts_centre	68	60	2	bank	39	20	2
bar	19	8	1	bicycle_parking	8	7	1
biergarten	11	12	1	brokerage	39	9	2
bus_station	588	737	3	cafe	17	10	1
car_rental	70	43	2	car_wash	43	48	2
childcare	101	130	3	cinema	75	43	2
clinic	61	32	2	community_centre	52	74	2
conference_centre	401	519	3	courthouse	459	201	3
dentist	17	12	1	doctors	324	568	3
embassy	68	38	2	fast_food	25	24	1
fire_station	52	27	2	fountain	24	22	1
fuel	25	27	1	library	102	124	3
marketplace	325	228	3	music_rehearsal_place	33	15	1
nightclub	32	9	1	nursing_home	97	47	2
parking	182	309	3	pharmacy	65	38	2
place_of_worship	60	62	2	police	137	124	3
post_office	39	11	2	pub	25	25	1
public_building	280	236	3	recycling	28	20	1
restaurant	22	16	1	school	740	1280	3
social_centre	30	21	1	social_facility	356	801	3
stripclub	50	10	2	studio	268	307	3
swimming_pool	16	9	1	swingerclub	27	4	1
theatre	174	191	3	toilets	7	5	1
training	72	94	2	veterinary	21	7	1

Note: The mean and standard deviation values were reduced by a 20x factor. The categories for cosine vectors are 0 - 35, 36 - 100, 100+.

Table 13. The correlation results computed using similarity values based on *amenity area*.

cluster size	amenity type	datapoints (train/test)	cosine	rank_cosine	emd	rank_emd
8	time spent	agg/agg	-0.53	-0.52	0.30	0.17
8	area	agg/agg	0.05	0.11	0.14	0.22
16	time spent	agg/agg	-0.16	-0.11	0.10	0.05
16	area	agg/agg	-0.09	-0.03	0.09	0.07

Note: All models above were build, trained, and tested on aggregated datapoints.

- (ii) a feature corresponding to the emd Gaussian value for that cluster, loosely interpreted as the "total accumulated time spent value" for that cluster, in case of *time spent*, or the "total accumulated area" among all amenities in that cluster, for the amenity area; mathematically it is expressed in eq. (9)

$$feature(emd) = x \cdot f(x) \quad (9)$$

$$\text{where } x = \begin{cases} \text{duration in minutes,} & \text{if sim} = \textit{time spent} \\ \text{area in square meters,} & \text{if sim} = \textit{amenity area.} \end{cases}$$

and f is the constructed emd Gaussian function that equals the number of amenities in the cluster for any value x .

We subsequently build *extended machine learning models* that additionally contain the features above. Since these features are identical for all datapoints in a cluster, the model needs to be trained on multiple clusters. Therefore, models on $n - 1$ out of n clusters will be build and tested on the remaining cluster. These *all-but-one* or *total models* will be constructed for all n combinations of $n - 1$ clusters and their averaged test errors will compared to total models that contain the normal features. To build the models, the same methods and parameters are used as described before.

The resulting test errors are shown in table 14 for various number of clusters. In three out of four cases, the addition of the two features did not help with finding better parking occupancy values for the test clusters, at least using the gradient boost model we used in our experiments. For the 16-cluster case, the total model returns a superior result. In the one case, the simple total models achieved a better prediction performance.

Table 14. Total models extended with cosine and emd features compared to total models with the previous feature set.

cluster size	model	test error average
4	xgb	14.92
4	xgb total	16.18
8	xgb	18.12
8	xgb total	19.58
16	xgb	18.19
16	xgb total	18.09
32	xgb	18.65
32	xgb total	20.38

7. Further Possible Variations

To further investigate parking occupancy prediction given the assumptions in this work, there are several improvements or alternative approaches that can be experimented with:

- (i) **Use parking data from other locations.** In the present work, several pieces of data could not be integrated because of merging issues, i.e., the location units among multiple datasets did not coincide. Traffic, events, weather, etc. might improve estimation results and hence the final estimations for clusters without parking data. Other sources for parking occupancy data can be found for the cities of Cologne ²⁹, Zurich ³⁰, Santa Monica ³¹. In Germany, Deutsche Bahn provides an API to obtain data from parking around train stations ³². Data pertaining to street occupancy is, however, hard to find. At the time of writing, open data portals mostly provide the location of parking lots, parking meters, parking prices and opening times, if applicable.

- (ii) **Gather better *time spent* information** Accessing data on the time durations people spend in various amenities is restricted by Google by not providing an API for it. Other social media services, such as Foursquare offers information on how busy an amenity is by means of the number of check-ins that people send from their phones^c ^d. The API is however subject to costs and we have not tried it. By gathering more information automatically than we managed to collect manually, the evaluating similarity values based on the *time spent* information will likely increase in accuracy and so will the resulting correlation values.
- (iii) **Use more OSM (meta)data.** The mathematical representations in the present work are relying on public amenities available from OpenStreetMap. OSM has great potential as a collaborative map service but it currently still lacks many pieces of information that could be useful (in comparison with Google's *time spent* information, for example). Data such as opening hours, if it would be widely available, would be interesting to include in the mathematical representations, which would then take into account the number of public amenities that are available at a certain point in time. Overall, more and finer city data, together with an appropriate representation and similarity function could eventually improve the occupancy estimations for clusters without parking data.
- (iv) **Experiment with other clustering methods** Our approach is dependent on the fact that the computed K-Means clusters are of approximately the same size. Ideally, this would mean that the number of amenities of a certain *time spent* value is virtually equal between clusters. But since for both *cosine similarity* and *emd* it is not the absolute number of amenities for a certain *time spent* value that matters, instead the relative number between certain *time spent* values, it may be that clusters of clearly opposite sizes are very similar. All the more is the reason to experiment with other clustering methods, such as DBSCAN and OPTICS. Both focus on detecting clusters based on density of neighborhoods, which may translate in practice into separating regions in the city that are more sparse, e.g., large office areas, from regions that are very dense, e.g., residential or old town areas.
- (v) **Apply semi-supervised machine learning.** Another relevant machine learning approach in this case is based on organizing the city areas as an undirected graph. The vertices represent the clusters with their respective occupancy data, while the edges between them are assigned similarity values. Initially, only a part of the vertices have the occupancy value known, i.e., the clusters with parking data, while the rest has undetermined occupancy, i.e., the clusters without parking data. At each step, the value for a vertex whose value is undetermined is being computed by considering the occupancies of the linked vertices and

^c<https://developer.foursquare.com/docs/api/endpoints>^d<https://developer.foursquare.com/docs/api/venues/details>

their corresponding similarity values.

- (vi) **Accounting for completely filled parking spaces in the vicinity.** Currently, we only take information from nearby amenities into account for the model. However, one could imagine that when nearby areas have completely full parking areas, then cars will ‘spill over’ to the area under investigation. A model having to account for this type of factors would, however, be a lot more complicated to express.

8. Conclusion

We presented our work on approximating street parking occupancy in cities. Given the fact that parking sensors cannot capture all parking areas and assuming that parking demand can be modeled using urban features, we proposed an alternative solution to the ones previously developed for this problem. The data used was composed of parking data from the SFpark project, that investigated parking pricing in San Francisco, and OpenStreetMap. The two data sources were merged and afterwards split into small city areas that separated the regions that contained parking data and other regions that did not contain parking data. We built mathematical representations for amenities, any form of facility or building in a city, and their *time spent* information, visiting duration that have been gathered using smartphones. This led to computing similarity values between city areas by applying functions such as *cosine similarity* and *earth mover’s distance* to the mathematical objects. The K-Means algorithm was used to cluster the city areas, while four methods were employed to train models for the clusters: *decision trees*, *support vector machines*, *multilayer perceptrons* and *extreme gradient boosting*. The occupancy estimations for clusters *without* parking data were defined in terms of model estimations from clusters *with* parking data and the corresponding cluster similarity values. The estimations are expressed as intervals which extend the model prediction values by the magnitude of the similarity values.

The data source for our work, the SFpark project, gathered parking data for more than 2 years starting in 2011 in San Francisco and now offers it for free usage. The city data was collected from OpenStreetMap as amenity information, and from Google Places as stay duration values. Both sources are open and free of charge. Over 30 types of public amenities were found in the San Francisco blocks and data from over 470 Google Places sources was collected.

Following our experiments, the best machine learning model for our problem setting turned out to be extreme gradient boosting. We used the clusters *with* parking data for the evaluation of the similarity values and calculated correlation coefficients between the similarity values and the estimation errors, using both absolute values and ranks. The best correlation were reached for the 100m *merge distance* for 8 clusters. In the same configuration, both *cosine similarity* and *emd* reached their best results from all the experimented configurations.

We further investigated the test error and correlation values when models were

built and tested on all datapoints instead of aggregated datapoints. The aggregation factor has been shown to play a role in finding the best estimations.

Overall, we touched upon several aspects that influence parking occupancy and provide a ground for further experimentation, which can further improve parking occupancy estimations when transferring models from monitored to unmonitored regions.

References

1. M. N. Smith, The number of cars worldwide is set to double by 2040 (2016), <https://www.weforum.org/agenda/2016/04/the-number-of-cars-worldwide-is-set-to-double-by-2040>.
2. D. Shoup, Free parking or free markets (2001), <https://www.accessmagazine.org/spring-2011/free-parking-free-markets/>.
3. M. K. Zavitsas, I. Kaparias, Transport problems in cities (2012), https://trimis.ec.europa.eu/sites/default/files/project/documents/20120402_173932_45110_D%201.1%20-%20Transport%20problems%20in%20cities%20-%20v3.pdf.
4. INRIX, Searching for parking costs americans 73 billion us dollars a year (2017), <http://inrix.com/press-releases/parking-pain-us/>.
5. A. Ionita, A. Pomp, M. Cochez, T. Meisen and S. Decker, Where to park?: Predicting free parking spots in unmonitored city areas, in *Proceedings of the 8th International Conference on Web Intelligence, Mining and Semantics WIMS '18*, (ACM, New York, NY, USA, 2018), pp. 22:1–22:12.
6. A. Ionita, Extending estimation of parking occupancy to untracked city areas using city background information, master's thesis (December 2017).
7. T. Lin, H. Rivano and F. Le Mouël, A survey of smart parking solutions, *IEEE Transactions on Intelligent Transportation Systems* **18** (Dec 2017) 3229–3253.
8. T. S. Lin, *Smart Parking: Network, Infrastructure and Urban Service*, theses, INSA Lyon, (INSA Lyon, December 2015).
9. B. Xu, O. Wolfson, J. Yang, L. Stenneth, S. Y. Philip and P. C. Nelson, Real-time street parking availability estimation, in *2013 IEEE 14th International Conference on Mobile Data Management* **1**, IEEE, (IEEE, June 2013), pp. 16–25.
10. Z. Chen, J. C. Xia and B. Irawan, Development of fuzzy logic forecast models for location-based parking finding services, *Mathematical Problems in Engineering* **2013** (2013).
11. A. Nandugudi, T. Ki, C. Nuessle and G. Challen, PocketParker: Pocketsourcing parking lot availability, in *Proceedings of the 2014 ACM International Joint Conference on Pervasive and Ubiquitous Computing UbiComp '14*, (ACM, New York, NY, USA, 2014), pp. 963–973.
12. A. Koster, A. Oliveira, O. Volpato, V. Delvequio and F. Koch, Recognition and recommendation of parking places, in *Advances in Artificial Intelligence – IBERAMIA 2014*, eds. A. L. Bazzan and K. Pichara Springer, (Springer International Publishing, Cham, 2014), pp. 675–685.
13. T. Rajabioun, B. Foster and P. Ioannou, Intelligent parking assist, in *21st Mediterranean Conference on Control and Automation IEEE*, (IEEE, June 2013), pp. 1156–1161.
14. T. Rajabioun and P. A. Ioannou, On-street and off-street parking availability prediction using multivariate spatiotemporal models, *IEEE Transactions on Intelligent Transportation Systems* **16**(5) (2015) 2913–2924.
15. T. Tiedemann, T. Voegelé, M. Krell, J. Metzen and F. Kirchner, Concept of a data

- thread based parking space occupancy prediction in a Berlin pilot region **29**, (AAAI, 2015).
16. F. Richter, S. Di Martino and D. C. Mattfeld, Temporal and spatial clustering for a parking prediction service, in *2014 IEEE 26th International Conference on Tools with Artificial Intelligence* (IEEE, Nov 2014), pp. 278–282.
 17. San Francisco Municipal Transportation Agency, SFpark - open data (2011–2013), <http://sfpark.org/how-it-works/open-data-page/>.
 18. A. O. Kotb, Y.-C. Shen, X. Zhu and Y. Huang, iParker – a new smart car-parking system based on dynamic resource allocation and pricing, *IEEE Transactions on Intelligent Transportation Systems* **17**(9) (2016) 2637–2647.
 19. J.-H. Shin and H.-B. Jun, A study on smart parking guidance algorithm, *Transportation Research Part C: Emerging Technologies* **44** (2014) 299–317.
 20. S. Mathur, T. Jin, N. Kasturirangan, J. Chandrasekaran, W. Xue, M. Gruteser and W. Trappe, ParkNet: Drive-by sensing of road-side parking statistics, in *Proceedings of the 8th International Conference on Mobile Systems, Applications, and Services MobiSys '10*, (ACM, New York, NY, USA, 2010), pp. 123–136.
 21. M. Caliskan, A. Barthels, B. Scheuermann and M. Mauve, Predicting parking lot occupancy in vehicular Ad Hoc networks, in *2007 IEEE 65th Vehicular Technology Conference - VTC2007-Spring* IEEE, (IEEE, April 2007), pp. 277–281.
 22. A. Klappenecker, H. Lee and J. L. Welch, Finding available parking spaces made easy, *Ad Hoc Networks* **12** (2014) 243–249.
 23. P. Szczurek, B. Xu, O. Wolfson, J. Lin and N. Rishe, Learning the relevance of parking information in vanets, in *Proceedings of the Seventh ACM International Workshop on Vehicular InterNetworking VANET '10*, (ACM, New York, NY, USA, 2010), pp. 81–82.
 24. San Francisco Municipal Transportation Agency, SFpark (2011–2013), <http://sfpark.org>.
 25. San Francisco Municipal Transportation Agency, SFpark – pilot project evaluation summary (2011–2013), http://sfpark.org/wp-content/uploads/2014/06/SFpark_Eval_Summary_2014.pdf.
 26. A. Ionita, Parking prediction web application (2017), <https://datalab.rwth-aachen.de/parking-prediction/>.
 27. Google, Google my business (2017), <https://www.google.com/business/>.
 28. O. Community, OpenStreetMap (2004), <https://www.openstreetmap.org/>.
 29. Stadt Köln and DKAN, Öffene Daten Köln (2015), https://www.offenedaten-koeln.de/dataset/taxonomy/term/52/field_tags/Transport%20und%20Verkehr-52?query=park&sorting=changed%7CDESC.
 30. Stadt Zürich, Stadt Zürich - open data (2015), <https://data.stadt-zuerich.ch/dataset/parkleitsystem>.
 31. City of Santa Monica, Santa Monica – open data (2014), <https://data.smgov.net/Transportation/Parking-Lot-Counts/ng8m-khuz>.
 32. DB BahnPark GmbH, DB - Parkplätze API (2016), <http://data.deutschebahn.com/dataset/api-parkplatz>.