

A comparison of probabilistic forecasting methods for extreme NO₂ pollution episodes

Sebastián Pérez Vasseur

*Artificial Intelligence Department
Universidad Nacional de Educación a Distancia — UNED
c/ Juan del Rosal, 16, Madrid, Spain*

José L. Aznarte¹

*Artificial Intelligence Department
Universidad Nacional de Educación a Distancia — UNED
c/ Juan del Rosal, 16, Madrid, Spain*

Abstract

Keywords: probabilistic forecasting, air quality, quantile regression, nitrogen dioxide, Madrid

1. Introduction

2. Probabilistic forecasting with quantile regression

As mentioned above, the prediction from most regression models is a point estimate of the conditional mean of a dependent variable, or response, given a set of independent variables or predictors. However, the conditional mean measures only the center of the conditional distribution of the response, and if we need a more complete summary of this distribution, for example in order to estimate the associated uncertainty, quantiles are in order. The 0.5 quantile (i.e., the median) can serve as a measure of the center, and the 0.9 quantile marks the value of the response below which reside the 90% of the predicted points. Recent advances in computing have inducted the development of regression models for predicting given quantiles of the conditional distribution. The technique is called quantile regression (QR) and was first proposed by Koenker in 1978 [?] based on the intuitions of the astronomer and polymath Rudjer Boscovich in the 18th century. Elaborating from the same concept of estimating conditional quantiles from different perspectives, several statistical and CI models that implement this technique have been developed: from the original linear proposal to multiple or additive regression, neural networks, support vector machines, random forests etc.

Quantile regression has gained an increasing attention from very different scientific disciplines [?], including

financial and economic applications [?], medical applications [?], wind power forecasting [?], electric load forecasting [?], environmental modelling [?] and meteorological modelling [?] (these references are just examples and the list is not exhaustive). To our knowledge, despite its success in other areas, quantile regression has not been applied in the framework of air quality.

Thus, as we can estimate an arbitrary quantile and forecast its values, we can also estimate the full conditional distribution, which will entail us to the results presented in Section 4.

Also, probabilistic forecasting is an advantage as we need to predict when the target will be above a certain threshold (180). So instead of having a Yes/No Answer, we are calculating the probability of the target being the above the threshold.

Among the array of methods that allow to estimate and forecast data-driven conditional quantiles, in this study we have chosen k-neighbors quantile regression, quantile regression forests and quantile XGBoost. We will compare the different algorithms through the CRPS metric for the distribution and the RMSE, MAE, Correlation and Bias for the quantile 50.

3. Data description and experimental design

3.1. Training Data

There are 20 pollution stations around the city and they are organized in different zones as you can see in the picture below:

All of the 24 stations of Madrid's monitoring system capture hourly data for NO₂. They are spatially distributed according to European regulations and are classified into

Email address: jlaznarte@dia.uned.es (José L. Aznarte)

¹This work has been partially funded by Ministerio de Economía y Competitividad, Gobierno de España, through a *Ramón y Cajal* grant (reference: RYC-2012-11984).

- NO2 past levels
- O3 past levels
- calendar variables that represent the status of the day:
bank holiday, laboral, type of holiday ...
- weather data: Precipitation, Temperature, Humidity,
Wind
- ECMWF numerical pollution prediction

3.2.1. Calendar Fields Summary

Horizontal bar chart showing the correlation of various variables with the variable 'Festivo'. The x-axis represents the correlation coefficient, ranging from -600 to 900. The y-axis lists the variables. The bars are color-coded: red for positive correlation, green for negative correlation, and blue for zero correlation.

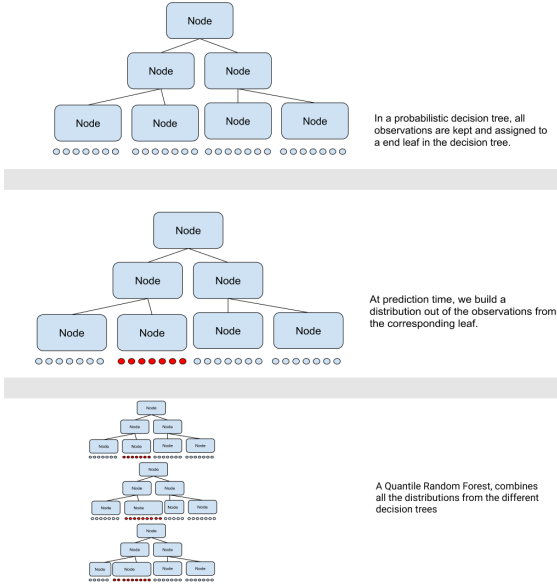
Variable	Correlation Coefficient (approx.)
r Festivo EqNoc	50
r Festivo FesNov	-20
r Festivo InmCns	10
r Festivo NavAnu	80
ar Festivo PriVer	100
Festivo SemSan	100
r Festivo VigAgo	150
lav NocheBuena	20
tlav NocheVieja	-50
ndar OprRetorno	70
rSalida PriNoLab	-300
prSalida Vispera	900
ndar PuenteLab	-300
ensiva fin_curso	30
siva inicio_curso	-100
ensiva navidad	-500
no_lectivo otros	-400
aciones navidad	-400
s.semana_santa	10
aciones verano	400

3.2.2. Seasonality Extraction

- Every 12 hour Seasonality
- Yearly Seasonality
- Daily Seasonality
- Every 4 year seasonality
- Weekly Seasonality

2

Figure 2: Random Forest Steps



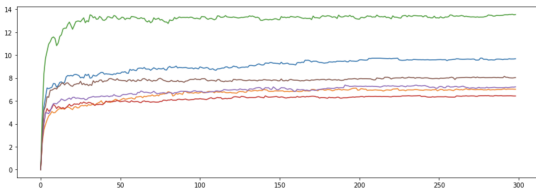
3.2.3. Previous Values

As with any forecast technique based on machine learning, we add previous values to improve the accuracy of the analysis. Based on the seasonal analysis, we see it's interesting to add a week of values. We will not add more to keep a reasonable number of features as input.

3.3. k-Neighbors

The probabilistic k-Neighbors is based on the competition entry from [K-nearest neighbors for GEFCom2014 probabilistic wind power forecasting]. This algorithm is based on the standard k-Neighbor, where instead of aggregating the targets of the k nearest points to the input (by taking the mean or median for example), it builds a distribution from those neighbors.

We need to find an optimal number k for our model. We will use the CRPS of the predicted distribution to get the best k. As you can see in the chart, 50 seems to be the optimal number of neighbor.



3.4. Quantile Random Forest

We will use the quantile random forest as described in [1]. Quantiles are built from the observations in the training set. We take the observation points that belong to the same leaf than the input and we build a distribution from those points. See figure for explanation:

For a detailed discussion on quantile regression forests, see [?].

3.5. Gradient Boosted Tree

Gradient Boosted Trees is a technique that consist on growing trees based on the compromise of a cost function and a regularization functions. This cost function is usually used to forecast the mean of the signal. But we are modifying the cost function in order to forecast the quantile of the function.

as an illustration of the concept (for a detailed discussion of quantile regression, refer to [?]), given a set of vectors (x_i, y_i) , in point forecasting we are usually interested in what prediction $\hat{y}(x) = \alpha_0 + \alpha_1 x$ minimizes the mean squared error,

$$E = \frac{1}{n} \sum_i \epsilon_i = \frac{1}{n} \sum_i [y_i - (\alpha_0 + \alpha_1 x)]^2. \quad (1)$$

This prediction is the conditional sample mean of y given x , or the location of the conditional distribution. But we could be interested in estimating the conditional median (i.e., the 0.5 quantile) instead of the mean, in which case we should find the prediction $\hat{y}(x)$ which minimizes the mean absolute error,

$$E = \frac{1}{n} \sum_i \epsilon_i = \frac{1}{n} \sum_i |y_i - (\alpha_0 + \alpha_1 x)|. \quad (2)$$

The fact is that, apart from the 0.5 quantile, it is possible to estimate any other given quantile τ . In that case, instead of (2), we could minimize

$$E = \frac{1}{n} \sum_i f(y_i - (\alpha_0 + \alpha_1 x)) \quad (3)$$

where

$$f(y - q) = \begin{cases} \tau(y - q) & \text{if } y \geq q \\ (1 - \tau)(q - y) & \text{if } y < q \end{cases}, \quad (4)$$

with $\tau \in (0, 1)$. Equation (3) represents the median when $\tau = 0.5$ and the τ -th quantile in any other case.

By forecasting different quantiles, we can forecast the CDF of the time series. The main drawback is that we need to build a model per quantile. And since quantiles are calculated separately, we can have quantile crossing, i.e. the non monotonicity of the predicted CDF.

In order to solve that, we will apply the technique from [?]. This chart shows the predicted 10 and 90 percentiles of the first 100 values.

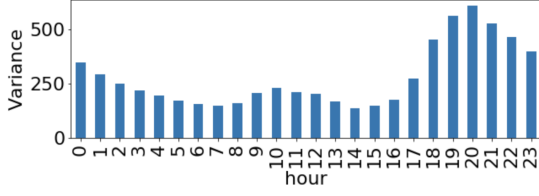
3.6. Nitrogen dioxide data

3.7. Weather data

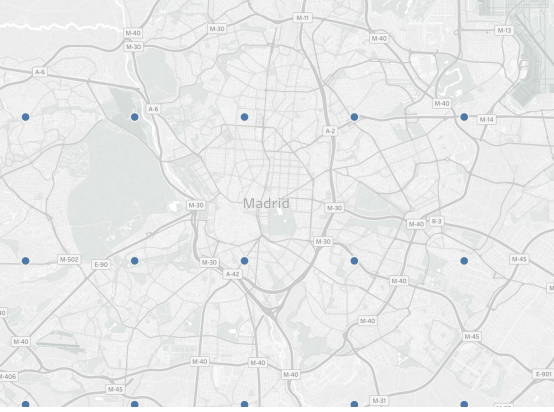
3.8. ECMWF numerical pollution prediction

The European Centre for Medium-Range Weather Forecasts implements the Copernicus Atmosphere Monitoring Service. This service provides CAMS delivers a daily production of near-real-time European air quality analyses and forecasts with a multi-model ensemble system. As you can see in the picture, the scope of the forecast is european

Figure 3: Variance of NO2 per hour



and does not have the needed granularity to forecast the levels of NO2 in a station.



3.9. Experimental design

We will train the models with data prior to 2017 and we will test our models with 2017 data. We will always test with predictions done at 10:00, as this is the time the forecast will be done and the alert will be decided or not.

Also, we will create different models for each horizon (how many hours in advance we do the forecast)

3.10. Evaluation of probabilistic forecasts

3.10.1. Variance of the Time series per hour

The time series have very different variance for each hour of the day. This is important when we evaluate our model, since we are always forecasting at 10:00.

3.10.2. Metrics

We will evaluate the 50 percentile through the RMSE, MAE, Bias and Corr and the whole forecasted CDF through the CRPS.

3.10.3. CRPS

The CRPS measures the accuracy of a probabilistic forecast.

4. Results and discussion

4.1. Reference models

In the first experiment, we used quantile regression to compute point-forecasts of the expected value (median) for one-day ahead predictions of NO₂ concentrations.

Table 1 shows Arima outperforms the other Machine Learning algorithms

Table 1: CRPS Error for the different methods at different horizons

horizon	Xarima	KNN	RF	XGB
1	5.88	9.30	7.01	7.92
12	14.92	20.22	18.02	16.04
13	14.55	16.91	15.41	13.96
14	12.95	14.62	13.86	12.24
20	6.82	8.79	8.57	10.04
37	15.29	23.20	20.91	18.42
45	8.50	9.69	8.78	9.14
55	7.51	12.11	11.58	16.59

Figure 4: RMSE Error per horizon and Method

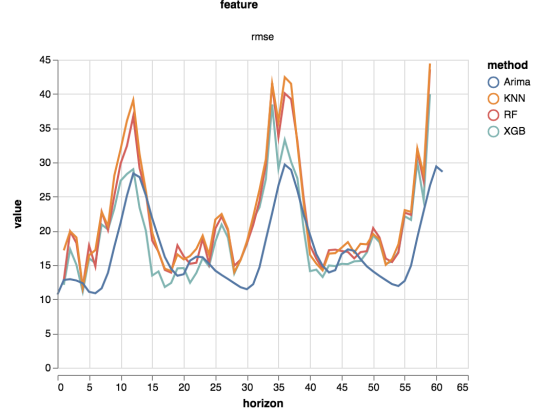


Figure 5: MAE Error per horizon and Method

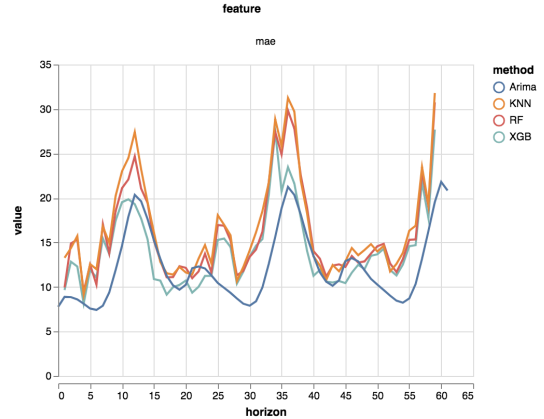
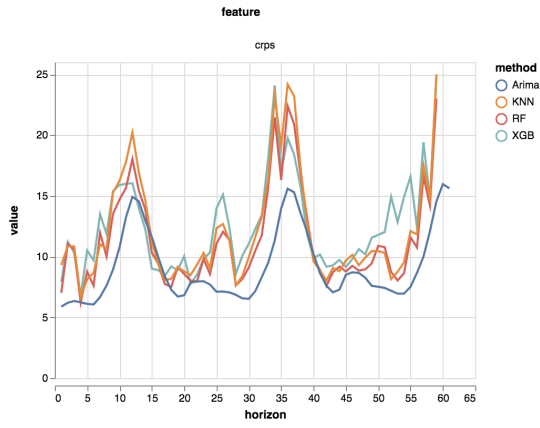


Figure 6: CRPS Error per horizon and Method



4.2. Probabilistic forecasting of extreme values

4.3. Forecasting the probability of alerts

5. Conclusions

References

- [1] D. of Random Forest, quantregforest.
URL <https://stat.ethz.ch/~nicolai/quantregforests.pdf>