

# IML Term project paper

Julia Palorinne, Pyry Silomaa, Sara Sippola

05.12.2022

## Introduction

In this project we trained a classifier on a data set of atmospheric measurements. The task is to predict whether new particle formation (NPF) happens or not on a given day based on the atmospheric data.

## Preprocessing of the data

### Initial data analysis

The training data consists of several variables measured on 464 non-consecutive days. The variables are daily means and standard deviations of measurements such as carbon dioxide concentration, solar radiation and air temperature. Some of the variables are of the same phenomenon measured at different heights.

Many classifiers are affected by correlation and colinearity between variables. As expected, we found that many of the variables describing the same phenomenon are correlated, as are variables related to radiation (see figure below). We take a more detailed look into this in the section *Correlation between parameters*.

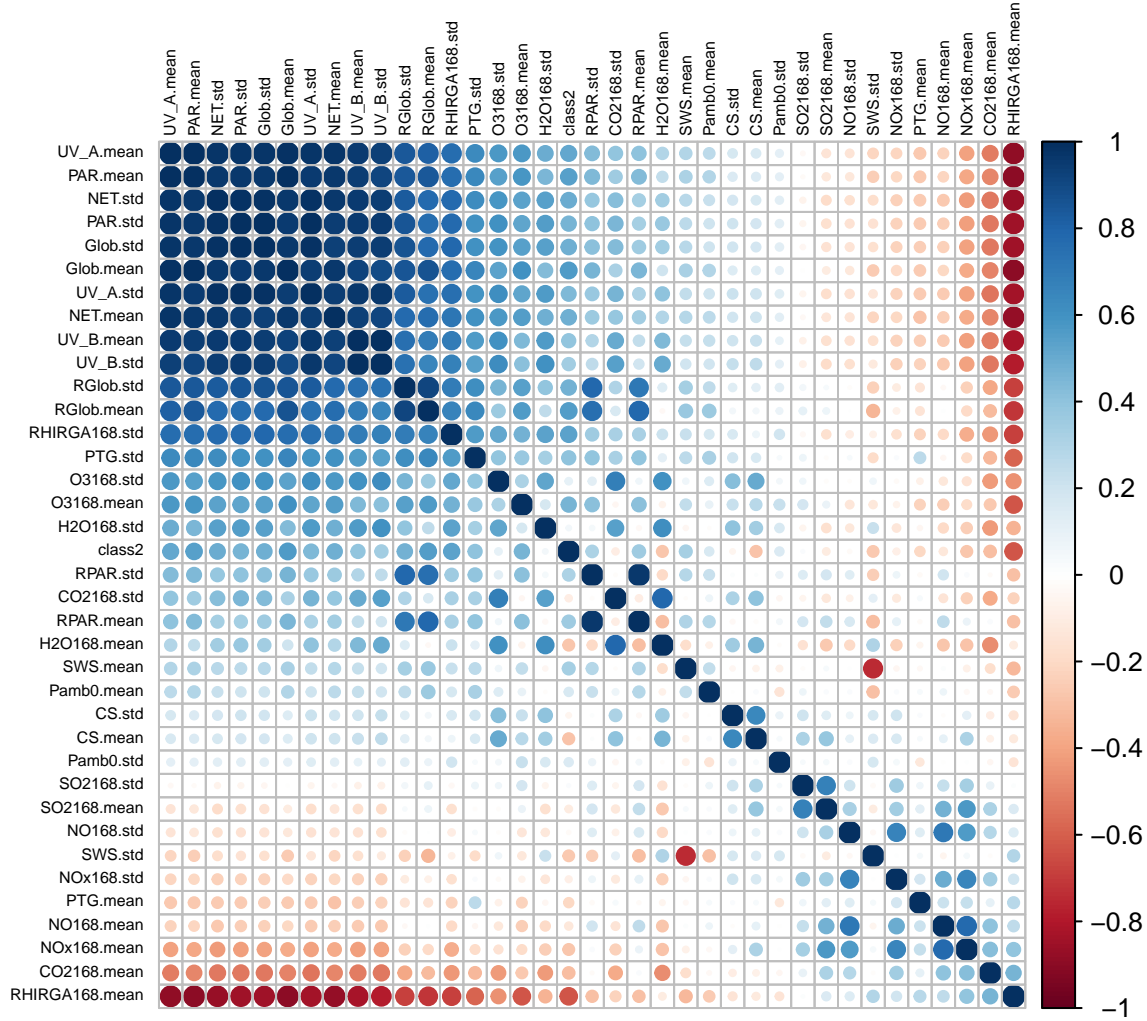
We also familiarized ourselves with the data through the [smear.avaa.scs.fi](https://smear.avaa.scs.fi) webpage that offers further visualizations details about the data and the measurement site, and with the variable names and details available at [wiki.helsinki.fi](https://wiki.helsinki.fi).

(The data is over non-consecutive days, ok, but are there any significant gaps e.g. a month?)



Table 1: Correlation (H2O)

	H2O168.mean	H2O336.mean	H2O42.mean	H2O504.mean	H2O672.mean	H2O84.mean
H2O168.mean	1.0000000	0.9998966	0.9997062	0.9997158	0.9994631	0.9998894
H2O336.mean	0.9998966	1.0000000	0.9993506	0.9999302	0.9997498	0.9996330
H2O42.mean	0.9997062	0.9993506	1.0000000	0.9990202	0.9986416	0.9999330
H2O504.mean	0.9997158	0.9999302	0.9990202	1.0000000	0.9998589	0.9993631
H2O672.mean	0.9994631	0.9997498	0.9986416	0.9998589	1.0000000	0.9990316
H2O84.mean	0.9998894	0.9996330	0.9999330	0.9993631	0.9990316	1.0000000



The correlation plot for the remaining parameter shows us that the remaining highly correlated parameters are all radiation-related.

## **Classifier**

Description of considered machine learning approaches

Chosen classifier, pros and cons of this particular classifier for this application

## **Results**

Classifier performance (numerical)

Insights, conclusions, discussion etc.