

Capstone Proposal

Machine Learning Engineer Nanodegree

Christian Wolff

October 27, 2017

Proposal

Domain Background

In recent years, bike-share programs in cities around the world have grown in popularity. [Wikipedia](#)¹ currently lists 330 bike-sharing systems across the world. Their ultimate goal is to create a more sustainable transportation landscape that provides new mobility options for short trips and to improve connectivity to other modes of transportation.

As bike-sharing programs grow in popularity, cities begin to realize additional benefits which should mean plenty of related opportunities for years to come. But there are also challenges to overcome. One of the challenges are the running costs of said programmes, mostly from maintenance of bikes and stations, redistribution of bikes by vans during the day and the necessary IT infrastructure.

In this project, we will take a closer look at [Hubway](#)², Boston's regional bike sharing system. As of September 2017, it gives members access to more than 1,700 bikes at 170 stations. The station network is designed to offer a one-way transportation option and an alternative to public transport like bus or train.



A Hubway station

Problem Statement

Bike-sharing services can only be successful when they meet the demand for mobility. In this regard, two important factors for customer satisfaction are

- high availability of bikes at all stations,
- and the possibility to return a bike to any station.

But cities are one of the most complex systems known to us, and far from being homogenous or in equilibrium. Thus, some bike-sharing stations are frequented more often than others, with changing distributions over the day and week. This means, that bikes have to be continuously redistributed to meet the demand.

Hubway tackles this problem on a real-time basis, as we can read on their website:

Hubway employs 4-5 rebalancing vans, each with a payload of 20-25 bikes, used to redistribute bicycles between 6am and 10pm, 7 days a week. Bicycles are redistributed dynamically in response to real-time data.

¹https://en.wikipedia.org/wiki/List_of_bicycle-sharing_systems#Cities

²<https://www.thehubway.com>

From this statement follows the question:

Can Hubway predict the demand for bikes and plan the necessary rebalancing ahead of time?

I assume that proper planning of rebalancing will decrease running costs and improve customer satisfaction. The goal of this project is the implementation of a workflow to predict the daily demand of bikes at the stations, so that rebalancing can be planned ahead of time, instead of relying on real-time rebalancing. The problem will be stated as a regression task, with the number of bike departures and arrivals at each station as target variables.



A Hubway rebalancing van

Datasets and Inputs

From the very beginning of their bike-sharing program, Hubway has collected usage data. In 2012, they ran a public challenge to visualize and analyze the data. The results of the challenge are available on [Hubways' official challenge website](http://hubwaydatachallenge.org)³.

Hubway publishes trip data every quarter on its official [system data page](https://www.thehubway.com/system-data)⁴. The main data source for this project will be data from January 2015 through November 2016, containing over 2.3 million bike trips collected at almost 200 stations. The set contains fine-grained data of single trips, e.g. start and end timestamp of each trip, together with the respective station.

Field definitions in Hubway trip data	
tripduration	Time of trip in seconds.
starttime	Start timestamp of trip with date and time, in EST.
stoptime	End timestamp of trip with date and time, in EST.
start station id	ID of start station.
end station id	Station ID of end station.
bikeid	ID of bicycle used.
usertype	Subscription type. "Registered" is a user with membership. "Casual" is a user without membership.
birth year	Birth year of the user, self-reported by member
gender	Gender of the user, self-reported by member

The number of arrivals and departures of bikes at each station per hour can be calculated by aggregation of the trips data. The trips data has to be cleansed first by removing any unusual trips (very short or long trips) and all features that can not be aggregated by hour (e.g. gender).

I expect the trip dataset to be correlated with the following three factors, which will be used to augment it:

- Local weather conditions, mostly temperature and precipitation. This includes effects of the calendar seasons. Hourly historic weather summaries are available from [DarkSky](https://darksky.net)⁵.
- We can expect a baseline of regular commuters during work days, and spikes of casual usage during weekends and public holidays. Information about public holidays in Boston can be acquired from the [Massachusetts state government](https://www.mass.gov/)⁶.
- Local events, like festivals and marathons increase the demand for mobility. These can be taken into account by using the methodology of event labeling, as proposed by [Fanaee-T and Gama](https://doi.org/10.1007/s13748-013-0040-3)⁷.

³<http://hubwaydatachallenge.org>

⁴<https://www.thehubway.com/system-data>

⁵<https://darksky.net>

⁶<https://www.mass.gov/>

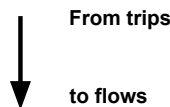
⁷<https://doi.org/10.1007/s13748-013-0040-3>

Solution Statement

The maximum demand at a station for a given day can be evaluated as follows:

1. From the trips dataset, compute the flow of bikes per hour. The flow at a station is defined as the number of bikes that are returned to the station, minus the number of bikes rented from the station, during a given hour.
2. Derive the total number of bikes for each hour from the flow. Assume that the station has no bikes at exactly midnight. The total number of bikes at a given hour is then the rolling total of the flow until this hour.
3. It is clear that the total number of bikes can be negative under the assumption that the station is empty at midnight. By taking the minimum of the running total of the day, we can evaluate how many bikes are actually needed at the start of the day, so that the station is never depleted.

Trip ID	Start time	Start station	End time	End station
1	7/28/2011 12:11:00	2	7/28/2011 12:23:00	8
2	7/28/2011 16:08:00	5	7/28/2011 16:20:00	4
3	7/28/2011 19:18:00	8	7/28/2011 19:46:00	1



	In	Out	Flow	Total
Midnight to 6am	0	1	-1	-1
6am to 12am	0	3	-3	-4
12am to 6pm	5	3	2	-2
6pm to midnight	6	2	4	2

← We assume, that the station is empty at midnight.

← The minimum total over the day is -4. This means, that we actually need to put 4 bikes at this station at midnight.

(No causal relation between the tables.
Flow data is for one day and one station.
Hourly data condensed into four time periods for illustration purposes.)

Transforming trips into flows

The maximum demand for each station and day can then be augmented with daily weather summaries and information about public holidays and weekends. The resulting dataset can be used to train different regression models or do unsupervised clustering for more data exploration.

So far I assumed that the rebalancing is done once each day at midnight. Because stations can also 'overflow', rebalancing is probably needed several times a day. I will try to extend the project to a point where it is possible to plan a good balancing of bikes, such that each station is close to half it's maximum capacity at any given time during the day.

Benchmark Model

As far as I know, Hubway has not yet implemented a predictive system for the given problem. This means that there is no existing benchmark model we can directly compare to. Instead, I will train a simple regression model with default parameters, like Gaussian Naive Bayes or a single decision tree. The performance of the simple model on the training and test sets will then serve as benchmark for the more sophisticated approach.

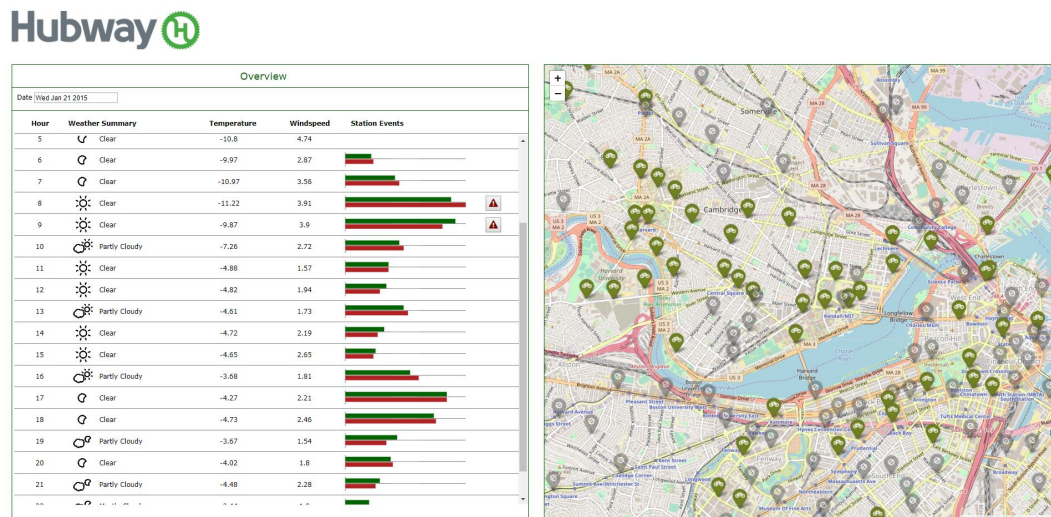
Another benchmark we could use is cost efficiency, by comparing the current costs for real-time rebalancing to the expected cost from the proposed machine learning approach. These costs can be approximated by the distance the rebalancing vans have to travel each day. Since the trip data does not include information about the historical rebalancing tours, calculating this benchmark is beyond the scope of this project.

Evaluation Metrics

The trips dataset encompasses the time from 2015/01/01 to 2016/11/30. For the training and evaluation of the predictive model, the dataset will be split at a certain date into two sets. The set with the older data will be used for training the model, the younger set for testing. This ensures that the model can predict future developments from the past.

The target value of the model will be the demand for bikes at a station for given day. The root of the mean squared error over all stations and dates of the test set will be used to quantify the success of the model.

Furthermore, the historical and predicted flow data will be visualised in an interactive web application. The application shall give the user an overview of bike movements for a selected date, and the option to drill down into a single station to reach a detailed view. The overall design of the application shall give insight into historical data and also show real-time predictions of future events.



Web application mockup

Project Design

Programming language and libraries

- Python 3.
- scikit-learn. Machine Learning library for Python.
- XGBoost. Open Source Gradient Boosting implementation for Python.
- Seaborn. A data visualisation library for Python.
- HTML, CSS, JavaScript. The basic web programming languages.
- flask. A web server framework written in Python.
- Leaflet. A JavaScript framework for geographical maps.

Project steps

1. Requisition of all necessary data.
2. Data cleansing and transformation.
 - Remove all trips with extreme trip durations, e.g. less than a minute or more than six hours.
 - Remove all user related features which can not be aggregated by hour.
 - Calculation of hourly flow values by aggregation. This should yield number of arrivals, departures and total flow per hour.
 - Join data sets: aggregated flow, weather and public holidays.
3. Data exploration with matplotlib and seaborn. Visualize correlations of trip numbers with weather features, using e.g. boxplots and jointplots.
4. Training of a benchmark model. The benchmark regression model will be trained and tested with the same datasets used for the sophisticated model, using default parameters without further refinement.
5. Model training and optimisation of model parameters. Different regression models will be evaluated, with a focus on ensemble models. XGBoost was chosen because of two reasons: 1) It gives a wide variety of different target objectives (Poisson count, Gamma regression, Tweedie Regression) and 2) it performs well on large datasets. Optimisation of parameters will be done with a grid search.
6. Visualize predicted flow and historical flow side by side in a web application. The user shall be able to select a specific date and get an overview which highlights potential rebalancing events for each hour of the day. An interactive map shall enable drill down into single stations.