

~~Title of This Paper~~

## AIR POLLUTION MEASUREMENTS PREDICTION

~~AUTHOR~~ YU LI

ABSTRACT. In this competition you are predicting the values of air pollution measurements over time, based on basic weather information (temperature and humidity) and the input values of 5 sensors. The three target values to you to predict are: target-carbon-monoxide, target-benzene, target-nitrogen-oxides.

### CONTENTS

|                               |   |
|-------------------------------|---|
| 1. <del>Introduction</del>    | 2 |
| 2. <del>Preliminaries</del>   | 2 |
| 2. <u>Data Description</u>    | 2 |
| 3. <del>Method</del>          | 2 |
| 3. <del>Conclusions</del>     | 4 |
| <del>Acknowledgement</del>    | 5 |
| 3. <u>Feature Engineering</u> | 5 |
| 4. <u>Model Training</u>      | 6 |
| 5. <u>Result</u>              | 6 |

---

*Date:* (None).

*2020 Mathematics Subject Classification.* Artificial Intelligence.

*Key words and phrases.* Machine Learning, Pollution Prediction.

## 1. INTRODUCTION

At a high level, what is the problem area you are working in and why is it important? It is important to set the larger context here. Why is the problem of interest and importance to the larger community?

This paragraph narrows down the topic area of the paper. In the first paragraph you have established general context and importance. Here you establish specific context and background.

"In this paper, we show that ...". This is the key paragraph in the intro-- you summarize, in one paragraph, what are the main contributions of your paper given the context you have established in paragraphs 1 and 2. What is the general approach taken? Why are the specific results significant? This paragraph must be really good.

You should think about how to structure these one or two paragraph summaries of what your paper is all about. If there are two or three main results, then you might consider itemizing them with bullets or in test. In this competition you are predicting the values of air pollution measurements over time, based on basic weather information (temperature and humidity) and the input values of 5 sensors. The three target values to you to predict are:

- e.g., First ... target-carbon-monoxide
- e.g., Second ... target-benzene
- e.g., Third ... target-nitrogen-oxides

If the results fall broadly into two categories, you can bring out that distinction here. For example, "Our results are both theoretical and applied in nature. (two sentences follow, one each on theory and application)"

Keep this at a high level, you can refer to a future section where specific details and differences will be given. But it is important for the reader to know at a high level, what is new about this work compared to other work in the area.

"The remainder of this paper is structured as follows..." Give the reader a roadmap for the rest of the paper. Avoid redundant phrasing, "In Section 2, In section 3, ... In Section 4, ... " etc.

Test citation [?], and [?] or ?].

This is for, and this is for.

Number: , , , and

We have, the range:  $\cdot 1/2$ .

For, as shown below:

$$a = b \times \sqrt{ab}$$

## 2. PRELIMINARIES

## 2. DATA DESCRIPTION

Before model training, data needs to be analyzed to determine the required features. Here is the statistics of training data and test data:

## 3. METHOD

Experiment and Analysis

 (None)-(None) ((None))

## Precision Comparison on Event Detection Methods

TABLE 1. Train Data Description

| <u>Elements</u>                   | <u>OR-Event-Detection-Number</u>   |
|-----------------------------------|------------------------------------|
| <u>datetime</u>                   | <del>AC-Event-Detection-7111</del> |
| <u>degC</u>                       | <del>TC-Event-Detection-408</del>  |
| <u>relative – humidity</u>        | 762                                |
| <u>absolute – humidity</u>        | 5451                               |
| <u>sensor1</u>                    | 3882                               |
| <u>sensor2</u>                    | 3882                               |
| <u>sensor3</u>                    | 3882                               |
| <u>sensor4</u>                    | 3882                               |
| <u>sensor5</u>                    | 3882                               |
| <u>target – carbon – monoxide</u> | 95                                 |
| <u>target – benzene</u>           | 405                                |
| <u>target – nitrogen – oxides</u> | 3268                               |

TABLE 2. Test Data Description

| <u>Elements</u>                       | <u>Number</u>         |
|---------------------------------------|-----------------------|
| <del>precision-datetime</del>         | <del>0.83-2247</del>  |
| <u>degC</u>                           | <del>0.69-280</del>   |
| <u>relative – humidity</u>            | <del>0.46-653</del>   |
| <del>recall-absolute – humidity</del> | <del>0.68-1915</del>  |
| <u>sensor1</u>                        | <del>0.48-1758</del>  |
| <u>sensor2</u>                        | <del>0.36-1816</del>  |
| <del>F-score-sensor3</del>            | <del>0.747-1833</del> |
| <u>sensor4</u>                        | <del>0.57-1877</del>  |
| <u>sensor5</u>                        | <del>0.4-2017</del>   |

In order to understand the change trend of data, the data is visualized and analyzed based on the visualization results.

3. CONCLUSIONS



FIGURE 1. Target Overall Situation

It can be seen from the figure1 that the values of the three target pollutants in August each year will be lower, gradually rising from September, and significantly higher than the level before August, so it is necessary to take the month as a feature of the model.

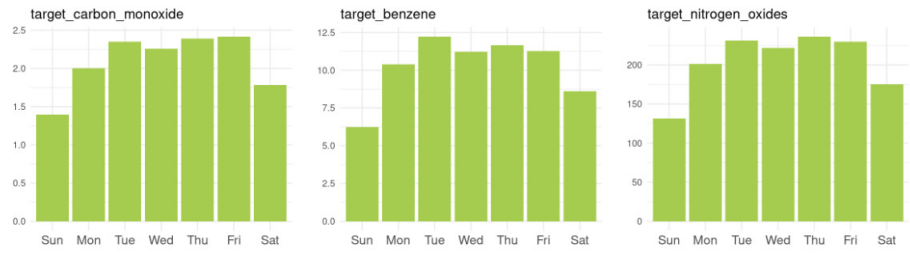
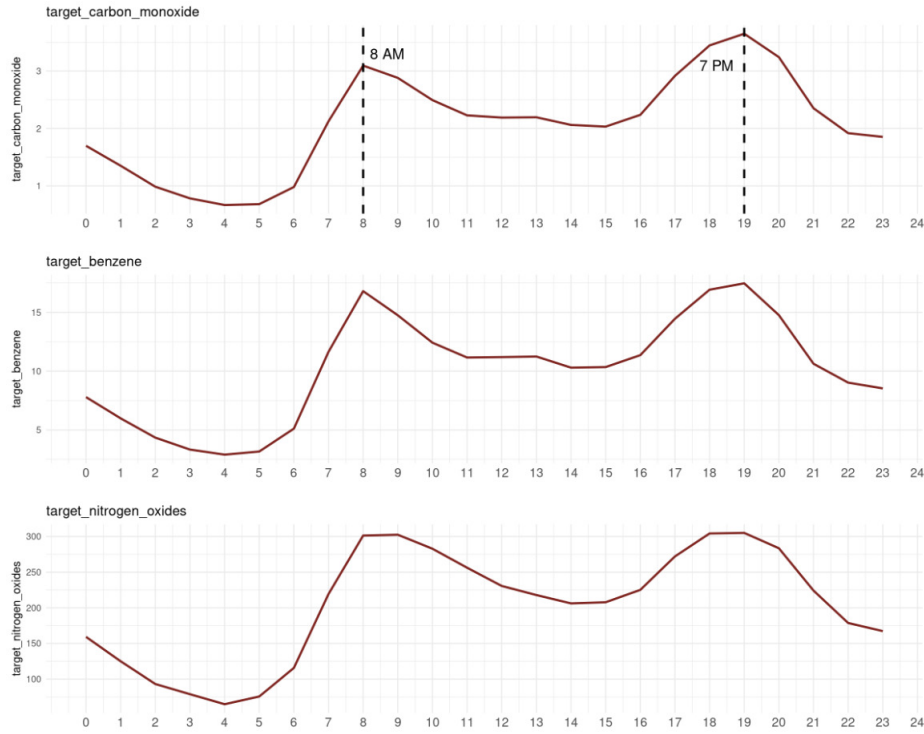


FIGURE 2. Target Weekly Situation

It can be seen from the figure2 that the content level of each pollutant at the weekend of each week will decrease significantly, so it is necessary to take whether this day is a weekend as a feature of the model.

## ACKNOWLEDGEMENT

FIGURE 3. Target Daily Hourly Change

It can be seen from the figure3 that the level of each pollutant is the lowest at 5:00 a.m. every day, and then gradually rises to 8:00 a.m. to reach the first peak, and then gradually falls to 4:00 p.m., and then rises to 7:00 p.m. to reach the second peak, and then continues to decline, so it is necessary to take time as a feature of the model.

3. FEATURE ENGINEERING

According to the analysis of training data, the following features are used for model training:

- absolute-humidity
- deg-C
- relative-humidity
- sensor1-5
- month
- week
- is-weekend
- hour

~~The authors would like to thank ...~~



(None)-(None) ((None))

4. MODEL TRAINING

Data fitting using LGBMRegressor, the algorithm is easy to use. It only needs to put the set features and three prediction targets into the model for training, but there is no parameter optimization, which has a certain impact on the training results.

5. RESULT

- Use RMSLE(Root Mean Squared Logarithmic Error) to evaluate the results.

$$RMSLE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\log(\hat{y}_i + 1) - \log(y_i + 1))^2}$$

- Private Score:0.33979
- Public Score:0.387

(A. 1) SCHOOL OF COMPUTER SCIENCE,, NANJING UNIVERSITY OF SCIENCE AND TECHNOLOGY,  
JIANGSU 246000, CHINA

Email address, A. 1: yli@tulip.academy