

Creating and Evaluating Stochastic Regression Models on the Basis of Heterogeneous Sensor Networks

Bachelors Thesis, Supervisors - Dr. Johannes Riesterer, Dr. Sebastian Lerch

Stanislav Arnaudov | 7. November 2018

TECO - DAS TELECOOPERATION OFFICE



Given is a heterogeneous network of sensors – the network contains “good” **and** “bad” sensors.

- good – calibrated and high precision of the measurements
- bad – uncalibrated and low precision of the measurements

Interesting Questions:

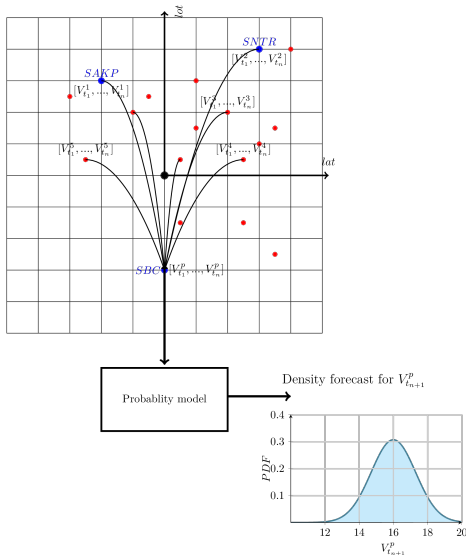
- Can we use the bad sensors in order to predict the values of the good ones?
- Can we identify the weak parts of the network?

Our approach to the problems:

- Use of stochastic regression models.
- Formal evaluation of the models through the use of proper scoring rules, verification rank histograms and predictive performance checks.
- Analyzing the relevance of each predictor for the final prediction through the use of a feature importance technique.

- Heterogeneous Network of air pollution sensors in Stuttgart.
 - LU-BW (*Landesanstalt für Umwelt, Messungen und Naturschutz Baden-Württemberg*) - **3 Sensors** of high quality.
 - Station names: SBC, SNTR, **SAKP**
 - *luftdaten.info* – Public data from cheap DIY sensors.
- Considered period: the year of 2018.
- Challenges:
 - Noisy data.
 - Holes in the data.
 - Station SAKP provides air pollution values only for the last four months of the considered year.
 - The DIY sensors provide values for each minute of the day, the LUBW sensors - for thirty-minute intervals

General View



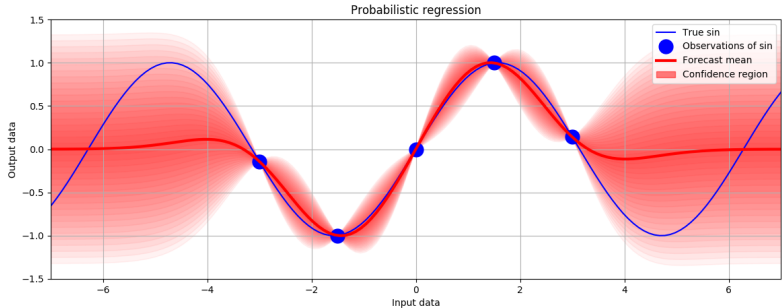
Goals:

- Investigate how the considered models perform on the data.
- Show that the use of stochastic regression models is a feasible approach to predicting air pollution values in heterogeneous networks of sensors.
- Show that the unreliable parts of the network can be identified through the use of feature importance technique*.

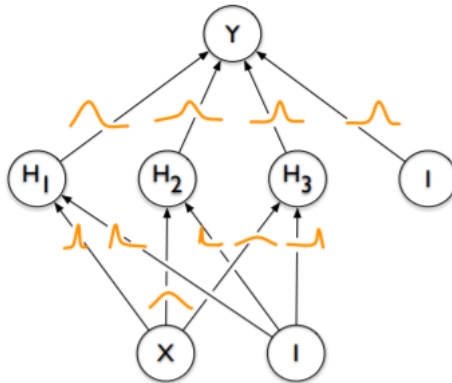
Not goals:

- Build the best possible model for predicting air pollution values.

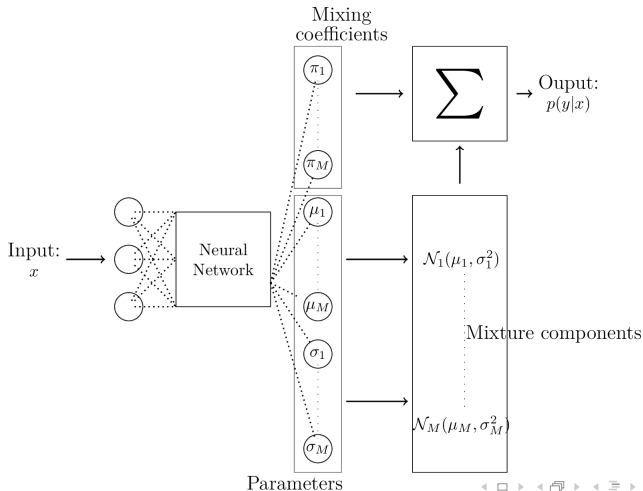
- Regression: feature vector \mapsto a real value
- Stochastic Regression: feature vector \mapsto a probability distribution



■ Bayesian Neural Networks



■ Mixture Density Networks



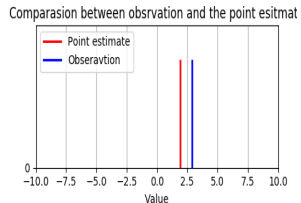
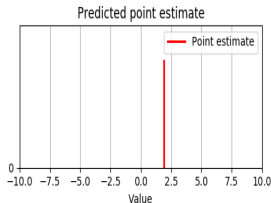
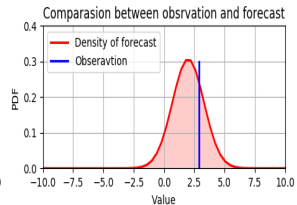
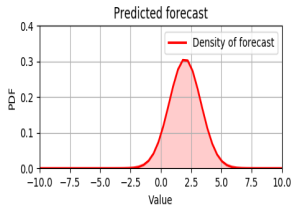
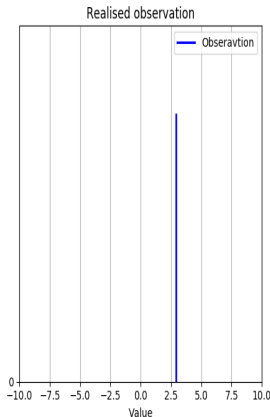
Empirical model

- Past observations are treated as values of a random variable. From that a distribution for a predicted future value is implied.
- We use this model as a baseline.

We hope to achieve better results than the empirical model with the other two models.

Evaluation of probability distributions

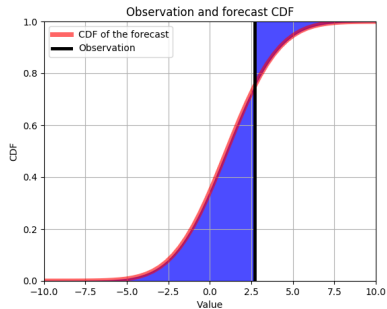
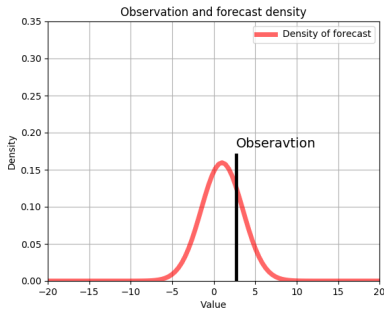
- We do not compare point estimate with a realized real value but rather a probability distribution with a real value.



Continuous Rank Probability Score

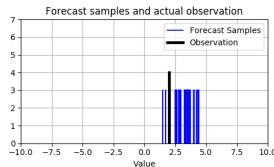
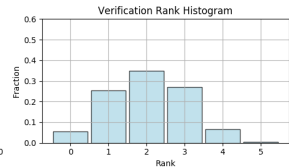
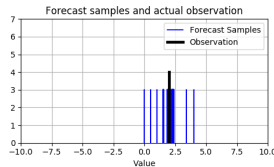
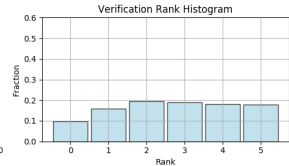
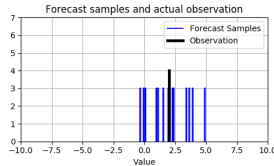
- CRPS compares a distribution with an observation, where both are represented as cumulative distribution functions.

$$CRPS(F, y) = \int_{-\infty}^{\infty} (F(x) - \mathbb{1}\{y \leq x\})^2 dx$$



Verification Rank Histograms

- Does the observation behave like a random sample from the forecast distribution?



- Assess the **relative** importance of all features used by a model.
 - **Important:** Importance is relative only to the model itself.
- Random shuffling of each predictor/feature in the test set one at a time and observing the increase/decrease in mean the CRPS value compared to the unpermuted features.

■ Approach

Based on the data and the examined types of models there are several orthogonal considerations to make when it comes to training:

- Target LUBW-Station – SAKP, SBC, SNTR
- Integration period – one day, twelve hours, one hour
- Use of the other two LUBW station – yes or no
- Predicted air pollution value – PM10 or PM2.5
- Used model - BNN or MDN (two types of MDNs)

We wanted to investigate every possible combination.

■ Challenges:

- BNNs are very hard to train.
- A lot of possible ways to train several different models.

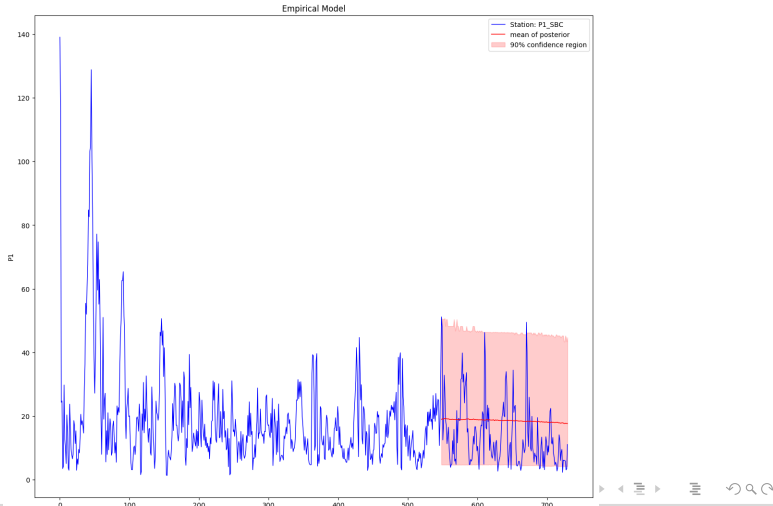
Results

Results presented here: representative sample of all results.

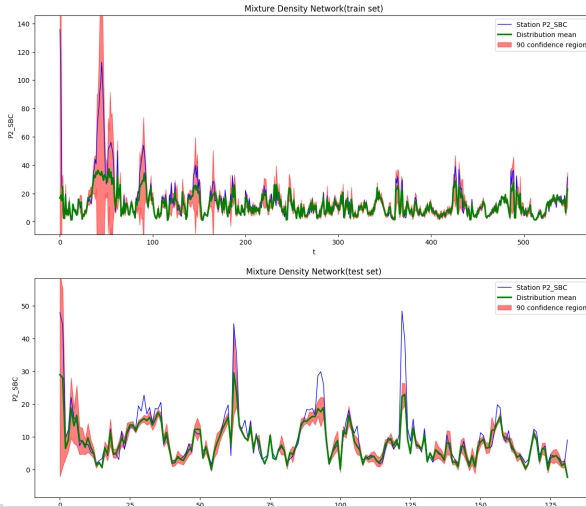
Results presented here: representative sample of all results.

- Plots showing how exactly the built models predict the values of their train and test sets.

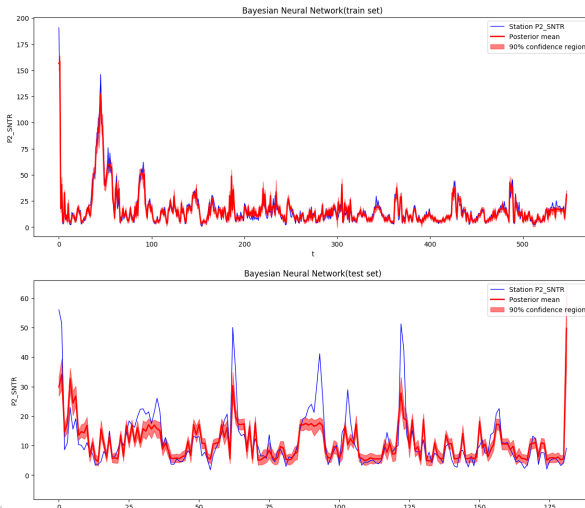
Empirical Model, twelve hour averaged data



MDN, twelve hour averaged data



BNN, twelve hour averaged data



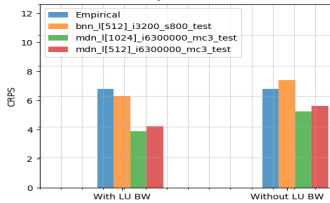
Results presented here: representative sample of all results.

- Plots of the models modeling their respective test and train sets.

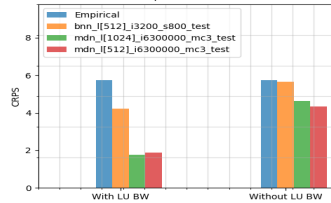
Results presented here: representative sample of all results.

- Plots of the models modeling their respective test and train sets.
- Plots comparing the built models with respect their CRPS score across LUBW Station, use of LUBW data and predicted air pollution value.

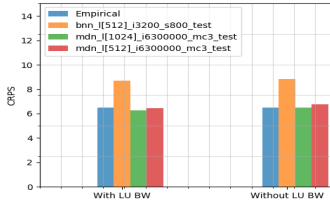
Station: SBC, Predicted value:PM10



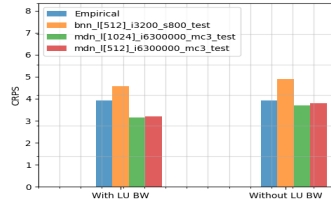
Station: SBC, Predicted value:PM2.5



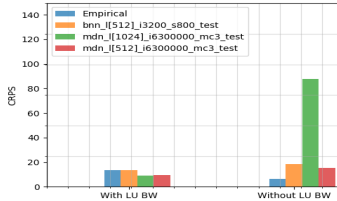
Station: SAKP, Predicted value:PM10



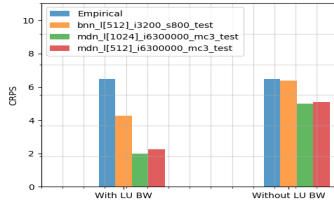
Station: SAKP, Predicted value:PM2.5



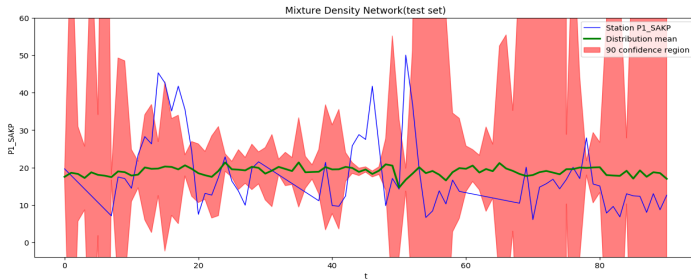
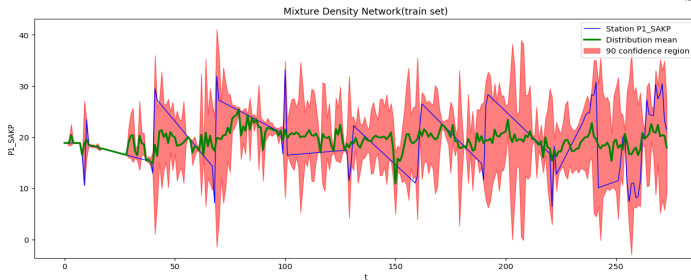
Station: SNTR, Predicted value:PM10



Station: SNTR, Predicted value:PM2.5



Results



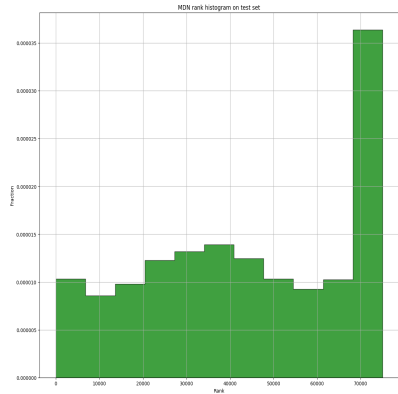
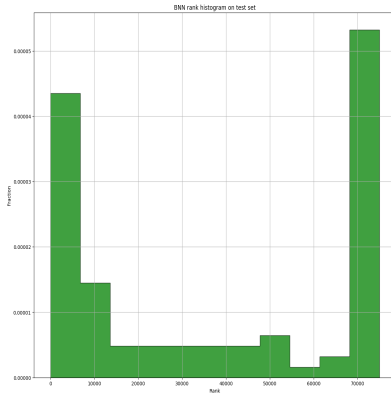
Results presented here: representative sample of all results.

- Plots of the models modeling their respective test and train sets.
- Plots comparing the built models with respect their CRPS score across LUBW Station, use of LUBW data and predicted air pollution value.

Results presented here: representative sample of all results.

- Plots of the models modeling their respective test and train sets.
- Plots comparing the built models with respect their CRPS score across LUBW Station, use of LUBW data and predicted air pollution value.
- Two rank histograms reflecting the results of all built models.

BNN (one day average) and MDN (one hour average), Rank Histograms

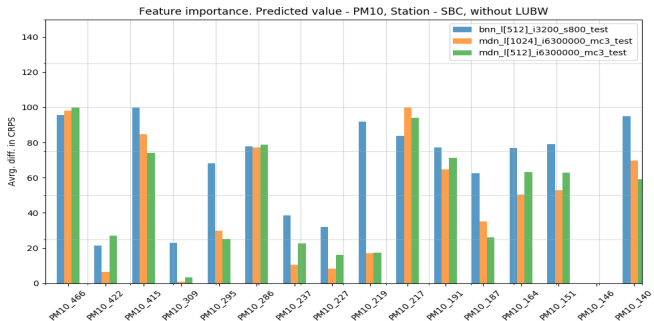
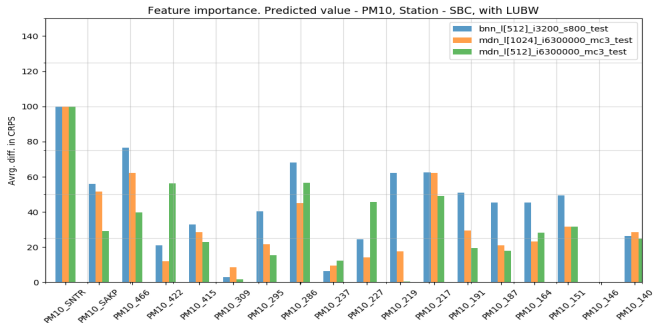


Results presented here: representative sample of all results.

- Plots of the models modeling their respective test and train sets.
- Plots comparing the built models with respect their CRPS score across LUBW Station, use of LUBW data and predicted air pollution value.
- Two rank histograms reflecting the results of all built models.

Results presented here: representative sample of all results.

- Plots of the models modeling their respective test and train sets.
- Plots comparing the built models with respect their CRPS score across LUBW Station, use of LUBW data and predicted air pollution value.
- Two rank histograms reflecting the results of all built models.
- Feature importance data.



The results show:

- Stochastic regression models are a viable approach for predicting values in the investigated sensor network.
- Few sensors have consistently low feature importance for the built models.
- There is clearly room for improvement of the models.

Thank you for your attention.

Questions?