

# Term Project - Group 80

Sari Ropponen, Outi Boman, Loic Dreano

2022-12-09

## Description of data

The training data npf\_train.csv and testing data npf\_test\_hidden.csv were downloaded. The data includes 104 features. The number of observations in training data is 464 and 965 in testing data.

Table 1: Summary of the dataset

	npf_train	npf_test
Measurements	464	965
Variables	104	104

The features include a lot of daily measurements taken in Hyytiälä forestry field station. Some of the features like temperature T and CO2 are measured in different heights. The height is indicated in the name of the feature, for example, T84.mean is the mean temperature at 8.4 meters above the mast base.

## Preprocessing data

The summary of the first four features in training data is: — I changed this to be the 4 first because only they are needed to do the next steps of data cleaning.

```
##          id          date          class4          partlybad
## Min.      : 1.0    Length:464    Length:464    Mode :logical
## 1st Qu.:116.8    Class :character  Class :character  FALSE:464
## Median :232.5    Mode  :character  Mode  :character
## Mean      :232.5
## 3rd Qu.:348.2
## Max.      :464.0
```

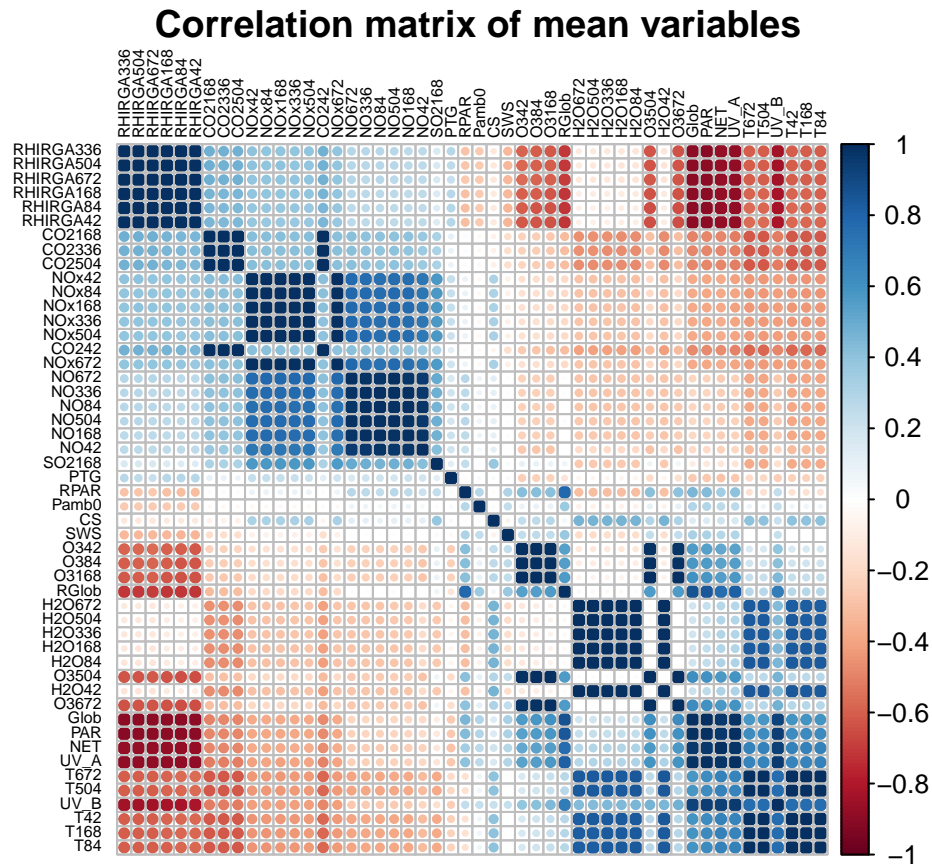
The column “date” was set to be the row names. Columns “id”, “date” and “partlybad” were removed from the data. Because the value of the logical variable “partlybad” is FALSE for all the observations, it doesn’t give any information.

A qualitative variable “class2” is added to the training data. It gets either value “event” or “nonevent” according to “class4”. Variable “class4” indicates the type of the event if it has happened, values “Ia”, “Ib” or “II”, or “nonevent” if no event has happened during the day.

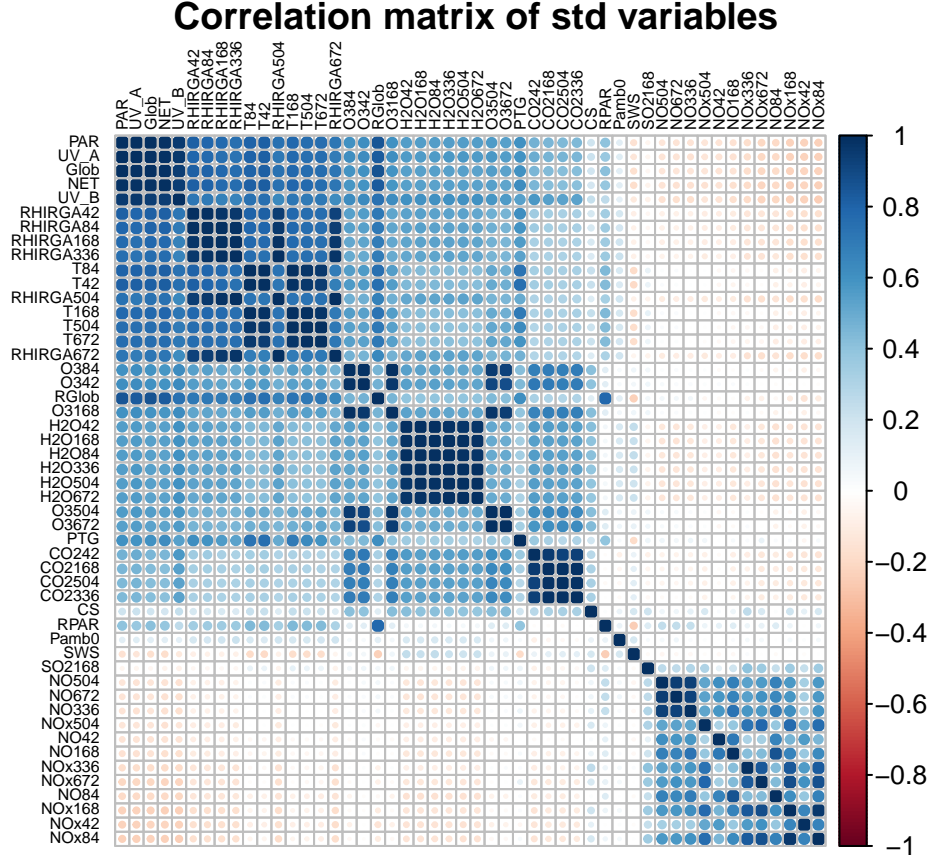
The task is to build a binary classifier which predicts if an NPF event will happen or not during the day according to the observed measurements.

## Data exploration

Because the data includes same measurements at different heights it is expected that there are correlation between variables. Correlations between different mean values are shown in the matrix below:



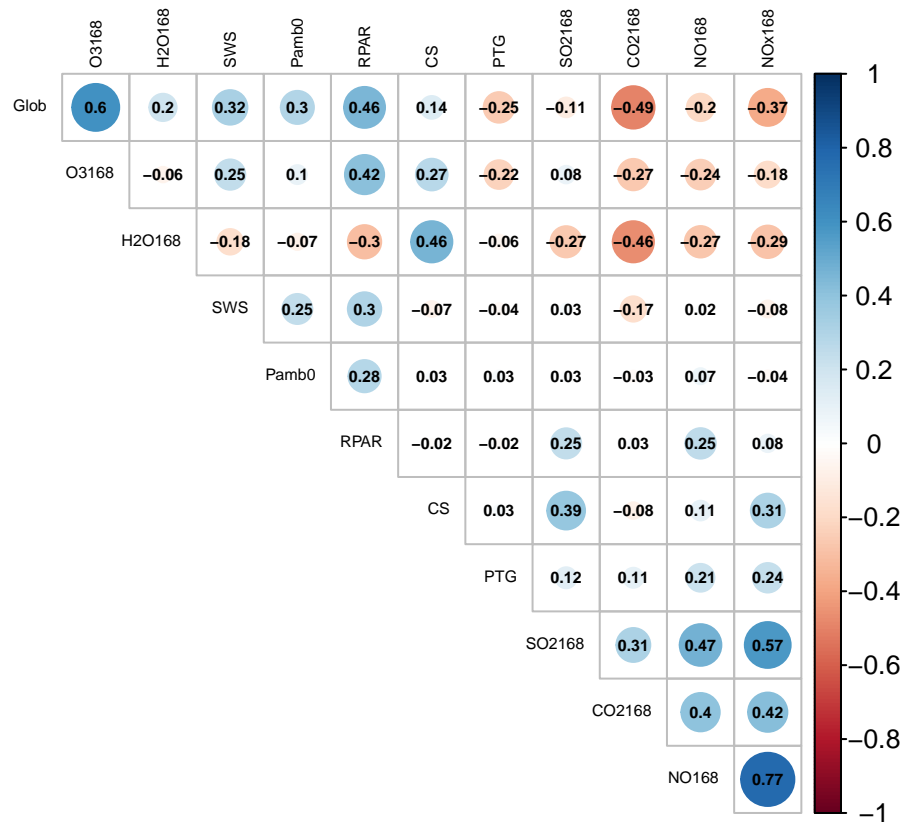
The same picture for different standard deviation values is



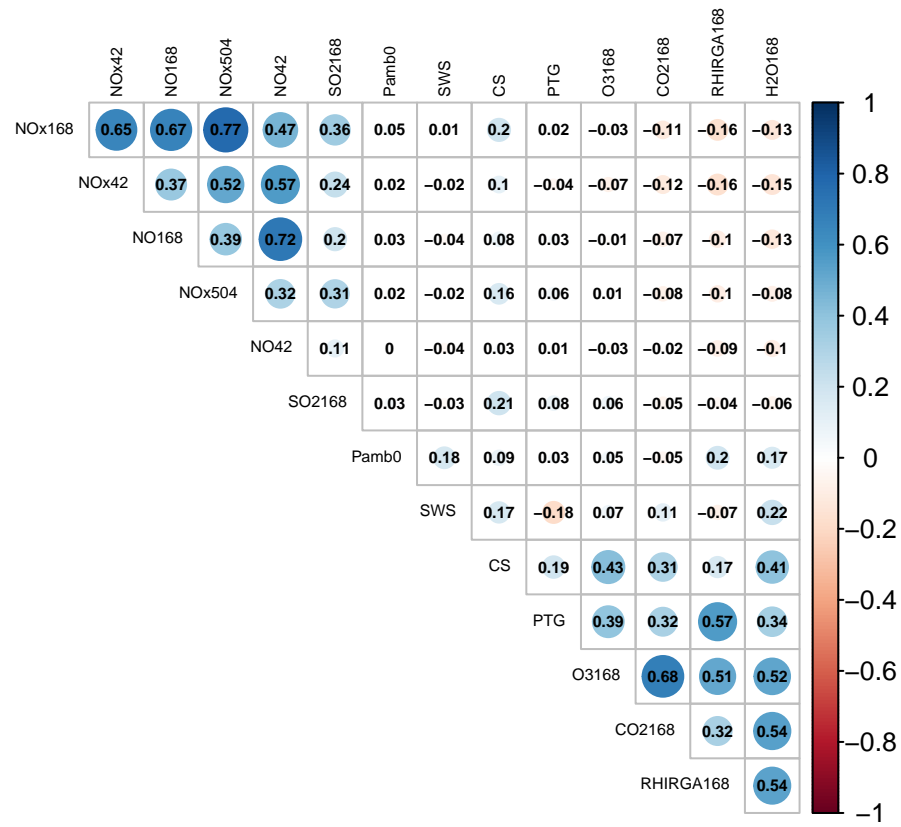
As we can see, there are a lot of highly correlated variables, keeping them all would not improve the predictivity of our models (since they carry the same information). In order to increase the power of our models to identify independent variables that are statistically significant and to make them simpler to interpret (simpler model in general) we will remove the highly correlated variables. To do so the variables are clustered together based on their correlation; every pair of variables which have an absolute correlation greater than 0.8 are clustered together. Then a random variable from each cluster is kept for further analysis.

The correlations for mean and standard deviations after cleaning the data are shown in the pictures below as well as a table of variables left in the data (table 2). Table 3 shows that the data includes only 26 variables. — should we put the means and std in the same correlation matrix to see if there are correlations between them and to save space?

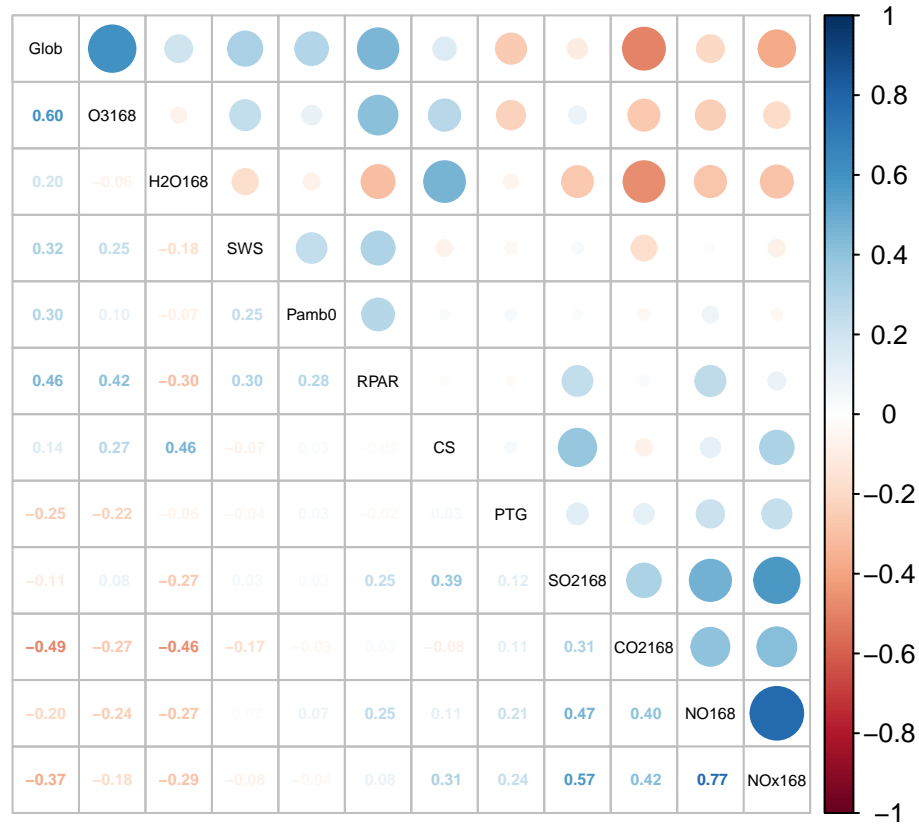
## Correlation matrix of mean variables for the clean dataset



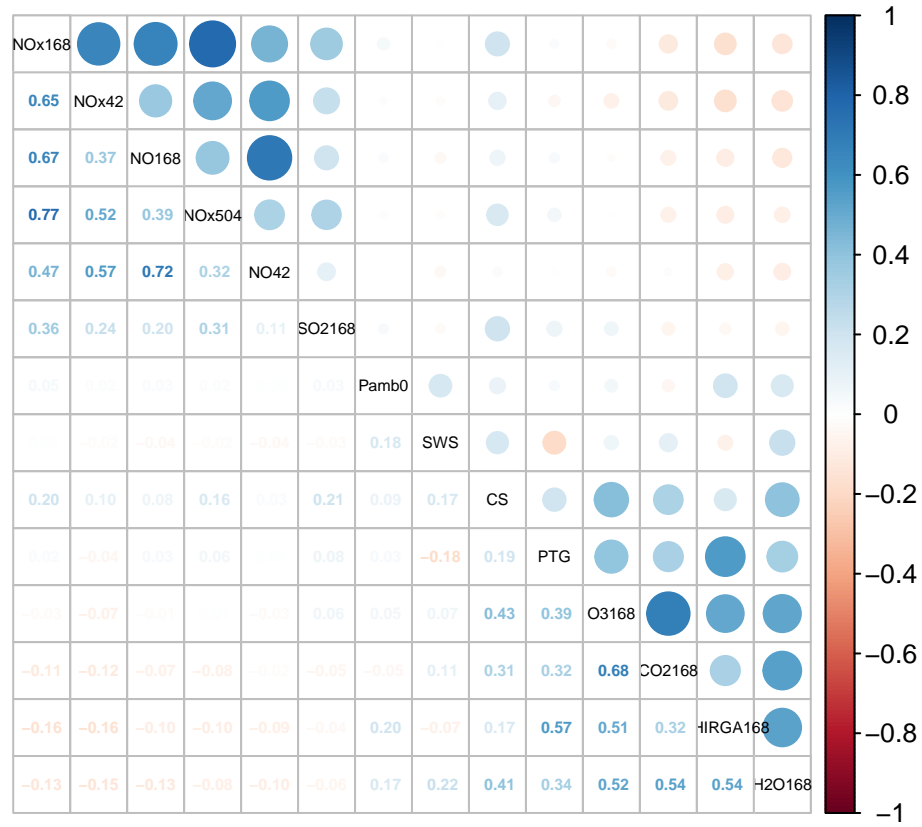
## Correlation matrix of std variables for the clean dataset



**Correlation matrix of mean variables for the clean dataset**



## Correlation matrix of std variables for the clean dataset



## Correlation matrix of variables for the clean dataset

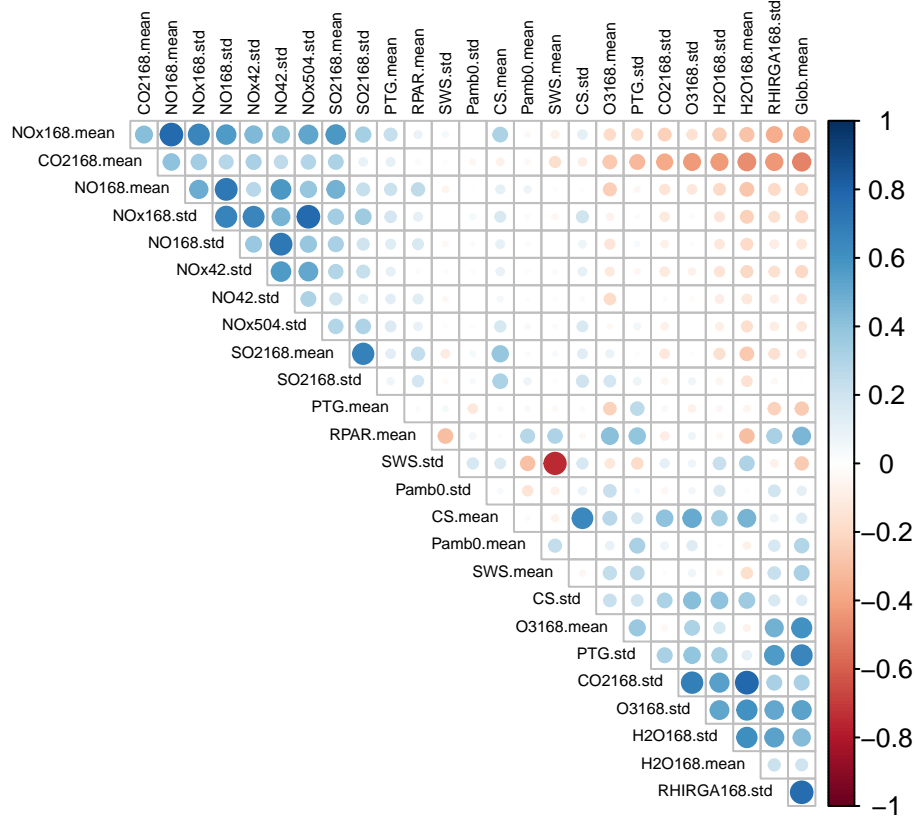


Table 2: List of kept variables

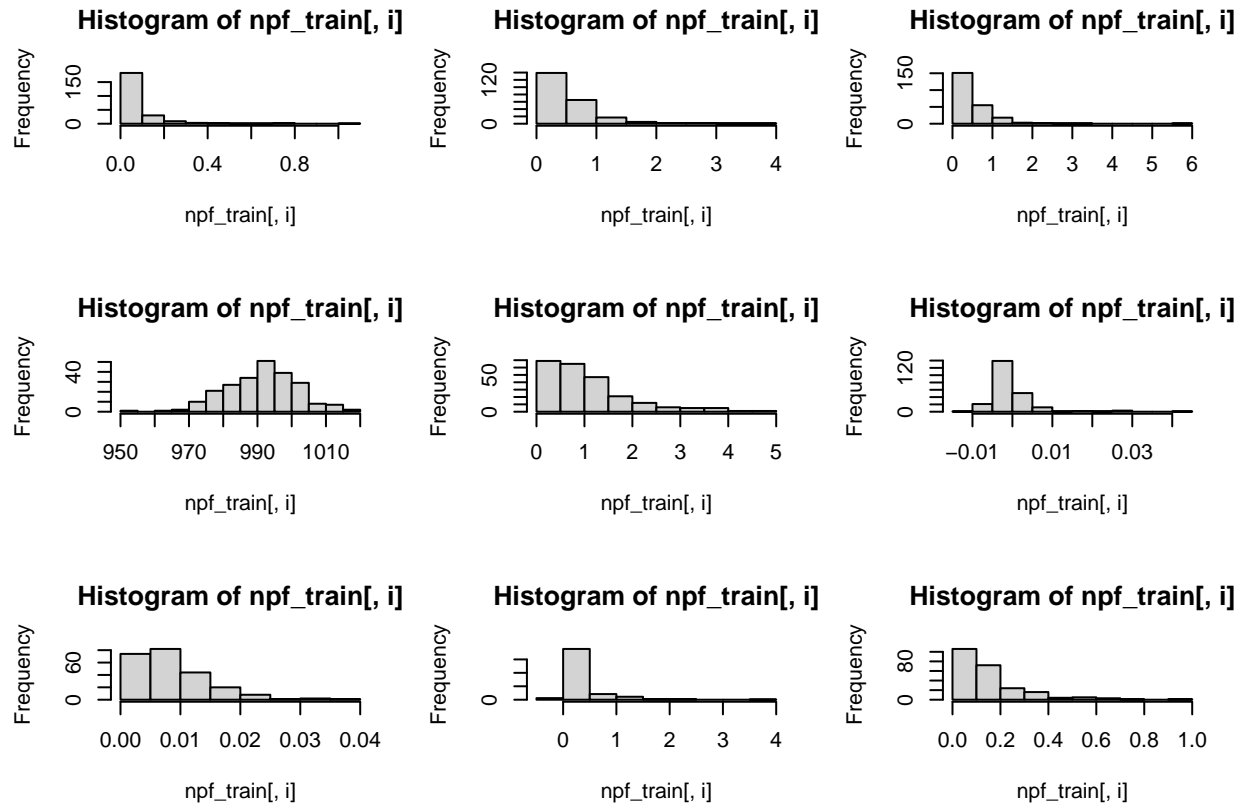
NO42.std	SO2168.mean	CO2168.std	NOx168.std
NOx42.std	SO2168.std	H2O168.mean	O3168.mean
NOx504.std	SWS.mean	H2O168.std	O3168.std
Pamb0.mean	SWS.std	NO168.mean	RHIRGA168.std
Pamb0.std	CS.mean	NO168.std	RPAR.mean
PTG.mean	CS.std	NOx168.mean	Glob.mean
PTG.std	CO2168.mean	NOx168.std	

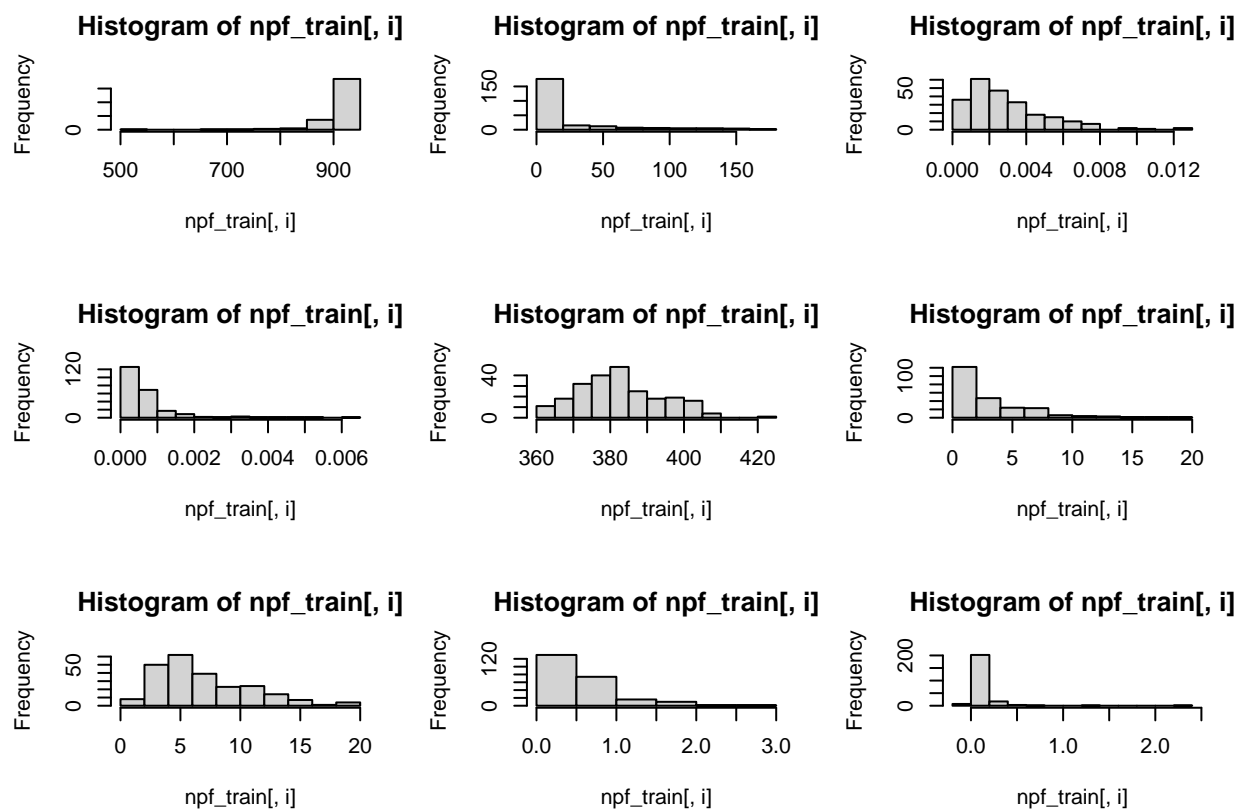
Table 3: Summary of clean dataset

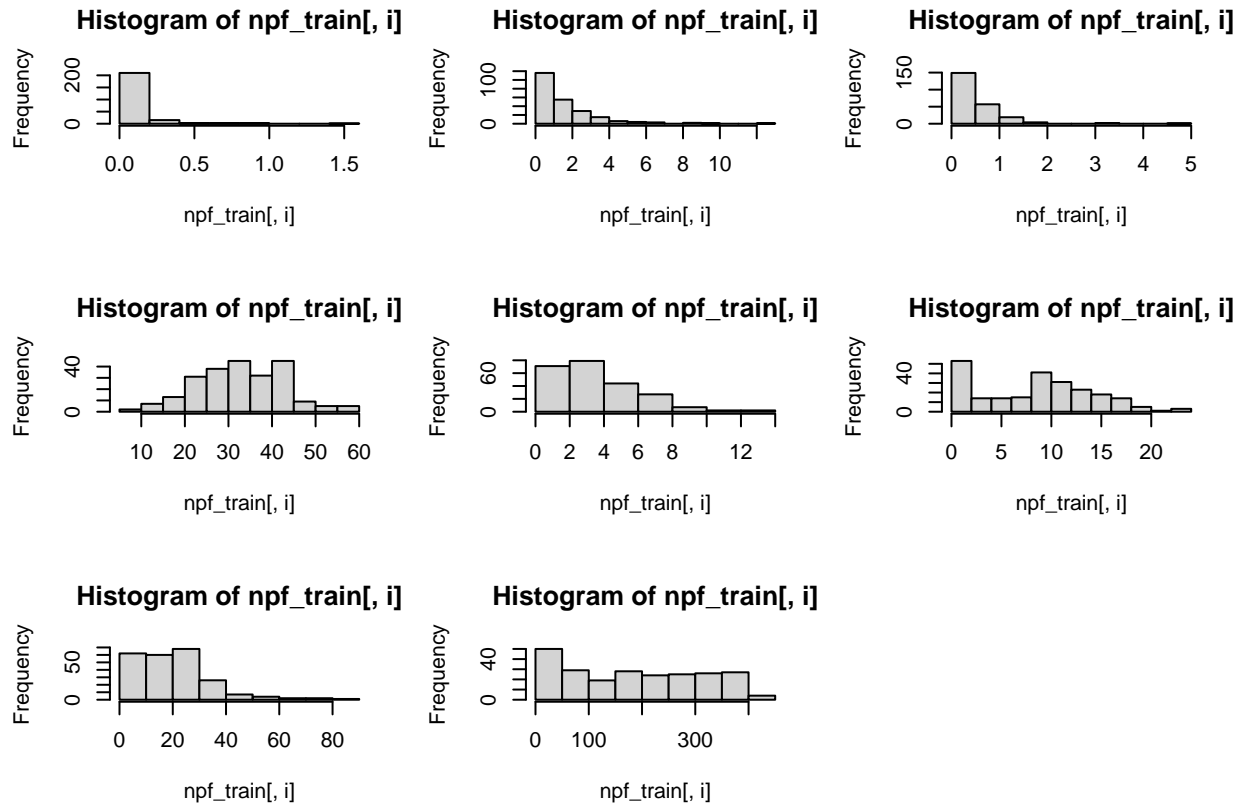
	npf_train	npf_test
Measurements	464	965
Variables	26	26

— I think we should put here some plot or graphical summary of the variables left after cleaning to get an overview of the values (for example to justify scaling in PCA). — Is there any good “summary picture” we could use? I tried histograms but there are too many of them...









## Performance measures

To compare different classifiers two measures are used: accuracy and perplexity. Accuracy is the proportion of the observations that has been classified correctly. Perplexity is a rescaled variant of log-likelihood. If perplexity equals to 1 the classifier predicts always the probability of an observation to an actual class. Perplexity of 2 corresponds to coin flipping.

For Random Forest and decision tree perplexity is not calculated but instead ROC curves are looked at. — Did we agree on this? :) ROC should still be added to code below

If the method predicts a probability of an event, observation is classified as “event” if the estimated posterior probability is more than 0.5. Otherwise the observation is classified as “nonevent”.

Performance measures are calculated using validation method and 10-fold Cross-Validation. For validation method the training data is randomly divided into 2 equally large data sets, training data to fit the model and validation data to estimate the accuracy and perplexity.

## Investigation of features

To find the most important variables logistic regression with Lasso, decision tree (“basic”), random forest and Principal Component Analysis (PCA) are used to cleaned data. The accuracy of the approaches are also calculated as well as perplexity of logistic regression with Lasso.

- 1) Logistic regression with Lasso, lambda selected by Cross-Validation:

Logistic regression is a discriminant classifier which assumes that the log odds is linear in variables. When combined with Lasso a subset selection of variables are done by adding so called penalty term to residual sum of squares which is minimized in parameter estimation process. The bigger the estimated coefficients of the variables are the bigger the penalty. The penalty term forces some of the coefficients to zero. The amount of penalty depends on a coefficient called lambda. The value of lambda is selected by cross-validation so that the value of the test error is minimized. The lambda is

```
## [1] 0.009059896
```

The estimated coefficients are

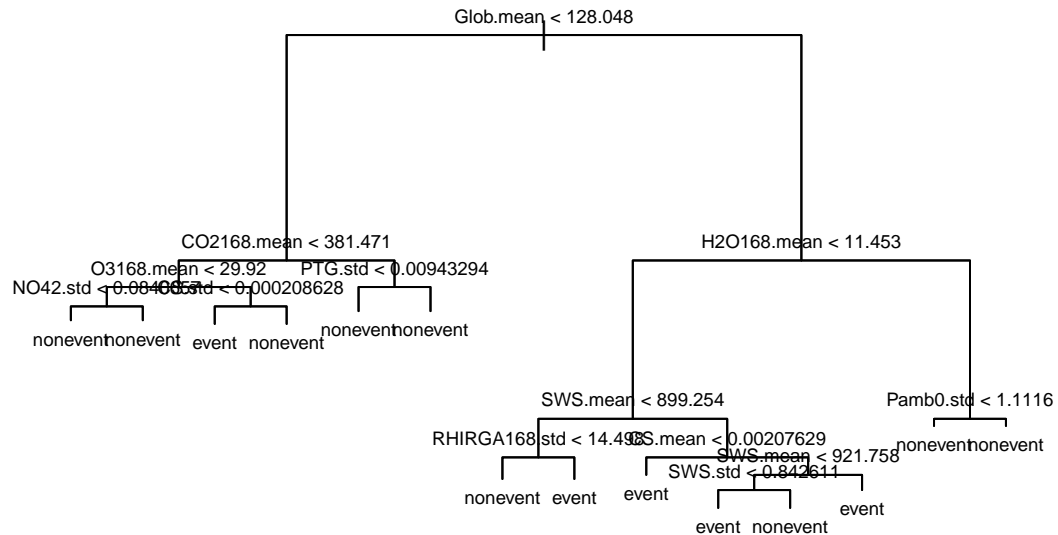
```
## 27 x 1 sparse Matrix of class "dgCMatrix"
##              s0
## (Intercept)  5.178482e-01
## N042.std     7.355450e-02
## N0x42.std    .
## N0x504.std   .
## Pamb0.mean   .
## Pamb0.std    2.284548e-01
## PTG.mean     .
## PTG.std      4.088522e+01
## S02168.mean  .
## S02168.std   4.954213e-01
## SWS.mean     1.451465e-02
## SWS.std      .
## CS.mean      -7.942649e+02
## CS.std       3.210184e+02
## C02168.mean  -4.745895e-02
## C02168.std   7.771859e-02
## H20168.mean  -2.422650e-01
## H20168.std   .
## N0168.mean   .
## N0168.std    2.096218e+00
## N0x168.mean  .
## N0x168.std   .
## O3168.mean   1.449695e-01
## O3168.std    9.487847e-02
## RHIRGA168.std 6.634250e-02
## RPAR.mean    .
## Glob.mean    7.590128e-03
```

As we can see, some of the variables have a zero coefficient meaning that Lasso has done variable selection.

2) A “normal” decision tree selects the following variables with the misclassification as follows

```
##
## Classification tree:
## tree(formula = class2 ~ ., data = npf_train)
## Variables actually used in tree construction:
## [1] "Glob.mean"      "C02168.mean"    "O3168.mean"     "N042.std"
## [5] "CS.std"         "PTG.std"        "H20168.mean"    "SWS.mean"
## [9] "RHIRGA168.std" "CS.mean"        "SWS.std"        "Pamb0.std"
```

```
## Number of terminal nodes: 14
## Residual mean deviance: 0.2277 = 49.64 / 218
## Misclassification error rate: 0.05172 = 12 / 232
```



Accuracy of the tree can be calculated from the confusion matrix. Predicted classes vs. actual classes for validation data:

```
##
## tree.pred  nonevent  event
##  nonevent      104     21
##   event        13     94
```

Predicted classes vs. actual classes for training data: — this can actually be seen from the results above (summary(tree.npf)) - do we want to leave this confusion matrix away?

```
##
## tree.pred_train nonevent event
##      nonevent      113      10
##      event         2      107
```

The accuracies are respectively

```
## [1] 0.948
```

```
## [1] 0.853
```

3) Random Forest with 5 variables (square root of the number of features) used in each run

Confusion matrix for training data

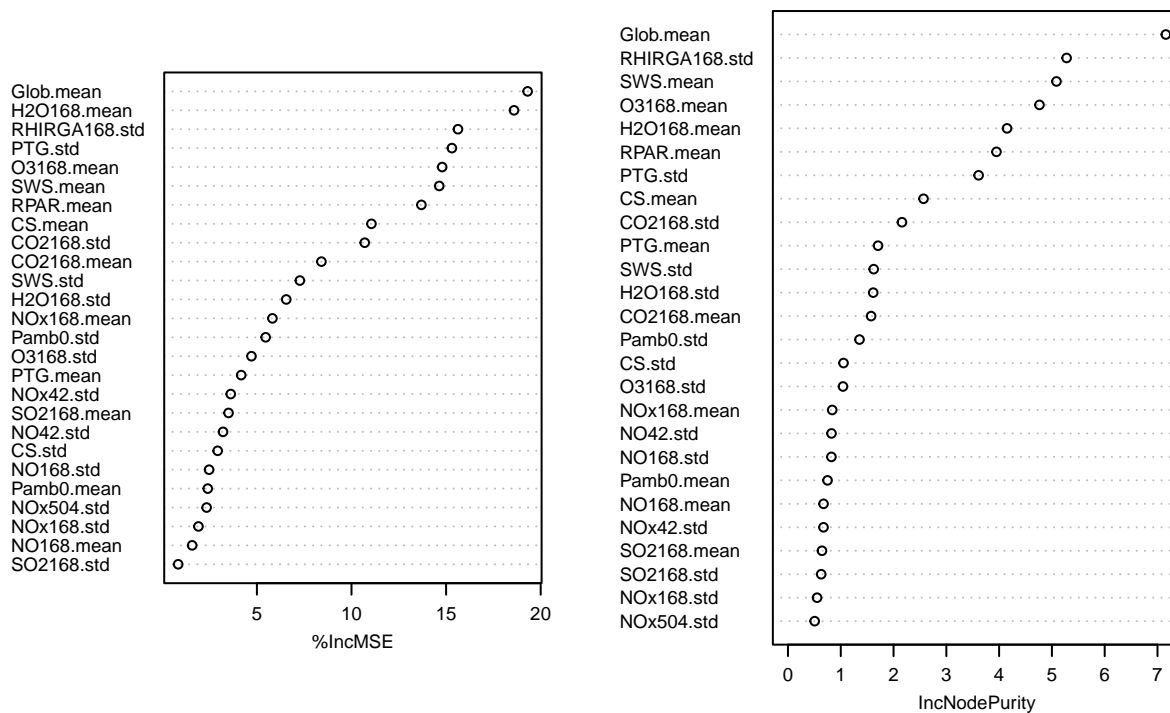
```
##
##      nonevent  event
##  0         115      0
##  1          0     117
```

Confusion matrix for validation

```
##
##      nonevent  event
##  0         102     18
##  1          15     97
```

The importance of the variables are:

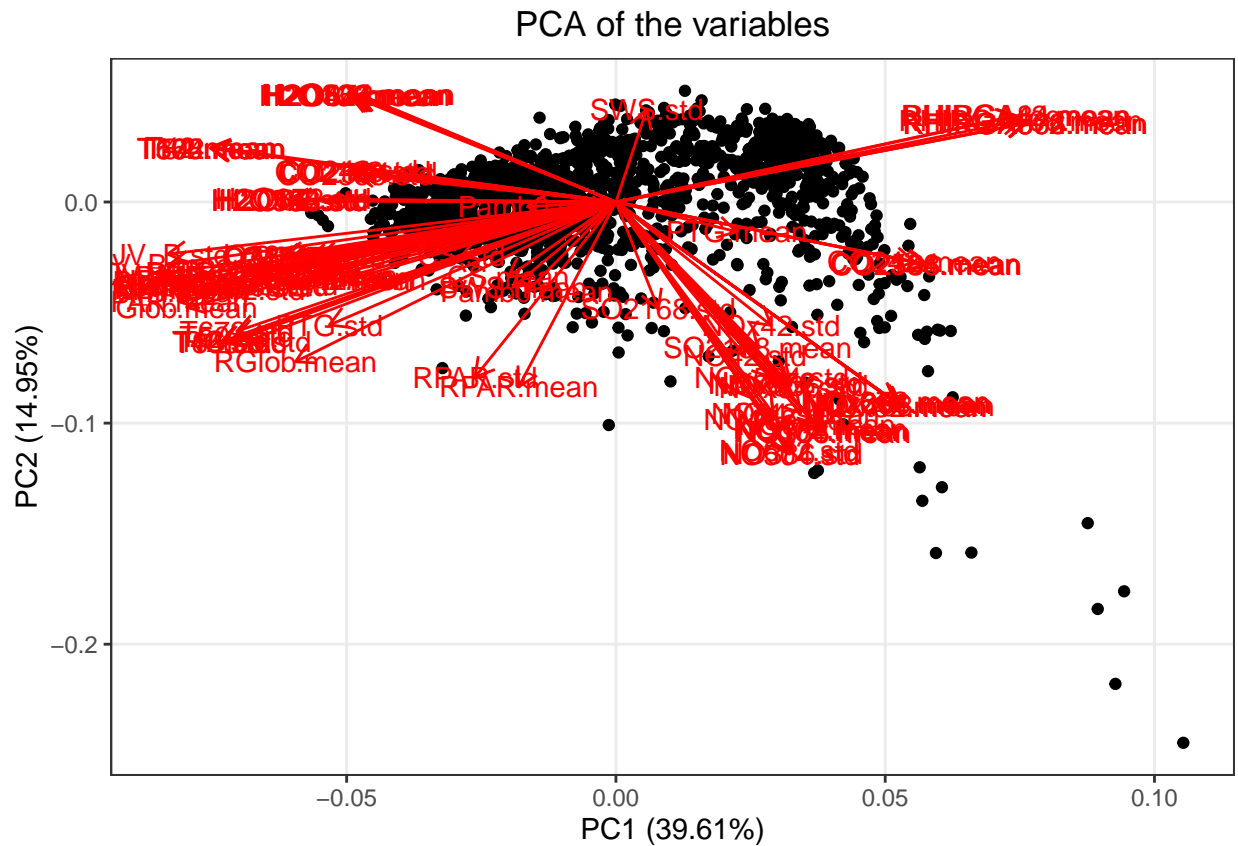
## Random Forest



```
## [1] 1
## [1] 2
## [1] 3
## [1] 4
## [1] 5
## [1] 6
## [1] 7
```

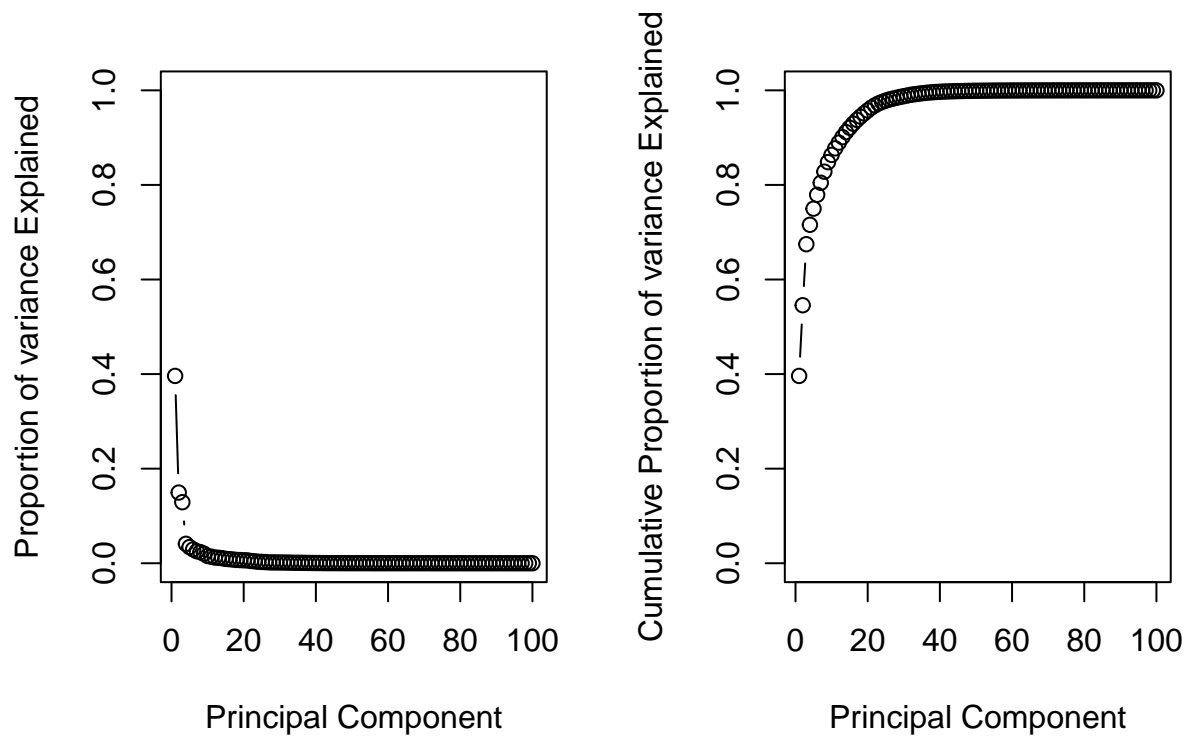
```
## [1] 8
## [1] 9
## [1] 10
## [1] 11
## [1] 12
## [1] 13
## [1] 14
## [1] 15
## [1] 16
## [1] 17
## [1] 18
## [1] 19
## [1] 20
```

- 4) PCA PCA is used to the original training data (npf\_train.csv with 464 observations) together with the original testing data (npf\_hidden.csv) where the variables are centered to have zero mean and scaled to have standard deviation one. The responses variables are removed from the original training data because we are using unsupervised learning method. The first two principal components are:



—Can we conclude from the picture above if there are any significant outliers?

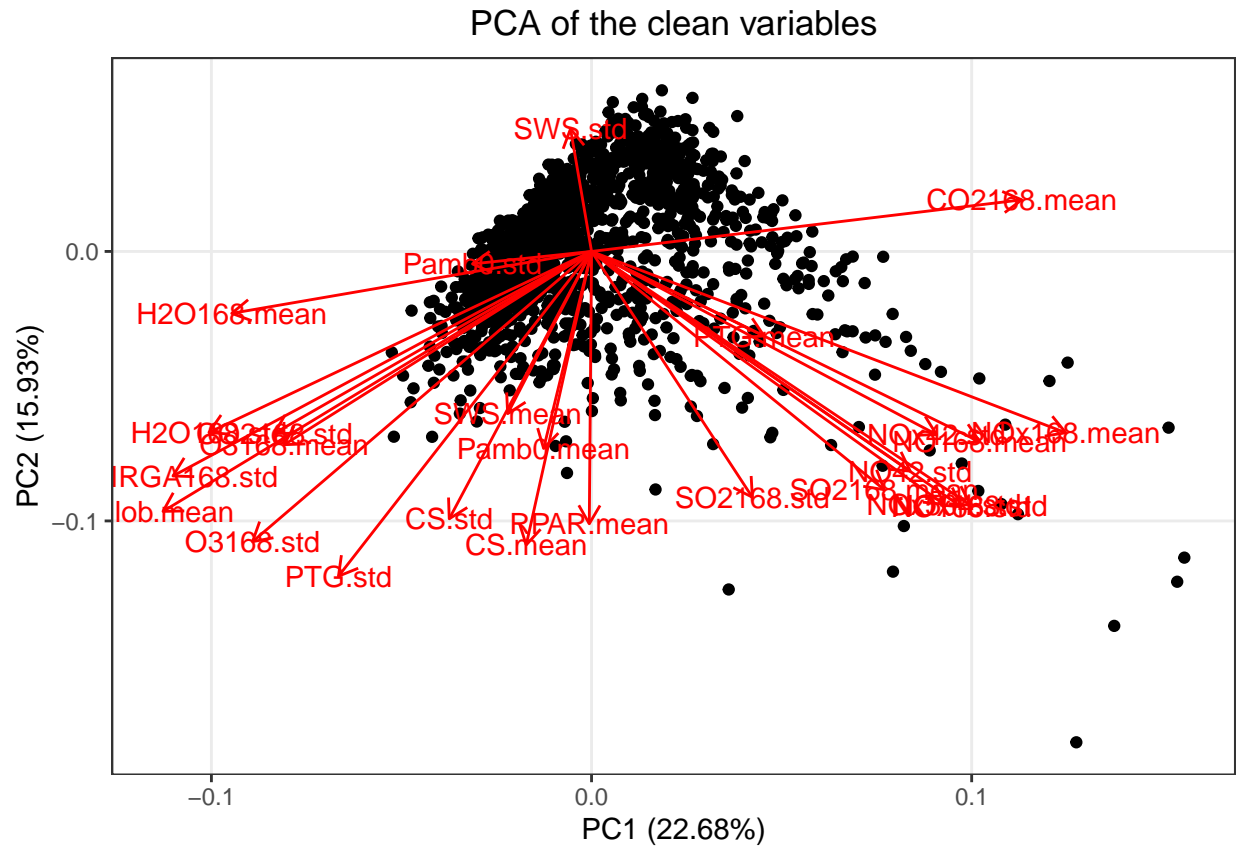
The proportion of variance explained (PVE) by each principal component and the cumulative PVE is shown



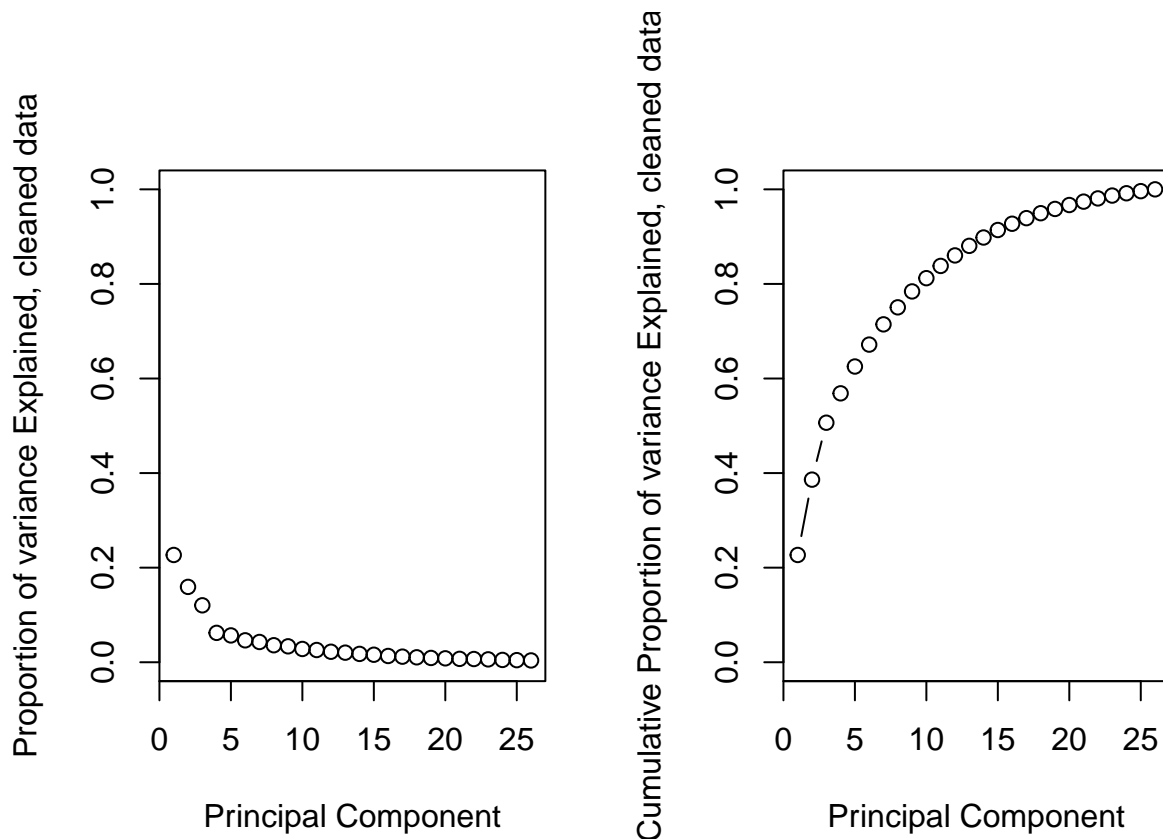
in the figures below:

PCA is also done for the data where variables with high correlation to other variables are removed (cleaned data). The first two principal components are:





The proportion of variance explained (PVE) by each principal component and the cumulative PVE is shown in the figures below:



According to the results above at least 15 first principal components should be used to explain about 90 % of the variance. — If you can conclude something else from the results, you can add some text :)

The performance of the classifiers investigated so far are:

##	Model	Train Accuracy	Validation Accuracy	CV Accuracy
## 1	Log reg Lasso CV	0.918		0.836
## 2	Tree	0.853		0.948
## 3	Random Forest	1		0.858
##	Train Perplexity	Validation Perplexity	CV Perplexity	
## 1	1.24	1.348	1.311	
## 2				
## 3	1.117	1.406	1.324	

## Conclusions and feature selection

— !!! This whole part can be removed if it's not needed. The text is not updated after cleaning the data so at least the conclusions should be updated. The models in the next session is set to use cleaned data instead of the selection previously done here. !!!

The accuracy is best for Random Forest, but also Logistic Regression with Lasso where lambda is selected using CV ("log reg CV") performs very well. — Should we investigate diagnostic plots for log reg? Should we give weight to well performed classifiers when selecting the variables?

All models use the following variables - RHIRGA: mean is more important than std according to tree and RF - H2O mean - O3 mean (Logistic Regression with Lasso, lambda = CV ("log reg CV") gives small value to std) - CO2: tree and RF uses only std, log regs give value to mean too - T std - SWS.mean

The following variables are used in all other models but RF - NO.std - Pamb0.mean or .std

The following variables are used in all other models but “normal” tree - CS.mean (.std only in log reg CV)

These variables are used in some models - UV\_B.std is used only on log reg CV - UV\_A.mean is used only in RF - RGlob.mean and RGlob.std (RF, mean also in log reg 0.1) - PTG.std (tree, RF) - NET.mean - RPAR.mean - SO2

Because the measures in different heights/levels are highly correlated, we select just one of them. — which ones? Should we put more weight to the selections in random forest and log reg CV? — Should be produce boxplots, histograms and/or scatterplots to selected variables

We select the following variables to be used in models: — selected for testing the same variables than in Log Reg Lasso CV but only once if there are multiple values with different heights ———

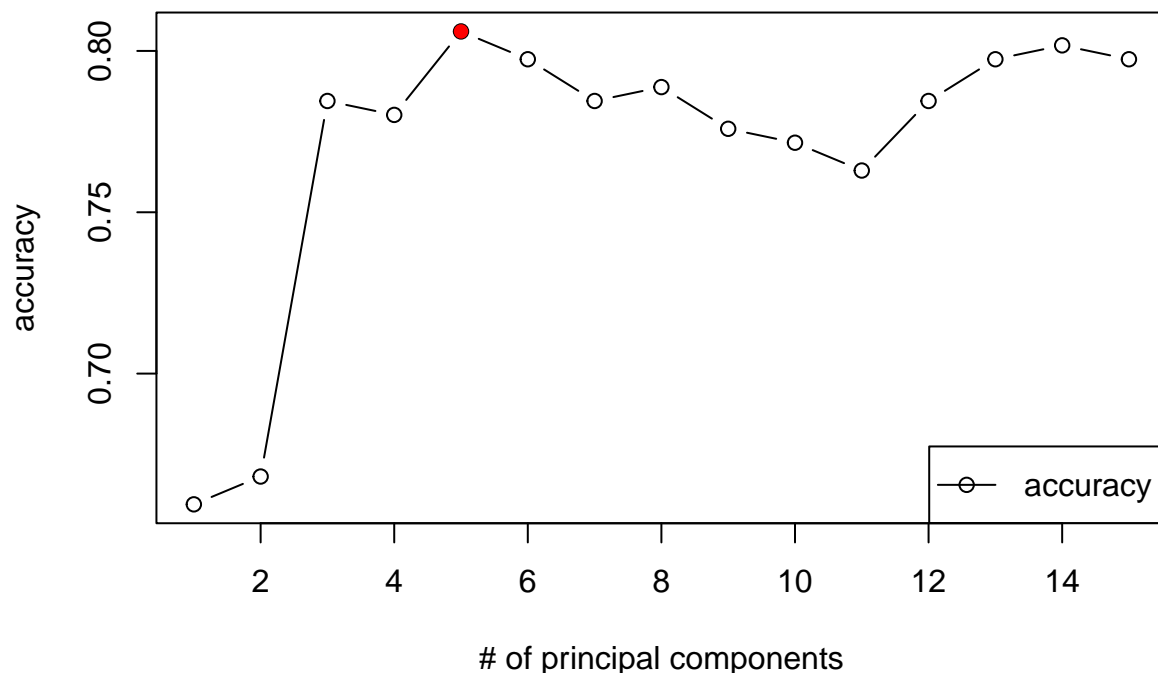
Let's check if there are any correlation left:

## Model selection

The models tested in addition to the ones already tested are

- 1) Dummy: — write a description of the method and the code itself. . .
- 2) Naive Bayes

Accuracy of Naive Bayes when using the first 15 principal components to reduce the dimensionality of the data. The highest accuracy is marked with red.



The highest accuracy is received when using 5 first PCs so this is added to be one of the possible methods.

- 3) Logistic regression — with or without interactions? now it's without
- 4) k-NN K nearest neighbor is tried with different values of k: 1, 5, 10, 15, 20 and 50. Accuracy on validation set for each of them are respectively

```
## [1] 0.780 0.836 0.836 0.828 0.823 0.780
```

When  $k = 5$  (or 10), the accuracy is highest so 5-NN is added

Performance measures for the methods are:

##	Model	Train Accuracy	Validation Accuracy	CV Accuracy
## 1	Log reg Lasso CV	0.918	0.836	0.886
## 2	Tree	0.853	0.948	
## 3	Random Forest	1	0.858	0.881
## 4	Naive Bayes	0.853	0.746	0.759
## 5	Naive Bayes with PCA	0.875	0.806	0.819
## 6	Logistic regression	0.94	0.866	0.89
## 7	5-NN	0.892	0.836	0.5

##	Train Perplexity	Validation Perplexity	CV Perplexity
## 1	1.24	1.348	1.311
## 2			
## 3	1.117	1.406	1.324
## 4	2.106	Inf	6.232
## 5	1.46	1.471	1.626
## 6	1.184	1.411	1.352
## 7			