

Introduction to Learning and Intelligent Systems – Spring 2015

taubnert@student.ethz.ch
junkerp@student.ethz.ch
kellersu@student.ethz.ch

March 15, 2015

1 Project Regression – Team “awesome”

1.1 logscore

Since we can not use logscore directly as a distance function during regression, we need to transform our data x, y . Suppose we want to search a function $f(x)$. Then to minimize $\text{logscore}(f(x), y)$ we minimize the two-norm $\|f'(x) - y'\|_2$ instead. Looking at the definition of *logscore*, we see that this can be accomplished by choosing $f'(x) = \log(1 + f(x))$ and $y' = \log(1 + y)$. The function f can then be reconstructed by $f(x) = \exp(f'(x)) - 1$.

1.2 Regressors

We used a number of different regressors. Most of them we understand how they work but we also used the *RandomForestRegressor* which we don't understand at all.

In the end, we compared a simple linear regression, a ridge regression, a k-nearest-neighbours regression, a lasso regression with the random forest regression. We concluded that we can do almost as good as the random forest regression.

1.3 Features

Different heuristics lead us to use different basis functions for our features. Because the data is about train usage and we have a timestamp provided, we assumed that there will be some periodicity observable. This assumption let us add the fourier and also the discrete cosine transformation as base-functions.

We deliberately didn't use all of the available data to fit. We concluded that time in the minute frequency and the weather data D corresponds mostly to noise.

1.4 Scores

Linear

Ridge

Lasso

K-NN

Random Forest