

Hochschule Ulm



Masterproject

Data Science

Analytics and Prediction for Rental Bike Usage in London

Contributors:

Anass Khaldi, Guillaume Goni, Kathi Rodi, Pascal
Riedel

Expert: Prof. Dr. von Schwerin

Expert: Prof. Dr. Herbort

Expert: Prof. Dr. Goldstein

Wednesday 3rd July, 2019

Contents

1 Project Presentation and Organization	2
1 Introduction	2
2 Task Description (Semester 1)	2
3 Task Description (Semester 2)	4
4 Collaboration Technologies	5
4.1 IceScrum	5
4.2 Box	6
4.3 GitHub	6
4.4 Skype	6
5 Responsibilities	6
2 Postal Code Database (Nominatim and Graphhopper)	8
1 Nominatim	8
1.1 Installation	9
1.2 Database import	12
1.3 Usage	15
1.4 Notebooks	16
2 Graphhopper	23
2.1 Deployment	23
2.2 Usage	24
3 Database	26
3.1 SQL vs NOSql	26
3.2 Installation	27
3.3 Datamodel	27
3.4 Initialize Database	28
4 Distance Database	29
4.1 Algorithm	29
4.2 2-digits Database	30

5	Django Web Application	31
5.1	Mockup	31
5.2	Implementation	32
3	Bike Rental in London	35
1	Evaluation of Hadoop Distributions	35
1.1	Hadoop Distributions Overview	35
1.2	Definition of criteria	40
1.3	Weighted evaluation matrix	42
1.4	Performance Comparison with micro benchmarks . . .	45
1.5	Installation of sample Hadoop Distribution	46
1.6	Run MapReduce Jobs	51
1.7	HDP Issues	53
1.8	Recommendation	55
1.9	Hadoop Rollout	57
2	Data Profiling	62
2.1	Data Preparation	62
2.2	Conclusion of the file sizing problem	70
2.3	CyclingRoutes	70
2.4	Plotting Data	71
2.5	Conclusion	72
3	Prediction on Daily Basis	73
3.1	Data Profiling Part II	73
3.2	Feature Engineering Kings Cross	78
3.3	Data Prediction with MLPRegressor	81
3.4	Data Prediction with Spark	86
3.5	Modeling (Polynomial Regression)	108
4	Prediction on Hourly Basis	111
4.1	Data Preparation	112
4.2	Learning and Prediction	119
5	Additional Features	120
5.1	Data Preparation	123
5.2	Data Evaluation	125
4	Conclusion	127
1	Conclusion	127
C	List of Abbreviations	134
D	Attachments	137

Chapter 1

Project Presentation and Organization

1 Introduction

No matter in which area you move there is a huge chance that data science will be there too. It is one of the fast moving trends during the last years which is caused by the unimaginable amount of data we produce everyday. We could gain a lot of knowledge out of this data, but therefore we must analyze it. That is when data science comes into play. But with the help of data science we can not just analyze the past we can also predict the future.

Let's say we have a bike rental station which provides to rent bikes for a certain amount of time. We track the duration of the rental, the time of rental, the routes the bikes rode and lot more information. By processing all this data it is possible to predict the future usage of the bikes and the most frequented routes to specific times. That is what the main goal of this project represents. This project report deals with tasks which can be summarized by *preparatory steps* to reach this final goal. The project duration is divided into two semesters. The tasks of the first part will be discussed in this paper.

2 Task Description (Semester 1)

The first part can be roughly divided into two parts:

1. Data Engineering

- Data Profiling
- identify potential difficulties
- store data set

2. Virtual Infrastructure

- install and configure Nominatim
- install and evaluate different Hadoop distributions

The project should be implemented in an agile manner, therefore we will follow the development process of **Scrum**. This allows flexibility during the project in terms of specific tasks and topics.

To gain a better understanding of the project and its goals we created a „Big Picture“ of the applied technologies and the relations among themselves which can be seen in figure 2.1.

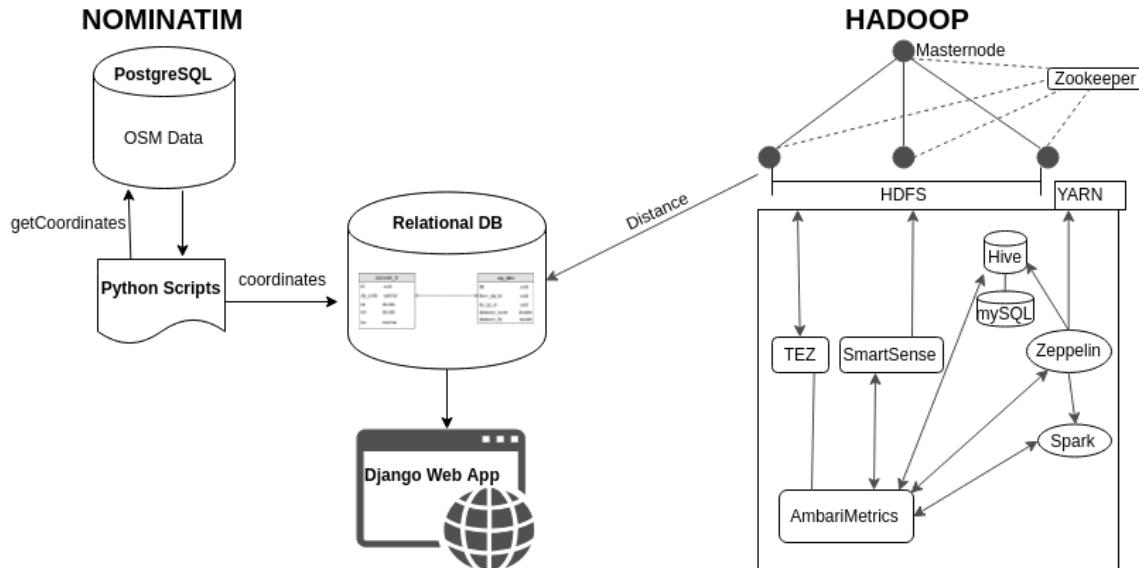


Figure 2.1: Overview of project architecture

According to this picture the project can be classified more detailed into four parts:

1. Nominatim

2. Hadoop
3. Database
4. WebApp

The implementation of those four parts will be documented in the following chapters.

3 Task Description (Semester 2)

The first part of the second semester is dedicated to finish the remaining tasks of the previous semester. This includes finish the Django Web Application as well as implementing an appropriate algorithm to fill the distance database.

Moreover the actual goal of the second semester is to implement prediction algorithms in order to forecast the usage of rental bikes in London on a daily as well as on a hourly base. Therefore data profiling tasks need to be done in order to investigate the provided data. Furthermore additional features like weather data should be researched and prepared to be added to the feature matrix. At the same time several algorithms need to be tested to find out the best fitting one to our needs. For this purpose two different libraries (Scikit-Learn and Spark) should be analyzed and tested.

A big part depicts of the data profiling part where data should be processed for further prediction usage. This includes data cleansing as well as data preparation. Once the data is ready a model can be implemented and first predictions can be made on a daily base. After that, the algorithms need to be evaluated and improved. In order to do so a suitable rating as well as meaningful plots need to be found, to get a better idea of the prediction results.

Moreover the same need to be done for the prediction on a hourly base. Which also needs specific data preparation steps.

After first predictions the features should be evaluated to gain more knowledge of which features are meaningless and which improve the prediction accuracy. A Further task is to find proper additional features which could also improve the prediction. Therefore research need to be done as well as further data preparation steps and evaluation tasks.

Another goal is to configure a Hadoop Cluster at University which should provide more computing capacity in order to make predictions out of very large data sets.

4 Collaboration Technologies

Since we want to ensure efficiency which includes that everyone has access to required data, can exchange data and that data can be versioned, we decided to work with several technologies which will be explained further in the following sections.

4.1 IceScrum

As mentioned in the chapter before the development process should be done by using Scrum methods. Those methods includes *Sprint Plannings*, *Daily Scrums*, *Retrospectives* and *Sprint Reviews*. To apply the methods properly it was very helpful to use *IceScrum* which is a web application in order to use Scrum. It provides storyboards on which tasks can be accepted. Furthermore it helps to keep track of progress and remaining tasks by burndowncharts as well as reports. We used IceScrum mainly for the sprint planning. To be precise for writing, estimating, prioritizing tasks, maintain our backlog and to distribute the incurred tasks as well as keeping track of our remaining story points and our progress.

IceScrum offers furthermore different roles which should be satisfied in Scrum. There is the *Team* which is responsible to deliver the committed delivery in time and with the defined quality. In our case the Team consists of Anass Khaldi, Guillaume Goni and Pascal Riedel. Within the team there is one *Scrum Master*, this role was taken by Kathi Rodi. The Scrum Master is responsible to ensure the team keeps to the values and practices of Scrum, sort of like a coach. He removes impediments, facilitates meetings and work with product owners as well. This leads us to the next role in Scrum the *Product Owner* which is the project's key stakeholder. The Product Owner does not get to determine how much work happens in the sprint cycles, or alter the goals for that sprint. Product Owners must be available to the team, and engage actively with it. In this project the Product Owners are Prof. Dr. von Schwerin, Prof. Dr. Herbort and Prof. Dr. Goldstein.

4.2 Box

Besides IceScrum we also used *Box* for collaborating. Our Box includes for example the documentation of our daily scrums. We created a template which is used to log every daily meeting. According burndowncharts and reports will be stored there as well. Furthermore we used the box for recording our working hours and to exchange documents.

4.3 GitHub

To group our code base and have versioned on it, we created a GitHub account which is available at <https://github.com/dataBikeHsUlm>. For the sake of clarity we created three repositories ordered by the main topics:

- WebApp: for the web server
- NominatimLibrary: for the Jupyter notebooks
- MapReduce: for the scripts for Hadoop

4.4 Skype

We carried out our daily scrums every monday at 4 pm as well as every thursday at 1.30 pm. On thursdays we have lectures together the whole day, hence the daily scrums on thursdays takes place at the university. Since it wasn't possible for the whole team to meet mondays as well, we decided to hold those daily scrums via Skype.

5 Responsibilities

5.1 First semester

Within Scrum the team is self-organizing and cross-functional. This means that every team member can choose tasks fitting his knowledge base as well as his preferences and interests. Hence everyone was free to suggest and accept the tasks he liked most. It turned out that the team leads a very good self-organization because everyone of us has a specific interest and knowledge according to the tasks of the project. Therefore Anass Khaldi was responsible to install the Nominatim server and it's components. In the second

sprint his main part was to research Django and how we can display maps on our web application. In the third sprint it was his responsibility to install graphhopper and to fix the according issues as well as issues of Hadoop. The main part of the Jupyter notebook to research nominatim functions like compute centroids, compute distances by route and doing all in bulk was done by Guillaume Goni. In the second sprint he created all GitHub repositories for storing our work. He wrote the script to initialize the database as well as the one to fill the database. Therefore he researched for best ways to collect zip codes and find the solution in the PostgresDB. Moreover he prepared a tutorial about Django for the rest of the team members. As mentioned before Kathi Rodi took over the role of the Scrum Master, hence it was her responsibility to ensure that the scrum artifacts were kept as well as correspond with the Product Owners. Apart from her tasks as a Scrum Master she was responsible to provide a box for collaborating with all needed documents like templates for daily scrums, sprint plannings, reviews, retrospectives and tracking the working hours. Furthermore she created the presentation for every sprint. Moreover she solved several tasks like create Jupyter notebooks for geocode, reverse geocoding and compute distances as the crow flies, create the ER model, install the MySQL database and to research for Graphhopper. In the end it was her responsibility to write a project report and add the task descriptions of the team members to it.

Pascal Riedel was responsible for research, document and compare the different Hadoop distributions. With the help of his weightened comparison table we could take the decision for one specific distribution. It was also his part to install the HDP as well as to implement a Map Reduce job. During the first project phase he solved several issues according to HDP. Furthermore he was responsible for the data profiling part and to create associated Zepelin notebooks.

5.2 Second semester

Kathi Rodi continued her task of Scrum Master as last semester by organising sprints and meetings, as well as preparing reports and presentations.

Pascal Riedel continued taking care of Hadoop by fixing running issues on our cluster and by setting up another one on dedicated machines. He also worked on the management of the Nominatim server by cleaning the Nominatim PostgreSQL database to import only Europe data and configured Apache to run with Django, so Django doesn't run in develop mode any more.

Finally on the configuration part, he fixed Graphhopper, finished preparing it and worked on an API to use it effectively. For these tasks, he was helped by Anass Khaldi. This done, he could research a better analysis on the bike rental on the whole London network with a map of all the stations and the paths used by the users. Pascal used PySpark on the Hadoop cluster to prepare the raw bike rental dataset for prediction. He first prepared a daily-based dataset joined with weather usage and holidays, and later worked on a hourly-based one.

Kathi Rodi and Anass Khaldi mainly worked on the data learning and prediction parts. They prepared the data and researched different learning algorithms with different parameters for a more accurate prediction. They also investigated the relevance of the different features and better ways to measure accuracies. Kathi also extracted and prepared extra features that could be added to the dataset to improve it (demographic data, events, income and political opinions). Guillaume Goni finished preparing and parsing these new features and integrating them into our existing datasets.

Guillaume Goni first focused on finishing the first part of the project concerning the postal codes and distance calculation by writing scripts to extract the postcodes from different sources and format/filter them and another script to calculate the distance between each pair of two-digits regions. Finally he completed the Django view to run the algorithm and query the database for distance information. Afterwards, Guillaume helped the team in diverse areas, for instance he adapted the learning algorithms configured by Kathi and Anass on the hourly-based dataset and displayed the predictions for comparison and helped Kathi with the extra features.

Chapter 2

Postal Code Database (Nominatim and Graphhopper)

1 Nominatim

Nominatim is a tool to search OpenStreetMap (OSM) data by name and address and to generate synthetic addresses of OSM points (reverse geocoding). Nominatim provides geocoding based on OpenStreetMap data. It uses a PostgreSQL database as a backend for storing the data. There are three basic parts to Nominatim's architecture: the data import, the address computation and the search frontend. The data import stage reads the raw OSM data and extracts all information that is useful for geocoding. This part is done by **osm2pgsql**, the same tool that can also be used to import a rendering database. It uses the special gazetteer output plugin in **osm2pgsql/output-gazetter.[ch]pp**. The result of the import can be found in the database table **place**. The address computation or indexing stage takes the data from **place** and adds additional information needed for geocoding. It ranks the places by importance, links objects that belong together and computes addresses and the search index. Most of this work is done in Procedural Language/-PostgreSQL (PL/pgsql) via database triggers and can be found in the file **sql/functions.sql**.

The search frontend implements the actual API. It takes queries for search and reverse geocoding queries from the user, looks up the data and returns the results in the requested format. This part is written in php and can be found in the **lib/** and **website/** directories.

When we received the Nominatim server, we created the non-administrator user account for our use and we added the SSH public keys of all the team members.

1.1 Installation

These instructions expect that you have a freshly installed Ubuntu 18.04. Make sure all packages are up-to-date by running:

```
sudo apt-get update -qq
```

Now you can install all packages needed for Nominatim:

```
sudo apt-get install -y build-essential cmake g++  
libboost-dev libboost-system-dev \  
libboost-filesystem-dev libexpat1-dev zlib1g-dev  
libxml2-dev \  
libbz2-dev libpq-dev libproj-dev \  
postgresql-server-dev-10 postgresql-10-postgis-2.4 \  
postgresql-contrib-10 \  
apache2 php php-pgsql libapache2-mod-php php-pear php-  
db \  
php-intl git
```

If you want to run the test suite, you need to install the following additional packages:

```
sudo apt-get install -y python3-setuptools python3-dev  
python3-pip \  
python3-psycopg2 python3-tidylib phpunit php-cgi  
pip3 install --user behave nose  
sudo pear install PHP_CodeSniffer
```

System Configuration

The following steps are meant to configure a fresh Ubuntu installation for use with Nominatim. You may skip some of the steps if you have your OS already configured. Nominatim will run as a global service on your machine. It is therefore best to install it under its own separate user account. In the following we assume this user is called **nominatim** and the installation will be in **/srv/nominatim**. To create the user and directory run:

```
sudo useradd -d /srv/nominatim -s /bin/bash -m  
nominatim
```

You may find a more suitable location if you wish. To be able to copy and paste instructions from this manual, export user name and home directory now like this:

```
export USERNAME=nominatim  
export USERHOME=/srv/nominatim
```

Never, ever run the installation as a root user. Make sure that system servers can read from the home directory:

```
chmod a+x $USERHOME
```

Setting up PostgreSQL

Tune the PostgreSQL configuration, which is located in `/etc/postgresql/9.5/main/postgresql.conf`. You might want to tune your PostgreSQL installation so that the later steps make best use of your hardware. You should tune the following parameters in your `postgresql.conf` file.

```
shared_buffers (2GB)  
maintenance_work_mem (10GB)  
work_mem (50MB)  
effective_cache_size (24GB)  
synchronous_commit = off  
checkpoint_segments = 100 # only for postgresql <= 9.4  
checkpoint_timeout = 10min  
checkpoint_completion_target = 0.9
```

The numbers in brackets behind some parameters seem to work fine for 32GB RAM machine. Adjust to your setup. For the initial import, you should also set:

```
fsync = off  
full_page_writes = off
```

Don't forget to reenable them after the initial import or you risk database corruption. **Autovacuum** must not be switched off because it ensures that the tables are frequently analyzed.

Setting up the Webserver

The **website** directory in the **build** directory contains the configured website. Include the directory into your web browser to serve php files from there.

Make sure your Apache configuration contains the required permissions for the directory and create an alias:

```
<Directory "/srv/nominatim/build/website">
Options FollowSymLinks MultiViews
AddType text/html
.php
DirectoryIndex search.php
Require all granted
</Directory>
Alias /nominatim /srv/nominatim/build/website
```

The path of the build directory (**/srv/nominatim/build**) should be replaced with the location of your build directory. After making changes in the apache config you need to restart apache. The website should now be available on `http://localhost/nominatim`. Restart the PostgreSQL service after updating this config file.

```
sudo systemctl restart postgresql
```

Finally, we need to add two PostgreSQL users: one for the user that does the import and another for the webserver which should access the database for reading only:

```
sudo -u postgres createuser -s $USERNAME
sudo -u postgres createuser www-data
```

You need to create an alias to the website directory in your apache configuration. Add a separate nominatim configuration to your webserver:

```
sudo tee /etc/apache2/conf-available/nominatim.conf <<
EOFAPACHECONF
<Directory "$USERHOME/Nominatim/build/website">
Options FollowSymLinks MultiViews
AddType text/html
.php
DirectoryIndex search.php
```

```
Require all granted</Directory>
Alias /nominatim $USERHOME/Nominatim/build/website
EOFAPACHECONF
```

Then enable the configuration and restart apache:

```
sudo a2enconf nominatim
sudo systemctl restart apache2
```

Building and Configuration

Get the source code for the release and change into the source directory:

```
cd $USERHOME
wget https://nominatim.org/release/Nominatim-3.2.0.tar.bz2
tar xf Nominatim-3.2.0.tar.bz2
cd Nominatim-3.2.0
```

The code must be built in a separate directory. Create this directory, then configure and build Nominatim in there:

```
mkdir build
cd build
cmake ..
make
```

You need to create a minimal configuration file that tells nominatim where it is located on the webserver:

```
tee settings/local.php << EOF
<?php
@define('CONST_Website_BaseURL', '/nominatim/');
EOF
```

After those steps Nominatim is ready to use.

1.2 Database import

The following instructions explain how to create a Nominatim database from an OSM planet file and how to keep the database up to date. It is assumed that the Nominatim software itself is already successfully installed.

Flatnode files

If you plan to import a large dataset (e.g. Europe, North America, planet), you should also enable flatnode storage of node locations. With this setting enabled, node coordinates are stored in a simple file instead of the database. This will save you import time and disk storage. Add to your settings/local.php:

```
@define( 'CONST_Osm2pgsql_Flatnode_File' , '/path/to/
flatnode.file' );
```

Replace the second part with a suitable path on your system and make sure the directory exists. There should be at least 40GB of free space.

Initial import of the data

Important: first try the import with a small excerpt, for example from **Geofabrik**. Download the data to import and load the data with the following command:

```
./ utils/setup.php --osm-file <data file> --all [--osm2pgsql-cache 28000] 2>&1 | tee setup.log
```

The **--osm2pgsql-cache** parameter is optional but strongly recommended for planet imports. It sets the node cache size for the **osm2pgsql** import part (see -C parameter in osm2pgsql help). As a rule of thumb, this should be about the same size as the file you are importing but never more than 2/3 of RAM available. If your machine starts swapping reduce the size. Computing word frequency for search terms can improve the performance of forward geocoding in particular under high load as it helps PostgreSQL' query planner to make the right decisions. To recompute word counts run:

```
./ utils/update.php --recompute-word-counts
```

This will take a couple of hours for a full planet installation. You can also defer that step to a later point in time when you realise that performance becomes an issue. Just make sure that updates are stopped before running this function.

Updating the data

There are many different possibilities to update the Nominatim database. The following section describes how to keep it up-to-date with **Pyosmium**.

For a list of other methods see the output of `./utils/update.php --help`. It is recommended to install Pyosmium via pip. Run (as the same user who will later run the updates):

```
pip install --user pyosmium
```

Nominatim needs a tool called **pyosmium-get-updates**, which comes with Pyosmium. You need to tell Nominatim where to find it. Add the following line to your **settings/local.php**:

```
@define('CONST_Pyosmium_Binary', '/home/user/.local/bin/pyosmium-get-changes');
```

The path above is fine if you used the **-user** parameter with pip. Replace user with your user name.

Next the update needs to be initialised. By default Nominatim is configured to update using the global minutely diffs. If you want a different update source you will need to add some settings to **settings/local.php**. For example, to use the daily country extracts diffs for Ireland from geofabrik add the following:

```
// base URL of the replication service
@define('CONST_Replication_Url', 'https://download.geofabrik.de/europe/ireland-and-northern-ireland-updates');
// How often upstream publishes diffs
@define('CONST_Replication_Update_Interval', '86400');
// How long to sleep if no update found yet
@define('CONST_Replication_Recheck_Interval', '900');
```

To set up the update process now run the following command:

```
./utils/update.php --init-updates
```

It outputs the date where updates will start. Recheck that this date is what you expect. The **--init-updates** command needs to be rerun whenever the replication service is changed.

The following command will keep your database constantly up to date:

```
./utils/update.php --import-osmosis-all
```

Note that even though the old name „import-osmosis-all“ has been kept for compatibility reasons, Osmosis is not required to run this - it uses Pyosmium behind the scenes.

Nominatim will respond with the following data formats:

- JSON
- JSONv2
- GeoJSON
- Geocode JSON
- XML

1.3 Usage

As we want to use python for working with Nominatim we use the python package **geopy** to interact with the server.

Up to now, we were using the default public Nominatim server [34]. As we want to use our own Nominatim server we researched for connecting the geopy library to it.

First, the Nominatim server needed to be configured with Apache, the configuration that was done earlier wasn't working properly so we firstly had to fix that. The server can be accessed at the following url: <http://i-nominatim-01.informatik.hs-ulm.de/nominatim/>.

After that HTTPS isn't enabled on the server. Of course, we can set it up later but for now we need to disable it for geopy, it is done by setting:

```
geopy.geocoders.options.default_scheme = "http"
```

Third, our server is not very powerful and being quite new, it has not cached enough queries. This can make some queries take a long time to complete and according to that the geopy query timeouts. Therefore we had to increase this limit. The final connection is shown in figure 1.1.

```
1 from geopy.geocoders import Nominatim
2
3 DOMAIN = "i-nominatim-01.informatik.hs-ulm.de/nominatim/"
4 TIMEOUT = 100000
5
6 geopy.geocoders.options.default_scheme = "http"
7 locator = Nominatim(user_agent="Test", timeout=TIMEOUT, domain=DOMAIN)
```

Figure 1.1: Connection of Nominatim with geopy

1.4 Notebooks

In attempt to gain more knowledge on how it is possible to work with Nominatim in python we created several Jupyter notebooks. The different methods we need on the Nominatim API are grouped in a python library on the project's GitHub repository which can be accessed under <https://github.com/dataBikeHsUlm/NominatimLibrary>. It needs **geopy** as a dependency and **osmapi** for **pyroutelib2** which can be installed via pip:

```
pip3 install geopy osmapi
```

Geocoding

Geocoding an address means we input an address and as an output we expect the corresponding coordinates. Therefore we have to extract the coordinates as follows:

```
def locate (address):
    """Gets the coordinates to a given address

    Args:
        address(str): String List of following
                      arguments:house_number,road,town,city,
                                 county,state_district,state,postcode
                                 ,country,country_code

    Returns:
        Associated address and coordinates
    """
    location = geolocator.geocode(address)
    print(location.address)
    print(location.latitude, location.longitude)
```

We could also extract coordinates not by searching for a whole address but by postalcodes:

```
def locateCords(postcode):
    """Gets the coordinates to a given address

    Args:
        postalcode(str): String of postalcode
```

```

    Returns :
        Associated coordinates as latitude and
        longitude
"""
location = geolocator.geocode(postcode)
lat = location.latitude
lon = location.longitude
print "The Latitude of ", postcode, " is: ", lat, " the
        Longitude is ", lon
return (lat, lon)

```

The results are the corresponding coordinates.

Reverse Geocoding

Vice versa we can apply reverse geocoding to compute the corresponding address to given coordinates.

```

def reverseLocate (coordinates):
    """Gets the address to given coordinates

```

```

    Args:
        coordinates(str): String List of
            coordinates

```

```

    Returns :
        Associated address
"""

```

```

location = geolocator.reverse(coordinates)
print(location.address)

```

The result was compared with Google Maps which is shown in figure 1.2.

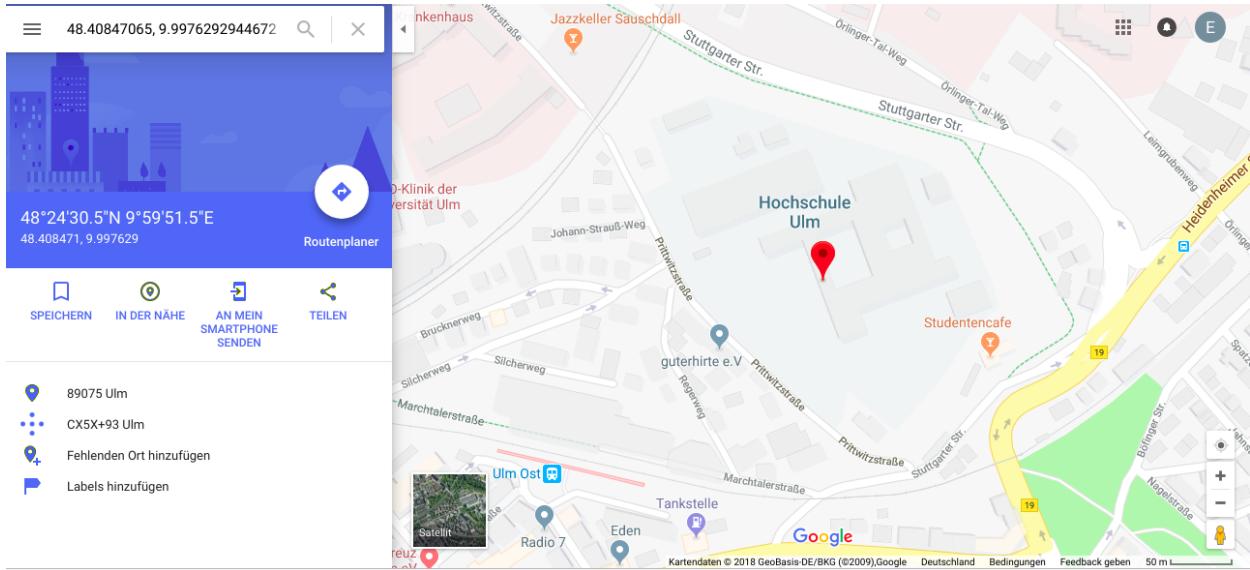


Figure 1.2: Reverse Geocoding

Querying Centroids

As already mentioned geopy offers methods to access addresses and coordinates. We used this as follows for querying centroids:

```
location = locator.geocode("Baden-Wuerttemberg",
                           Deutschland")
print(location.address)
print(location.latitude, location.longitude)
```

We received the following result:

```
Baden-Wuerttemberg, Deutschland
48.6296972 9.1949534
```

We entered the given coordinates for the centroid of Baden-Wuerttemberg (**48.6296972, 9.1949534**) in Google Maps and made a route between this centroid and the one given by Google Maps as is shown in figure 1.3.

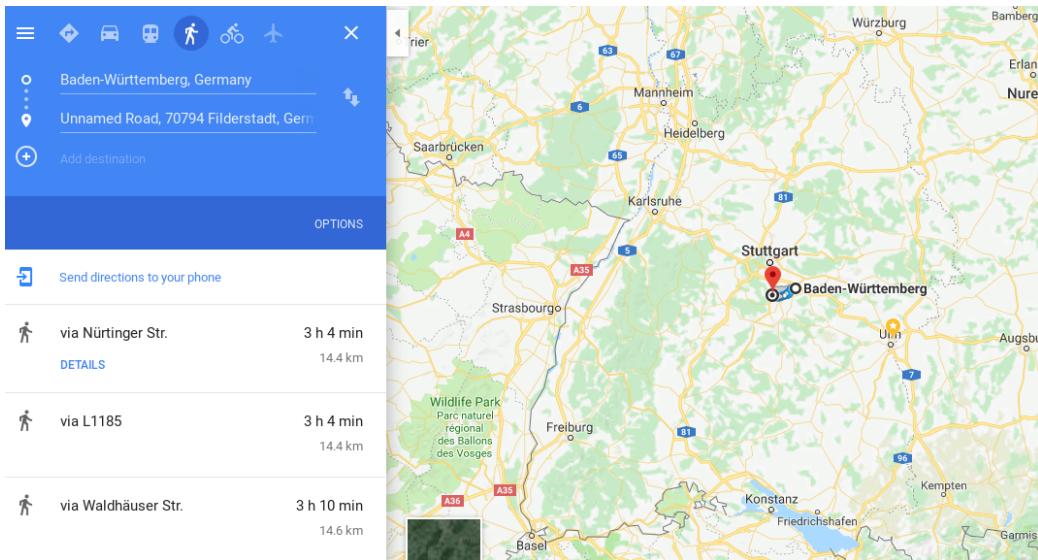


Figure 1.3: Centroid of Baden-Wuerttemberg shown in Google Maps

We can see that there is only 14.4km between the two points. If we consider Baden-Wuerttemberg is roughly 200km wide, we have a difference of 7.2%. The values used here are obviously not very accurate but we can still see that the given centroid is located roughly in the center of the requested region. The difference with Google Maps can maybe be explained by a difference of precision in the borders.

Distances as the crow flies

To compute the distance between two locations as the crow flies we can use the **distance()** method which is provided by geopy. Therefore we have to transfer coordinates into points:

```
def locatePoint (postcode):
    """Gets the coordinates to a given address

    Args:
        postalcode (str): String of postalcode

    Returns:
        Associated coordinates as a Point
```

```

"""
location = geolocator.geocode(postcode)
lat = location.latitude
lon = location.longitude
p = location.point
return p
locatePoint("80336")

Subsequently we can compute the associated routes:

def getDistance(start, end):
    """Get distance of two postalcodes

    Args:
        start(str): String of start postalcodes
        end (str): String of end postalcode

    Returns:
        Distance in km
    """
    a = locatePoint(start)
    b = locatePoint(end)
    d = distance.distance(a, b).km
    print("The distance between Ulm and Munich is: ", d, "km")
    return d

getDistance("89075", "80336")

```

The result was compared to Google Maps and can be seen in figure 1.4.

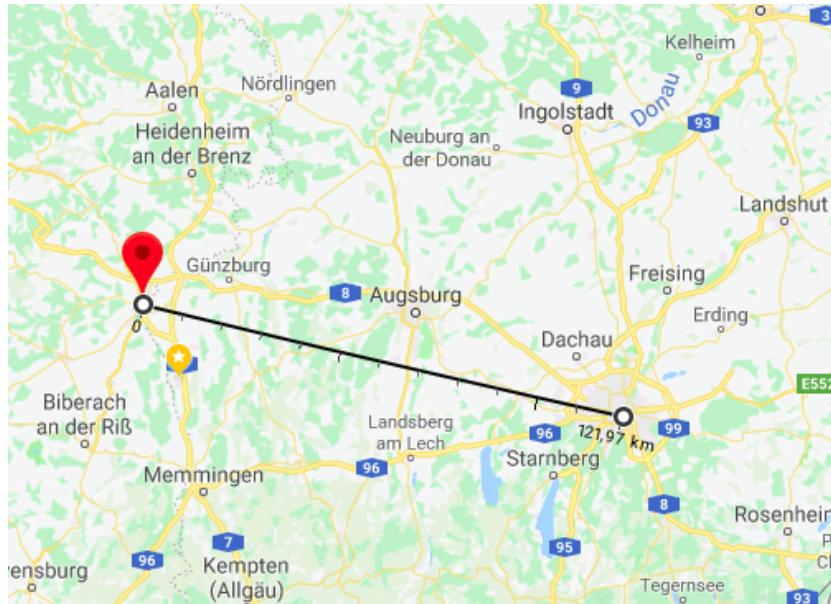


Figure 1.4: Compute Distances as the Crow Flies shown in Google Maps

Distances by Road

There is not much library providing a routing algorithm for python and working with OpenStreetMap. But we found **pyroutelib2** [36]. It does not seem to be supported therefore we had to make a few modifications to the original code to make it work.

The modified version can be found in the NominatimLibrary GitHub repository.

The dependency osmapi is required.

First, we get the coordinates of the starting and finishing points:

```
location = locator.geocode("Prittwitzstrasse, Ulm")
a = (location.latitude, location.longitude)
location = locator.geocode("Albert-Einstein-Allee, Ulm")
b = (location.latitude, location.longitude)
from pyroutelib2.loadOsm import LoadOsm
from pyroutelib2.route import Router
# By default, it uses the open API, this can be changed
# directly in the file
```

```

# Here , we use bicycle to calculate routes .
data = LoadOsm("cycle")
router = Router(data)
# This gets the node ids of the two points
node_a = data.findNode(a[0] , a[1])
node_b = data.findNode(b[0] , b[1])
# `doRoute` calculates the route and returns a list of
# coordinates tuples
result , route = router.doRoute(node_a , node_b)
if result == 'success':
# Do something ...
    pass
else:
    print("Error calculating the route .")
Finally, to get the distance, we compute the distance between each couple of
points:
from geopy import distance
lats = []
lons = []
if result == 'success':
for i in route:
    node = data.rnodes[ i ]
    lats.append( node[0] )
    lons.append( node[1] )
else:
    print("Error calculating the route .")
distance_route = 0
for i in range(len(lons)-1):
    p1=(lats[i],lons[i])
    p2=(lats[i+1],lons[i+1])
    distance_route += distance.geodesic(p1,p2 , ellipsoid='
        GRS-80').km
print("Total_distance_on_the_route : %.2fkm" %
      distance_route)

```

As a result we get from Prittwitzstrasse to Albert-Einstein-Allee, a distance by route of 5.87 km. Google Maps gives a similar result of 5.2 km as shown in figure 1.5.

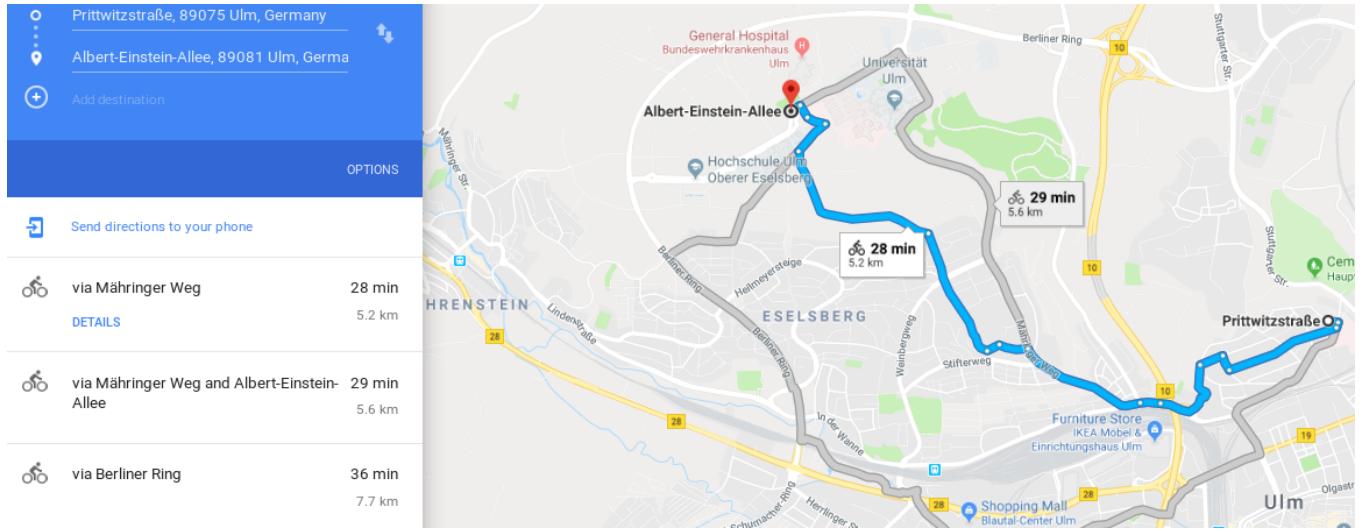


Figure 1.5: Distance shown in Google Maps

For some in-city routes, the library is pretty fast, but for longer routes (Stuttgart → Muenchen), **pyroutelib2** took nearly an hour. Consequently, this is not a possible solution to calculate all the distances between postcodes. We will then have to turn to another way (external to python) of calculating the routes for the project.

2 Graphhopper

Graphhopper is an open-source routing library and server written in Java which provides a web interface called Graphhopper Maps as well as a routing API over HTTP.

2.1 Deployment

For simplicity you could just start jetty from maven and schedule it as background job:

```
graphhopper.sh -a web -i europe-germany-berlin.pbf -d
--port 11111
```

Then the service will be accessible on port 11111. For production usage you have a web service included. Use **-c config.yml** in the script to point to it.

Increase the **-Xmx/-Xms** values of your server e.g. for world wide coverage with a hierarchical graph do the following before calling **graphhopper.sh**:

```
export JAVA_OPTS="-server -Xconcurrentio -Xmx17000m -Xms17000m"
```

Graphhopper is able to handle coverage for the whole OpenStreetMap road network. It needs approximately 22GB RAM for the import (CAR only) and 1 hour (plus 5h for contraction). If you can accept slower import times this can be reduced to 14GB RAM - you'll need to set **datareader.dataaccess=MMAP**. Then 'only' 15GB are necessary. Without contraction hierarchy this would be about 9GB. With CH the service is able to handle about 180 queries per second (from localhost to localhost this was 300qps). Measured for CAR routing, real world requests, at least 100km long, on alinux machine with 8 cores and 32GB, java 1.7.0_25, jetty 8.1.10 via the QueryTorture class (10 worker threads).

Especially for large heaps you should use **-XX:+UseG1GC**. Optionally add **-XX:MetaspaceSize=100M**. Avoid swapping e.g. on linux via **vm.swappiness=0** in **/etc/sysctl.conf**. If you want to use elevation data you need to increase the allowed number of open files. Under linux this works as follows:

- sudo vi /etc/security/limits.conf
- add: * - nofile 100000 which means set hard and soft limit of "number of open files" for all users to 100K
- sudo vi /etc/sysctl.conf
- add: fs.file-max = 90000
- reboot now (or sudo sysctl -p; and re-login)
- afterwards ulimit -Hn and ulimit -Sn should give you 100000

2.2 Usage

Our endpoint is [https://graphhopper.com/api/\[version\]/vrp](https://graphhopper.com/api/[version]/vrp). Graphhopper can provide up to 6 APIs:

- Routing API

- Route Optimization API
- Isochrone API
- Map Matching API
- Matrix API
- ZGeocoding API

Routing API

Get distances between two points:

```
curl "https://graphhopper.com/api/1/route?
point=51.131,12.414&point=48.224,3.867&vehicle=car&
locale=de&key=[YOUR_KEY]"
```

Example in our case:

```
https://localhost:11111/1/route
?point=51.131,12.414&point=48.224,3.867&vehicle=car&
locale=de&key=3160e710-58ed-45da-
bb39-09d383b1c5b2
```

Pricing

The Graphhopper routing engine is open source under the permissive Apache License and is therefore free to use for anything. You could even integrate it in your products, modify Graphhopper and sell this, without noticing or contributing back. Although it is encouraged to contribute back so that your feature gets maintained for free by graphhopper service. Also you can host Graphhopper on your own servers for 'free' and do whatever you want with it. The Graphhopper Directions API that we host falls under usage terms and always requires an API key. It was decided to make it free for development purposes and open source projects, both with a limit of currently 500 queries per day. So, the free usage of the API in a company internally would not be allowed. But there are custom packages possible.

One Routing API request costs one credit. Every 10 via-points cost one more credit. E.g. 11 via-points cost two credits, 21 via-points costs three

credits and so on. And if you specify **optimize=true** the credits will be multiplied by 10 i.e. one requests costs 10 credits for 1 to 10 locations, 20 credits for 11 to 20 locations and so on. Changing the parameter algorithm costs additionally two credits, i.e. calculating the alternative route between two points costs $1+2=3$ credits. For more info about the other APIs: Graphhopper API Docs

3 Database

Since we want to compute distances between the centroids of two postalcode areas we need a database where those information can be stored.

3.1 SQL vs NOSql

First, it was necessary to find a suitable database. For this purpose we compared the Not only Sql (NoSQL) and the Structured Query Language (SQL) approaches with each other.

SQL happens to be the more structured, right way of storing data, like a phone book. For a relational database to be effective, you will have to store your data in a very organized fashion. SQL databases remain popular because they fit naturally into many vulnerable software stacks, including LAMP and Ruby-based stacks. These databases are widely supported and well understood, which could be a major plus point if you run into problems.

The main problem with SQL is scaling it as your database grows. You see, even though scalability is usually tested in production environments, its often lower than NoSQL databases.

For dealing with massive amounts of unstructured data and data requirements which arent clear at the outset, you probably dont have the luxury of developing a relational database with a clearly defined schema. You get much more flexibility than its traditional counterparts, with non-relational databases. Picture non-relational databases as file folders, assembling related information of all types.

Since we already know the scope of our database and therefore don't want to scale it up and due to the lack of time and the good knowledge of team members in SQL, we decided to make use of the relational database **MySQL**.

3.2 Installation

Since MySQL is the most popular open-source relational database management system, we decided to install this system on our server. First we installed the MySQL package with the following command:

```
sudo apt install mysql-server
```

Once the installation is completed, the MySQL service will start automatically. To check whether the MySQL server is running or not, we type:

```
sudo systemctl status mysql
```

After we have made sure that the database is running, we create a user account:

```
CREATE USER IF NOT EXISTS 'database_user'@'localhost'  
IDENTIFIED BY 'user_password';
```

Copy

Afterwards a new database could be created and initialized.

3.3 Datamodel

To implement an appropriate database which offers fast responses to queries it was indispensable to design a datamodel which fulfills these requirements. The figure 3.1 shows the output of the datamodelling process.

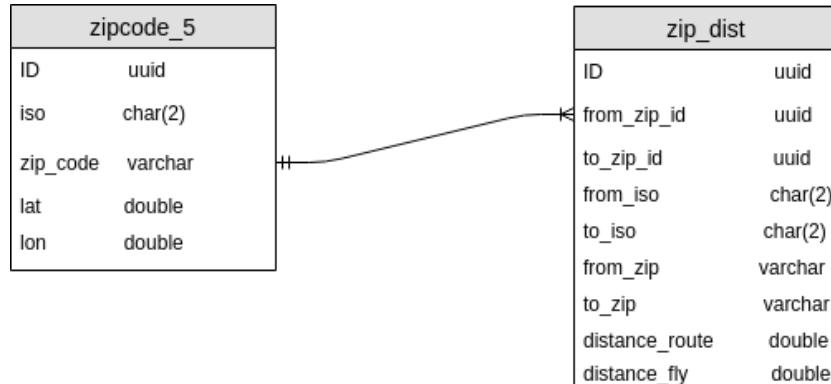


Figure 3.1: ER model for the MySQL Database

The model consists of two entities which are connected over a one-to-many relationship. The datamodels for the postcodes and the distances were directly written in the Django source code. That way, Django directly provides a python object to work with the tables. It is even possible to write methods to manipulate this data directly from it. The implementation is done in the file: `model.py`. For further information on working with the database in Django, see our corresponding GitHub repository.

3.4 Initialize Database

Geonames

To be able to fill the postalcode table, we need a list of all postcodes per country. The Nominatim API doesn't directly provide a function to do that, however the *GeoNames*, which is a geographical database, provides a list of all postalcodes which are free to download [12]. We implemented a first version of a script using this list and started filling the database. Geonames was incomplete, some countries (like Greece) were missing and the script had problems finding some cities. Therefore we had to find another, more reliable way to get the postcodes.

Nominatim Database

We decided to directly look into the Nominatim PostgreSQL database. There is in fact a table named `location_area` containing the column postcodes and the related `country_code`. Some rows in this table have empty `country_code` or `postcode`, so we need to filter them out. The resulting SQL query looks as followed:

```
1 SELECT DISTINCT country_code,postcode \
2   FROM location_area \
3  WHERE country_code IS NOT NULL \
4    AND postcode IS NOT NULL \
5 ORDER BY country_code, postcode;
```

Then, for every postcode, we need to query the Nominatim API to get the centroid coordinates, but querying only with the country code and the postcode is too ambiguous and can lead to wrong centroids. To fix that,

we used the **country_name** table that lists countries with names and ISO codes. We joined this table to the previous one:

```
1 SELECT DISTINCT location_area.country_code,postcode,country_name.name \
2   FROM (location_area LEFT JOIN country_name \
3     ON location_area.country_code = country_name.country_code) \
4 WHERE location_area.country_code IS NOT NULL \
5   AND postcode IS NOT NULL \
6 ORDER BY country_code, postcode;
```

Final Script

To ensure that we don't left out any postal codes, we use both sources to fill the database. In the table, 443493 postcodes from 51 countries were inserted with this script. Using the resulting data, the script to fill the zipcode table was re-written, the database cleaned and filled with the new version of the script. The script can be found at the GitHub repository WebApp. To run the script, run the accompanying shell script **fill_db.py**, it will install the dependencies before running the python code.

4 Distance Database

4.1 Algorithm

If we computed the distance between all cities in Europe, it would take a very long time and file the hard drive consequently. To reduce the effort was decided to group each postcode by its region (most of the time represented by the two first digits). Then we compute the distance "as the crow flies" and by road between each of these regions and store them into the database.

When requesting the distance between two postcode, we process as such:

1. transform both postcode into coordinates
2. calculate the distance as the crow flies (named $d_{direct,crow}$)
3. determine the two 2-digit regions associated
4. get the distance "as the crow flies" ($d_{region,crow}$) and by road ($d_{region,road}$) between the two regions from the database
5. calculate the ratio $ratio = d_{region,road}/d_{region,crow}$

- the final distance estimate is $d_{direct,crow} * ratio$

The calculation is summarized by figure 4.1.

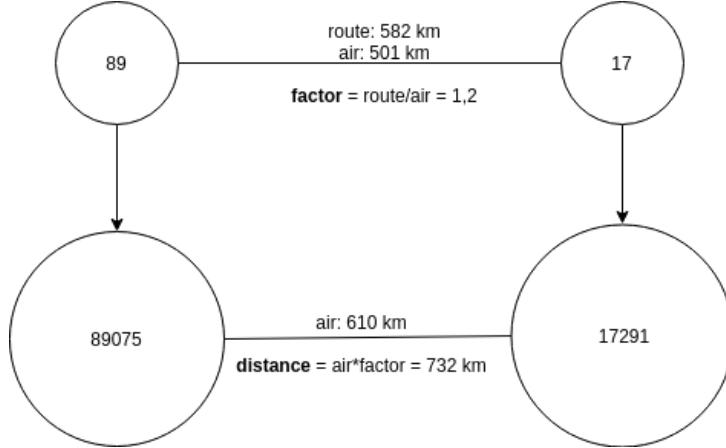


Figure 4.1: Theoretical example of distance calculation

4.2 2-digits Database

List of 2-digits Regions

The countries' regions being often indicated as the two first digits of the postcodes, we can simply create our groups based on that. However, we must take into account that there are some exceptions; for instance, Luxembourg prefixes its postcodes with “LE-”. Some other countries also introduce non-alphanumerical characters, such as: spaces, hyphens, dots, ...

We also restrict this operation to European countries only, we can do that by filtering postcodes by their coordinates: Europe being found between -11 West and 41 East, and 35 to 71 North.

Calculation of a Region's Centroid

To calculate the centroid for each region, we can compute the average coordinates of all postcodes of that region. However, we will have to use Graphhopper with these coordinates and Graphhopper can only work when coordinates point to a road/street/... To test if a coordinate is working for Graphhopper, we simply query it from our coordinate to an already checked place.

First, we test Graphhopper for the average centroid. If it can't be used, we go through the closest postalcodes to find one accepted by Graphhopper. Finally, if we still don't have a good point, we spiral around the average centroid until we find a good point or we reach 3 degrees radius (3 degrees in longitude is around 245km in Germany). We discard every other regions that fail these three tests.

Filling the Database

For each couple of region centroids calculated above, we calculate the distance "as the crow flies" thanks to the GeoPy function and the distance by road thanks to Graphhopper API and store those two results in the MySQL database. We obtain, in the end, 3399528 distances in the database in 2607 regions.

The script can be found at the GitHub repository WebApp.

5 Django Web Application

The aim of this task was to develop a web application to provide distance calculation based on our Nominatim as well as our Graphhopper server. The calculation should be done by postalcodes as well as iso codes. As a result it should return the air distance as well as the distance by route of the related postalcodes.

5.1 Mockup

After the requirements were clear, we designed a mockup (see figure 5.1) in order to define the layout of the web application.

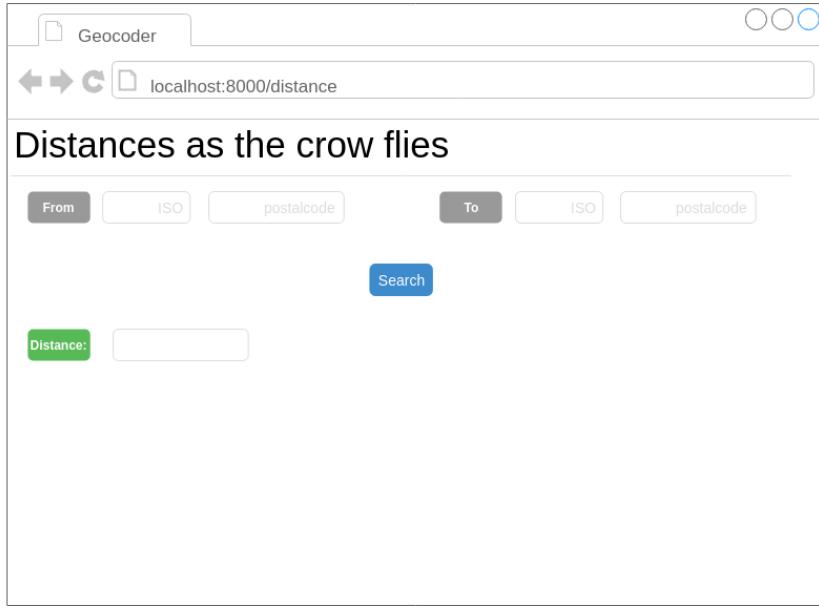


Figure 5.1: Mockup for Django web app

The mockups shows four input fields where the user can enter the postalcode and the corresponding iso code of his start and his starting point as well as his destination. The „Search“-button triggers the calculation and represents the result in the input field with the label „Distance“below. In this mockup only the distance calculation as the crow flies was taken into consideration. However in the final web application we included the distance calculation by route as well.

5.2 Implementation

For our frontend web application, we will use **Django** as it works in python and is the most used.

View

This mockup was then implemented as an HTML file. In order to get a quick result with a sophisticated and responsive design we decided to include the Bootstrap toolkit. This open source framework provides themes, templates and code snippets for a uniform layout and rapid processing. Therefore we

installed bootstrap on our Nominatim server. The view consist mainly of four parts:

1. index.html
2. style.css
3. urls.py
4. views.py
5. forms.py

The *index.html* file is responsible for the general structure of the web application. The *style.css* file is mainly responsible for the layout of the footer which is customized and not just adapted by bootstrap.

In the *urls.py* we defined which HTML file should be displayed for which url. Since the web application consist of one page only we set the *index.html* as the entry point by calling the web application.

Django provides a very secure form handling which is normally a very complex business including editing the form with a convenient interface, send it to the server, validate and clean the input and then save or pass it to further processing. With Django's form functionality all this work is automated and also do it more securely then we would probably be able to do. Therefore we created the file *forms.py* where we directly specified the input fields with their corresponding length as well as if they are required or not.

```
class CrowForm(forms.Form):  
    from_zip = forms.CharField(max_length=5, required=  
        True)  
    from_iso = forms.CharField(max_length=2, required=  
        True)  
    to_zip = forms.CharField(max_length=5, required=True  
        )  
    to_iso = forms.CharField(max_length=2, required=True  
        )
```

The input of this form is then processed within the *views.py* file, since here the post and get requests are handled. The *views.py* file is responsible to handle post ad get requests. This file calls also the responsible methods of the datamodel in order to calculate the distances. Moreover it returns

the results to the input fields „Direct“ and „By route“ which can be seen in figure 5.2.

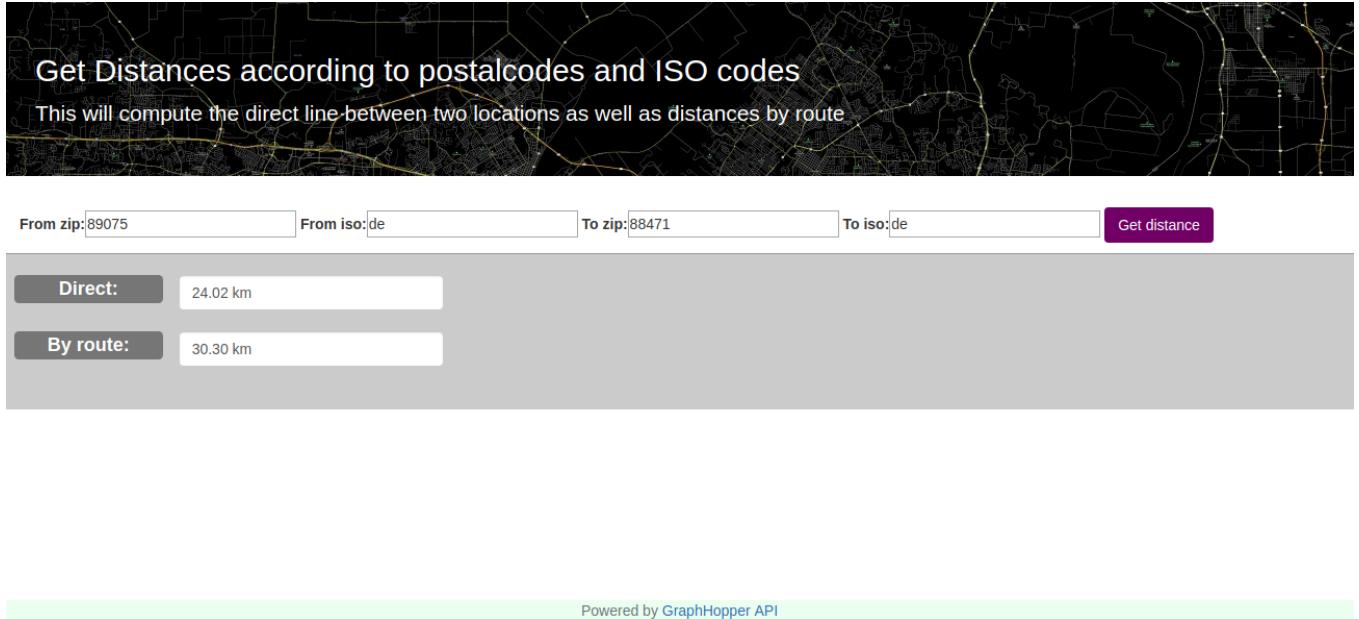


Figure 5.2: Final layout of the Django web application

Calculation

As described in Section 4.1, we first need to determine the corresponding region for the two given postal codes and country ISO codes and extract from the distance database the two distances to calculate the air/road ratio. We use Nominatim to convert the postcodes to coordinates to calculate their “as the crow flies” distance and we finally apply the regional ratio to it.

If the two postcodes are in the same region, we could apply an average ratio to correct a little the difference with the by road distance. We could also directly use Graphhopper again as these are two close points and shouldn't require much resources to compute.

Chapter 3

Bike Rental in London

1 Evaluation of Hadoop Distributions

This chapter shows a general conducted comparison of Hadoop distributions which are dominating the big data market. Since it is necessary for our data science project to compute in a parallel and efficient fashion, such a comparison helps us to identify the right distribution. Firstly, an overview of the identified distributions with its features and pricing tiers will be showed. Afterwards we will declare some important criteria for our project based on an extensive research of internet blogs and manufacturer websites. Then a weighting of these criteria will be done. These criteria will be used to evaluate the different Hadoop distributions. The resulting evaluation matrix shows that Hortonworks might be the best distribution for our data science problem. However, a detailed look at the performance criterion shows us that MapR is the most powerful in five out of six cases. For this performance comparison, evaluation results of various MapReduce jobs found on the Internet are used. Hence, in a pure performance comparison MapR would win over Hortonworks, but in an overall comparison Hortonworks convinces more. To round off the evaluation, a test implementation of Hortonworks was carried out on the four virtual machines procured for this purpose where some simple MapReduce jobs are carried out.

1.1 Hadoop Distributions Overview

For the first part of the evaluation, a search for widespread Hadoop distributions was conducted [28, 46, 21, 6, 24]. As the table 3.1 below shows, seven

noteworthy distributions were found.

Hadoop Distribution	Features	Pricing Gear	Hadoop Ecosystem
Apache Hadoop 3.1.1 hadoop.apache.org	<ul style="list-style-type: none"> • No extra tools • Blank installation of Apache Hadoop (HDFS, YARN and Hadoop Web UI) 	Open Source 100%	Has to be installed manually [2]
Cloudera CDH (v5.15.0) www.cloudera.com	<ul style="list-style-type: none"> • Oldest distribution • Very polished • Comes with good (and proprietary) tools to install and manage a Hadoop cluster: <ul style="list-style-type: none"> – Cloudera Manager for managing and monitoring clusters [5]. – Cloudera Search (Free-Text) on Hadoop [7]. – Apache Crunch framework for MapReduce Pipelines [7]. 	Freemium (Cloudera Manager require license). Also source code is not fully available and enterprise edition has 60 trial-day [9]. Final costs may depend on cluster size [6].	Accumulo, Flume, HBase, HCatalog, Hive, HttpFS, HUE, Impala, Kafka, KMS, Mahout, Oozie, Pig, Sentry, Snappy, Spark, Sqoop, Parquet, Whirr, ZooKeeper [5, 7]

Hadoop Distribution	Features	Pricing Gear	Hadoop Ecosystem
Hortonworks Data Platform (HDP3.0.1) www.hortonworks.com	<ul style="list-style-type: none"> • Newer distributions • Tracks Apache Hadoop closely • Comes with standard and open source tools for managing clusters: <ul style="list-style-type: none"> – Improved HDFS in HDP 2.0: automated fail over with a hot standby and full stack resiliency [14]. – Ambari for managing and monitoring clusters. – Includes almost all Hadoop extensions from the Apache foundation [14]. 	Open Source, optional enterprise paid support [9]	Atlas, HBase, HDFS, Hive Metastore, HiveServer2, Hue, Spark, Kafka, Knox, Oozie, Ranger, Storm, WebHCat, YARN, SmartSense, ZooKeeper [15]
MapR 6.1 www.mapr.com	<ul style="list-style-type: none"> • Uses own HDFS (MapR-FS) • Integrates own database systems (MapR-DB seven times faster than HBase) • Mostly used for big data projects • Free from Single Point of Failures • Offers Mirroring and Snapshotting • Might be the fastest Hadoop distribution • MapR supports backward compatibility across multiple versions of projects • Supports a global event replication for streaming at IoT scale (MapR-ES) 	Freemium [26]:Community Edition contains no high availability, no disaster recovery and no global replication for streaming data.	AsyncHBase, Cascading, Drill, Flume, HBase Client and MapR Database, Binary Tables, Spark, HCatalog, Hive, HttpFS, Hue, Impala, MapR Event Store For Apache Kafka Clients and Tools, Myriad, OpenStack, Manila, Oozie, Pig, Sentry, Spark, Sqoop, MapR Object Store with S3- Compatible API[25]

Hadoop Distribution	Features	Pricing Gear	Hadoop Ecosystem
Intel www.hadoop.intel.com	<ul style="list-style-type: none"> • Partnered with Cloudera [29] • Encryption support • Hardware acceleration added to layers of stack to boost performance • Admin tools to deploy and manage Hadoop 	Premium (90 days trial)	Offers same services as Cloudera.
Pivotal HD www.gopivotal.com	<ul style="list-style-type: none"> • Partnered with Hortonworks [29] • Fast SQL on Hadoop • Proprietary Software [38]: <ul style="list-style-type: none"> – Pivotal DataLoader – USS (external file system) – Spring Data – Pivotal ADS-HAWQ (parallel SQL query engine) 	Premium	Offers same services as Hortonworks.
IBM Open Platform www.ibm.com	<ul style="list-style-type: none"> • Partnered with Hortonworks [29] • Proprietary tools: <ul style="list-style-type: none"> – IBM Big SQL allows concurrent process of Hive, HBase and Spark and other sources using a single database connection [29]. – IBM BigInsights v4.2 provides Text-Analytics module [19]. • Highly compatible to other IB products. 	Free for non-commercial purposes, optional enterprise paid support [20]	Offers same services as Hortonworks.

Table 3.1: Overview of Hadoop distributions

It is worth to mention that each Hadoop distribution in table 3.1 comes up with a minimum of services from the Apache foundation: Yet Another Resource Negotiator (YARN), Hadoop Distributed File System (HDFS), Hive, Pig, HBase, Kafka, Storm, Mahout, HCatalog... But they may use a proprietary implementation of them e.g. MapR uses own HDFS instead of Apache HDFS implementation. It is also worth to note that a few cloud providers are offering Hadoop cluster services over their platforms. For instance, Microsoft Azure provides with HDInsight a full manageable Hadoop respectively Spark Cluster [30]. The user could profit from fast provisioning and also from less costs since no on-prem hardware cluster infrastructure is required. Therefore, Hadoop (or better Spark) by Cloud Computing is also a noteworthy option for our data science project in the future, although it may be over sized at the moment. It is worth to mention that each Hadoop distribution in table 3.1 comes up with a minimum of services from the Apache foundation: YARN, HDFS, Hive, Pig, HBase, Kafka, Storm, Mahout, HCatalog etc. But they may use a proprietary implementation of them e.g. MapR uses own HDFS instead of Apache HDFS implementation. It is also worth to note that a few cloud providers are offering Hadoop cluster services over their platforms. For instance, Microsoft Azure provides with HDInsight a full manageable Hadoop respectively Spark Cluster [30]. The user could profit from fast provisioning and also from less costs since no on-prem hardware cluster infrastructure is required. Therefore, Hadoop (or better Spark) by Cloud Computing is also a noteworthy option for our data science project in the future, although it may be over sized at the moment.

Nevertheless, Intel, ERM/Pivotal and IBM partnered with Hortonworks or Cloudera 3.1. These vendors therefore have the same extensions as their partners offer, but also can provide exclusive tools for users (e.g. IBM Big SQL [29]). If we take this into account, we find that there are mainly 3 different Hadoop distributions: Cloudera with CDH, Hortonworks with its Hortonworks Data Platform (HDP) and MapR with its own HDFS. In addition to the distributions shown in 3.1, there are many other Hadoop distributions such as Altiscale from SAP or Elastic MapReduce from Amazon. Both run only in a cloud environment. The Forrester Wave shows also the three global players Cloudera, Hortonworks and MapR (see 1.1). Since cloud solutions do not play a role for this use case, Microsoft, Google or Amazon will not be considered.

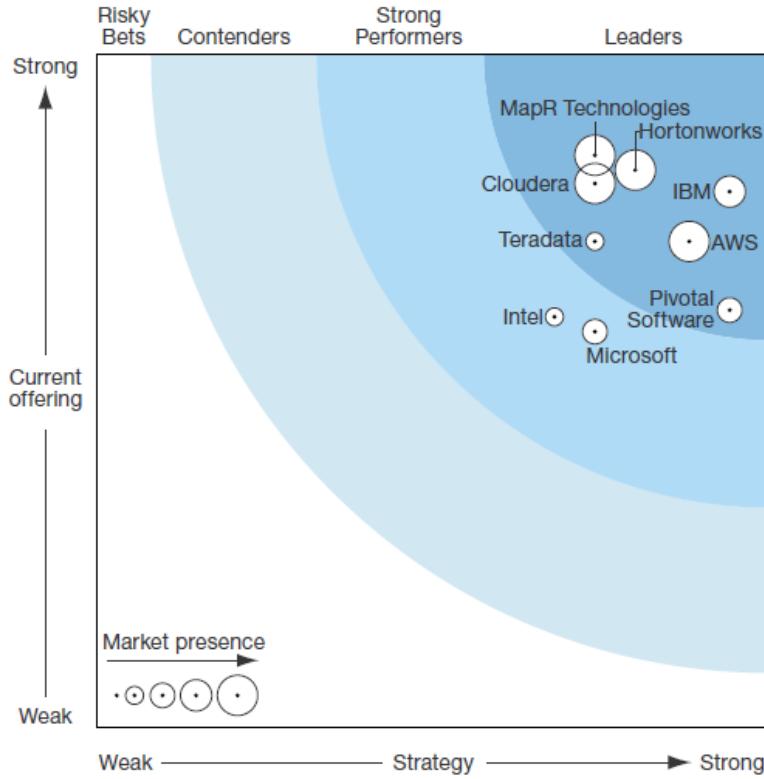


Figure 1.1: Forrester Wave of Hadoop Distributions [46]

As we can see from 1.1 the three mentioned Hadoop distributions are building a cluster and Cloudera even overlaps with MapR slightly. That means, they are providing similar products and enjoy a similar market position. From the Forrester Wave (1.1) we can derive that MapR seems to have the best offering. However, since we want to use one of them for our analytics project we have to dig a step deeper and make a decision based on well-defined criteria.

1.2 Definition of criteria

Hadoop distributions are rated using selected criteria, with scores ranging from 1 (very poor) to 5 (very good). Based on the research carried out, the following 20 clear criteria can be determined:

Criterion	Description
Batch Processing	Batch process jobs can run without any end-user interaction or can be scheduled to start up on their own as resources permit
Cloud Support	Cloud compatibility of Hadoop distributions
Cluster Scalability	Creating new clusters ad-hoc and set up cluster size (e.g. 4,8 or 12 nodes)
Data querying possibilities	Ways to query Hadoop with SQL (e.g. Hive, Stinger, Spark SQL,...)
Database Support	Additional databases in Hadoop (e.g. PostgreSQL, MongoDB, HBase, Impala etc.)
Disaster-Discovery	There is a fallback option available in case of unexpected Hadoop errors
Ease of Use	The complexity of the usage of the Hadoop distribution
High-Availability	Self-healing across multiple services or single failure recovery (i.e. fail-over clusters)
Install complexity	Is there a simple installation routine or does it require lot of manual configuration (e.g. Multi node Cluster...)?
Licensing and Pricing	Cost model of Hadoop distributions
ML support	Providing interface for enlarged machine learning tasks (e.g. Spark MLlib...)
Operating System Support	Guaranteed and certified OS compatibility
Performance	Cluster performance in case of parallelized MapReduce Jobs
Programming Language Support	Support of data science programming languages (Python and R)
Streaming Analytics	Vendor offering real-time analytics capabilities (e.g. IoT platform hub..)
Support of secondary services	Services on top of Hadoop (e.g. Zookeeper, Spark...)
Third party module integration	Usage of additional components from other vendors
User Support/Community	How fast is the response in the community?
Expertise	Experience with Hadoop (life time of company, ...)

Table 3.2: Definition of criteria

The weight scale ranges from 1 (trivial) to 4 (very important), with the focus on our Santander Bicycle project. This means that cloud support,

for example, is a crucial criterion in many use cases, but has only little importance for our project which leads to a trivial weight. A percentage weighting is not applied because 20 criteria would result in a fine-granular distribution (which is not good to read at all). Therefore, absolute values are used for the comparison 1.2 in the next section.

1.3 Weighted evaluation matrix

As already mentioned in Section 1.1, the three major Hadoop distributions are: Cloudera CDH 5.15.0, Hortonworks HDP 3.0.1 and MapR 6.1.0. With regard to the defined comparison criteria, a weighted decision matrix can be mapped. At the same time, this matrix represents the starting point for the decision of a Hadoop distribution. It cannot be denied that a certain degree of subjectivity is included in the evaluation (table 3.2). In addition, only free editions are compared. For example, if a feature only exists in the premium version, this feature will be rated one (worst) since we dont want to spend money.

Criteria	Weight	KO	MapR	Hortonworks	Cloudera CDH
Batch Processing	3		5	5	5
Cloud support	1		3	4	5
Cluster Management Complexity	4		3	3	4
Cluster Scalability	2		4	4	3
Data querying possibilities	3		4	4	5
Database Support	4		5	3	5
Disaster-Discovery	2		1	4	1
Ease of Use	2		3	3	4
High-Availability	2		1	4	5
Install complexity	2		3	2	3
Licensing and Pricing	2		1	5	2
ML support	4	x	5	5	5
Operating System Support	3	x	2	3	2
Performance	4	x	5	4	4
Programming language support	3	x	5	5	5
Streaming Analytics	1		5	5	5
Support of secondary services	3		3	5	4
Third party module integration	2		2	4	1
User Support & Community	2		2	4	3
Expertise	2		4	5	3
Total Score			179	209	199
Legend:	(bad) 1	(is okay) 2	(mediocre) 3	(good) 4	(ideal) 5

Figure 1.2: Weighted Comparison Table of Hadoop Distributions

The comparison in figure 1.2 clearly shows that there are almost no significant differences between the selected distributions. Although there are a few deviations such as disaster-discovery or pricing model, the overall distribution of points is relatively the same. From the comparison matrix it can be deduced that Hortonworks HDP achieves the highest score and is probably the best option for our project at the moment. It is possible that after the first test phase (see chapter 1.5), it turns out that the distribution is not convenient after all which may require to update the evaluation matrix. Therefore figure 1.2 can be considered as a continuous iterative updateable weighted comparison matrix that will likely be updated on further sprints. Another representation of the evaluated Hadoop distributions, a spider chart may be appropriate as it makes it easier to read and to identify outliers or

similarities even quicker as with a raw table.

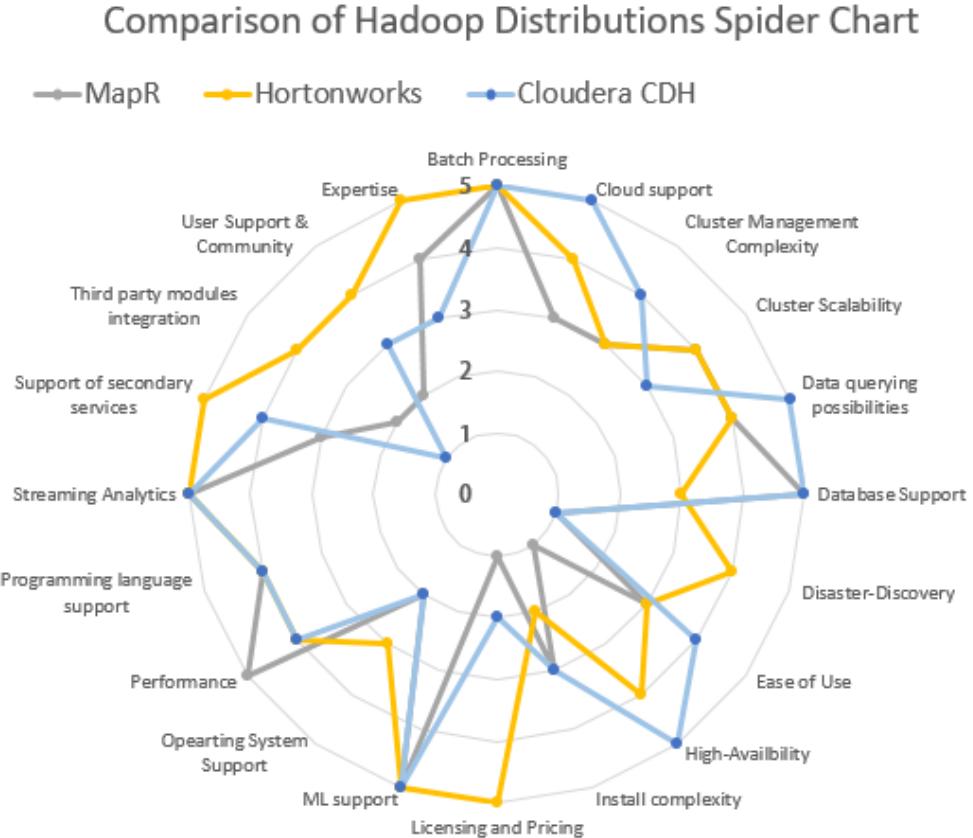


Figure 1.3: Spider Chart of compared Hadoop Distributions

The KO-criteria from figure 1.2 are all satisfied by the selected Hadoop vendors. One criterion that stands out is the license model, which was rated with 1 point for MapR and 2 points for Cloudera. This is due to the fact that the price models of both vendors are not transparent. So Cloudera offers a free community edition but to use the Cloudera Manager (the real strength of Cloudera) you have to pay again. Only Hortonworks offers a complete Open Source package with except the Hortonworks business support is fee required. However, Hortonworks has a very active community as well. MapR comes off worst with the criterion Support of secondary services, because there is no Spark and HBase support in the "Converged Community Edition" [26], but both are important Hadoop components for our project. So if we wanted

to use MapR, we would have to use the paid version. There are also some other differences in figure 1.3 such as Disaster-Recovery which is only on Hortonworks HDP completely free. An interesting aspect of the comparison is the performance criterion where MapR has the most points. Since this criterion is a KO-criterion that means it is a very important component for our project, it makes sense to examine the performance comparison between the distributions in a more detailed fashion.

1.4 Performance Comparison with micro benchmarks

For a performance comparison it is good to know how fast they compute when running in a concurrent mode. For this task MapReduce Jobs like WordCount or Data File System I/O (DFSIO) Read/Write might be helpful. Following figures are extracted from a sophisticated evaluation work by Alturos [1].

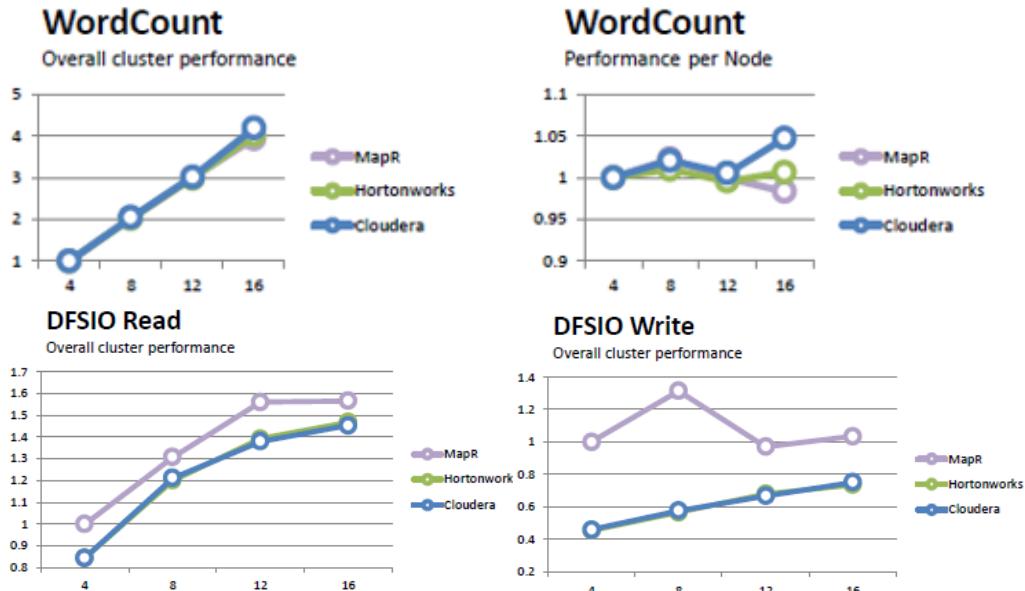


Figure 1.4: Micro benchmarks on DFSIO / WordCount MapReduce Jobs

For the performance comparison of figure 1.4 Hadoop cluster each one with 16 nodes have been established in order to measure computation time. An interesting point is that all three Hadoop distributions almost achieve the same overall speed on the famous WordCount MapReduce job. But if we

look at the performance per node then there is a small difference between them. This could be due to a certain error rate in the execution of the job. If the same test were repeated, the results would be negligibly different. The DFSIORRead/Write (figure 1.4) job shows MapR is definitely faster than Hortonworks and Cloudera. The reason for this might be the fact that MapR uses its own HDFS which is according to the vendor seven times faster than the original one from Apache. This is the reason why it gets five points in the evaluation matrix (see table 3.2) whereas other ones only reach four points. Hortonworks and Cloudera seems to have the same performance on the DFSIO job because both using Apache HDFS.

In overall, Hortonworks HDP wins the competition for the moment, because they have the longest experience and are completely compatible with secondary services. Based on the evaluation matrix from chapter 1.3 and the performance measurements, it can be deduced that MapR is the most powerful Hadoop distribution on the market today, but when considering the other criteria, Hortonwork's Hadoop is simply more convincing. Especially the complete Open Source guarantee at Hortonworks is a decisive criterion for the choice of this distribution.

Beside from the comparison, we would generally recommend to use rather Spark than Hadoop since Spark is around 100 times faster due to in-memory processing. Also Spark provides the most important libraries (ML, Streaming, etc) for data science and works well on top of Hadoop (thanks YARN). It can even be installed in a standalone manner, even though it doesn't benefit from distributive multi mode computing.

1.5 Installation of sample Hadoop Distribution

Hortonworks offers a configurator on their website that can be used to check products, operating systems, databases, browsers, JDKs and supported processors. It is noticeable that Ubuntu is only supported up to 16.04 LTS. However, our VMs are already v.18.4 LTS. This means Hortonworks does not officially support our installed OS [16]. Also Cloudera CDH 5.15 [8] and the newest MapR 6.1 distribution [27] only support Ubuntu 16.04 LTS (Xenial). Anyway, it is still worth a try to set up HDP Cluster on our virtual machines since the vendors are not explicitly warning Ubuntu 18.04 is not supported. So it may work. For the installation of HDP 3.0.1 the Ambari Wizard [13] will be used.

First of all, we have to configure the **etc/hosts** since every node in the clus-

ter must be able to communicate with each other. The hosts file should look like this:

1	127.0.0.1	localhost
2		
3		
4	141.59.29.111	i-hadoop-01.informatik.hs-ulm.de i-hadoop-01
5	141.59.29.112	i-hadoop-02.informatik.hs-ulm.de i-hadoop-02
6	141.59.29.113	i-hadoop-03.informatik.hs-ulm.de i-hadoop-03
7	141.59.29.114	i-hadoop-04.informatik.hs-ulm.de i-hadoop-04

Figure 1.5: etc/hosts configuration file

After that step, it is necessary to add public key authentication with SSH. The master node (**i-hadoop-01**) should login to its worker nodes without using a password. We can achieve this goal by generating a private / public key pair and distribute the public key to the 3 worker nodes by using **ssh-copy-id** command. Also for using Ambari the Hadoop user need sudo execution without password. Adding an additional line to **visudo** configuration file should solve that issue. At next, all nodes need to have Java installed. Since Ubuntu 18.04 doesn't come up with a default Java JDK we install Oracle JDK 1.8 manually on each node. After that it is time to download the Ambari repository file to a directory on our Hadoop master host. The Ambari v. 2.7.1.0 is used for installation. With the command **apt-get install ambari server** the server will be installed on i-hadoop-01. Afterwards, we can start the server. Now we can go to the Ambari surface via following link: <http://i-hadoop-01.informatik.hs-ulm.de:8080>. The default login credentials are **admin admin**. After successful login, the real installation process begins.

Installer

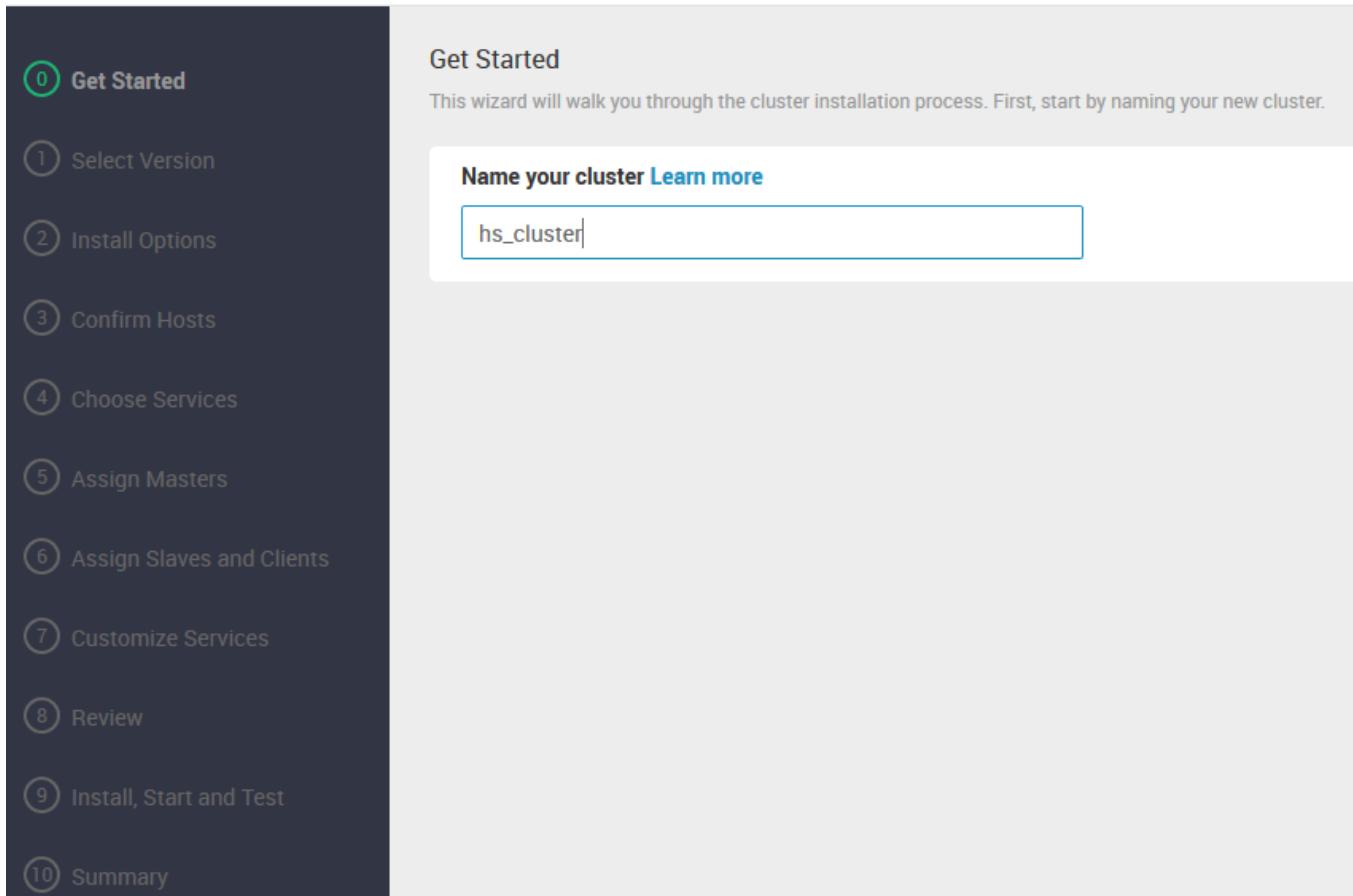


Figure 1.6: Ambari Wizard Installer

Next, we click through the installation routine until the point appears where the cluster nodes are defined. At this point the fact that we are using Ubuntu 18.04 will cause some troubles. If we adding the hosts Ambari is trying to check some dependencies and is running also some check routines. It will fail to add the worker nodes to the cluster since OS is not supported (termination condition). The solution to this problem would be a downgrade of the OS but we dont have the privileges to do this. So we manually changed the `/etc/issue`, `etc/lsb-release` and `etc/os-release` file to Ubuntu 16.04 version. Of course, a backup of the original ones has been created as well. After this workaround, the check condition will not fail because Ambari con-

sider our OS as Ubuntu 16.04. Obviously, that's a dangerous operation, so we change it back to original state after installation has been completed. In the next step we are choosing our Hadoop services that may be relevant for our data science project. For a start we choose following services: HDFS, YARN, MapReduce2, Tez, Hive, ZooKeeper, Ambari Metrics, SmartSense, Spark2 and Zeppelin Notebook. Interesting to see is that the proprietary service SmartSense is the only one that cannot be deselected. We have to install it, even if we don't want to use it. That's certainly not very user-friendly. However, it is simple to add further Apache services on running Ambari but more difficult to uninstall them so we do not enable all possible Hadoop services from the beginning. In addition, some services are not of interest at the moment, hence they would only consume storage without having a practical effect.

After some additional configuration steps we can choose which services should run on master or worker nodes. Also we have to decide which node should be data node and node manager. In this cluster all four VMs are data nodes as well as node managers. A database connection is also mandatory at this step. Depending on selected services there are several database connections. For instance, the service Hive needs a working SQL database connection. We can choose between MySQL (MariaDB), PostgreSQL and Oracle. For a first test, MySQL is used as Hive database. We can change the database connection setting on Ambari at any time. If the configuration has been finished, a summary with all decisions taken is shown. That's the last possibility to change a fundamental Hadoop setting. It is worth to check all configurations carefully. Afterwards there is no go back option and the uninstall routine of Ambari is complicated and might lead to tremendous work of removing operations. That is also negative criteria of HDP and Ambari. As soon as the installation completes we can switch to the Ambari monitor. Unfortunately, almost all selected services were offline or could not start properly. A deeper investigation showed that there were some permission problems with HDFS and the Hadoop user. The HDFS folder was only for root user but not for user „Hadoop“. With the command **chown -R hdfs:hadoop /hadoop/hdfs** the correct permissions could be granted.

Another problem was that YARN could not start the name nodes due to connection refused error messages. As **netstat -tupln — grep 50070** shown the service was only running for localhost address. After changing the binding address to **0.0.0.0** on HDFS, the communication between the cluster nodes was working and the name nodes could properly start. However, after fixing

this issue, the most services were running. The community of Hortonworks was quite helpful in this case. We got an answer within four hours whereas for getting an answer in Clouderas community took three times longer.

The data node of the master node was not running correctly. It could be started from Ambari and after few seconds the data node went offline. A look at the error log of this node showed that the cluster id was already used by another node. Removing the HDFS data folder of the data node solved this issue because afterwards it generates a new cache file with a new cluster id.

The following screenshot shows a (successfully) completed Ambari and HDP installation on our VMs.

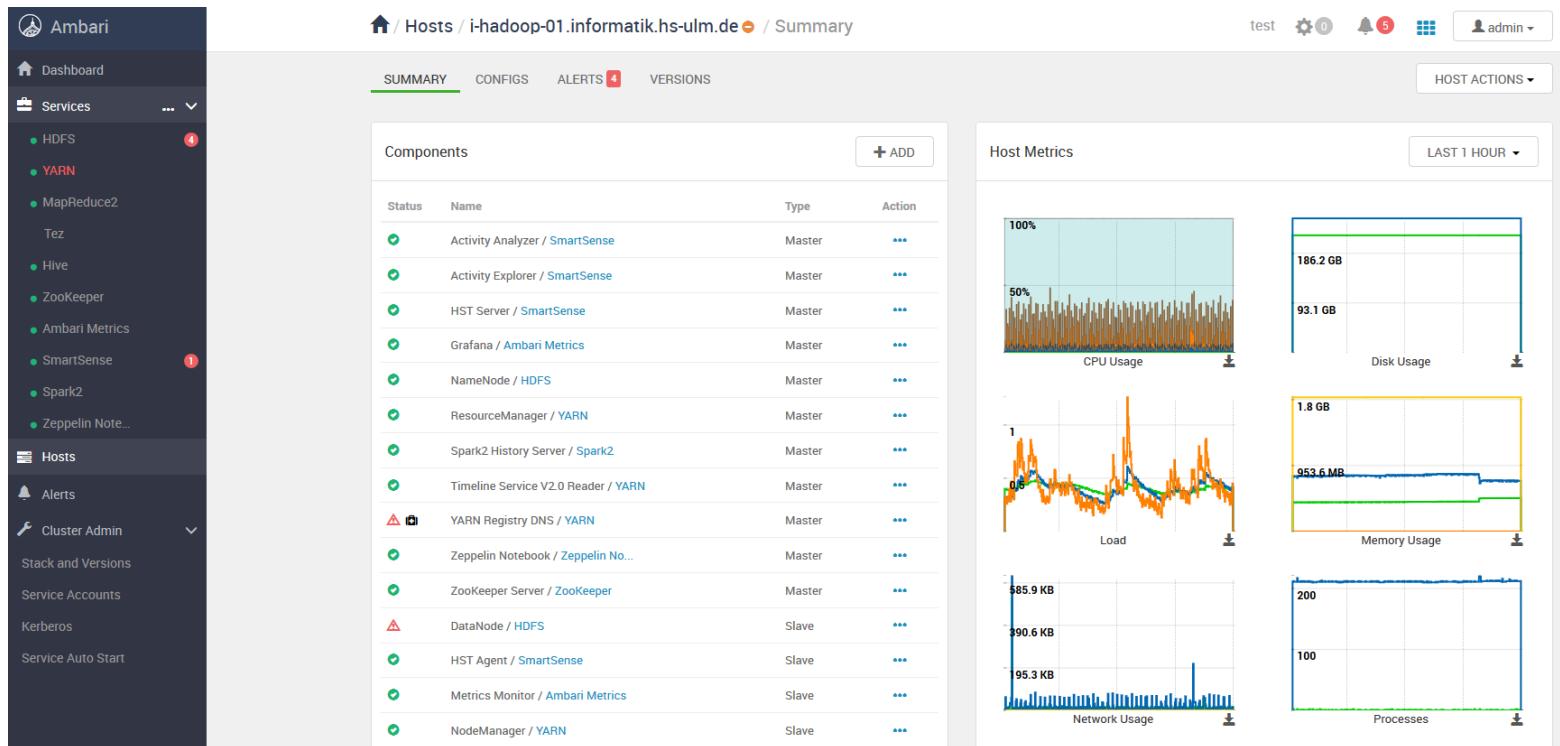


Figure 1.7: Ambari Surface, Host metrics on master node

Ambari has been set up and most Hadoop services are running (see Figure 1.7). Ambari shows lot of statistics that might be helpful for cluster management. We can control each service and either turn it off or set it into a maintenance mode. The ladder one encapsulates the service without influenc-

ing other services (high availability). Overall, Ambari is a powerful Hadoop cluster management tool that has lot of control options and supervising tools but also lacks in removing services and usability (e.g. no Ubuntu 18.04 support, only Python 2.7.x support, etc). The fact, that there is no official uninstaller from Hortonworks makes HDP with Ambari a risky installation. On the other hand, it is completely open source and has a strong community. Furthermore, HDP 3.0.1 comes up with a bunch of useful Hadoop services and well documentation so that we will likely stay on HDP.

Now it is time to run some MapReduce jobs on our newly created Hadoop cluster.

1.6 Run MapReduce Jobs

With four workers (master is also a worker node) a significant speed advantage over single computing should be achieved. First, with a MapReduce job, the accuracy of the decimal places of PI is determined using the quasi Monte Carlo method. This MapReduce job is also offered by the Apache foundation (Source). In the console of the master node the MapReduce job can be started with YARN as follows:

```
yarn jar /usr/local/hadoop/share/hadoop/mapreduce/
    hadoop-mapreduce-examples-
3.1.1.jar pi 16 1000
```

The *PI job* calculates 16 maps where each map process contains 1000 samples. After a short time, the job is finished and the following output is achieved:

```
Job Finished in 39.463 seconds
Estimated value of Pi is 3.14250000000000000000
```

Figure 1.8: PI Quasi Monte Carlo method executed as Cluster job

The same test has been performed with two worker nodes (i.e. two of the four workers were temporarily disabled via Ambari). After that the job took much longer (about 56 seconds). This shows that the Hadoop cluster works and already has a performance advantage over single computing. The next test was a simple WordCount job in Java that calculates the frequency of unique words in an input file distributed across the cluster. A MapReduce job consists of at least a mapper method, a reducer method and a driver method which is often the start method. Before executing the new job one

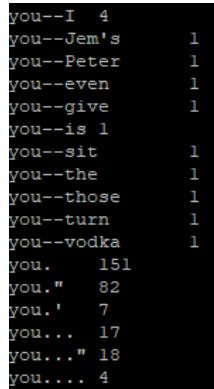
has to put the file of interest on Hadoops HDFS. This is simply done as shown below:

```
hadoop fs -put big.txt
```

The put command will upload the file into the HDFS (i.e. it is stored fault-tolerant and replicated 3 times). After that, the following command runs a WordCount job on this text file:

```
hadoop jar hadoop-mapreduce-examples-3.1.1.jar  
wordcount wordcountFile big.txt
```

Note that now „Hadoop“ is used, which starts also the YARN service. There is no difference between executing this with the Hadoop or YARN command. (Although latter is newer and should be used). The output of the „reduced“file looks like figure 1.9.



A black and white screenshot of a terminal window displaying the output of a WordCount job. The output lists words and their counts, separated by spaces. The words are listed in descending order of frequency. The counts are in the right column. The output is as follows:

you--I	4
you--Jem's	1
you--Peter	1
you--even	1
you--give	1
you--is	1
you--sit	1
you--the	1
you--those	1
you--turn	1
you--vodka	1
you.	151
you."	82
you.'	7
you...	17
you..."	18
you....	4

Figure 1.9: Output of WordCount job

Quick to recognize that the words were separated and counted by spaces. In fact, the input text was split afterwards and each word was mapped as "key value pair", where the key in this case is the word. The original text file was 6.4 MB, while the newly created file is only 0.93 MB in size. This shows another potential benefit of the cluster, namely the possible separation of files to minimize the total size of the initial data set. The running time was around 10 seconds. The same job has been executed locally on a commodity Windows client in NetBeans as well and it took around 8 minutes. There is a plenty of further typical MapReduce jobs (e.g. terasort, randomwriter, sudoko, etc), but for an evaluation these two jobs should already show that the cluster is working apparently faster than a single node.

1.7 HDP Issues

Due to the workaround in Chapter 1.5, a few services could not start properly. Also, permissions for some subfolders were set incorrectly, hence permission denied errors occurred when running the services. During the test phase some other problems were discovered, which are described below. Table 3.3 also shows the solution that was found to resolve these issues.

Issue	Solution
No starting NameNodes, Connection refused.	Go to HDFS→ Filter at https 0.0.0.0 and localhost address on master node http 127.0.0.0:50070 → 0.0.0.0:50070 Hive could not start on Node 2 since database was missing (mysql) Permission on HDFS was set to root and not to user hadoop
Wrong Permissions on HDFS	tail hadoop-hdfs-namenode chown -R hdfs:hadoop /hadoop/hdfs
Hive Database not found on i-hadoop-02	<ul style="list-style-type: none"> • Start mysql • Show databases; • create database hive; • CREATE USER 'hive'@'localhost' IDENTIFIED BY 'hive'; • GRANT ALL PRIVILEGES ON . TO 'hive'@'localhost'; • CREATE USER 'hive'@'%' IDENTIFIED BY 'hive'; • GRANT ALL PRIVILEGES ON . TO 'hive'@'%'; • GRANT ALL PRIVILEGES ON . TO 'hive'@'localhost' WITH GRANT OPTION; • GRANT ALL PRIVILEGES ON . TO 'hive'@'%' WITH GRANT OPTION; • FLUSH PRIVILEGES;
Bind address in mysql-conf (0.0.0.0):	<ul style="list-style-type: none"> • netstat -tupln — grep 3306 • Nano /etc/mysql/mysql.conf • Change bind-address to 0.0.0.0 • Systemctl mysql restart • Service mysql restart
YARN DNS Server not starting, connection refused	<ul style="list-style-type: none"> • Change Port 53 to other Port • Set right permissions on YARN folder

Issue	Solution
Master Node worker is not starting (Conflicting cluster id)	<ul style="list-style-type: none"> • Stop service • sudo wget archive.apache.org/dist/zeppelin/zeppelin-0.8.0/zeppelin-0.8.0-bin-all.tgz • sudo tar -xf zeppelin-0.8.0-bin-all.tgz • Move interpreter to /usr/hdp/current/zeppelin-server/interpreter/ • sudo chown -R zeppelin:zeppelin /usr/hdp/current/zeppelin-server/interpreter/ • Start service
Umask is changing from 0022 to 0002	<ul style="list-style-type: none"> • sudo gedit /.bashrc • Add line umask 022 • Logout and login
FAILED [2018-12-0813:36:00.609] java.io.IOException: mkdir of /hadoop-/yarn/local/file-cache/60_tmp failed	<ul style="list-style-type: none"> • YARN permissions • sudo chown yarn yarn • sudo chmod 750 yarn

Table 3.3: Problems and solutions of HDP

The described problems and their solutions of table 3.3 show that many setting options should be checked carefully during the Ambari installation routine (e.g. if a MySQL database exists). Furthermore, it can be assumed that most of these issues can be explained by the „officially“ unsupported operating system version. However, after the fix, all Hadoop services in Ambari had a green status point.

1.8 Recommendation

The evaluation showed that Hortonworks with the Hadoop distribution HDP is best suited for our goals. Ambari provides an efficient, comprehensive cluster management tool for monitoring and controlling the Hadoop cluster. In

addition, an active community and complete open source freedom can convince more than its competitors. Although all three Hadoop distributions may have their individual advantages and disadvantages (see evaluation matrix in figure 1.2), HortonWorks was able to convince us due to the reasonably good documentation and the easy to follow Ambari installation routine. It should be noted, however, that the Cloudera Manager installation wizard also has good documentation and a comparable installation effort. However, we liked the Ambari interface more than the CDH one. In addition, Cloudera's pricing is not entirely transparent and third-party module integration is rather difficult. For example, since CDH 5.11 there has been a proprietary „data science workbench“ service replacing Jupyter and IPython notebooks. HortonWorks, on the other hand, offers Zeppelin notebooks that resemble the famous Jupyter notebooks. Apart from that, Jupyter could also be installed as a third module during the test phase of HDP, even though such installations are naturally risky and are not officially supported by HortonWorks.

It is interesting that MapR, HortonWorks and Cloudera offer Docker images for immediate testing. For example, MapR offers a VMware image free of charge. However, all these packages have the disadvantage that they only have a single node cluster and limited functionality. For example, I could not select any services or create new ones in MapR and the function to add more worker nodes to the existing cluster was also not available. Therefore, such images or sandboxes are ideal for a quick check of the user interface and some standard functions, but rather insufficient for a complete range of functions or a productive system.

In conclusion it can be said that all three Hadoop distributions would meet the requirements of a big data project or our project perfectly, but HortonWorks could simply convince more with the overall score than the more commercially oriented counterparts. Ultimately, the performance also depends on the hardware of the individual cluster nodes. Due to a low physical RAM, MapReduce jobs even with MapR's high performance HDFS would be processed slowly by YARN and would probably lead to *MemoryOutOfExceptions*. Because Hadoop and every single service on top needs resources and these can quickly bring a commodity computer with standard equipment to its knees. In any case, the HDP cluster is ready for big data tasks and will be extremely helpful to us in the further course of the project. With Apache Spark, even better performance values might be achieved due to its „in-memory“ engine. The test cluster with the four VMs also showed which

special features and requirements you should pay attention to (e.g. operating system version).

1.9 Hadoop Rollout

For the Hadoop rollout ten fat clients each with 32 GB RAM and 2 TB HDD as well as 128 GB SSD have been used. The nodes were connected with a 48port Gigabit switch and a routable gateway, so that all nodes were in a single local domain. For various reasons (e.g. security) the installation of a hypervisor (Proxmox) has proven to be successful, with which any number of VMs can be generated and configured (figure 1.10).

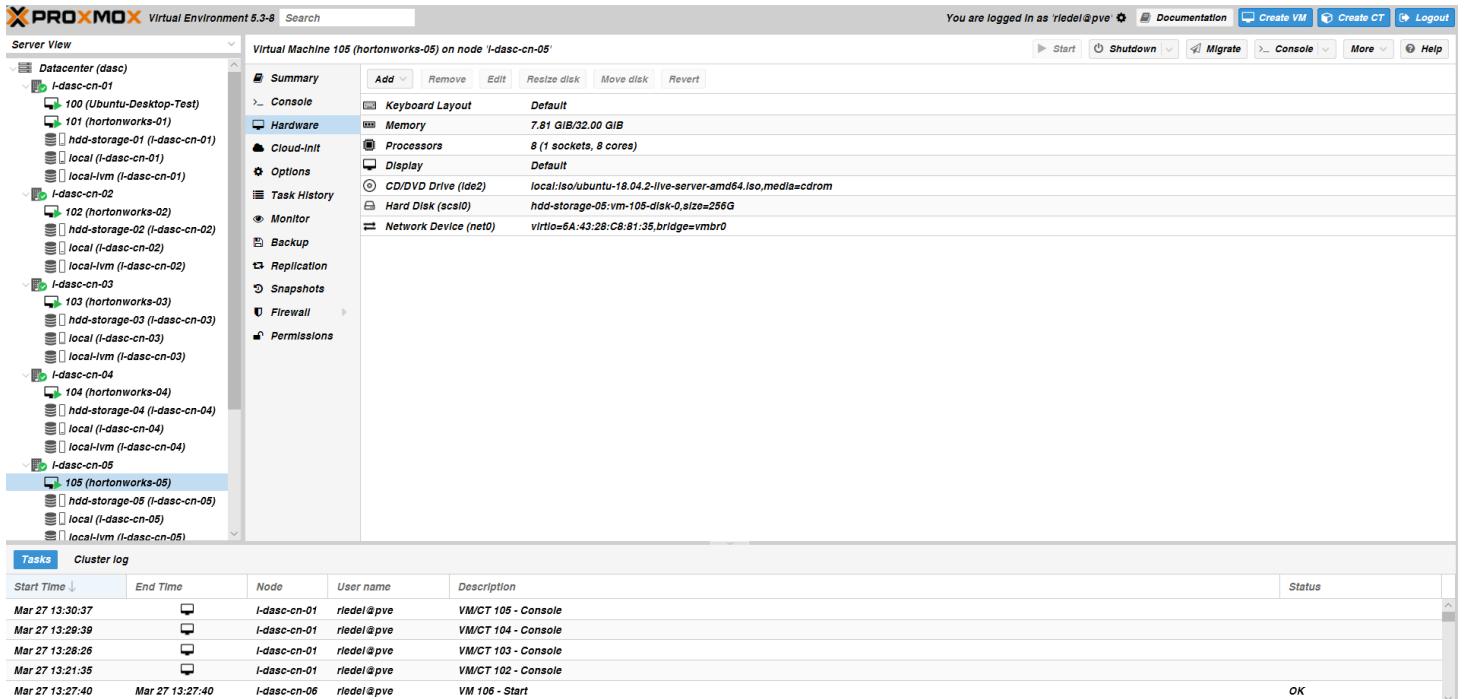


Figure 1.10: Proxmox surface

A VM with Ubuntu Server 18.04.02 LTS was installed for each node. Each VM got an allocated memory of 256 GB and a reserved memory for the Logical Volume Manager (LVM). Furthermore, sparse files were allowed and VirtIO SCSI interfaces were configured for the individual VMs, which should guarantee maximum VM performance ([37]).

For the first rollout attempt, six cluster nodes were configured for Hadoop. The remaining four nodes were later added to the existing Hadoop cluster via Ambari surface. The chosen Hadoop stack was HDP from Hortonworks, as already explained in the Hadoop evaluation part of this work. In addition, Webmin was installed on the master node, which offers extensive monitoring services on the node. For example, the current CPU load, memory usage, kernel information and much more can be displayed. Figure 1.11 shows the surface of Webmin. Although monitoring services are also offered with Ambari Metrics, they run on top of Hadoop and can only be called when Ambari is running as well ([17]).

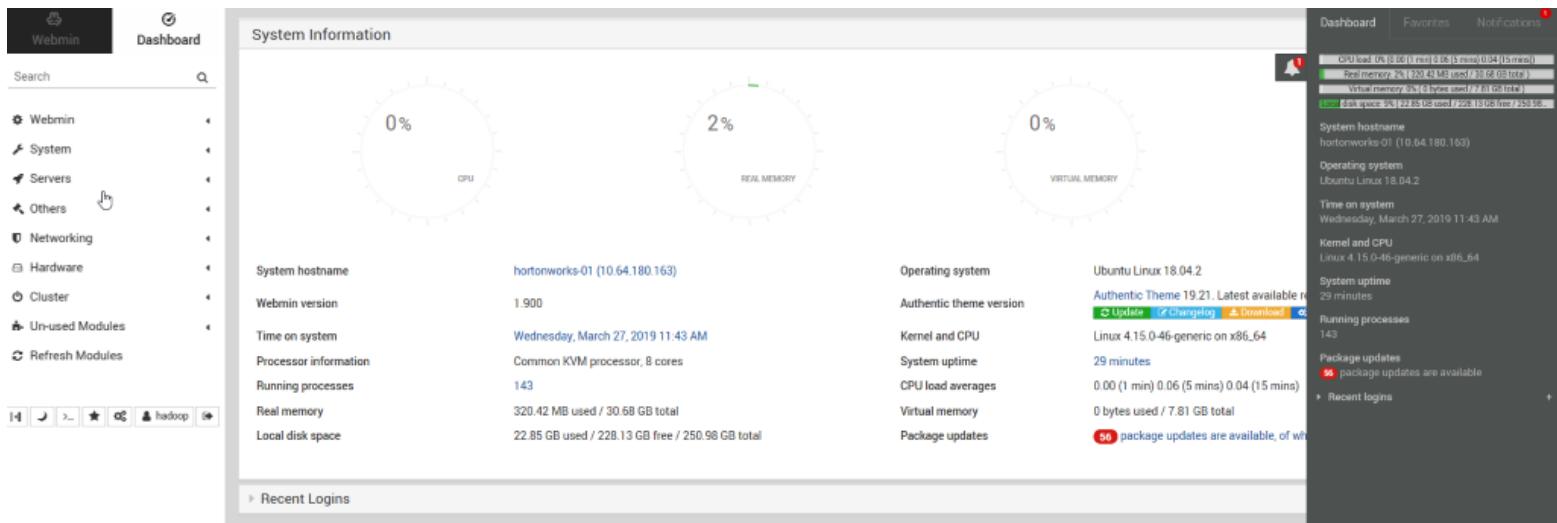


Figure 1.11: Webmin dashboard

Before the Ambari Wizard can install HDP, some pre-configurations have to be done on each VM. On the one hand, „ulimit“ must be set to at least 10000, because Ambari installs several thousand dependencies. On the other hand, a password-less ssh authentication is necessary so that the master node can connect to its worker nodes without entering a password. In addition, the master node must be able to execute „sudo“ commands without entering a password. This can be done by editing the „visudo“ file and adding „username ALL=(ALL) NOPASSWD:ALL“. Another necessity is to add the IP addresses and hostnames of all cluster nodes under „/etc/hosts“. This

must also be done on each node individually. The following table shows the current host configurations of a worker respectively slave node:

IP-Address	List of hostnames
10.64.180.163	hortonworks-01.dasc.cs.thu.de hortonworks-01
10.64.83.106	hortonworks-02.dasc.cs.thu.de hortonworks-02
10.64.79.161	hortonworks-03.dasc.cs.thu.de hortonworks-03
10.64.227.154	hortonworks-04.dasc.cs.thu.de hortonworks-04
10.64.159.100	hortonworks-05.dasc.cs.thu.de hortonworks-05
10.64.204.57	hortonworks-06.dasc.cs.thu.de hortonworks-06

Table 3.4: etc/hosts of a worker node

Databases (Hive, Ranger, Druid...) are created by the Ambari Wizard as PostgreSQL databases automatically. Java as well as a corresponding JDBC connector are required for the individual services to execute database statements. The node „hortonworks-01“ (see table 3.4) is both a master and a worker node. On the master node the Ambari Wizard installs Ambari server. The Ambari Agents, which are required for communication in the cluster, are installed on each worker node. Afterwards, the Ambari Wizard can be called under the following URL:

hortonworks-01.dasc.cs.thu.de:8080.

The installation routine guides the user through all necessary steps. It is important that the Ambari agents are installed on the individual worker nodes, otherwise the Ambari Wizard cannot add the nodes to the cluster. The most important step is probably the selection of the HDP services. Similar to the evaluation step with the VMs described in section 1.5, the identical service packages were selected, i.e. YARN+MapReduce2, Tez, Hive, HBase, Pig, ZooKeeper, Ambari Metrics, Spark2, Zeppelin and SmartSense. But also some new services were added for the production cluster: Storm, Accumulo, Infra Solr, Atlas and Kafka. Since the cluster is to become a long-lived high performance cluster, it might be reasonable to rollout the playground for

distributed streaming platforms like Kafka or Storm which can be used for streaming analytics use-cases.

Selection of services that are running on top of Hadoop is an important part of the Hadoop cluster setup process. Following services have been chosen from the HDP stack:

all 6 virtualized nodes are DataNodes (workers) and each have a YARN NodeManager (which takes care of the resource distribution and monitoring of a node).

Furthermore, all master components run on „hortonworks-01“. However, the cluster is fault-tolerant, so if the master node fails, the worker „hortonworks-02“ becomes the active master. This is made possible by the secondary NameNode service that runs on another worker node.

Once the individual accounts have been created for the different services, a final briefing is performed by Ambari before the cluster is started. A useful feature of Ambari is to download the complete configuration as a JSON template [43]. This makes another Hadoop installation much easier because the template can be reused.

This time, the installation process was done without any problems. HortonWorks (HDP 3.1.0) supports recently Ubuntu 18.04.02 LTS which makes tweaking of the operating system superfluous. Initializing and starting Hadoop services may take a few hours, depending on the size of cluster. The productive cluster now runs on Ambari 2.7.3.0 and HDP 3.1.0. By default, Zeppelin comes with a Python 2 kernel. However, it is possible to switch to the Python 3 kernel (IPython with Python 3.x).

The Ambari interface is intuitive to use. A first look at Ambari Metrics showed that the services worked properly and all workers in the cluster were active. Only YARN Registry DNS did not seem to start due to a connection refused error because Hadoop relies heavily on a functioning DNS server [18]. However, changing the YARN RDNS Binding Port from „53“ to „5300“ solved the problem.

Remark: the same issue happened in the evaluation part where a port conflict prevented a successful start of the Ambari server.

The physical hardware configuration of the cluster consists of 10 fat nodes,

as figure 1.12 shows.



Figure 1.12: Physical Hadoop cluster

These nodes were locally connected with a switch and a gateway so that all nodes could communicate with each other. Unfortunately, access outside the intranet was not possible, as the necessary infrastructure measures on the part of the data center are still pending. In principle, however, it would be possible to access the cluster from outside using VPN and a reliable authentication method.

First tests with the six configured Hadoop nodes could be carried out successfully. These tests were based on the Zeppelin notebooks from the previous data profiling chapter 2, which already worked successfully in the virtual cluster. Compared to the virtual cluster this time the execution was much faster, because more RAM was available and more workers (six! instead of four) were used.

Thus the cluster (figure 1.12) is in an operational and ready configured state. Of course, it is possible that a service may fail on its own or no longer run properly over time. In the evaluation phase, for example, it was shown that the YARN Timeline service fails more frequently. Usually, however, a restart of the corresponding service via the Ambari interface is sufficient.

Most Hadoop services also run autonomously, i.e. a corrupt service cannot block other running services (exception: HDFS). With the new ready to use Hadoop cluster, further data profiling action of the bicycle data can now be performed in the cluster.

2 Data Profiling

Data profiling is the process of reviewing source data, understanding structure, content and interrelationships, and identifying potential for data projects. For the project, the Santander Bicycle data will be profiled more closely [44]. For this I mainly used Zeppelin notebooks and the HDP cluster with four worker nodes initialized in chapter 1.5. Zeppelin works similar to Jupyter notebooks and also supports magic commands. However, Zeppelin stores the notebooks in .json format, while Jupyter notebooks (python) uses .ipynb. A corresponding import of Zeppelin notebooks into Jupyter is therefore not possible and should be considered when working with one of them. Furthermore, as already mentioned, HDP only supports python 2.7.x, so all Python code in Zeppelin is python 2.7.x as well. (The Nominatim server and Django already use Python 3.7.x).

2.1 Data Preparation

The Santander bicycle data is available on the public Transport for London TfL website. The structure of these data is as followed:

- CycleCounters: has information about speed, bike details (length, wheels, axles, class...), direction, timestamps and serial number.
 - Blackfriars May and June and July
 - Embankment May and June and July
- CycleParking: contains proprietary map material
- CycleRoutes: contains geographical spatial data (GeoJSON) and .KML files
- Usage-stats: information about rental stations, duration, start / end time

With regard to the folder structure shown above, CycleRoutes and Usage-stats are decisive. In „CycleCounters“, apart from the maximum driven speed (56 mph!) with a borrowed Santander Cycle, no noteworthy information is contained. Therefore, only the CycleRoutes and Usage-stats data are considered in the first step of Data Profiling. The CycleParking folder could contain useful map material, but I only know from Google Earth that can read .kml files and that it is neither compatible with Geopandas nor with Nominatim. Most raw data are available either as .csv files or as .xlsx files. One problem that arose when reading the usage-stats was that the folder contained 156 csv sheets and each file was about 30 MB in size. Reading and concatenating as Pandas DataFrame unfortunately completely occupied the maximum free memory of the master node, so the python job could not be finished and was in an endless run mode. To work around the problem, we started a Spark job with the PySpark API on the cluster instead of the Panda DataFrame. The spark job had created a DataFrame from 156 raw files within 12 seconds (i.e. it read 38.122.372 lines of raw text!). What we get here is a so called Resilient Distributed Dataset (RDD) or Spark DataFrame which, in contrast to the Pandas DataFrame, is redundant and distributed over the Cluster DataNodes. With the following PySpark method an RDD can be generated:

```

1  %spark2.pyspark
2
3  rdd_usage = spark.read.csv("/user/hadoop/usage-stats/*.csv", header=True)

```

Of course the files of the folder „usage-stats“ need to be stored on HDFS before it can be read in with Spark. There are several ways of doing this, we used just a Zeppelin cell with magic command shell (%sh) to execute shell based commands within the notebook. With that one can use:

```

execute hdfs dfs -put
/home/hadoop/TFL_Cycling_Data_raw/usage-stats
/user/hadoop/usage-stats

```

to push the files on the datanodes of HDP. Remember the default replication size of Hadoop is 3, i.e. the files are replicated on three worker nodes but not on fur, which should still be resilient enough. The next decisive step is to prepare the usage-stats data and clean it up a bit. If we read in the data

as described above, our Spark DataFrame has the following columns:

Rental Id	Duration	Bike Id	End Date	EndStation Id	
				EndStation Name	
Start Date		StartStation Id		StartStation Name	

Since we only want to plot the bicycle stations at this time, the Rental ID columns and the date columns have been removed. This has the advantage that the DataFrame is already smaller and we don't carry any unnecessary ballast. Furthermore, NaN values have been replaced by mandatory „0“ and duplicates have been removed. The TFL also offers numerous .xml based live feeds, from which we can get the coordinates (longitude and latitude) to the rental stations. We can also use this to fetch the coordinates:

```
1  #spark2.pyspark
2  import requests
3  from xml.etree import ElementTree as ET
4  import pandas as pd
5
6  #Getting the coordinates tuples of the rental stations
7  site = "https://tfl.gov.uk/tfl/syndication/feeds/cycle-
8  hire/livecyclehireupdates.xml"
9  response = requests.get(site)
10 root = ET.fromstring(response.content)
11
12 #fetch relevant information
13 id_list = [int(root[i][0].text) for i in range(0, len(root))]
14 name_list = [root[i][1].text for i in range(0, len(root))]
15 lat_list = [float(root[i][3].text) for i in range(0, len(root))]
16 lon_list = [float(root[i][4].text) for i in range(0, len(root))]
17 capacity_list = [int(root[i][12].text) for i in range(0, len(root))]
18
19 #zip tuples
20 all_locs = pd.DataFrame(list(zip(name_list, id_list, lat_list,
21                                     lon_list, capacity_list)), columns =
22 ["name", "id", "lat", "lon", "capacity"])
23
24 #zip
25 all_locs.to_csv('/home/hadoop/TFL_Cycling_Data_raw/rental_stations_save
26 d.csv', header=True, index=None)
27 print(all_locs.shape)
28
29 #786 unique rental stations
30 all_locs.head()
```

This gives us not only the longitude and latitude, but also the maximum bicycle capacity of any station in London. This information could still be

quite useful. As you can see, we have 786 unique bicycle stations in London (which are quite a lot). Remark: `Livecyclehireupdates.xml` contains also following attributes:

- installed
- locked
- installDate
- nbBikes
- nbEmptyDocks
- nbDocks

Given that there are 786 stations across London (at least at the time of writing), this allows for $786 * 785 = 617.010$ possible journey combinations if we ignore those that start and end at the same station. The next step is some data cleansing steps such as sorting, renaming and converting data types. PySpark offers similar operations for the RDDs as for Pandas DataFrames. However, some pandas such as „`iloc()`“ cannot be transferred to Spark RDDs. Therefore Spark offers the option to directly execute SQL queries in DataFrames. Anyway, the fetched data from the live feed must still be merged with the `rdd_usage` Spark DataFrame. For this purpose we first saved the Pandas DataFrame locally and then uploaded it to the HDFS. This checkpoint is then read in again with PySpark as RDD. Merging works similar to SQL or/and python. To do this, simply use the „`join`“ command to link the two DataFrames:

```
1  %spark2.pyspark
2  from pyspark.sql.functions import *
3
4  ##Merging
5  #Merge cycle_usage RDD with RDD from livecyclehireupdates for
6  StartStations
7  ta = unq_rentals.alias('ta')
8  tb = rdd_rental_locs.alias('tb')
9  unq_rentals = ta.join(tb, ta["StartStation Id"] == tb.id)
10
```

In this case we also used alias (similar to the alias of SQL) in order to shorten the joining function name. The rest of the merging stuff is automatically done

by the Spark API. After all the data munging steps we got following schema of our combined and prepared usage-stats table:

```
StartStation Id | EndStation latitude | StartStation Id |  
Duration | StartStation longitude | StartStation  
Address | StartStation capacity | EndStation Address |  
EndStation latitude | EndStation longitude | EndStation  
capacity |
```

One can note, that there are also the columns „StartStation ID Used“ and „EndStation ID Used“ which stands for frequency of rented bikes on a single station. This gives us an indication of how popular the station is in comparison with other ones. It might be also used as a weighting factor for plotting on maps which will be described in just a moment. With this, we can begin to plot statistics. But at first let compare the size of the new file with the total size of all loaded cycle-usages. The total size was around 5.4 GB (which could not be simply loaded into a Pandas DataFrame) whereas the new DataFrame has only 0.9 GB size. The reduction of the size can be explained by the fact that we removed some unused columns and dropped duplicates. This means also we are now able to use Pandas because to load 1 GB into a local memory should be possible for most common computers. In fact, the Hadoop test environment has shown its functionality and legitimation for the using in big data projects.

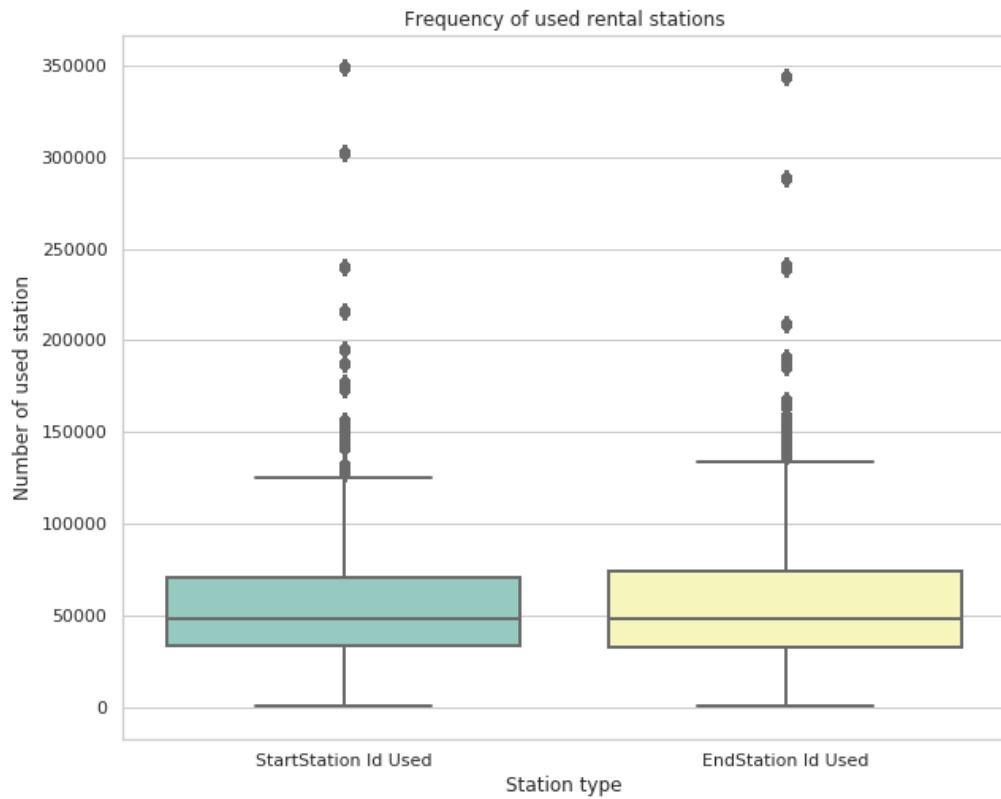


Figure 2.1: Frequency of used StartStations and EndStations

As we can see from Figure 2.1, the use of end stations slightly outweighs the starting stations, which merely means that more different end stations than different start stations were used. Apparently there are some very popular rental stations like the outliers from figure 2.1 show. It is worth noting that the data were all collected from 2016 - 2018. (The year 2014-2016 are not included in this case). The next plot should serve us to show the density center of the rental stations locations. For this a joint kernel density estimation plot could be helpful. We are plotting longitude against latitude of each single StartStation. The illustrations looks then like this:

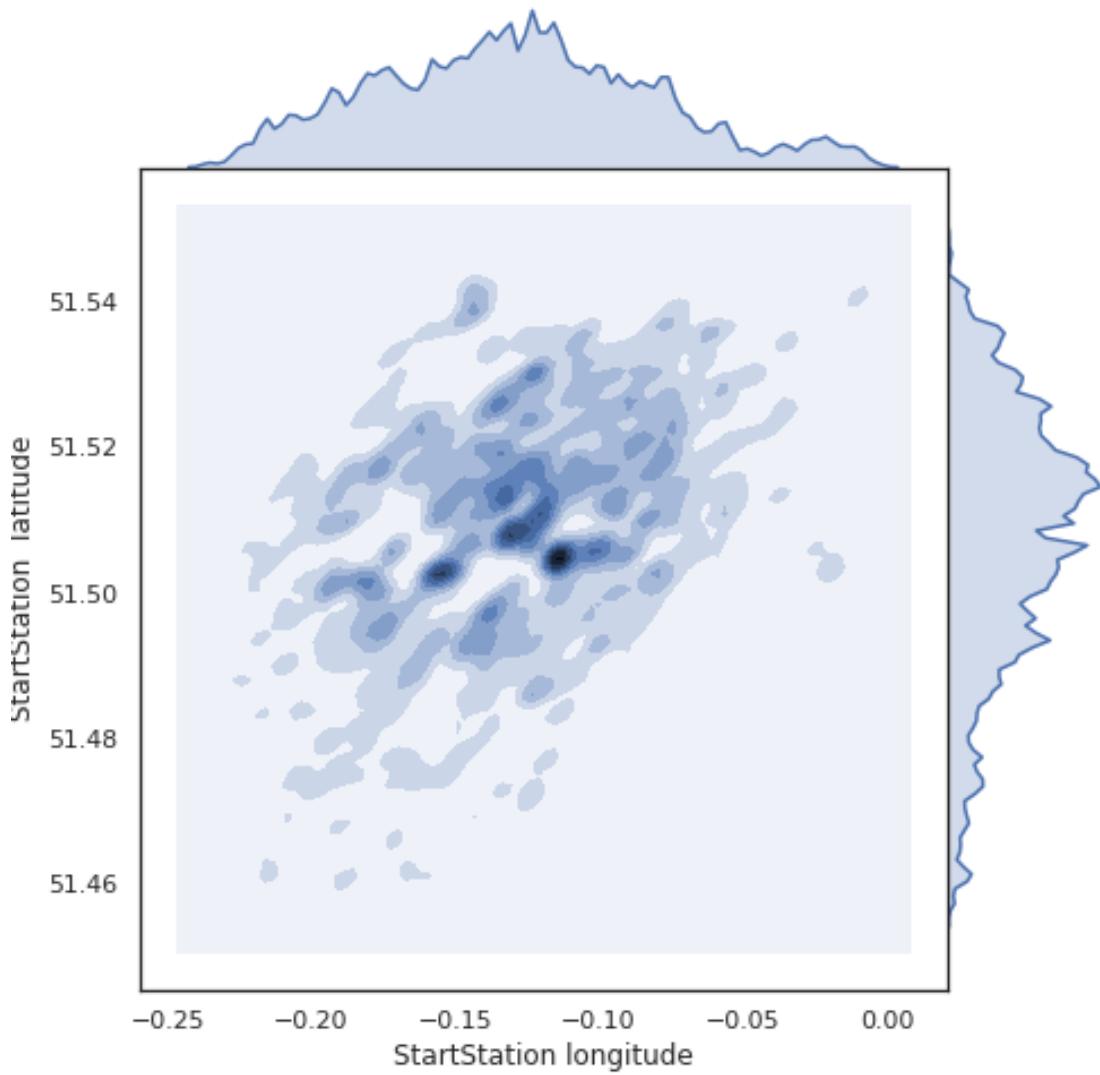


Figure 2.2: Frequency of used StartStations and EndStations

The centroid or center of this plot is indicated by the dark spots. There are the most rental stations as also the frequency distributions (i.e. histograms) on the right border and on the top are clearly showing. Apart from that center of gravity, on the right corner there are also rental stations placed (guessing between longitude: -0,01 and latitude 51,6) which indicates that

these rental stations must be placed in an outer borough of London. At a next, it would be useful to show the rental stations on a map based on their usage. As already mentioned we have 786 unique rental stations so that we can drop the duplicates for this case. In fact, we have saved the number of used stations on a separate column. At this step we prefer to use Jupyter on a local machine since Zeppelin on the cluster is not supporting IPython Display which is required by third-party packages such as folium. So we loaded the prepared cycle-usage data into a local Pandas dataframe. Now we can plot a heatmap for the rental stations:

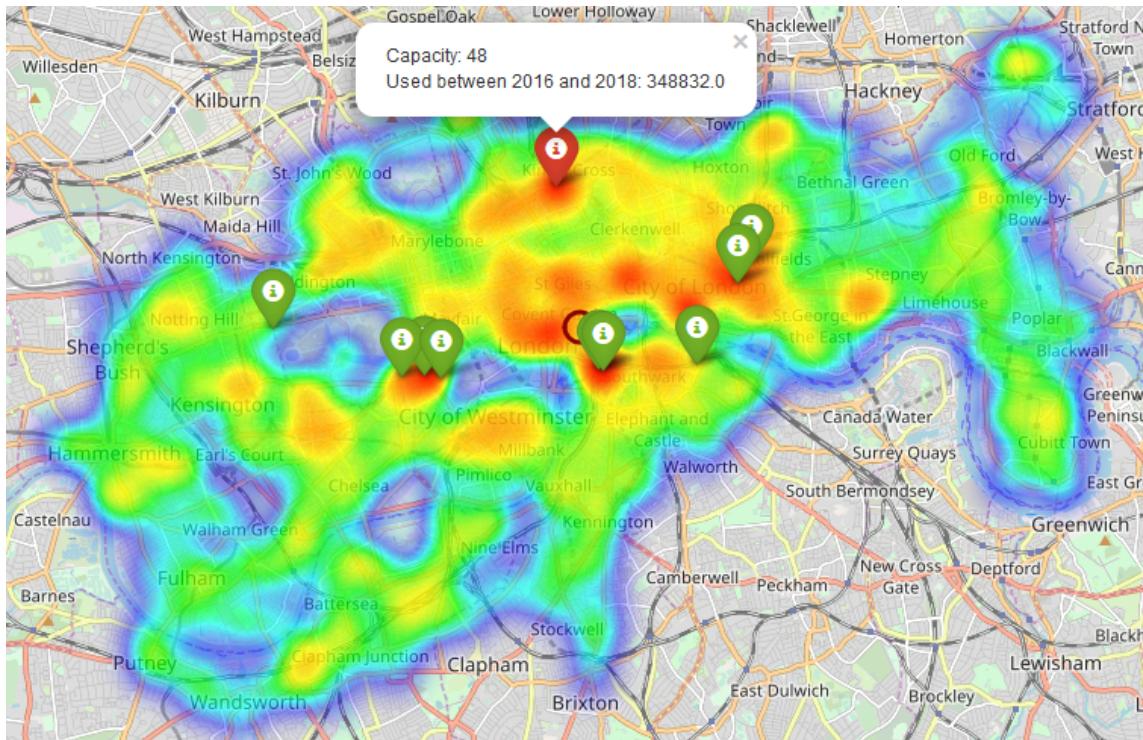


Figure 2.3: Top 10 used Santander rental stations

The red marker in figure 2.3 indicates the most often used bicycle station in London in last two years. This is not surprisingly because close to the biking station there is London train station. Probably many travelers who are arriving by train are using a Santander cycle. However, this interactive map has been saved and uploaded on a webserver on the Hadoop cluster (link) As we can see, there are many rental stations close together. For a

first impression this should be enough as one can already see the capacity of each single Santander rental bike station by clicking on the marker and also the location of each station is visible on the map. As an addition one can see the frequency of each station by looking at the sizes of the heat spots or/and reading marker details.

2.2 Conclusion of the file sizing problem

1. The sizing problem can be well handled by Spark
2. The HDP is working and storing the cycle-usage files on HDFS
3. For plotting tasks packages like folium or seaborn cannot be used with Spark DataFrame
 - Converting back to Pandas DF might exceed local RAM
 - Using Jupyter notebook on local computer for plotting maps (e.g. folium)

2.3 CyclingRoutes

The Cycling routes from the folder of TFL [24] are geojson files and represent real geographical routes within London inner city. There are 5 *geojson* files in this folder available:

- *Central_London_Grid.json*
 - *Central London Grid*: is a set of safer, connected routes for cyclists across central London
- *Mini_Hollands.json*
 - Mini-Hollands have features that make cycling feel safer and more convenient.
- *Quietways.json*
 - *Quietways* are signposted cycle routes, run on quieter back streets to provide for those cyclists who want to travel at a more relaxed pace.

- *super-highways-quietways.json*
 - Combination of super highways and quietways
- *Cycle_Superhighways.json*
 - *Cycle_Superhighways* are on main roads and often segregate cyclists from other traffic.

The routes should be concatenated by a single GeoPandas GeoDataFrame. Doing so requires to read the geojson files with following scripts:

2.4 Plotting Data

Lets see if we can compare the frequency of used rental stations between StartStation and EndStation. We can not use swarm plots or any other kind of plots with fine granular details since „matplotlib“ or Python would try to draw over six million dots which will easily exceed the local RAM of the master node. Beside from that, we also would see too many points at the same location since lot of stations are drawn several times. A better plot in my opinion is to use a box plot as following figure shows:

```

1  %spark2.pyspark
2  import pandas as pd
3  import geopandas as gpd
4
5  pd.set_option('display.expand_frame_repr', True)
6  pd.options.display.max_columns = None
7
8  path ='/home/hadoop/TFL_Cycling_Data_raw/CycleRoutes/'
9  grid_data = ["Central_London_Grid.json",
10    "Mini_Hollands.json",
11    "Quietways.json",
12    "super-highways-quietways.json",
13    "Cycle_Superhighways.json"]
14
15
16 cycle_routes_df = gpd.GeoDataFrame(columns=["OBJECTID", "Label",
17 "Route_Name", "Status", "Shape_Length", "Tag", "geometry", "gridType"])
18
19 for i in grid_data:
20   print "Load " + i
21   df = gpd.read_file(path+grid_data[0])
22   #Extract and load json as gridType
23   df["gridType"] = str(grid_data[0])[:grid_data[0].index(".json")]
24   #Concatenate

```

```

25     cycle_routes_df = pd.concat([cycle_routes_df, df])
26
27
28 cycle_routes_df.dropna()
29
30 #Save it back as shape
31 cycle_routes_df.to_file(path+"CycleRoutes_Cleaned")
32
33 cycle_routes_df.tail()

```

So we get in only GeoPandas DataFrame which contains the 5 Cycling routes. In addition, the type of each route was stored in the column „gridType“, so that it can be filtered later for certain routes. For a first plot a Choropleth should show the different routes.

2.5 Conclusion

1. The data lacks in documentation and quality (Data Preparation & Cleansing necessary).
2. Some column names can be misunderstood (e.g. SITE-NUMBER is SITE_ID).
3. Data like SPEED or SPEED_MHP should be stored in one way.
4. Many columns like the axles are not meaningful and useful for our project. (e.g. Cyclecounters has 107 columns!)
5. Raw data seems not to be complete (e.g. only May, June, July)
6. *usage-stats* is very big. Necessity to use all usage-stats?

In the next sprint we will create similar plots as shown above, but this time with duplicates and times in mind. Because it is quite often the case that between two same bicycle stations several same entries exist at different times with different rental periods. For example, a use case would be to display the day/night change on the map, i.e. do more cyclists ride quietways or super-highways at night between rental stations?

3 Prediction on Daily Basis

Data profiling from the previous chapter has already provided initial information on bicycle data and rental stations in London (see chapter 2). We found out where most of the stations were close to each other (see figure 2.1) as well as the top 10 most used rental stations in London. For example, the map from figure 2.2 shows the station at Kings Cross, which has been the most used in the last 2 years with 348,832 borrowed bicycles. This is certainly due to the fact that there is a world-famous railway station and therefore many commuters prefer the bicycles on the spot. These findings and data serve as the basis for the next steps. One of these steps is the presentation of the bicycle usages of the individual stations over the last 4 years. The preferred route for each day should be displayed. Since this is de facto another application case for data profiling, the folium map from figure 2.1 is further enriched in subchapter 3.1 (Data Profiling Part 2). This is followed by the description of a feature engineering use-case for the "Kings Cross" station. Since it was discussed with the project stakeholders to carry out a daily based prediction, feature engineering also takes place on a daily basis. In this case, useful features are considered for the prediction of the model (e.g. holidays, weekends..., see subchapter 3.2). After this data preparation phase, various regression models are used to predict the bicycle usage of a rental station. These models are all trained with the same training data set and the error rate is compared. It turns out that the use of a neural network delivers the best results. The procedure and results of the prediction are described in more detail in section 3.3. In addition to the neural network, a polynomial regression is carried out (see section 3.5), which we considered as useful in consultation with the stakeholders. The polynomial regression shows that it has a comparable error rate similar to the neural network, but the calculation of higher dimensions turns out to be computational intensive, so that we finally decided on the neural network as regressor model.

3.1 Data Profiling Part II

As already indicated in Data Profiling Part 1 in chapter 2, the next step is to display the use of the routes at different times between the individual bicycle stations on a map. Since the Python package "folium" uses leaflet maps based on Javascript [33], the plotting of the routes on an hourly level is not performant, because too many poly lines have to be drawn and the map

can no longer be efficiently displayed in the browser. Therefore only the top 10% most used routes were plotted on the folium map. Another restriction was the aggregated granularity on a daily base. This means that the plotted map always showed the route usage for day x. With the folium plugin „TimestampedGeoJson“ time series data can be plotted in JSON format. An excerpt from the script „Cycleusage & Cycleroutes [allStations].ipynb“ shows the corresponding function call:

```

1 import folium as fo
2 from folium.plugins import HeatMap, MarkerCluster, TimestampedGeoJson
3
4 m = fo.Map(
5     location=(51.509865, -0.118092),
6     zoom_start=12,
7     prefer_canvas=True
8 )
9
10 #Init
11 marker_cluster = MarkerCluster().add_to(m)
12 liste = []
13 routes_complete.fillna(0, inplace=True)
14 freq_mean = routes_complete["frequency"].mean()
15 freq_std = routes_complete["frequency"].std()
16
17 for index, p in routes_complete[:top10].iterrows():
18     startX = p["lat"]
19     startY = p["lon"]
20
21     if (startX == 0):
22         fo.PolyLine(liste, weight=((freq - freq_mean)/freq_std),
23                     opacity=0.5, popup="Used between 2015 and 2018: " + str(int(freq)) +
24                     "x", color="blue").add_to(m)
25         liste = []
26     else:
27         freq = int(p["frequency"])
28         liste.append((startX, startY))
29
30 #Add markers (unique)
31 df_startStations[:len(df_startStations)].apply(lambda
32 row:fo.Marker(location=[row["StartStation latitude"], row["StartStation
longitude"]],popup=row['StartStation Address'] + '</br>' + "Capacity: "
33
34

```

```

35 + row['StartStation capacity'] + '</br>' + " Used between 2016 and
36 2018: " + str(int(row['StartStation Id Used']))
37
38
39 #Put dates and coordinates in required format
40 features = [
41     {
42         'type': 'Feature',
43         'geometry': {
44             'type': 'LineString',
45             'coordinates': coordinates[i],
46         },
47         'properties': {
48             'times': dates[i],
49             'style': {
50                 'color': 'red',
51                 'weight': weights[i][0],
52                 'opacity': 1
53             }
54         }
55     }
56     for i in range(len(dates))
57 ]
58
59 #Allows polylines over time (i.e. red lines), JSON format required;
60 TimestampedGeoJson({
61     'type': 'FeatureCollection',
62     'features': features,
63 },
64     auto_play=False,
65     loop=False,
66     loop_button=True,
67     date_options='YYYY/DD/MM',
68     duration='P2M',
69     add_last_point=False).add_to(m)
70
71 m.save(outfile= "routes_frequency.html")
72 m
73

```

Figure 3.1: Route usage over time plot

As the code from figure 3.1 shows, the poly lines on the map have been added iteratively. The weight parameter can be used to determine the thickness of the poly line. Since these should look as dynamic as possible on the map, the weight has been standardized. The fixed stations were initially

added, but no duplicate stations were plotted. A disadvantage of the plugin is that it needs the data in a JSON format. Therefore the coordinates for the single points of a poly line as well as the time series data had to be converted into a compatible format. As can be seen from the Python code, the coordinates must be given the type „LineString“. A LineString is defined as a sequence of uniquely assignable points. In this case, the longitude and latitude previously requested using the Graphhopper API on our Graphhopper server were sufficient. It is important that the parameter „date_options“ in „TimestampedGeoJson“ corresponds exactly to the date format as in the nested list dates, otherwise the time slider function on the map will not work properly. Applied to the bicycle data the following picture results for the time 26.06.2016:

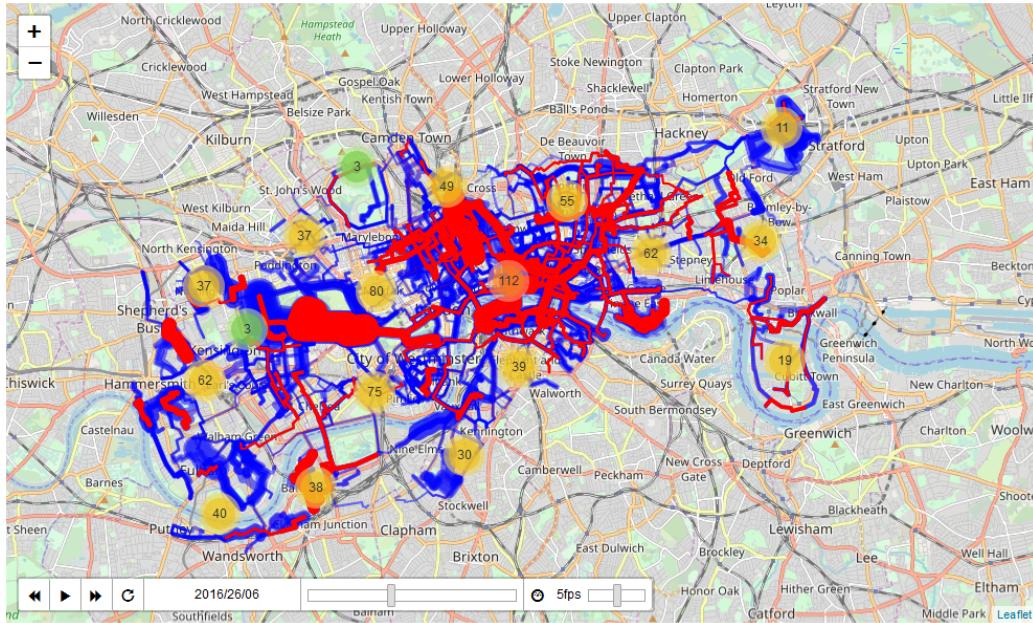


Figure 3.2: Used bicycle routes over time plot (high usage)

From figure 3.2 it can be seen that on this day there was a moderate use of Santander Bicycles in inner London. Especially the hubs such as Kings Cross or Hyde Park were obviously used very often with Santander Bicycles on this spring day. Furthermore, the map with the thickness of the poly lines shows how often this route was used in relation to the total use of all routes. The red colored routes from figure 3.2 are the actually used routes on this

particular day, while the blue routes are the „inactive“ ones. The bubbles with the numbers represent the respective rental stations. These have only been aggregated to provide a clearer representation. If the map is zoomed in (figure 3.2), the granularity is refined and the markers of the individual rental station locations are displayed. The color of these bubbles correlates with the number of aggregated stations in the vicinity. This method also has the advantage that it is easy to see where most rental stations are located. In fact, with 112 stations (see figure 3.2), the inner districts connected by the Waterloo Bridge and the London Blackfriars Bridge form the center of most stations which are close by. It should be noted, however, that new rental stations are constantly being added (even given up again!), and a new data extract could result in a different picture. Therefore this assumption is valid for the selected date from figure 3.2, but not for today or in two years. A useful feature that comes along with the folium plugin is the automatic data display sequence [33]. When someone clicks on the „Play“ button from figure 3.2, all time series data is automatically run through. The speed can be adjusted with the „fps“ slider. To give a counterexample, the following map shows the use of the routes on a winters day:

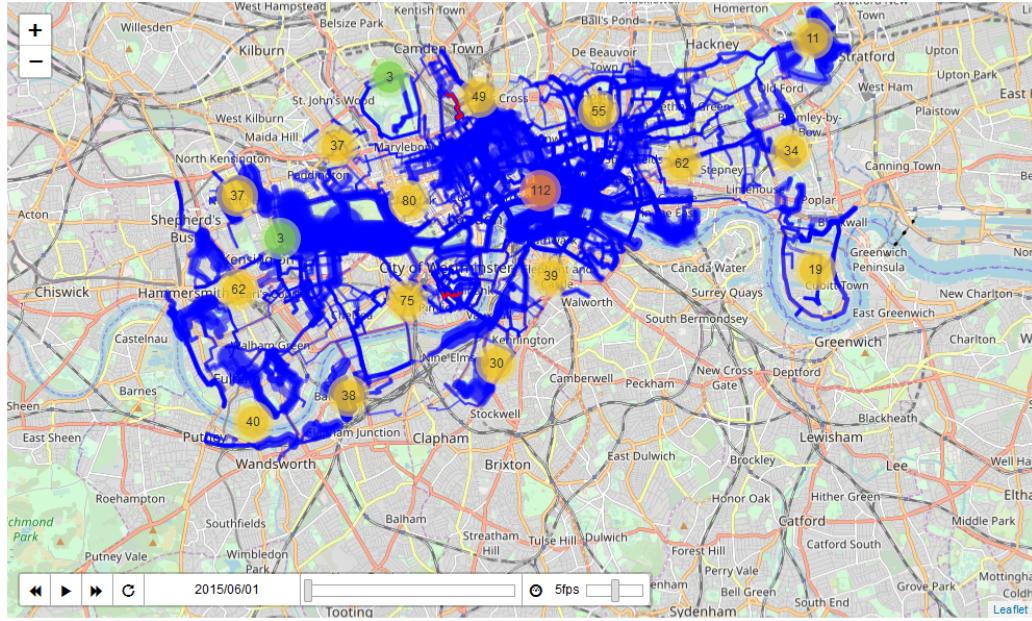


Figure 3.3: Used bicycle routes over time plot (low usage)

On 06.01.2015 obviously less bicycles were rented than in figure 3.2. This is due to the fact that at this time it was winter and on the other hand there were only 315 Santander bicycle rental stations all over London [23]. In the course of time, new stations were added again and again, resulting in a network which is monitored and managed by the local government and administrated by TfL. The interactive map is also available on the Hadoop Cluster. Overall, it can be observed that the usage of the rental stations has risen sharply on average over the last four years. This can be explained by the expansion of the network but also by the increased environmental awareness of the citizens. It is to be expected that further stations will be connected in the future, so that ideally there will be one rental station at each crossroad in the town.

3.2 Feature Engineering Kings Cross

Next, the features for the most frequently used station (Kings Cross) were prepared and aggregated on a daily basis. In addition, these features have been extended by more like *Holidays*, *Weekdays*, *Months*, *Seasons* and weather data. This enrichment of the features should help to achieve a higher accuracy of the learning model. For the weather data the weather API „Dark Sky“ was used. Like every available weather API the free use is limited. In the case of Dark Sky a maximum of 1000 API request calls per day can be executed with one key [32]. However, this limitation can be bypassed if multiple accounts are used. Since each account has a key, the keys can be collected and thus in fact significantly more requests can be executed. However, this cannot be applied to the existing data sets of the bicycle data, since the daily aggregated data of Kings Cross alone has already 376625 rows. Therefore a „dates“ file was created, which contains all days from 04.01.2015 - 14.04.2019 and thus covers all possible date values of the rental station data. This has the advantage that now only 1562 API requests are necessary and two keys are sufficient. The disadvantage of this method is that of course no precise weather data can be retrieved at every station. Therefore the coordinates for the weather data of the centroids (51.509865,-0.118092) of London were chosen. Dark Sky API also returns a JSON file as response, which can be searched for the desired information. The following function from the script „Cycleusage & Cycleroutes [allStations].ipynb“ shows the corresponding request call:

```

1  ### Read weather data for London with darksky api (1000 requests per
2  day)
3  df_dates = pd.read_csv("dates.csv")
4  df_dates['Start Date'].replace('.','/',inplace=True, regex=True)
5  df_dates.head()
6  #api_key2 = "f79a70fd65182047247a0506f998a96f"
7  api_key = "fbe9812e95cfc01d44c8d28e28d155e0"
8
9  df_dates["Daily Weather"] = ""
10 df_dates["Humidity"] = ""
11 df_dates["Windspeed"] = ""
12 df_dates["Apparent Temperature (Avg)"] = ""
13 df_dates["Hourly Weather"] = ""
14
15 for index, p in df_dates[0:1].iterrows():
16     timestamp = str(int(datetime.strptime(df_dates["Start Date"][index], '%d/%m/%Y'
17
18 ).replace(tzinfo=timezone.utc).timestamp())))
19     weather_resp =
20     requests.get('https://api.darksky.net/forecast/' + api_key + '/51.509865,-
21     0.118092,' + timestamp + '?exclude=currently,flags,minutely')
22     df_dates["Daily Weather"].iloc[index] =
23     weather_resp.json()['daily'][0]['data'].get('icon', 'No weather data')
24     df_dates["Humidity"].iloc[index] =
25     weather_resp.json()['daily'][0]['data'].get('humidity', 0)
26
27
28     df_dates["Windspeed"].iloc[index] =
29     weather_resp.json()['daily'][0].get('windSpeed', 0)
30     df_dates["Apparent Temperature (Avg)"].iloc[index] =
31     (float(weather_resp.json()['daily'][0].get('apparentTemperatureLow',
32     0)) + float(weather_resp.json()['daily'][0].get('apparentTemperatureHigh',
33     0))) / 2.0
34     df_dates["Hourly Weather"].iloc[index] =
35     weather_resp.json()['hourly'][0]['data']
36
37 df_dates.head()
38 df_dates.to_csv("dates and weather.csv", index=None, header=True)
39
40
41

```

Figure 3.4: Fetching the weather API with Python

As the code from figure 3.4 clearly shows, the following weather data were

considered as relevant: „Daily Weather, Humidity, Windspeed“and „Apparent Temperature (Avg)“. The feature „Daily weather“returns a string value, how general the weather was on that particular day (e.g. sunny, partly-cloudy...). Dark Sky makes this value dependent on the worst condition that occurred on that day [32]. That is, if it snowed for an hour, but the rest of the day was cloudy, „Daily Weather“will still show „snowy“because this condition is weighted higher.

In addition to the weather data mentioned above, the hourly weather data of each day was also queried as JSON lists, as otherwise every hour of each day would have to be queried separately, which would drastically increase the API calls. A drawback of this variant is that the dates and weather dataframe has a mixed structure. While the column „hourly weather“is a JSON format, the other columns have a normal structure. This problem is further addressed in the section „Data Preparation - Hourly Base“.

The processed weather data were next merged with the processed cycle usage dataframe. For the Holidays the Python package „holidays“was used, which returns the corresponding holidays to a selected location. Similarly, the features *Weekdays*, *Months* and *Seasons* were be generated with defined functions and afterwards merged back with the main dataframe.

Subsequently, the dataframe was supplemented by „past“and„future“data. This should improve the training of regression models, for example. A kind of „sliding window“was implemented for the past data. This means that these columns always contain the value of the previous day. With the future usage data, the daily number of borrowed bikes of the next day was also displayed. The column „Rented Bikes“(i.e. usage) corresponds to the class variable of interest. This has the consequence that the first and last day in the dataframe have missing values, because there is no data for this naturally. The prepared dataframe for Kings Cross on a daily basis looks after the preparation steps like this:

	Month	Season	Weekday	Holiday	Daily Weather	Daily Weather (Past)	Humidity	Humidity (Past)	Windspeed	Windspeed (Past)	Apparent Temperature (Avg)	Apparent Temperature (Avg) (Past)	Rented Bikes	Rented Bikes (Future)
1510	March	Spring	Friday	False	wind	wind	0.80	0.71	19.56	16.31	54.260	50.685	258	28
1511	March	Spring	Saturday	False	wind	wind	0.80	0.80	20.50	19.56	42.860	54.260	28	41
1512	March	Spring	Sunday	True	partly-cloudy-day	wind	0.74	0.80	13.56	20.50	38.345	42.860	41	327
1513	March	Spring	Monday	True	clear-day	partly-cloudy-day	0.74	0.74	7.60	13.56	46.570	38.345	327	323
1514	March	Spring	Tuesday	False	partly-cloudy-day	clear-day	0.83	0.74	6.26	7.60	50.315	46.570	323	33

Figure 3.5: Prepared data frame of Kings Cross (daily)

As one can see from figure 3.5, the past data was only included for the weather data whereas the future data was only added for the usage of the rental station. The string values were later be transformed into numerical values since machine learning methods always require number values instead of string values. Therefore a proper schema was defined, which is described in more detail in the modeling section.

With the processed bicycle data from 3.5, machine learning can already be used to predict, for example, the daily usage of bicycles (i.e. how many bicycles will be borrowed tomorrow starting from today?). The script „Feature Engineering Kings Cross.ipynb“ contains the preparation functions described above and can be found on our GitHub repository.

3.3 Data Prediction with MLPRegressor

Since we wanted to implement the prediction algorithm in Python we used the library *Scikit Learn* which provides a class *Multi-layer Perceptron Regressor (MLPRegressor)* with according functions. As we wanted to train a non-linear model we decided to use a Multi-layer Perceptron since it can learn a non-linear function approximator for either classification or regression. However the MLPRegressor uses backpropagation which is the most widely used algorithm for supervised learning with multi-layered feed-forward networks [39], this algorithm was implemented to train a prediction model for the rental bike usage on a daily base.

The outcome of the Data Profiling Part 2, as described in chapter ?? served as a basis for the prediction. As a first implementation we used the following data as an input for the feature matrix: *Month*, *Weekday*, *Day*, *Season*, *Daily Weather*, *Daily Weather Past*, *Humidity*, *Humidity Past*, *Windspeed*, *Windspeed Past*, *Apparent Temperature (Avg)*, *Apparent Temperature (Avg) Past*, *Rented Bikes*. Since we assume that weather conditions play an important role in bike usage we decided to add weather data. Another assumption was that on weekdays the frequency will rise, because a lot of people use rented bikes to travel to work. Furthermore we added the number of rentals of today in order to predict the ones of tomorrow. Therefore *Rented Bikes Future* served as our target variable which is to predicted.

Since the MLPRegressor works with data represented as dense and sparse numpy arrays of floating point values, data had to be encoded accordingly beforehand. Therefore we mapped all of the ordinal scaled data like *Month*, *Weekday*, *Season*, *Daily Weather* to numerical data, by implementing a dictionary which assigns numerical values to ordinal data. E. g. For the weekdays we used the following dictionary:

```
"Weekday": {"Monday": 1, "Tuesday": 2, "Wednesday": 3,
            "Thursday": 4, "Friday": 5, "Saturday": 6, "Sunday": 7
}
```

After the encoding we split the data into 80 % training data and 20% test data. Since the MLPRegressor is sensitive to feature scaling we normalized the data accordingly to the activation function. Since we used the *logistic* activation function which expects values between [0,1] we normalized the feature matrix as well as the target variable beforehand. Scikit Learn provides several scaling mechanism to rescale the data. After data was scaled we applied the MLPRegressor with *logistic* as an activation function and ten hidden layers with five neurons. After the prediction we denormalized the data back to its original state. In order to rate the accuracy of the prediction we computed the Root Mean Square Error (RMSE) which measures the difference between actual and predicted values of a model. The less the difference respectively the value, the better is the model. To get a better impression of how the prediction worked, the results were plotted as time series. The time series was plotted over for years, which causes the plot to be very large. In order to make it more user friendly, we implemented an interactive plot which gives the user the opportunity to zoom into specific parts or cut out some parts to

look at them more closely.

Normalization

In the first attempt we used the *MinMaxScaler* which rescales the data such that all values are in the range of [0,1]. This gave us a RMSE of 122.347. To get a better impression of what that means we plotted this result which can be seen in figure 3.8.

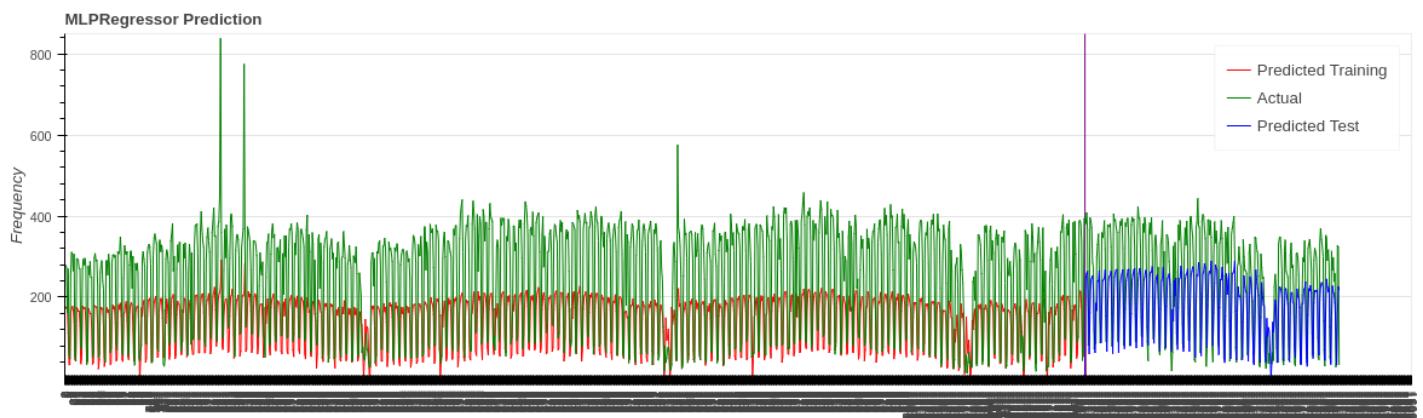


Figure 3.6: Result of MLPRegressor with MinMaxScaler

This result shows that the prediction lack accuracy. As we found out this was caused by the *MinMaxScaler*, since this scaler is very sensitive to the presence of outliers. Therefore we chose another scaler for normalization which is more prone to outliers. In order to find a more suitable scaler, experiments were made with all available scalers of Scikit Learn which met the requirements for the output of the scaling. The best result was achieved by the *QuantileTransformer* with a RMSE of 57.837, therefore we stuck with this normalization method.

Feature Evaluation

In order to improve accuracy, further experiments were carried out to evaluate the features. As mentioned earlier the previous predictions were made with a feature matrix of thirteen features. To evaluate which features improve the prediction and which are useless 18 test cases were carried out.

Each test case consists of different constellations of features. In the end the different RMSE values were compared, which can be seen in figure 3.7.

	Test	Da y	Month	Season	Week- day	Daily Weather	Daily Weather Past	Humid- ity	Humid- ity Past	Wind- speed	Wind- speed Past	A Temp	A Temp Avg	Rented Bikes	RMSE
MLP Regressor	1	1	1	1	1	1	1	1	1	1	1	1	1	1	57.837
	2	1	1	1	1	1	0	1	0	1	0	1	0	1	53.178
	3	0	1	1	1	1	1	1	1	1	1	1	1	1	56.406
	4	1	0	1	1	1	1	1	1	1	1	1	1	1	62.382
	5	1	1	0	1	1	1	1	1	1	1	1	1	1	56.827
	6	1	1	1	0	1	1	1	1	1	1	1	1	1	130.375
	7	1	1	1	1	0	1	1	1	1	1	1	1	1	57.200
	8	1	1	1	1	1	0	1	1	1	1	1	1	1	57.547
	9	1	1	1	1	1	1	1	0	1	1	1	1	1	58.855
	10	1	1	1	1	1	1	1	1	1	0	1	1	1	56.681
	11	1	1	1	1	1	1	1	1	1	1	1	0	1	226.706
	13	0	0	0	1	1	1	1	1	1	1	1	1	1	54.890
	14	1	1	1	1	0	1	1	1	0	1	0	1	1	59.239
	15	1	1	0	1	1	1	0	1	0	1	1	1	1	63.752
	16	0	1	0	1	1	1	0	1	0	1	1	1	1	51.414
	17	0	0	0	1	0	1	0	1	0	1	1	1	1	52.561
	18	0	0	0	1	0	1	0	1	0	1	0	1	1	54.866

Figure 3.7: Feature evaluation test cases

Figure 3.7 shows beneath the different test cases, included features, those are assigned to „1“ and excluded features which are assigned to „0“. The first RMSE value, colored in blue depicts the prediction accuracy with all 13 features. The red colored values show the constellations which are significantly worse than the original one and in turn the green values represent the values with an increased accuracy than the original one.

Result

This evaluation showed that the features *Weekday* and *Rented Bikes* are the most important ones and in turn the features *Day* and *Season* are less important. Based on our testing results we recommend to use the following nine features: *Weekday*, *Month*, *Past Data*, *Apparent Temperature (Avg)*, *Daily Weather* and *Rented Bikes*. Moreover experiments with scalers showed that the *QuantileTransformer* is more prone to outliers than others and is therefore the scaler of choice.

First plots were made with the library *Matplotlib* which turns out to be very

restricted and complicated to handle in order to create interactive plots. Therefore we recommend to use *Bokeh* which provides easy to implement interactive plots especially for large data sets.

Applying the recommendations regarding scaling and feature evaluation we received a RMSE value of 51.414 which is the best value achieved during the project phase. The overall result visualized with *Bokeh* can be seen in figure 3.8.

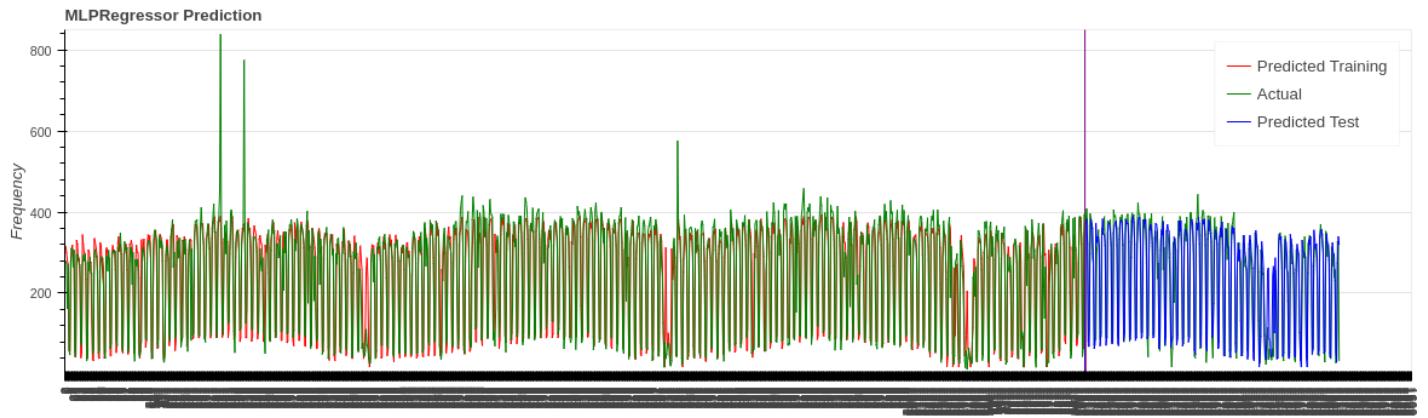


Figure 3.8: Accuracy with recommended features and scaling by the QuantileTransformer

Compared to figure 3.6 one can see significant improvements in accuracy. This model is based on the most used station, which is Kings Cross with a dataset of 1515 records. In order to confirm this model we applied it to the least used station in Farringdon Street, Holborn as well. The result can be seen in figure 3.9.

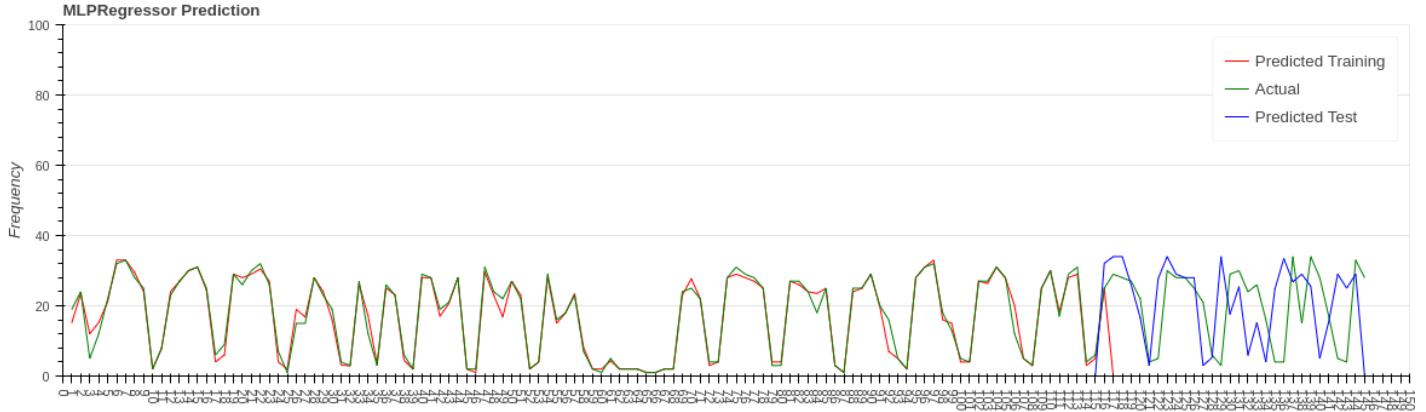


Figure 3.9: Accuracy with recommended features and scaling by the QuantileTransformer

The accuracy of the prediction of the least used station measures a RMSE of 9.311. This shows that the model works not only for highly frequented stations but also for stations like Farringdon Street with only 145 records.

3.4 Data Prediction with Spark

Apache Spark's Machine Learning Library (MLlib) is designed for simplicity, scalability, and easy integration with other tools. With the scalability, language compatibility, and speed of Spark, data scientists can focus on their data problems and models instead of solving the complexities surrounding distributed data (such as infrastructure, configurations, and so on). Built on top of Spark, MLlib is a scalable machine learning library consisting of common learning algorithms and utilities, including classification, regression, clustering, collaborative filtering, dimensionality reduction, and underlying optimization primitives. Spark MLlib seamlessly integrates with other Spark components such as Spark SQL, Spark Streaming, and DataFrames and is installed in the Databricks runtime. The library is usable in Java, Scala, and Python as part of Spark applications, so that you can include it in complete workflows. MLlib allows for preprocessing, munging, training of models, and making predictions at scale on data. You can even use models trained in MLlib to make predictions in Structured Streaming. Spark provides a sophisticated machine learning API for performing a variety of machine learning tasks, from classification to regression, clustering to deep learning [10].

MLlib is Sparks scalable machine learning library consisting of common learning algorithms and utilities that can be divided into two categories:

Supervised learning:

- Classification and regression
 - naive Bayes
 - linear models (SVMs, logistic regression, linear regression)
 - decision trees
 - ensembles of trees (Random Forests and Gradient-Boosted Trees)

Unsupervised learning:

- Collaborative filtering
 - alternating least squares (ALS)
- Clustering
 - k-means
- Dimensionality reduction
 - principal component analysis (PCA)
 - singular value decomposition (SVD)

Spark 1.2 includes a new package called spark.ml, which aims to provide a uniform set of high-level APIs that help users create and tune practical machine learning pipelines. It is currently an alpha component, and we would like to hear back from the community about how it fits real-world use cases and how it could be improved.

Selected Algorithms

Regression analysis consists of a set of machine learning methods that allow us to predict a continuous outcome variable (y) based on the value of one or multiple predictor variables (x). Since our goal is to make predictions on a continuous feature, the next step is to implement the algorithms provided

the library and compare their results.

Decision Tree Regression

A decision tree is built top-down from a root node and involves partitioning the data into subsets that contain instances with similar values (homogeneous). We use standard deviation to calculate the homogeneity of a numerical sample. If the numerical sample is completely homogeneous its standard deviation is zero [40].

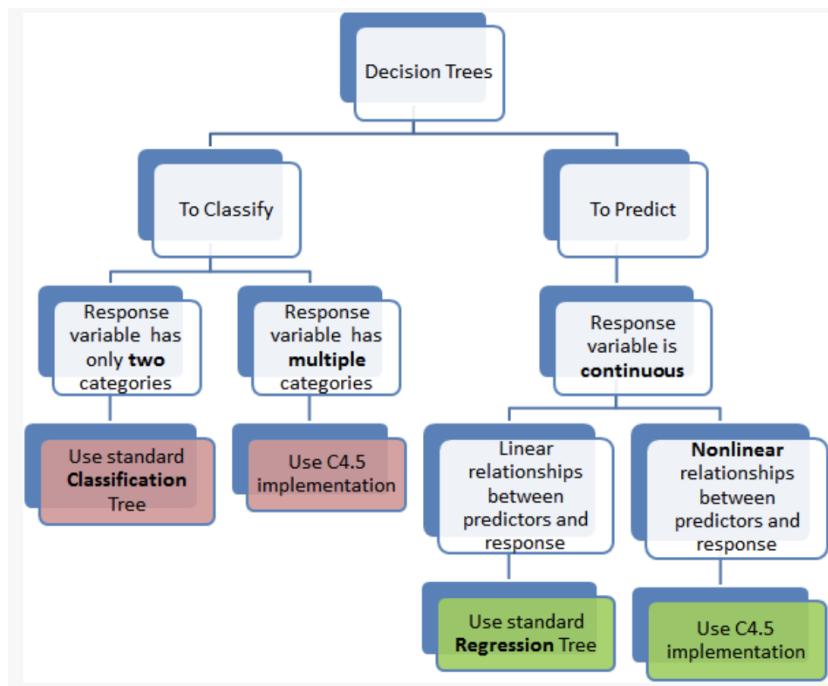


Figure 3.10: Decision tree usage diagram

In a standard classification tree, the idea is to split the dataset based on homogeneity of data. Lets say for example we have two variables: age and weight that predict if a person is going to sign up for a gym membership or not. In our training data if it showed that 90% of the people who are older than 40 signed up, we split the data here and age becomes a top node in

the tree. We can almost say that this split has made the data "90% pure". Rigorous measures of impurity, based on computing proportion of the data that belong to a class, such as entropy or Gini index are used to quantify the homogeneity in Classification trees [3]. In a regression tree the idea is this: since the target variable does not have classes, we fit a regression model to the target variable using each of the independent variables. Then for each independent variable, the data is split at several split points. At each split point, the "error" between the predicted value and the actual values is squared to get a Sum of Squared Errors (SSE). The split point errors across the variables are compared and the variable/point yielding the lowest SSE is chosen as the root node/split point. This process is recursively continued [3].

Regression trees are needed when the response variable is numeric or continuous. For example, the predicted price of a consumer good. Thus regression trees are applicable for prediction type of problems as opposed to classification. Keep in mind that in either case, the predictors or independent variables may be categorical or numeric. It is the target variable that determines the type of decision tree needed [3].

Random Forest Regression

A Random Forest is an ensemble technique capable of performing both regression and classification tasks with the use of multiple decision trees and a technique called Bootstrap Aggregation, commonly known as bagging. What is bagging you may ask? Bagging, in the Random Forest method, involves training each decision tree on a different data sample where sampling is done with replacement.

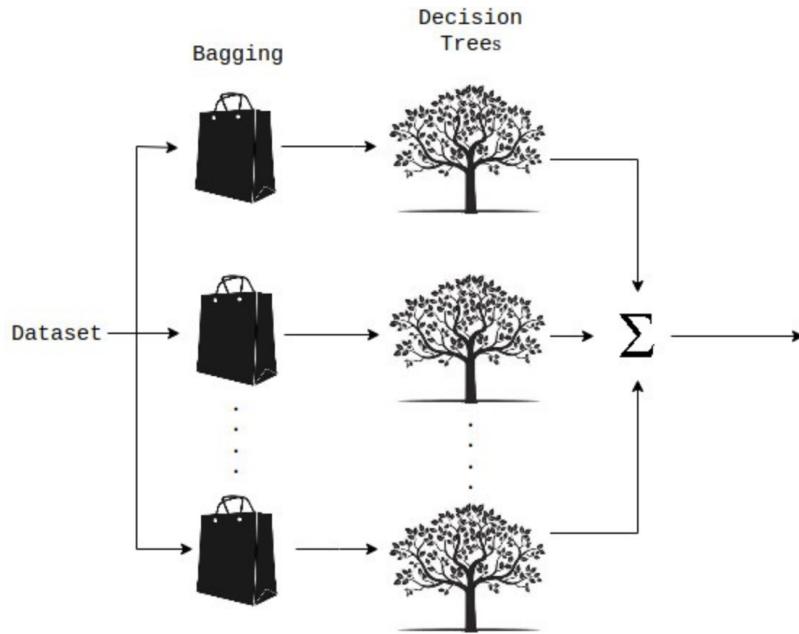


Figure 3.11: Random Forest Regression: Process

The basic idea behind this is to combine multiple decision trees in determining the final output rather than relying on individual decision trees. If you want to read more on Random Forests, I have included some reference links which provide in depth explanations on this topic.

Gradient Boosted Tree Regression

The idea of boosting came out of the idea of whether a weak learner can be modified to become better. Hypothesis boosting was the idea of filtering observations, leaving those observations that the weak learner can handle and focusing on developing new weak learners to handle the remaining difficult observations [4].

Gradient boosting involves three elements:

- A loss function to be optimized: The loss function used depends on the type of problem being solved. It must be differentiable, but many

standard loss functions are supported and you can define your own.

- A weak learner to make predictions: Decision trees are used as the weak learner in gradient boosting. Specifically regression trees are used that output real values for splits and whose output can be added together, allowing subsequent models outputs to be added and correct the residuals in the predictions. Trees are constructed in a greedy manner, choosing the best split points based on purity scores like Gini or to minimize the loss. Initially, such as in the case of AdaBoost, very short decision trees were used that only had a single split, called a decision stump. Larger trees can be used generally with 4-to-8 levels. It is common to constrain the weak learners in specific ways, such as a maximum number of layers, nodes, splits or leaf nodes. This is to ensure that the learners remain weak, but can still be constructed in a greedy manner.
- An additive model to add weak learners to minimize the loss function: Trees are added one at a time, and existing trees in the model are not changed. A gradient descent procedure is used to minimize the loss when adding trees. Traditionally, gradient descent is used to minimize a set of parameters, such as the coefficients in a regression equation or weights in a neural network. After calculating error or loss, the weights are updated to minimize that error. Instead of parameters, we have weak learner sub-models or more specifically decision trees. After calculating the loss, to perform the gradient descent procedure, we must add a tree to the model that reduces the loss (i.e. follow the gradient). We do this by parameterizing the tree, then modify the parameters of the tree and move in the right direction by (reducing the residual loss). Generally this approach is called functional gradient descent or gradient descent with functions [4].

Implementation

The first step is to apply the algorithm using all the features, then each time based on the features importance function, try to remove the least important columns and check for improvements in the RMSE value. The result of the testing process is described each step by a features importance plot. The results are as well listed on a table. The figures below show the correlation between the predicted data and the real future measured data. The first test series was implemented by regarding the most used station.

Decision Tree

We start by generating a model considering all the features with a RMSE of 71.35.

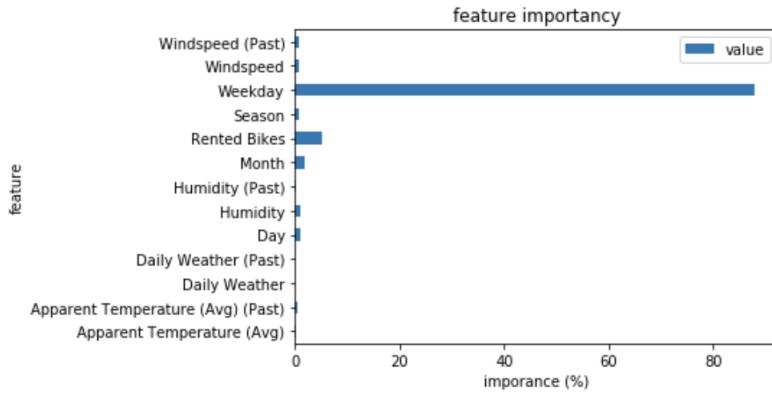


Figure 3.12: Test 1 features importance: 71.35

Figure 3.13 shows the removed columns and RMSE results on each test: The choice of removing features is explained by the features importance plots.

	TEST	Day	Month	Season	Weekday	Daily Weather	Daily Weather Past	Humidity	Humidity Past	Wind-speed	Wind-speed Past	Avg Temp	Avg Temp Past	Rented Bikes	RMSE	RMSE (on training)
Decision tree regression	1	1	1	1	1	1	1	1	1	1	1	1	1	1	71.35	47.66
	2	0	1	1	1	0	0	1	1	1	1	1	0	1	57.55	48.06
	3	0	1	1	1	0	0	1	1	1	1	1	0	1	57.55	48.06
	4	0	1	1	1	0	0	1	1	0	0	1	0	1	59.27	48.71
	5	0	1	0	1	0	0	1	1	0	0	1	0	1	56.98	49.09
	6	0	1	0	1	0	0	1	0	0	0	1	0	1	56.98	46.09
	7	0	1	0	1	0	0	0	0	0	0	0	0	1	54.43	50.78

Figure 3.13: Feature evaluation test cases

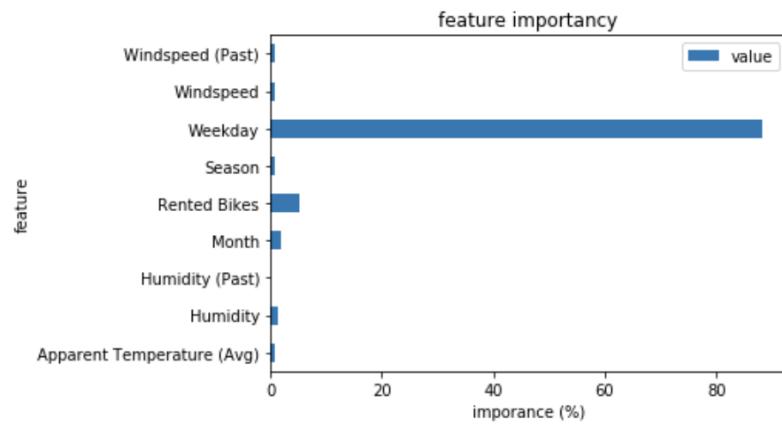


Figure 3.14: Test 3 features importance: 57.55

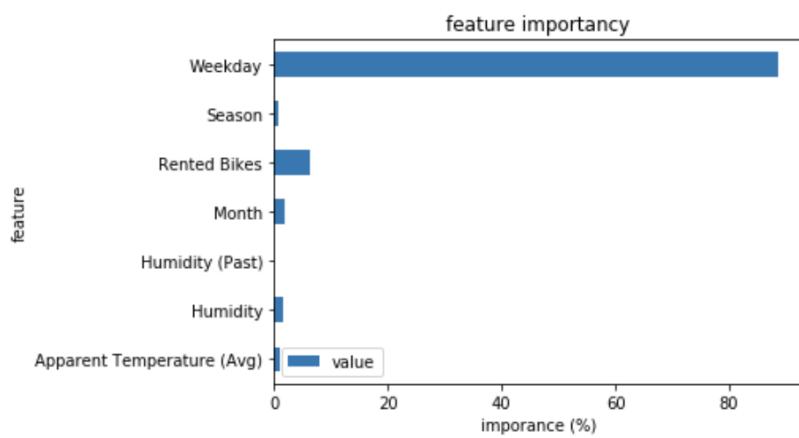


Figure 3.15: Test 4 features importance: 59.27

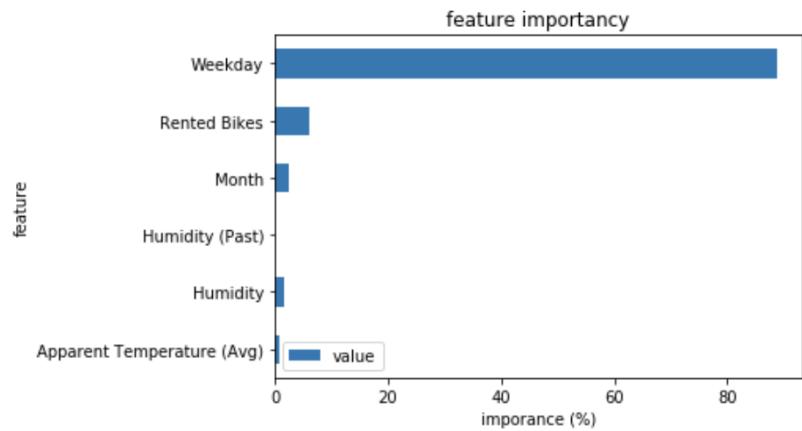


Figure 3.16: Test 5 features importance: 56.98

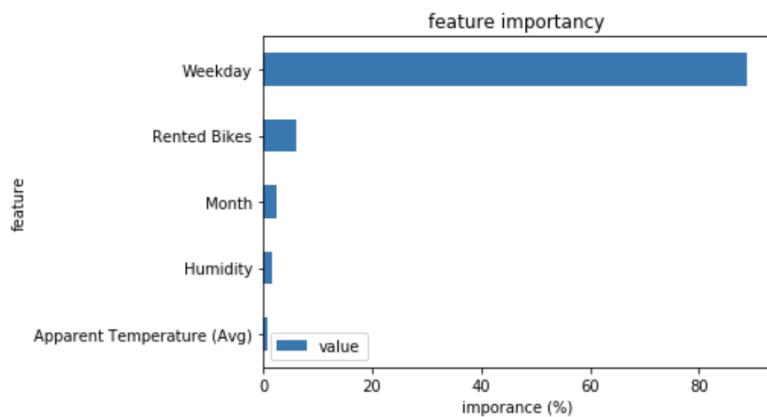


Figure 3.17: Test 6 features importance: 56.98

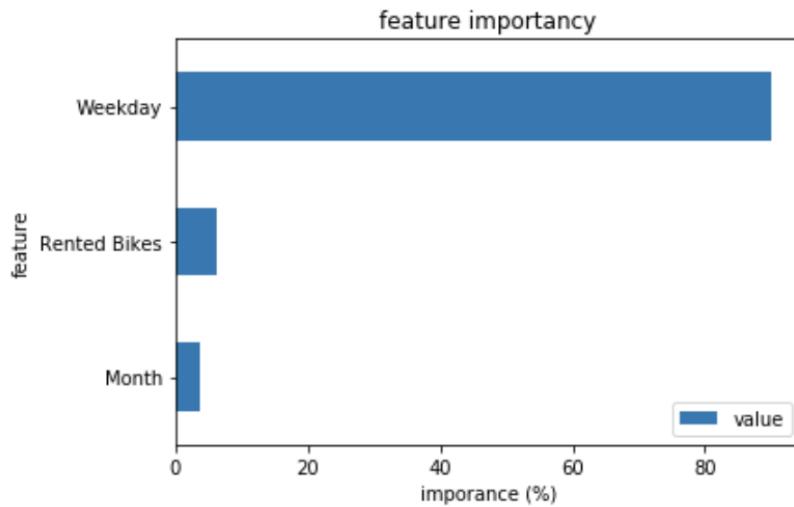


Figure 3.18: Test 7 features importance: 54.43

The final correlation results are shown in figure 3.19 with an RMSE of 54.43.

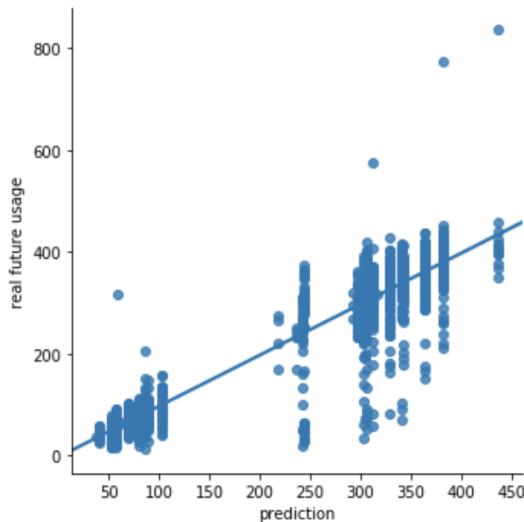


Figure 3.19: Data distribution of the actual and predicted data bike usage using Decision Tree

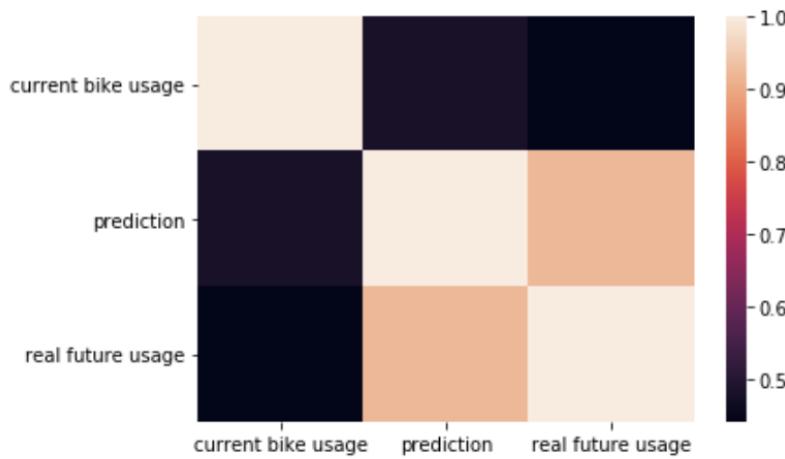


Figure 3.20: Heat map correlation plot showing interval value of the regression line (Decision Tree)

Random Forest

We start by generating a model considering all the features with a RMSE of 60.04.

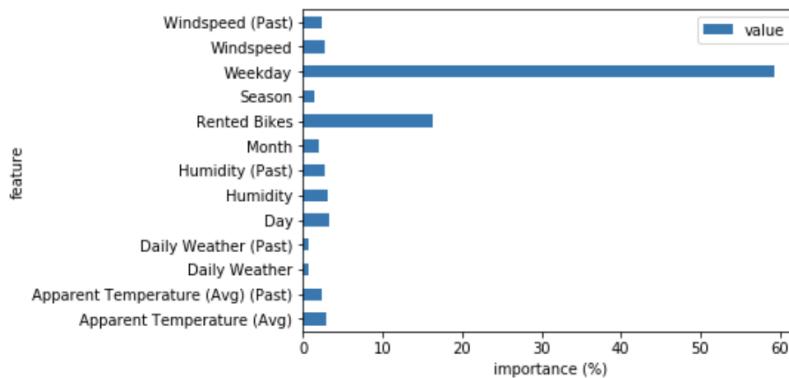


Figure 3.21: Test 1 features importance: 60.04

Figure 3.22 shows the removed columns and RMSE results on each test: The choice of removing features is explained by the features importance plots.

	TEST	Day	Month	Sea- son	We- ek- day	Daily Weat- her	Daily Weat- her Past	Hu- mi- dity	Hum- id- tyd-spd	Win- d-spd	Win- d-spd Past	A- Temp	A- Temp Avg	Rente- d Bikes	RM SE	RMS E(on trainin- g)
Random forest regressio n	1	1	1	1	1	1	1	1	1	1	1	1	1	1	60.0 4	22.30
	2	1	1	1	1	1	0	1	1	1	1	1	1	1	62.4 1	22.76
	3	1	1	0	1	1	0	1	1	1	1	1	1	1	60.9 8	22.42
	4	1	0	1	1	1	0	0	1	0	1	1	1	1	64.4 7	22.61
	5	0	0	0	1	0	0	0	1	0	0	0	0	1	63.0 8	44.46

Figure 3.22: Feature evaluation test cases

The process consists of removing after each test the features with the least importance.

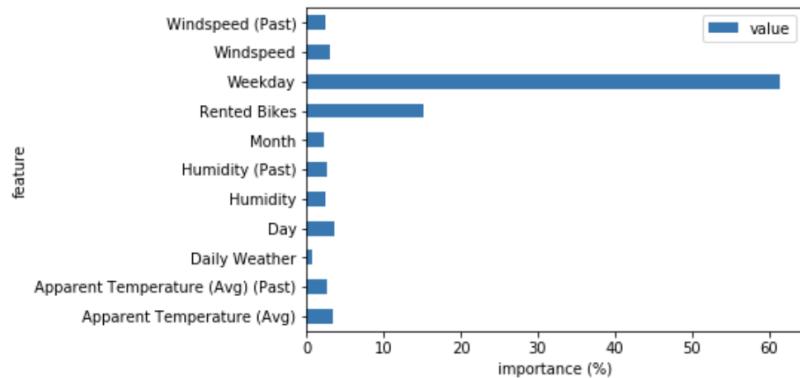


Figure 3.23: Test 2 features importance: 62.41

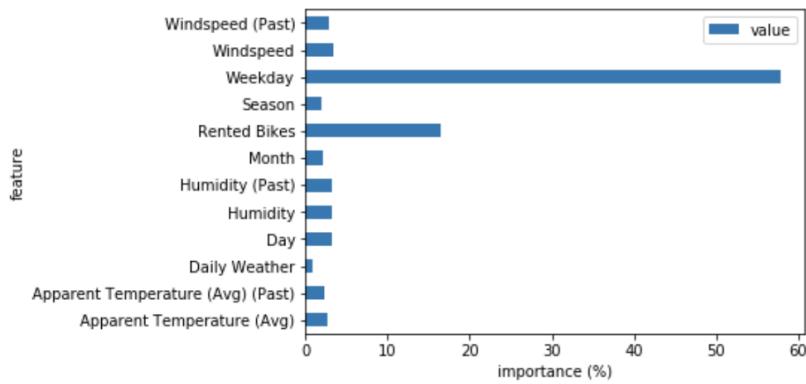


Figure 3.24: Test 3 features importance: 60.98

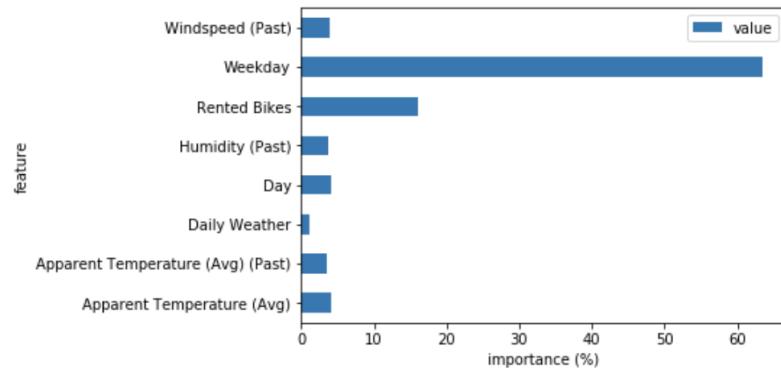


Figure 3.25: Test 4 features importance: 64.64

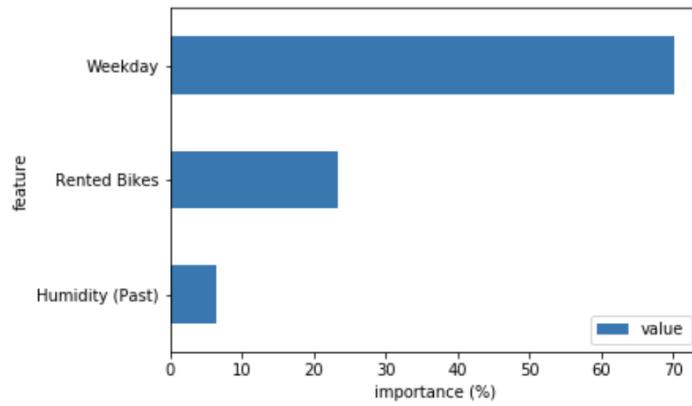


Figure 3.26: Test 5 features importance: 63.08

The final correlation results are shown in figure 3.27 with an RMSE of 54.43.

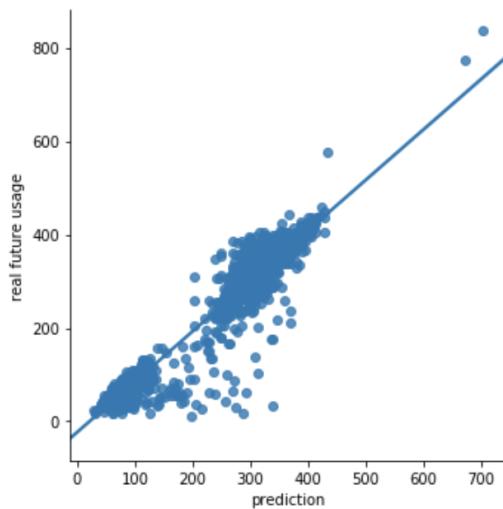


Figure 3.27: Data distribution of the actual and predicted data bike usage using Random Forest

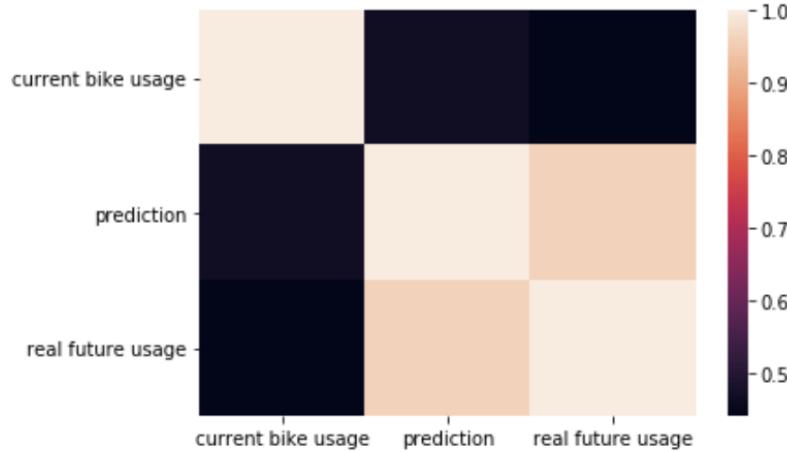


Figure 3.28: Heat map correlation plot showing interval value of the regression line (Random Forest)

Gradient boosted tree regression

We start by generating a model considering all the features with a RMSE of 81.30.

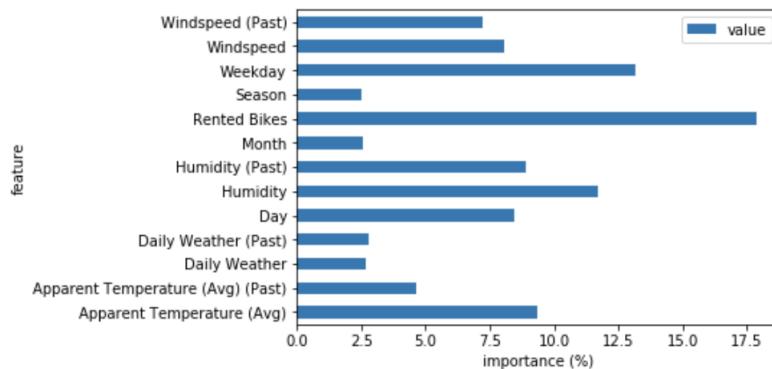


Figure 3.29: Test 1 features importance: 81.30

Figure 3.13 shows the removed columns and RMSE results on each test: The choice of removing features is explained by the features importance plots.

	TEST	Day	Month	Season	Week-day	Daily Weather	Daily Weather Past	Humidity	Wind-speed Past	Wind-speed Past	Avg Temp	Avg Temp Past	Rented Bikes	RMSE	RMSE (on training)
Gradient boosted tree regression	1	1	1	1	1	1	1	1	1	1	1	1	1	81.30	29.30
	2	1	1	0	1	1	0	1	1	1	1	1	1	88.48	30.79
	3	0	1	0	1	0	0	1	1	1	1	1	1	84.92	30.53
	4	1	0	0	1	0	0	1	1	1	1	1	1	76.71	30.35
	5	0	0	0	1	0	0	0	1	0	0	0	0	63.26	42.34

Figure 3.30: Feature evaluation test cases

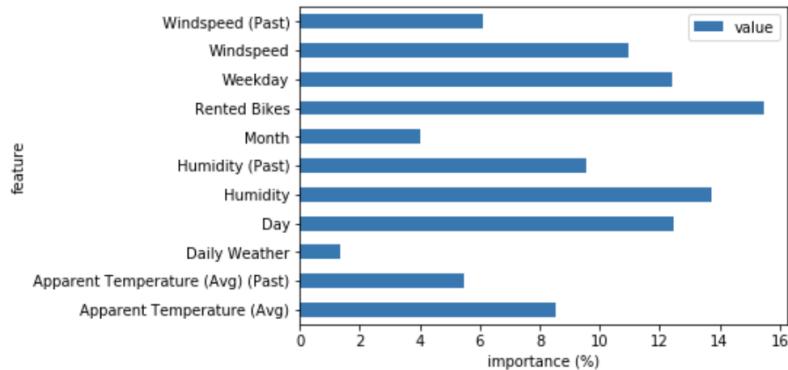


Figure 3.31: Test 2 features importance: 88.48

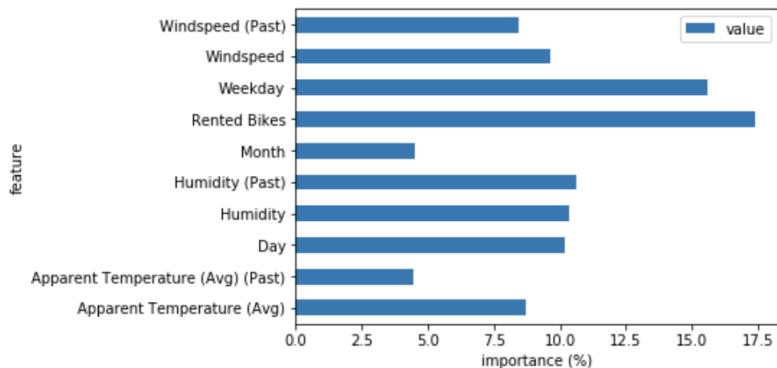


Figure 3.32: Test 3 features importance: 84.92

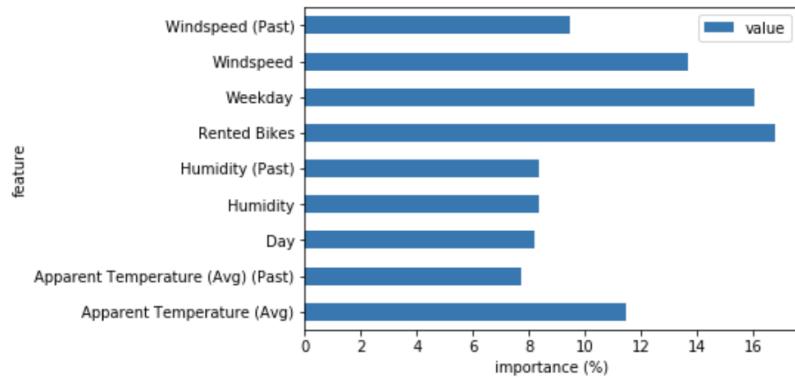


Figure 3.33: Test 4 features importance: 76.71

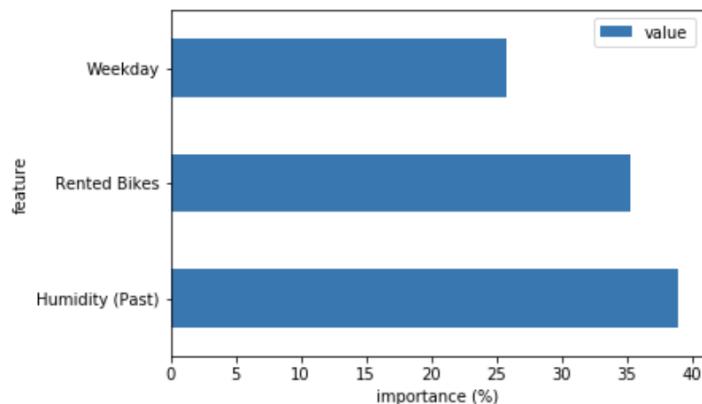


Figure 3.34: Test 5 features importance: 63.26

The final correlation results are shown in figure 3.37 with an RMSE of 63.26.

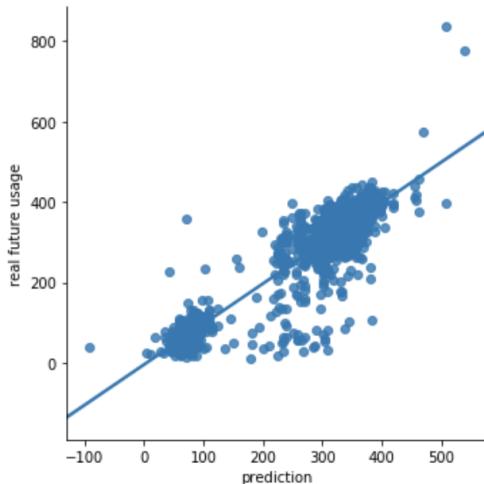


Figure 3.35: Data distribution of the actual and predicted data bike usage using Gradient boosted tree regression

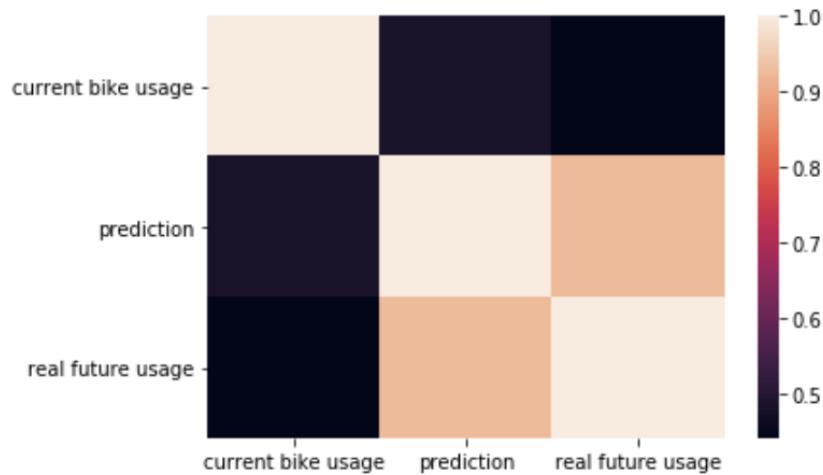


Figure 3.36: Heat map correlation plot showing interval value of the regression line (Gradient boosted tree regression)

The following test series was implemented by regarding the least used station with a data set consisting of 145 entries.

Decision Tree Regression

The final correlation results are shown in figure 3.37 with an RMSE of 11.72.

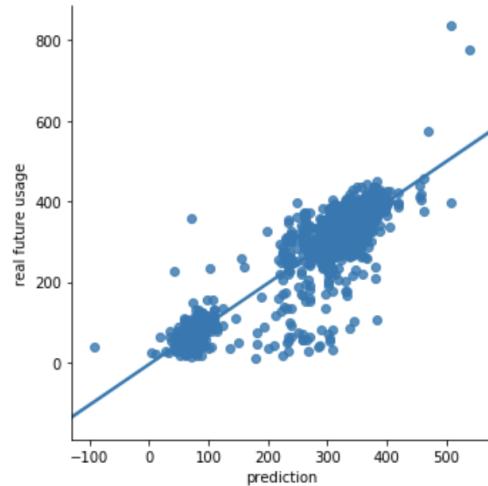


Figure 3.37: Data distribution of the actual and predicted data bike usage using Decision Tree Regression

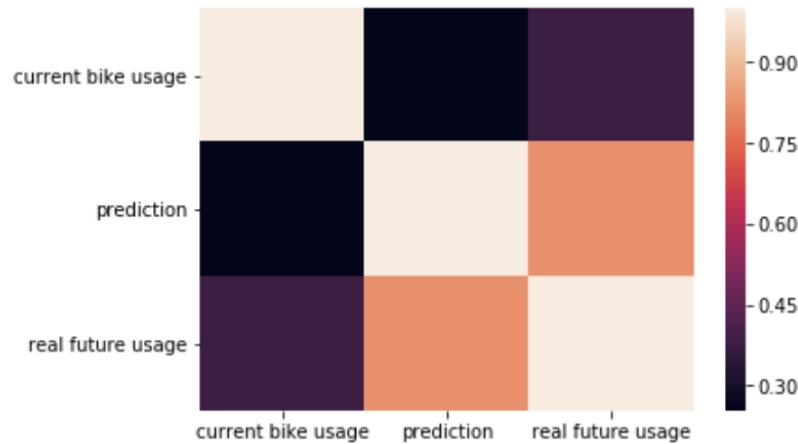


Figure 3.38: Heat map correlation plot showing interval value of the regression line (Decision Tree Regression)

Figure 3.39 shows the removed columns and RMSE results.

	TEST	Day	Month	Season	Weekday	Daily Weather	Daily Weather Past	Humidity	Humidity Past	Wind-speed	Wind-speed Past	Avg Temp	Avg Temp Past	Rented Bikes	RMS E	RM SE (on training)
Decision tree regression	1	1	1	1	1	1	1	1	1	1	1	1	1	1	11.72	2.46

Figure 3.39: Feature evaluation test cases

Random Forest

The final correlation results are shown in figure 3.40 with an RMSE of 10.69.

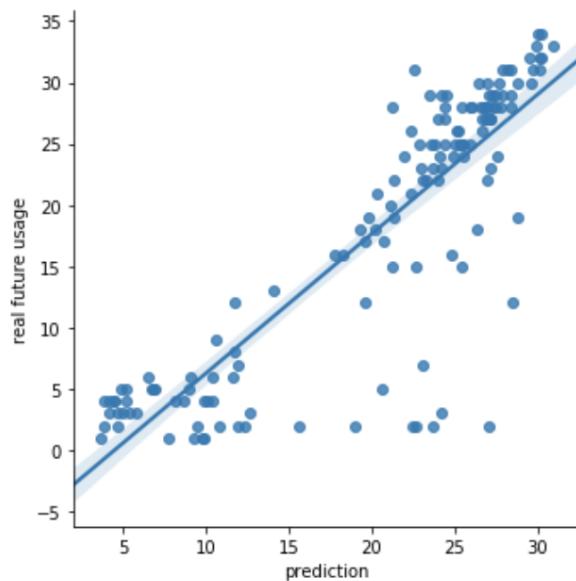


Figure 3.40: Data distribution of the actual and predicted data bike usage using Random Forest

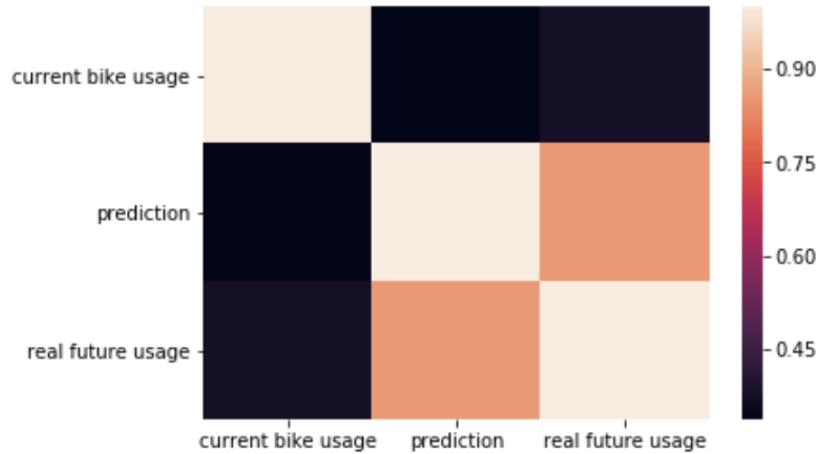


Figure 3.41: Heat map correlation plot showing interval value of the regression line (Random Forest)

Figure 3.39 shows the removed columns and RMSE results.

	TEST	Day	Month	Season	Week-day	Daily Weather	Daily Weather Past	Humidity	Humidity Past	Wind-speed	Wind-speed Past	Avg Temp	Avg Temp Past	Rented Bikes	RMSE	RMSE(on training)
Random forest regression	1	1	1	1	1	1	1	1	1	1	1	1	1	1	10.69	2.27

Figure 3.42: Feature evaluation test cases

Gradient boosted tree regression

The final correlation results are shown in figure 3.43 with an RMSE of 11.60.

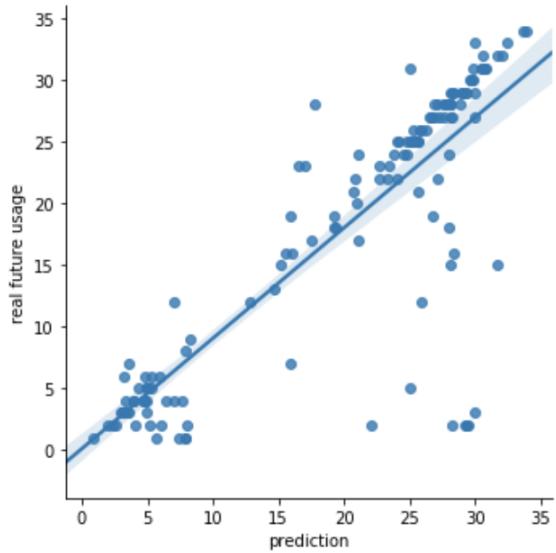


Figure 3.43: Data distribution of the actual and predicted data bike usage using Gradient boosted tree regression

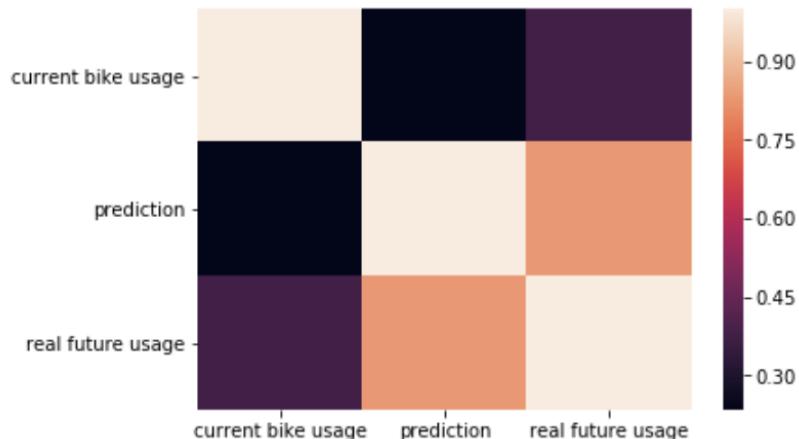


Figure 3.44: Heat map correlation plot showing interval value of the regression line (Gradient boosted tree regression)

Figure 3.45 shows the removed columns and RMSE results.

	TEST	Da-y	Da-n	Mont	Sea-son	We-ek-day	Daily-Weat-her	Daily-Weat-her-Past	Hu-mi-dity	Hu-mi-dity-Pas-t	Win-d-sp-eed	Win-d-sp-eed-Past	A-Tem-p-Avg	A-Tem-p-Avg-Past	Rente-d-Bikes	RMSE	RMS E(on training)
Random forest regression	1	1	1	1	1	1	1	1	1	1	1	1	1	1	10.69	2.27	

Figure 3.45: Feature evaluation test cases

3.5 Modeling (Polynomial Regression)

Since we don't have a linear relationship between the data, a linear regression will not be helpful. For example, the regress „Rented Bikes“ is not a linear correlation of temperature as even on rainy days there is slight chance that more people rent a bike than on a sunny day due to a special holiday. Therefore polynomial regression was a selected machine learning method for the prediction of rented bikes on a station.

Polynomial regression belongs to the regression forms [45]. In fact, it is just a modified version of a linear regression. This means the independent variable x and the dependent variable y is modeled as an n th digress (so-called polynomial) in x [45].

In a more formal way, the polynomial regression can be expressed as following:

$$Y = \beta_0 + \beta_1 * x + \beta_2 * x^2 + \beta_3 * x^3 + \dots + \beta_n * x^n$$

Where n is the degree of the regression.

With the Scikit-library in python a data scientist can import the function „PolynomialFeatures“ from „sklearn.preprocessing“ which transforms linear data into higher dimensional data. For example one could apply „poly = PolynomialFeatures(degree = 3)“ to get a polynomial regression in the third dimension. This should improve the accuracy as our underlying data has no linear relationships but maybe higher dimensional ones. Furthermore, the higher the degree the better the accuracy should be. Unfortunately the computation time is exponential. A degree of 4 already took several hours to perform and was only slightly better than a regression in the third dimension.

The RMSE error on a degree of 4 was around 48,4, which is in comparison to the other tested ones not really bad but maybe also not best one.

Figure 3.46 shows the different plots of each feature and the prediction (rented usage). It turns out that the feature Season, for example, has no substantial influence on the use of bicycles, whereby there was apparently an outlier in spring with 800 rented bicycles. It is also becoming apparent that bicycles are rented more frequently at low wind speeds than at high wind speeds. The average temperatures between 30 and 70 Fahrenheit are particularly high. Unfortunately, the addition of the past data has not caused much change in accuracy, as shown in figure 3.46.

Figure 3.47 shows a plot of the tested data (prediction) with the feature „Daily Weather“. The plot looks relatively good, except for a few single outliers at 2 (partly-cloudy-day), the prediction of the tested data matches the training data.

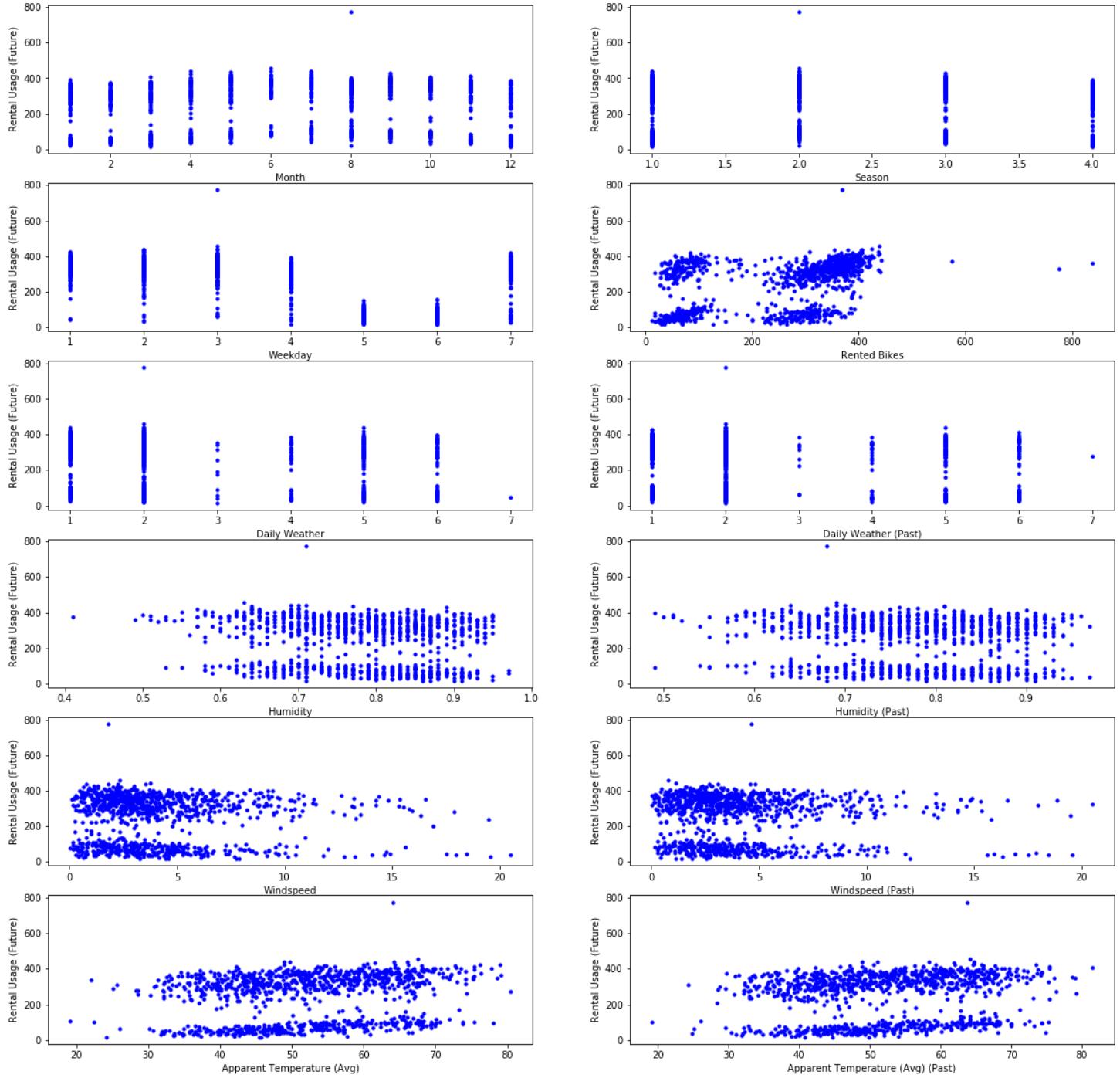


Figure 3.46: Plot of polynomial regression (features)

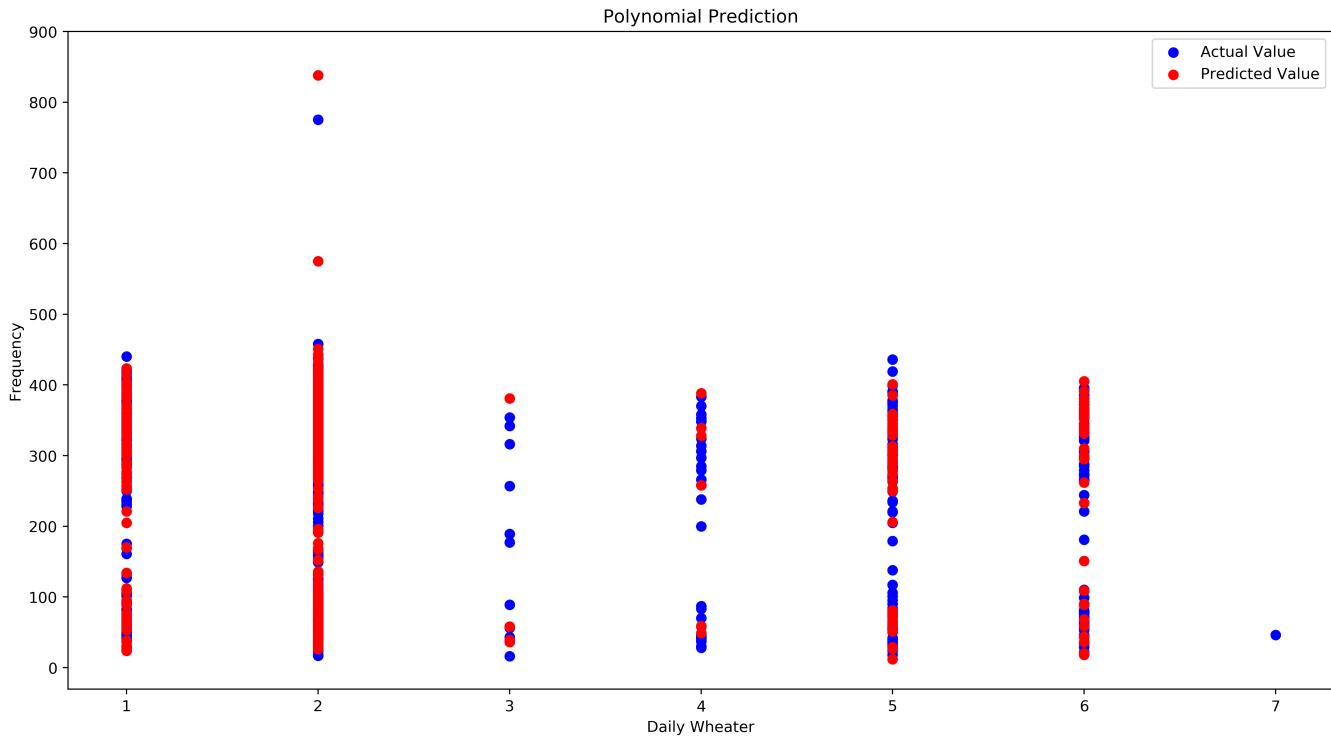


Figure 3.47: Polynomial Prediction of rental usage on daily weather

4 Prediction on Hourly Basis

Just like in chapter 3, data prediction is performed with the MLP regressor again. This time, however, the granularity is refined at an hourly level. This makes more training data possible. Furthermore, new correlations are created between the individual features. However, this also requires some changes in the data preparation script. Up to this stage of the project, the aggregation level at day level was sufficient. To be able to format the data from a daily format into an hourly format, both the weather data of the DarkSky API and the cycling usage data have to be re-processed, which is described in more detail in section 4.1. Since an hourly data prediction produces much more data than a daily prediction, this user story uses the built up Hadoop cluster as well as PySpark and Zeppelin. The tasks are typical tasks for a cluster-based calculation like Hadoop. Here, too, some parts are reused on the basis of previous scripts (e.g. the generated or calculated weather data).

4.2 deals with the training and prediction of the regression model based on the hourly prepared usage data. Subsequently it is described how the data set can be extended by additional features like "Age" or "local Events". The acquisition of these additional data was carried out by Ankur Mehra project team.

4.1 Data Preparation

Since the long term goal was to predict the usage on a hourly base, some further data transformation steps are necessary to achieve this goal. Furthermore the data records will increase significantly on an hourly granularity. Therefore the normal use of Anaconda and Jupyter on a local computer may be not sufficient due to low physical memory. An ideal use case for our newly installed Hadoop cluster! There PySpark can be used as already used under Data Profiling Part 1 in chapter 2 to manage „big data“transformations. Unfortunately the behavior and syntax of PySpark is sometimes a little more complicated than Pandas. For example, in PySpark it is not easily possible to iterate over rows since the data frame is distributed over the worker nodes and thus it only allows column wise operations. Moreover, Pandas operations such as iloc are not available in PySpark. But the API comes also with some advantages, e.g. it is quite performant on big data scale (i.e. it can easily perform several million of records) and it has an SQL approach. Functions like „select“, „where“and „filter“are syntactical close by to SQL as we know from MySQL and other database management systems.

The structure of the weather data is inconsistent due to "hourly weather" While the other columns only contain simple values, the column "hourly weather" contains nested JSON lists. This means that these nested lists must somehow become normal columns. This was a little more complicated than expected, but not impossible. Fortunately, PySpark allows one to define „schemas“that are used as a kind of blueprint by Spark to read the Spark data frame. With the following code one can already create normal columns from the JSON lists:

```

1 %spark2.pyspark
2 from pyspark.sql.functions import from_json
3 ...
4 #Defined schema
5 schema_hum = ArrayType(
6     StructType([StructField("time", StringType(), True),
7                 StructField("humidity", FloatType(), True)]))
8
9 rdd_weather = spark.read.csv("/user/hadoop/weather_new.csv",
10 header=True, sep=",")
11
12 #Add new column
13 rdd_weather = rdd_weather.withColumn(colChecker(i, "hum") + `i`,
14 from_json(rdd_weather["Hourly Weather"],
15 schema_hum)[k].getItem("humidity"))
16 ...

```

Figure 4.1: Transformation of humidity columns (excerpt)

The excerpt from 4.1 shows that „structTypes“ can be used to search the individual sub lists of „hourly weather“. The complete script is contained in the Zeppelin notebook „DFGeneration.json“ and can also be found on GitHub.

The described transformation creates a new column with the corresponding value for each hour. The script works dynamically. For example, the user can look at the data with two hours from today, which then looks like this:

Start Date	hum00	hum01	sum00	sum01	temp00	temp01	wind00	wind01
2015-01-04	0.93	0.94	partly-cloudy-night	partly-cloudy-night	37.24	36.35	0.94	0.56
2015-01-05	0.94	0.95	partly-cloudy-night		cloudy	37.62	37.97	0.37
2015-01-06	0.82	0.83		cloudy	45.99	45.87	1.44	1.58
2015-01-07	0.9	0.91	clear-night	clear-night	35.64	35.42	0.64	0.64
2015-01-08	0.9	0.9	partly-cloudy-night	partly-cloudy-night	47.21	50.5	6.27	5.51
2015-01-09	0.81	0.79	partly-cloudy-night	partly-cloudy-night	44.75	46.14	6.13	5.73
2015-01-10	0.82	0.85	partly-cloudy-night		cloudy	56.73	55.95	7.88
2015-01-11	0.76	0.79	clear-night	clear-night	35.06	35.14	4.56	4.27

Figure 4.2: Weather dataframe after transformation for two hours

This (figure 4.2), however, is more like a jumping window technique as it uses days as start date and not the single hours on each day. Neverthe-

less the structure of the data frame from figure 4.2 can be used as ground for applying sliding window method.

However, before the sliding window procedure could be implemented, the „Start Date“ from figure 4.2 had to be converted to an hourly format, whereby the individual hour columns had to be combined so that only one „current column“ existed. For example, „hum00“ and „hum01“ should become something like „Current Humidity“ which shows the humidity value of each hour. This means that further data preparation steps are necessary.

First, the exact timestamps of the bicycle data were read in from the provided TfL website and rounded off to an hourly level. For example, „23:38“ became „23:00“ Minutes and seconds were ignored. An arithmetic lap (half-round mode) did not make sense, because then there would be overlaps with the successor day when 23 o'clock is rounded up. However, this also means that a slight distortion must be assumed. For example, many bicycles could be rented at a bicycle station at 17:52 and significantly less after 18 o'clock. Then one would notice a peak at 17 o'clock and a low usage at 18 o'clock, although in fact more bicycles were rented around 18 o'clock.

Another problem was that not for every hour there was data about the usage of the bicycle stations. However, this data quality problem could be solved relatively easy. The missing hours can be determined by a user defined function. Since every day has 24 hours, this can be calculated manually as following (figure 4.3):

```

35 + row['StartStation capacity'] + '</br>' + " Used between 2016 and
36 2018: " + str(int(row['StartStation Id Used']))
37
38
39 #Put dates and coordinates in required format
40 features = [
41     {
42         'type': 'Feature',
43         'geometry': {
44             'type': 'LineString',
45             'coordinates': coordinates[i],
46         },
47         'properties': {
48             'times': dates[i],
49             'style': {
50                 'color': 'red',
51                 'weight': weights[i][0],
52                 'opacity': 1
53             }
54         }
55     }
56     for i in range(len(dates))
57 ]
58
59 #Allows polylines over time (i.e. red lines), JSON format required;
60 TimestampedGeoJson({
61     'type': 'FeatureCollection',
62     'features': features,
63 },
64     auto_play=False,
65     loop=False,
66     loop_button=True,
67     date_options='YYYY/DD/MM',
68     duration='P2M',
69     add_last_point=False).add_to(m)
70
71 m.save(outfile= "routes_frequency.html")
72 m
73

```

Figure 4.3: Adding missing hours on usage data

A prerequisite of this method from figure 4.3 is that at least one maximum value and one minimum value exist that serve as boundaries. If no data was collected on a day at a station, null values would appear. This does not happen, however, and even if it did, it would be removed at a later stage

as such data does not provide any advantage for training a machine learning model.

Another problem was how to transform the individual hour columns of the weather data so that only one column exists instead of e.g. 24 columns, whereby this new column should contain all values of the 24 columns arranged correctly. Fortunately the function „explode“ exists in PySpark and is provided via the class „pyspark.sql.functions“ With „explode“ an array can be passed, which then returns a new line for each element of the array [31]. The corresponding code excerpt looks as following:

```

1  %spark2.pyspark
2  import pyspark.sql.functions as F
3  from pyspark.sql.functions import lag, col, lead, first
4  from pyspark.sql.window import Window
5  from collections import Counter
6
7  def _combine(x,y):
8      """creates hourly interval for weather df"""
9      print(y)
10     d = str(x) + ' [5]:00:00'.format(y)
11     return d
12
13 def transformWeatherDF(df, past):
14     """add hours on key column and transform columoriented hours to
15     roworiented hours"""
16
17     #Remove daily columns
18     df = df.drop(*["Apparent Temperature (Avg)", "Daily Weather",
19     "Hourly Weather", "Humidity", "Windspeed"])
20
21     #use udf to generate hours on the dates
22     combine = F.udf(lambda x,y: _combine(x,y))
23
24     #fetch relevant columns and keep column datatype
25     cols_temp, dtypes = zip(*((c, t) for (c, t) in df.dtypes if c not
26     in ['Start Date'] and c.startswith("temp")))
27
28     ##temperature
29     #explode function for changing structure of dataframe (temperature)
30     kvs = F.explode(F.array([
31         F.struct(F.lit(c).alias("key"), F.col(c).alias("val")) for c
32     in cols_temp],)).alias("kvs")
33
34 ...
35     # go x hours back (temperature)
36     for i in range(past):
37         w = Window().partitionBy().orderBy(col("Start Date"))
38         df_temp = df_temp.select("*", lag("Current Temperature",
39         i+1).over(w).alias("temp"+`i`))
40 ...
41
42 ...

```

Figure 4.4: Further transformation of weather data (excerpt)

The code example from figure 4.4 is used to transform the temperature columns, which are combined into a single column as described above. Thus the weather data frame is prepared to the extent that every hour of a day

belongs to a specific hourly weather value. The lines 35 - 39 show the use of the PySpark function „lag“, that over a selected „window“ takes the starting boundary [31]. If one increments the index at this point, he will get the previous value from the window. Depending on how far one wants to go into the past, the window is moved over the data set and thus a sliding window is achieved. This can also be adapted for the future. With the function „lead“ the ending boundary of a „window“ is retrieved [31]. Again, the future sliding window is only applied for the usage (target variable) and not for the weather data or other columns.

The prepared hourly temperature data with sliding window looks finally as following:

	Start Date	Current Temperature	temp0	temp1	temp2
2015-01-04 00:00:00		37.24	null	null	null
2015-01-04 01:00:00		36.35	37.24	null	null
2015-01-04 02:00:00		35.61	36.35	37.24	null
2015-01-04 03:00:00		34.36	35.61	36.35	37.24
2015-01-04 04:00:00		33.43	34.36	35.61	36.35
2015-01-04 05:00:00		33.14	33.43	34.36	35.61
2015-01-04 06:00:00		32.86	33.14	33.43	34.36
2015-01-04 07:00:00		32.98	32.86	33.14	33.43
2015-01-04 08:00:00		33.09	32.98	32.86	33.14
2015-01-04 09:00:00		33.65	33.09	32.98	32.86
2015-01-04 10:00:00		34.03	33.65	33.09	32.98

Figure 4.5: Temperature dataframe after transformation for 3 hours (past)

In this example from figure 4.5 it is clearly recognizable that the initial values contain „null“ values, which is due to the fact that no weather data was available before 04.01.2015 (or at least they were not fetched from Dark Sky). The number of zero values increases with the number of selected hours into the past. The same applies to the last lines of the Spark data frame, where null values may appear for the „future usage“ columns. Therefore the first 20 and the last 20 lines of the data frame are removed at the end.

With the described PySpark transformations it was possible to create dynamic hourly data frames, which can then be used in the next step of the

modeling phase. A corresponding test file can be found on GitHub in the folder data preparation.

Holidays

To add the bank holidays are day-based, so they should be set on a 24 hour windows. However, fixing the bank holidays from midnight to midnight the next day might not be the best solution as people might consider the evening of a bank holiday the same way as they do for normal week day, while the evening of the day before will be more interesting. Indeed, when people want to go out in the evening and go to bed late, they will prefer doing that when they don't have to wake up early in the morning, thus they may be more likely to use a bike the end of the day before the bank holiday than of the day itself.

This is why it might be a good idea to test the accuracy of the prediction placing bank holidays from 6pm the day before to 6pm the day itself.

4.2 Learning and Prediction

For the learning part, we used the same algorithm as described in ???. We trained twenty-four models, one per future hour to predict. For each model, we only use the past and current data which we match to the future expected number of rented bike for the given following hour, no other future data is used.

We used the sliding-windows hourly based data over 48 hours (aggregating per slice of 6 hours after the first 24 hours). We have an accuracy of 0.800638 on the test dataset. On the graph below, you can see several randomly selected days in the testing set starting at midnight and predicting for the whole day. The title above each graph indicates the date displayed, the actual expected value is in red and the predicted one is in blue.

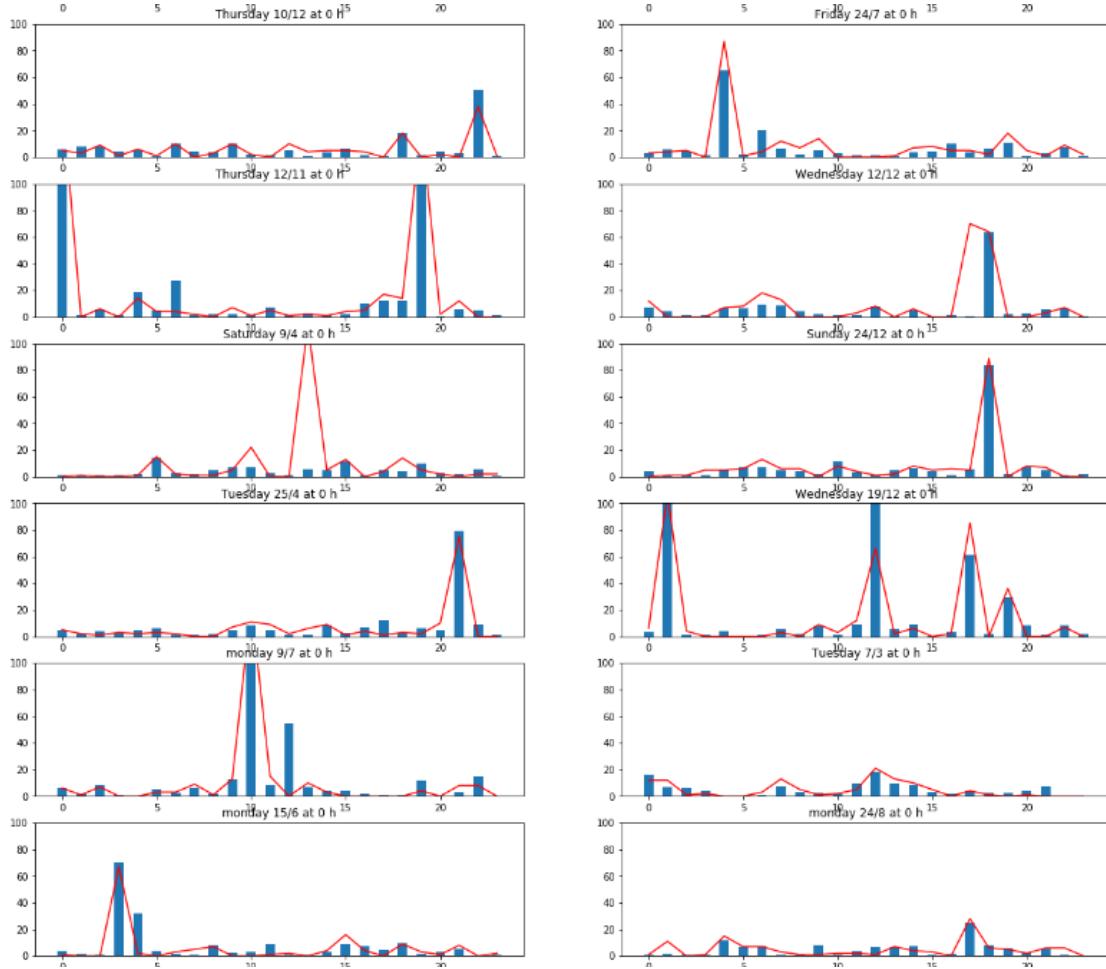


Figure 4.6: Comparison of expected (red) and predicted (blue) hourly rented bikes on test dataset

As we can see on the Figure 4.6 and considering the measured accuracy, we can see that we can predict quite well the amount of rented bikes for the following 24 hours based on the data of the previous 48 hours.

5 Additional Features

In order to make the predictions more precise we decided to improve the feature matrix with additional features. We took into consideration which

data could be meaningful for the frequency of bike rentals. We came to the conclusion that the following 4 features should be added:

- Age
- Events
- Earnings
- Political Attitude

We considered the first point because we assumed that there is a specific age group which uses rental bikes more than others. We conclude from that, if we find stations where a lot of people within this specific age group lives, we will have a higher frequency there. Therefore we investigated in the target group of Santander Rental Bikes and found out that their main target group depicts of people between the age of 16 and 54 [42]. Moreover approximately two third of their users are male. That is why we searched for data which gives us the total number of ages as well as the gender for each district in London where Santander provides rental station data.

Appropriate data was found at the *London Datastore* where data about London's population is freely provided as an excel document [11]. To use this data we first extracted only the district data we are interested in. Since there are ages from 1 to 99 available we created the same age groups like the ones Santander uses for their survey data and summed up the according values. Since we know that males using bikes more often we assigned the age group data to females as well as to males separately. With this task the data cleansing part was done. How the data was prepared further will be explained in the following chapter 5.1 Data Preparation.

Another feature we wanted to take into account is event data. We assume that when more events take place at a specific location, more people will visit this location and out of that more stations around this corresponding location will be used. For this purpose we instructed Ankur Mehra, who has done web scraping tasks before, to scrape appropriate data. The tasks was to find a source which provides event data of London over the last 4 years.

Event data was scraped from *LondonTown.com*[22]. The description of the

implementation can be found in the report of Ankur Mehra.

Santander distinguishes between *Casual Users* which use their bikes occasionally and just pay per ride and *Members* which have their own keys and pay beforehand to use bikes for one year. According to their last user survey they own more members than casual users. Moreover they found out that members have a higher average annual income than casual users, we investigated in income data as well. Our assumption was the districts in which more people with a higher income live are the districts where more potential members live and therefore where more frequency will take place.

The *Office for National Statistics* provides data of earnings in the UK from 2018, which we used as a source [35]. The data cleansing part consisted of extracting the data of the 15 districts we are interested in which can be seen in figure 5.1.

	2015	2016	2017	2018
District				
Camden	42161.6	43165.2	42406.0	45084.0
City of London	60351.2	62696.4	63642.8	68151.2
Hackney	36550.8	38854.4	39249.6	39124.8
Hammersmith and Fulham	36831.6	38818.0	42692.0	43097.6
Haringey	34008.0	32656.0	34470.8	35058.4
Islington	45099.6	47762.0	48984.0	48656.4
Kensington and Chelsea	35952.8	34039.2	38303.2	39993.2
Lambeth	41574.0	39301.6	40211.6	41646.8
Lewisham	33649.2	34403.2	34637.2	33919.6
Newham	32541.6	34377.2	36878.4	36296.0
Southwark	41449.2	40970.8	43004.0	46191.6
Tower Hamlets	57210.4	55302.0	60803.6	59144.8
Wandsworth	35620.0	36166.0	38204.4	40320.8
Westminster	44657.6	45614.4	49597.6	52364.0

Figure 5.1: Cleaned earnings data

The last additional feature we took into account was the political attitude. Since using bikes instead of cars or public transportation is healthier

and more environmental friendly, we assume that the most Santander users are interested in health, environmental protection and climate protection. Therefore we investigated in which parties have the most content according to those subjects, because we assume that voter of those parties will be more likely to use Santander bikes. Therefore districts with the most voters of according parties will show a higher frequency of rental bike usage.

Therefore we researched for specific content of election programs of the several parties within the UK. We found out that the *UK Independence Party* have mentioned least content regarding health and climate. Slightly more information was found in the election program of the *Liberal Democrats* as well as the *Conservatives*. The most content to health, climate and environment was found on the agenda of the *Labour Party* as well as the *Green Party* [41].

Therefore we assumed that districts with a high amount of Labour Party and Green Party voters will use rental bikes more often than other districts.

We searched for according data and were resourceful once again at the London Datastore. Here we found Data of the last four elections. We used the data of the election from 2016 and once again extracted all the data from the districts we are interested in.

Nevertheless those are all just assumptions. Of course there are several other reasons than trying to live healthier or doing something good for the environment, to use rental bikes. A lot of people will use it probably only because it is fun, cheaper and they may not have a driving license for cars.

5.1 Data Preparation

District to Stations

Nearly all of the additional data is stored by London district. To be used with our stations, we need to assign to every station a corresponding district. We first determined the centroid of the districts using Nominatim, which we compared to each station's coordinates to keep the closest district to each station. This enables us to link the *age*, *earnings* and *political attitude* datasets the stations.

Filling Missing Years

The data we extracted doesn't necessarily contain data for the same duration as our bike rental data. Thankfully, some data don't check much during a few years like *age*, *earnings* and possibly *political attitude* as well. With this supposition, we decided to fill missing periods with the closest data we have.

Events

First, the events are located with a postcode and a street, so we use Nominatim to extract their coordinates. For a few postcodes, we have to use Geonames as some are missing or no more used officially. Last, there is an event called "secret cinema" which happens at a location that is disclosed shortly before the rendez-vous time, so the postcode and street given don't exist.

Second, the date field follows many different formats, from a single date, to a list, or a range. Here is a list of different formats found:

- day month year
- day-of-week day month year
- day — day month year
- day-of-week day — day-of-week day month year
- day month — day month year
- day month to day month year
- day and day month year
- day, day and day month year
- day month, day month and day month year
- "15th May — 05th March 2016" (over nearly a year)
- "Various Venues"
- ...

We had to implement different RegEx parsing in order to parse all possibilities found in the database. Once parsed, we transformed the dates into a list of date (so “15th May — 15th June 2019” creates 32 days between those two dates included) to store each event each day it appears.

Finally, we created two features for each day in our original features dataset. For each day, we take the distance to the closest event and we count the number of events in a close perimeter (3 kilo-meters for instance).

5.2 Data Evaluation

Demographic, Earnings, Political Attitude

The demographic, earnings and political attitude data per district has been placed in a table shown in figure 5.2.

District	Nb inhabitants	Nb stations	Nb rented	r / station	r / person	person / station	Income 2018	% Con	% Lab	% Lib Dem	% Green	% UKIP
Newham	278090	2	3174	1587	0.011414	139045	36296.0	17.423339	65.307076	2.291541	4.473610	2.558457
Wandsworth	269043	65	26320	404	0.097828	4139	40320.8	40.797920	43.278400	3.791946	5.429276	1.620414
Hammersmith and Fulham	153076	32	33821	1056	0.220943	4783	43097.6	41.089339	39.768094	4.316411	5.916356	2.282932
Tower Hamlets	258943	39	85783	2199	0.331281	6639	59144.8	17.854048	61.271647	3.548958	7.121777	2.358437
Hackney	227081	53	140368	2648	0.618141	4284	39124.8	12.093692	66.721942	2.896786	10.193162	1.382321
Southwark	267121	45	258546	5745	0.967898	5936	46191.6	19.382631	54.966694	8.403446	8.316952	2.529683
Islington	202860	46	332246	7222	1.637809	4410	48656.4	16.990974	60.377152	4.623599	8.741961	2.371882
Camden	211902	69	592450	8586	2.795868	3071	45084.0	27.284130	51.053337	5.129850	8.448074	2.033980
Lambeth	279566	73	806646	11049	2.885351	3829	41646.8	21.426381	56.114814	5.388028	9.172221	1.449483
Westminster	217558	134	990930	7395	4.554785	1623	52364.0	44.371538	38.007135	4.247118	4.962441	2.113699
Kensington and Chelsea	134271	118	731817	6201	5.450298	1137	39993.2	55.599084	28.015772	4.075299	4.533198	1.836683
City of London	9177	110	1202336	10930	131.016236	83	68151.2	40.415105	37.647737	6.630153	7.004901	2.075526

Figure 5.2: Additional data per district

We can see that the district with the most rented bikes has also nearly the most station, but has the least amount of person living in it. This is the *City of London*. The City is a district used mostly as a workplace by Londoners and many of them live in other districts, so this might suggest that people are more likely to rent a bike near their workplace than near their living place. It might then be more useful to obtain the data based on where people work (and possibly also where they go shopping and have a meal) rather than where people live, but this kind of information isn’t usually available.

Furthermore, since it was decided to use a different learning model for each station, the *age*, *earnings* and *political attitude* datasets might not add much accuracy. They could, however, be used to predict the usage of a new station based on its location to determine whether it is relevant to proceed with construction or not.

Events

Concerning the events additional feature, it might be interesting to have more accurate information about the events: hour, duration, amount of people expected, kind of place/kind of event, ... Indeed, if the event takes place too late in the evening or involves drinking alcohol (if it is in a bar for example), people might not be wanting to use a bike. Also, if it doesn't expect many people to come, there won't be any noticeable change in the amount of bike rented. Additionally, some events last for nearly a year and might not create bike rental change during this whole period, rather at some precise dates but we don't have any details.

Chapter 4

Conclusion

1 Conclusion

This report presents in overall a model of a machine learning pipeline, this data science method can be applied in every domain since it provides a structured way to extract meaningful information out of a huge raw amount of data. We started as team of four members, we get initially ourselves acquainted with the different required technologies which were mostly unfamiliar to us. We afterwards set up the hardware, then processed, profiled and scraped the data. Finally, we applied machine learning regression algorithms to make predictions. More members joined our group afterwards, which made us think of new tasks and thus new ways to extend the project scope. This project was started from scratch and the results discussed previously were of course limited to our initial prediction goal. This means that it's still open to a lot of enhancements, especially from new members with fresh ideas about extracting new meaningful information out of the current data. This report presented for us a new way of collaboration in a group, for some of us it was more importantly an incentive to delve and specialize in this area in future works.

Bibliography

- [1] Altoros. Hadoop Distributions: Evaluating Cloudera, Hortonworks, and MapR in Micro- benchmarks and Real-world Applications. pages pp. 1 – 65, 2013.
- [2] Apache Foundation. Apache Hadoop, 2018.
- [3] Bala Deshpande. 2 main differences between classification and regression trees, 2011.
- [4] Jason Brownlee. A Gentle Introduction to the Gradient Boosting Algorithm for Machine Learning, 2016.
- [5] Cloudera. CDH Components, 2018.
- [6] Cloudera. Cloudera Pricing, 2018.
- [7] Cloudera. Installing CDH 5 Components, 2018.
- [8] Cloudera. Operating System Requirements, 2018.
- [9] D. Kumar. Since Cloudera and Hortonworks are 100% open source, can I use them freely as I would a Linux distribution?, 2016.
- [10] Databricks. Machine Learning Library (MLlib), 2019.
- [11] London Datastore. GLA Population and Household Projections.
- [12] GeoNames. Export Zip Codes.
- [13] Hortonworks. Apache Ambari Installation, 2018.
- [14] Hortonworks. Grundlagen der HDP-Sandbox, 2018.

- [15] Hortonworks. Starting HDP Services, 2018.
- [16] Hortonworks. Support Matrix, 2018.
- [17] Hortonworks. Apache ambari operations. available on: https://docs.hortonworks.com/hdpdocuments/ambari-2.6.2.2/bk_ambari-operations/content/ch_using_ambari_metrics.html. 2019.
- [18] Hortonworks. Check dns and nsqd. available on: https://docs.hortonworks.com/hdpdocuments/ambari-2.7.3.0/bk_ambari-installation/content/check_dns.html. 2019.
- [19] I. K. Center. IBM Knowledge Center - BigInsights features and architecture, 2018.
- [20] IBM Analytics. Hadoop-Testversionen, 2018.
- [21] G. Kirill. Comparing the top Hadoop distributions, 2014.
- [22] LondonTown. London Events 2019, 2019.
- [23] Ross Lydall. Boris johnson's bike hire scheme gets a 25m bonus from barclays. available on: <https://web.archive.org/web/20100913111233/http://www.thisislondon.co.uk/standard/article/23839406-boris-bike-hire-scheme-gets-a-pound-25m-bonus-from-barclays.do>. 2010.
- [24] M. Heo. Hortonworks vs. Cloudera vs. MapR, 2015.
- [25] MapR. MapR 6.1 Documentation. 2018.
- [26] MapR. MapR Editions, 2018.
- [27] MapR. Operating System Support Matrix (MapR 6.x). 2018.
- [28] B. Marr. Big Data: Who Are The Best Hadoop Vendors In 2017?
- [29] Merv Adrian. IBM Ends Hadoop Distribution, Hortonworks Expands Hybrid Open Source, 2017.
- [30] Microsoft. Manage Apache Hadoop clusters in HDInsight using Azure portal, 2018.

- [31] n.d. Apache pyspark documentation. available on: <https://spark.apache.org/docs/latest/api/python/pyspark.sql.html?highlight=explode>. 2019.
- [32] n.d. Dark sky weather api. available on: <https://darksky.net/dev/docs/faq>. 2019.
- [33] n.d. Folium package. available on: <https://python-visualization.github.io/folium/>. 2019.
- [34] Nominatim. OpenStreetMap Nominatim: Search, 2019.
- [35] Office for National Statistics. Employee earnings in the UK, 2019.
- [36] OpenStreetMap. PyrouteLib, 2017.
- [37] OVirt. Virtio-SCSI, 2019.
- [38] Pivotal. Overview of Apache Stack and Pivotal Components, 2018.
- [39] Martin Riedmiller and Heinrich Braun. A direct adaptive method for faster backpropagation learning: The rprop algorithm. In *Proceedings of the IEEE international conference on neural networks*, volume 1993, pages 586–591. San Francisco, 1993.
- [40] Sayad Saed. Decision Tree Regression, 2019.
- [41] Dominic Higgs Sandra Laville, John Sauven. How do the four main parties compare on the environment?, 2017.
- [42] Santander. Santander Cycles Customer Satisfaction and Usage Survey Casual Users Only: Quarter 2 2017/18. Technical report.
- [43] Jeff Sposetti. Ambari blueprints. available on: <https://cwiki.apache.org/confluence/display/ambari/blueprints#blueprints-step1:createblueprint>. 2017.
- [44] TFL. cycling.data.tfl.gov.uk, 2019.
- [45] Yongqiao Wang, Lishuai Li, and Chuangyin Dang. Calibrating classification probabilities with shape-restricted polynomial regression. *IEEE transactions on pattern analysis and machine intelligence*, 2019.

- [46] H. Kisker Yuhanna and D. Murphy M. Gualtieri N. Big Data Hadoop Solutions, Q1 2014. *The Forrester Wave*, 2014.

B Illustration Directory

2.1	Overview of project architecture	3
1.1	Connection of Nominatim with geopy	15
1.2	Reverse Geocoding	18
1.3	Centroid of Baden-Wuerttemberg shown in Google Maps . . .	19
1.4	Compute Distances as the Crow Flies shown in Google Maps .	21
1.5	Distance shown in Google Maps	23
3.1	ER model for the MySQL Database	27
4.1	Theoretical example of distance calculation	30
5.1	Mockup for Django web app	32
5.2	Final layout of the Django web application	34
1.1	Forrester Wave of Hadoop Distributions [46]	40
1.2	Weighted Comparison Table of Hadoop Distributions	43
1.3	Spider Chart of compared Hadoop Distributions	44
1.4	Micro benchmarks on DFSIO / WordCount MapReduce Jobs .	45
1.5	etc/hosts configuration file	47
1.6	Ambari Wizard Installer	48
1.7	Ambari Surface, Host metrics on master node	50
1.8	PI Quasi Monte Carlo method executed as Cluster job	51
1.9	Output of WordCount job	52
1.10	Proxmox surface	57
1.11	Webmin dashboard	58
1.12	Physical Hadoop cluster	61
2.1	Frequency of used StartStations and EndStations	67
2.2	Frequency of used StartStations and EndStations	68
2.3	Top 10 used Santander rental stations	69
3.1	Route usage over time plot	75
3.2	Used bicycle routes over time plot (high usage)	76

3.3	Used bicycle routes over time plot (low usage)	77
3.4	Fetching the weather API with Python	79
3.5	Prepared data frame of Kings Cross (daily)	81
3.6	Result of MLPRegressor with MinMaxScaler	83
3.7	Feature evaluation test cases	84
3.8	Accuracy with recommended features and scaling by the QuantileTransformer	85
3.9	Accuracy with recommended features and scaling by the QuantileTransformer	86
3.10	Decision tree usage diagram	88
3.11	Random Forest Regression: Process	90
3.12	Test 1 features importance: 71.35	92
3.13	Feature evaluation test cases	92
3.14	Test 3 features importance: 57.55	93
3.15	Test 4 features importance: 59.27	93
3.16	Test 5 features importance: 56.98	94
3.17	Test 6 features importance: 56.98	94
3.18	Test 7 features importance: 54.43	95
3.19	Data distribution of the actual and predicted data bike usage using Decision Tree	95
3.20	Heat map correlation plot showing interval value of the regression line (Decision Tree)	96
3.21	Test 1 features importance: 60.04	96
3.22	Feature evaluation test cases	97
3.23	Test 2 features importance: 62.41	97
3.24	Test 3 features importance: 60.98	98
3.25	Test 4 features importance: 64.64	98
3.26	Test 5 features importance: 63.08	99
3.27	Data distribution of the actual and predicted data bike usage using Random Forest	99
3.28	Heat map correlation plot showing interval value of the regression line (Random Forest)	100
3.29	Test 1 features importance: 81.30	100
3.30	Feature evaluation test cases	101
3.31	Test 2 features importance: 88.48	101
3.32	Test 3 features importance: 84.92	101
3.33	Test 4 features importance: 76.71	102
3.34	Test 5 features importance: 63.26	102

3.35	Data distribution of the actual and predicted data bike usage using Gradient boosted tree regression	103
3.36	Heat map correlation plot showing interval value of the regression line (Gradient boosted tree regression)	103
3.37	Data distribution of the actual and predicted data bike usage using Decision Tree Regression	104
3.38	Heat map correlation plot showing interval value of the regression line (Decision Tree Regression)	104
3.39	Feature evaluation test cases	105
3.40	Data distribution of the actual and predicted data bike usage using Random Forest	105
3.41	Heat map correlation plot showing interval value of the regression line (Random Forest)	106
3.42	Feature evaluation test cases	106
3.43	Data distribution of the actual and predicted data bike usage using Gradient boosted tree regression	107
3.44	Heat map correlation plot showing interval value of the regression line (Gradient boosted tree regression)	107
3.45	Feature evaluation test cases	108
3.46	Plot of polynomial regression (features)	110
3.47	Polynomial Prediction of rental usage on daily weather	111
4.1	Transformation of humidity columns (excerpt)	113
4.2	Weather dataframe after transformation for two hours	113
4.3	Adding missing hours on usage data	115
4.4	Further transformation of weather data (excerpt)	117
4.5	Temperature dataframe after transformation for 3 hours (past)	118
4.6	Comparison of expected (red) and predicted (blue) hourly rented bikes on test dataset	120
5.1	Cleaned earnings data	122
5.2	Additional data per district	125

C List of Abbreviations

YARN	Yet Another Resource Negotiator	39
-------------	---	----

SQL Structured Query Language.....	26
NoSQL Not only Sql.....	26
PL/pgSQL Procedural Language/PostgreSQL.....	8
OSM OpenStreetMap.....	8
HDP Hortonworks Data Platform.....	39
HDFS Hadoop Distributed File System	39
OSM OpenStreetMap.....	8
DFSIO Data File System I/O.....	45
MLPRegressor Multi-layer Perceptron Regressor	81
RMSE Root Mean Square Error	82
MLlib Machine Learning Library.....	86
ALS alternating least squares	87
SVD singular value decomposition	87
PCA principal component analysis.....	87

SSE Sum of Squared Errors.....	89
--------------------------------	----