

# Assignment 2

Group 4

2024-09-16

```
options(contrasts = c("contr.sum", "contr.poly"))
require("ggplot2")
require("dplyr")
require("ppcor")
require("caret")
require("tidyverse")
require("stringr")
require("lubridate")
require("tsibble")
require("ggfortify")
require("gridExtra")
require("reshape2")

library(ggcorrplot) # Easy cross-correlation plots
library(imputeTS)   # Time series missing value imputation

library(jsonlite)   # handle JSON data returned by Frost
library(tidyverse)  # unpack data from JSON format
library(lubridate)  # data manipulation with mutate etc, string formatting
library(tsibble)    # process date and time information
library(fpp3)       # autoplot() and gg_season() for time series
library(readr)      # to read the Frost client ID from file
```

## Task 1: Dimension reduction on air quality data

### Part A: Get

- Obtain data from <https://archive.ics.uci.edu/dataset/360/air+quality>.
- Provide a brief description of the data based on the information from the website.

```
airquality <- read.table("AirQualityUCI.csv", sep=";", dec=",", header= T)

airqual <- airquality %>%
  dplyr::select(-c("X", "X.1")) %>%
  na.omit() %>%
  mutate(timedate = dmy_hms(paste(Date, Time))) %>%
  dplyr::select(-c("Time", "Date")) %>%
  relocate(timedate) %>%
  as_tibble()

summary(airqual)
```

```

##      timedata          CO.GT.        PT08.S1.CO.
##  Min.   :2004-03-10 18:00:00  Min.   :-200.00  Min.   :-200
##  1st Qu.:2004-06-16 05:00:00  1st Qu.:  0.60  1st Qu.: 921
##  Median :2004-09-21 16:00:00  Median :  1.50  Median :1053
##  Mean   :2004-09-21 16:00:00  Mean   : -34.21  Mean   :1049
##  3rd Qu.:2004-12-28 03:00:00  3rd Qu.:  2.60  3rd Qu.:1221
##  Max.   :2005-04-04 14:00:00  Max.   : 11.90  Max.   :2040
##      NMHC.GT.        C6H6.GT.        PT08.S2.NMHC.        NOx.GT.
##  Min.   :-200.0   Min.   :-200.000  Min.   :-200.0   Min.   :-200.0
##  1st Qu.:-200.0   1st Qu.:  4.000  1st Qu.: 711.0   1st Qu.: 50.0
##  Median :-200.0   Median :  7.900  Median : 895.0   Median : 141.0
##  Mean   :-159.1   Mean   :  1.866  Mean   : 894.6   Mean   : 168.6
##  3rd Qu.:-200.0   3rd Qu.: 13.600  3rd Qu.:1105.0   3rd Qu.: 284.0
##  Max.   :1189.0   Max.   : 63.700  Max.   :2214.0   Max.   :1479.0
##      PT08.S3.NOx.    NO2.GT.        PT08.S4.NO2.        PT08.S5.03.
##  Min.   :-200     Min.   :-200.00  Min.   :-200     Min.   :-200.0
##  1st Qu.: 637    1st Qu.:  53.00  1st Qu.:1185    1st Qu.: 700.0
##  Median : 794    Median :  96.00  Median :1446    Median : 942.0
##  Mean   : 795    Mean   :  58.15  Mean   :1391    Mean   : 975.1
##  3rd Qu.: 960    3rd Qu.: 133.00  3rd Qu.:1662    3rd Qu.:1255.0
##  Max.   :2683    Max.   : 340.00  Max.   :2775    Max.   :2523.0
##      T            RH            AH
##  Min.   :-200.000  Min.   :-200.00  Min.   :-200.0000
##  1st Qu.: 10.900   1st Qu.: 34.10  1st Qu.: 0.6923
##  Median : 17.200   Median : 48.60  Median : 0.9768
##  Mean   : 9.778    Mean   : 39.49  Mean   : -6.8376
##  3rd Qu.: 24.100   3rd Qu.: 61.90  3rd Qu.: 1.2962
##  Max.   : 44.600   Max.   : 88.70  Max.   : 2.2310

head(airqual)

## # A tibble: 6 x 14
##   timedata          CO.GT.  PT08.S1.CO. NMHC.GT. C6H6.GT. PT08.S2.NMHC. NOx.GT.
##   <dttm>           <dbl>    <int>    <int>    <dbl>    <int>    <int>
## 1 2004-03-10 18:00:00    2.6     1360     150     11.9     1046     166
## 2 2004-03-10 19:00:00     2      1292     112      9.4      955     103
## 3 2004-03-10 20:00:00    2.2     1402      88       9      939     131
## 4 2004-03-10 21:00:00    2.2     1376      80      9.2      948     172
## 5 2004-03-10 22:00:00    1.6     1272      51      6.5      836     131
## 6 2004-03-10 23:00:00    1.2     1197      38      4.7      750      89
## # i 7 more variables: PT08.S3.NOx. <int>, NO2.GT. <int>, PT08.S4.NO2. <int>,
## #   PT08.S5.03. <int>, T <dbl>, RH <dbl>, AH <dbl>
```

Contains the responses of a gas multisensor device deployed on the field in an Italian city. Hourly responses averages are recorded along with gas concentrations references from a certified analyzer. Multivariate (15) and time series Has missing values.

The dataset has 15 columns:

1. Date: Measurement date
2. Time: Measurement time
3. CO(GT): Carbon Monoxide concentration in ppm (ground truth)
4. PT08.S1(CO): Sensor 1, Tin Oxide response related to CO
5. NMHC(GT): Non-Methane Hydrocarbons concentration in micrograms/m<sup>3</sup>
6. C6H6(GT): Benzene concentration in micrograms/m<sup>3</sup>
7. PT08.S2(NMHC): Sensor 2, response related to Non-Methane Hydrocarbons

8. NOx(GT): Nitrogen Oxides concentration in parts per billion (ground truth)
9. PT08.S3(NOx): Sensor 3, response related to NOx
10. NO2(GT): Nitrogen Dioxide concentration in parts per billion
11. PT08.S4(NO2): Sensor 4, response related to NO2
12. PT08.S5(O3): Sensor 5, response related to O3 (Ozone)
13. T: Temperature in Celsius
14. RH: Relative Humidity (%)
15. AH: Absolute Humidity (g/m<sup>3</sup>)

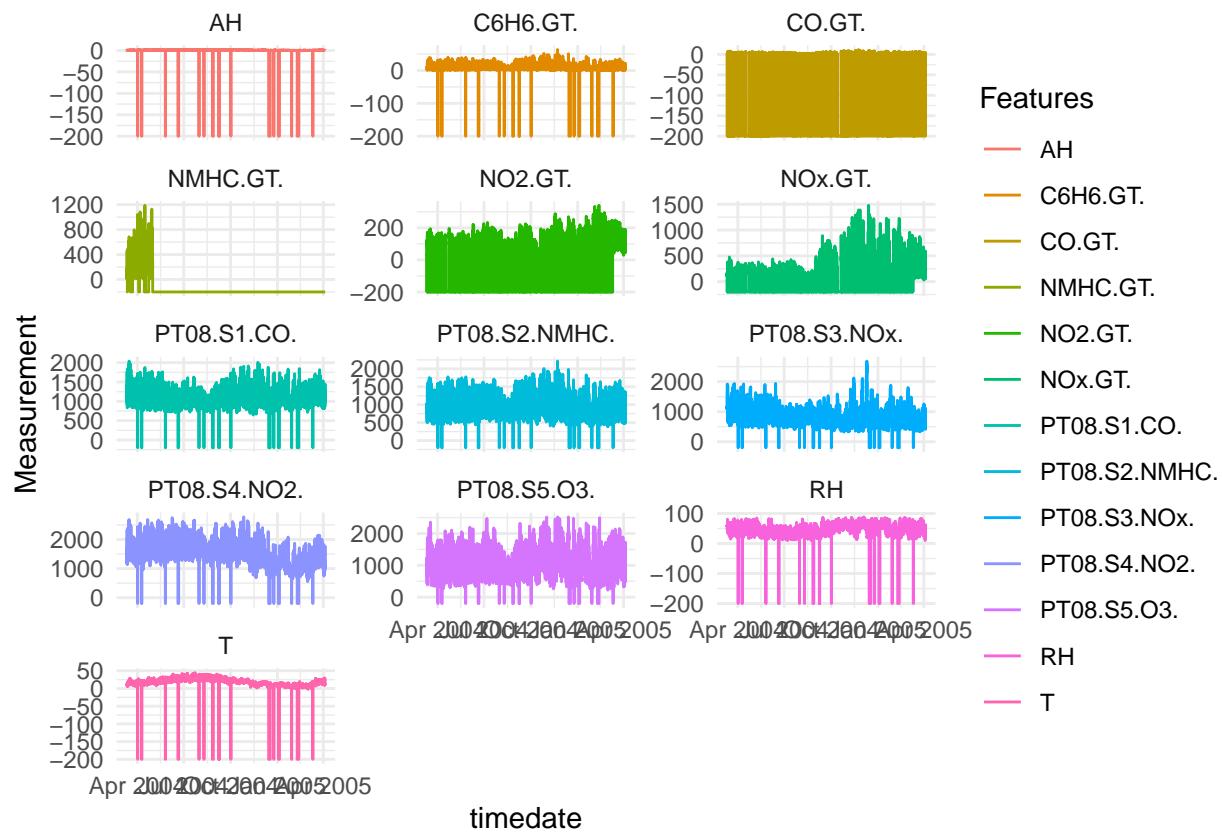
### Hints

- See options of `read.table()` for correct import

## Part B: Import and Visualize

- Load the data and convert to tsibble.
- Make sure dates and hours are converted into proper time objects
- Remove incomplete days at beginning and end of data
- Plot the data as is, preferably as multiple panels in a single plot
- Describe the data. What is most striking?

```
airqual %>%
  pivot_longer(!timedate, names_to="Features", values_to="Measurement") %>%
  ggplot(aes(x = timedate, y = Measurement, col = Features)) +
  geom_line() +
  theme_minimal() +
  facet_wrap(~ Features, scales = "free_y", ncol = 3)
```



*There seems to be multiple -200 in all the numerical (integer and Categorical) data that seems to be outliers or missing values. Should probably be removed or imputed from the dataset.* NHMC.GT. has a lot of NA from early on, should probably remove \* CO.GT., NO2.GT., NOx.GT. seems to have a lot of missing values, should investigate how severe the gaps are and if it can be remedied with imputations.

### Part C: PCA of data as is

- Perform PCA on the data as prepared in B

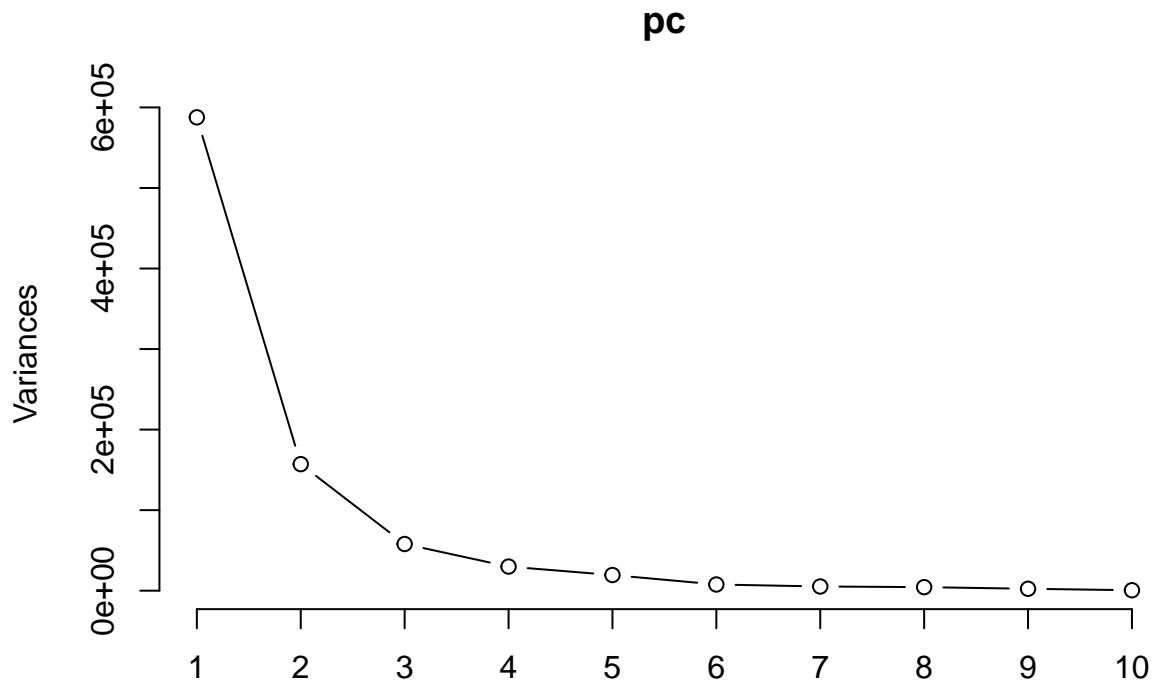
```
pc <- prcomp(airqual[,-1])
summary(pc)

## Importance of components:
##                PC1       PC2       PC3       PC4       PC5       PC6
## Standard deviation    766.6076 396.3514 240.73023 172.96207 139.37119 87.72389
## Proportion of Variance  0.6736   0.1801   0.06643   0.03429   0.02226   0.00882
## Cumulative Proportion  0.6736   0.8537   0.92012   0.95441   0.97668   0.98550
##                PC7       PC8       PC9       PC10      PC11      PC12
## Standard deviation    72.12259 66.31383 48.62607 23.26879 11.78504 2.29019
## Proportion of Variance  0.00596  0.00504  0.00271  0.00062  0.00016  0.00001
## Cumulative Proportion  0.99146  0.99650  0.99921  0.99983  0.99999  1.00000
##                PC13
## Standard deviation    0.7725
## Proportion of Variance 0.0000
## Cumulative Proportion 1.0000
```

- Create a screeplot and create biplots for 1st and 2nd and for 2nd and 3rd PCs

#### Screeplot

```
plot(pc, type = "l")
```



```

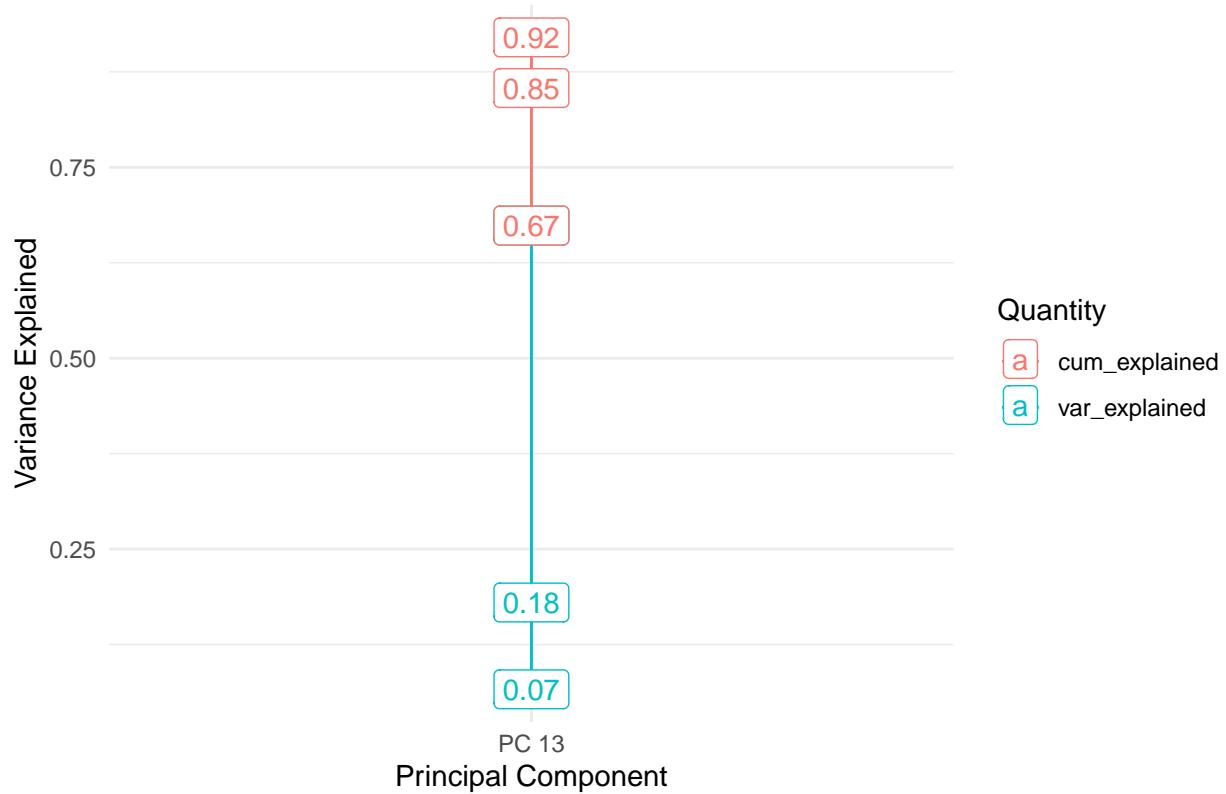
pc_v <- data.frame(PC = paste0("PC ", ncol(pc$x)),
                     var_explained = pc$sdev^2 / sum(pc$sdev^2)) %>%
  mutate(cum_explained = cumsum(var_explained))

pp <- pc_v[1:3,] %>%
  pivot_longer(!PC, names_to="Quantity", values_to="Explained") %>%
  ggplot(aes(x = PC, y = Explained, color=Quantity, group=Quantity)) +
  geom_line() + geom_point() +
  theme_minimal() +
  labs(title = "Variance Explained", x = "Principal Component",
       y = "Variance Explained")

pp + geom_label(aes(label = round(Explained, 2)))

```

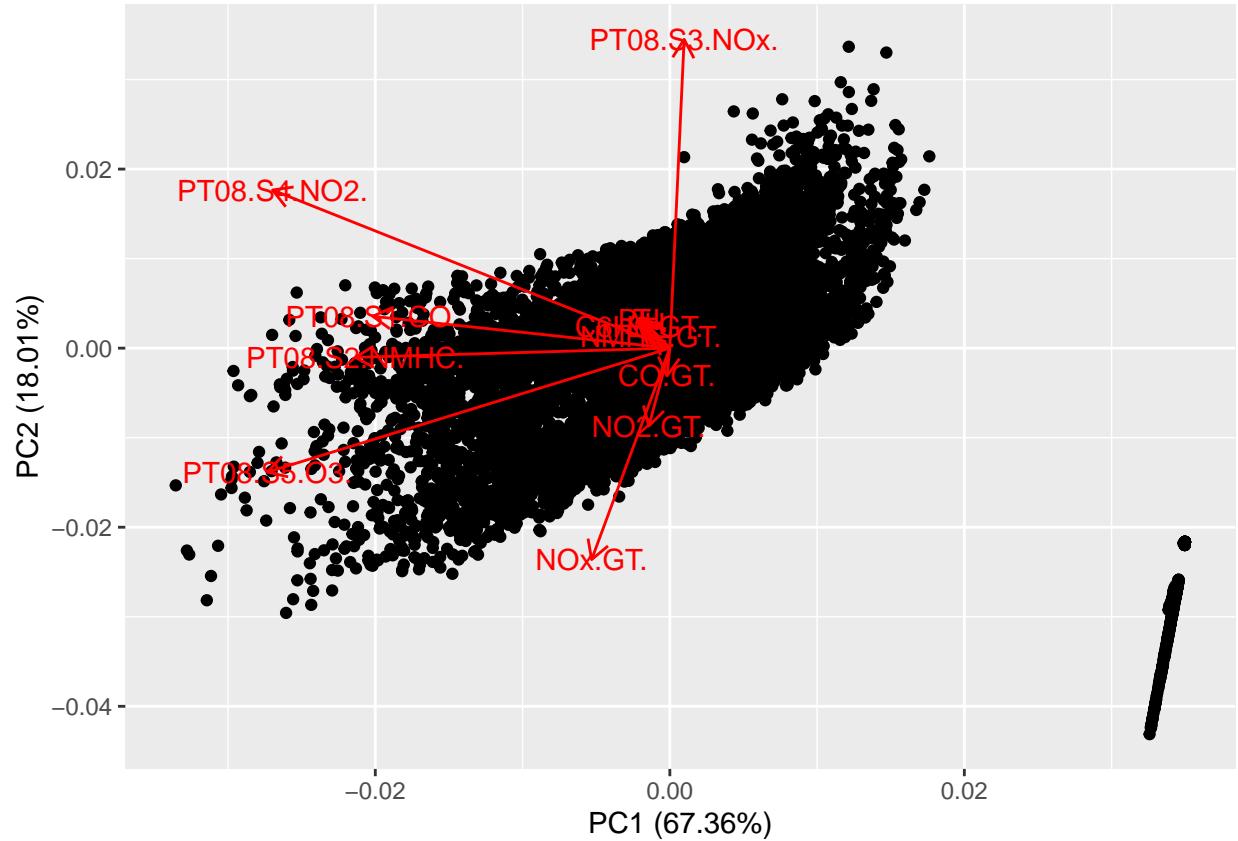
## Variance Explained

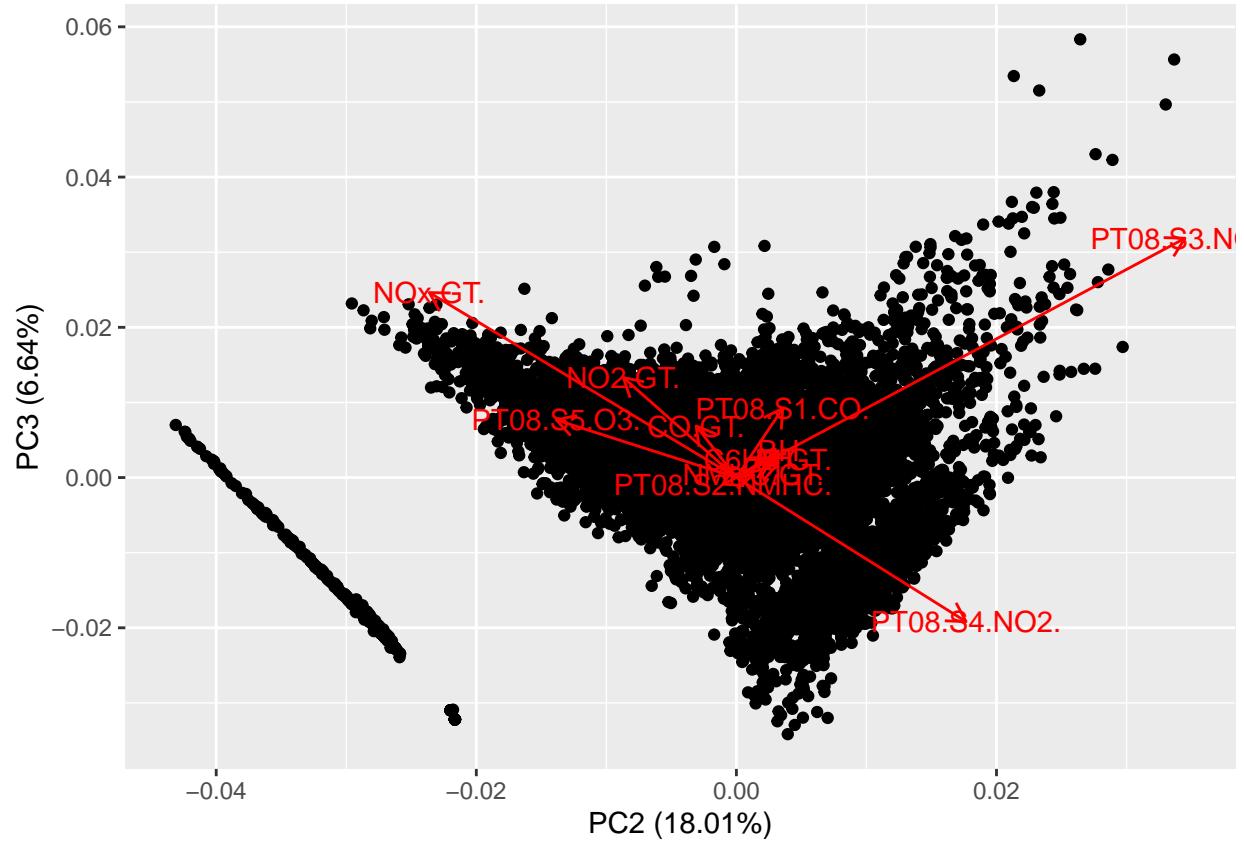


0.91 of variance is explained in the first 3 principal components

## Biplots

```
autoplot(pc, x=1,y=2, loadings=TRUE, loadings.label=TRUE)
```





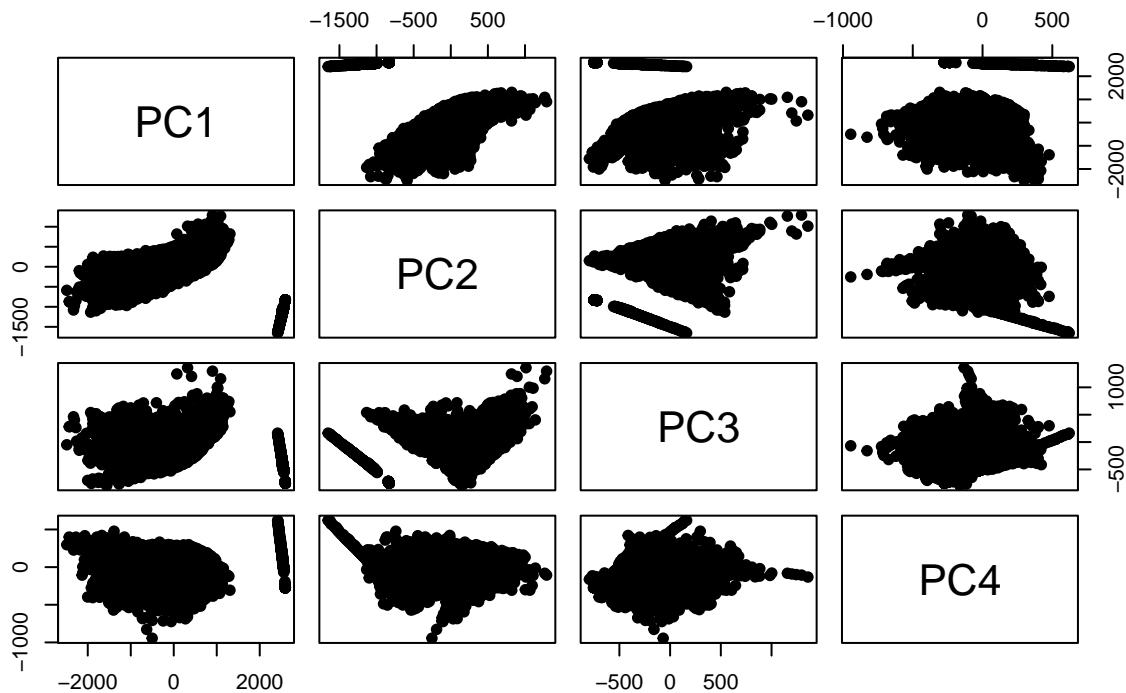
- Can clearly see the -200 outliers
- PRO8.S3.NOx takes a lot from PC2 than PC1
- NOX.GT. opposite
- Plot the scores for the PCs
- Comment on the results. Can you relate some features to your observations in part B?

Can clearly see that there is an line of outliers in almost every PC scoreplot combination, this is higly likely to be the -200 (NA) values

### Score plot

```
pairs(pc$x[,1:4], main = "PCA Score Plots", pch = 19)
```

## PCA Score Plots



### Hints

- `ggfortify` provides `autoplot()` for PCA results for ggplot-style biplots
- To plot the scores, you can use the same code as for plotting the original data

### Part D: Missing values

- Identify missing values in the time series

The website says that all -200 are NA

- Investigate to which degree missing values occur at the same time for multiple sensors

```
airqual_NA <- airqual
airqual_NA[airqual_NA == -200] = NA
missing_values <- colSums(is.na(airqual_NA)) %>% data.frame()
perc_NA_airqal <- (missing_values/nrow(airqual_NA))*100
```

```
missing_values
```

```
##
## timedate      .
## CO.GT.       1683
## PT08.S1.CO.   366
## NMHC.GT.     8443
## C6H6.GT.     366
## PT08.S2.NMHC. 366
## NOx.GT.      1639
```

```

## PT08.S3.NOx.    366
## NO2.GT.        1642
## PT08.S4.NO2.    366
## PT08.S5.03.    366
## T               366
## RH              366
## AH              366

perc_NA_airqal

##
##          .
## timedate      0.00000
## CO.GT.       17.98653
## PT08.S1.CO.   3.91151
## NMHC.GT.     90.23191
## C6H6.GT.     3.91151
## PT08.S2.NMHC. 3.91151
## NOx.GT.      17.51630
## PT08.S3.NOx.   3.91151
## NO2.GT.       17.54836
## PT08.S4.NO2.   3.91151
## PT08.S5.03.   3.91151
## T               3.91151
## RH              3.91151
## AH              3.91151

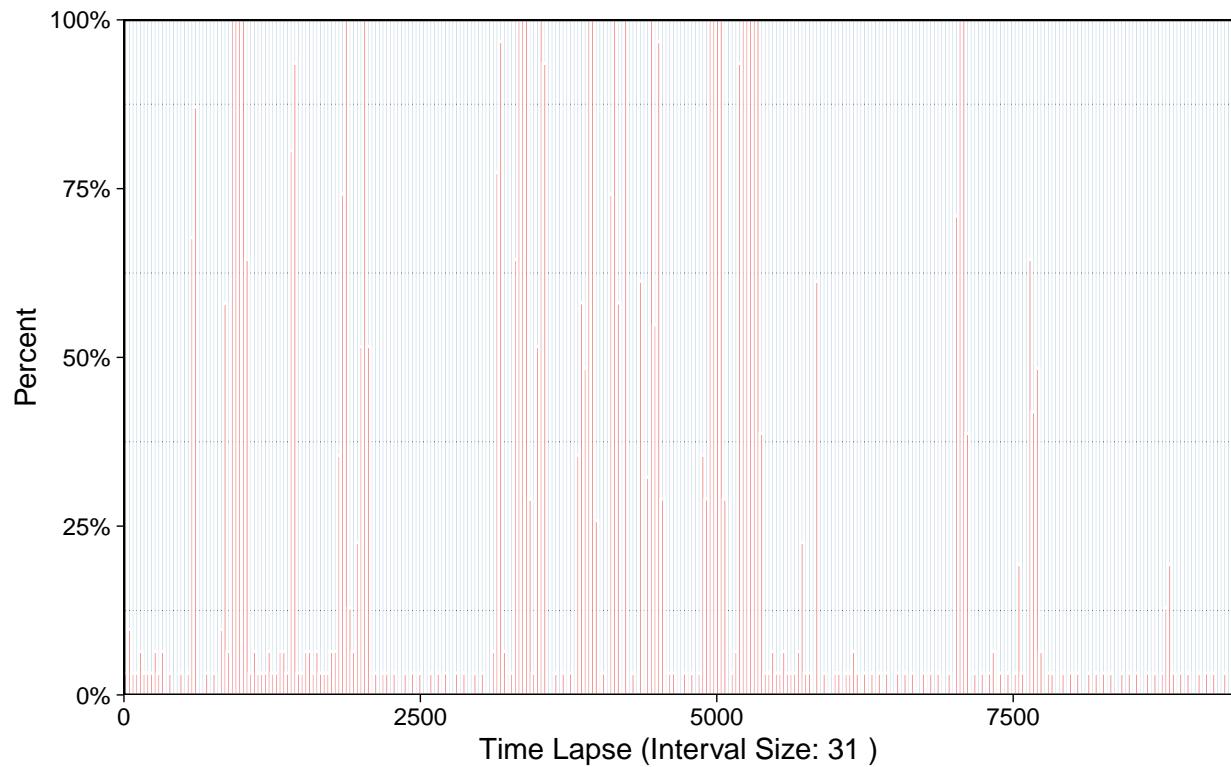
```

CO.GT. , NOx.GT. and NO2.GT. seems to have significant amount of missing values, investigate how distributed the gaps are to get an impression how well imputations will work.

```
ggplot_na_distribution2(airqual_NA$CO.GT., interval_size = 31)
```

## Missing Values per Interval

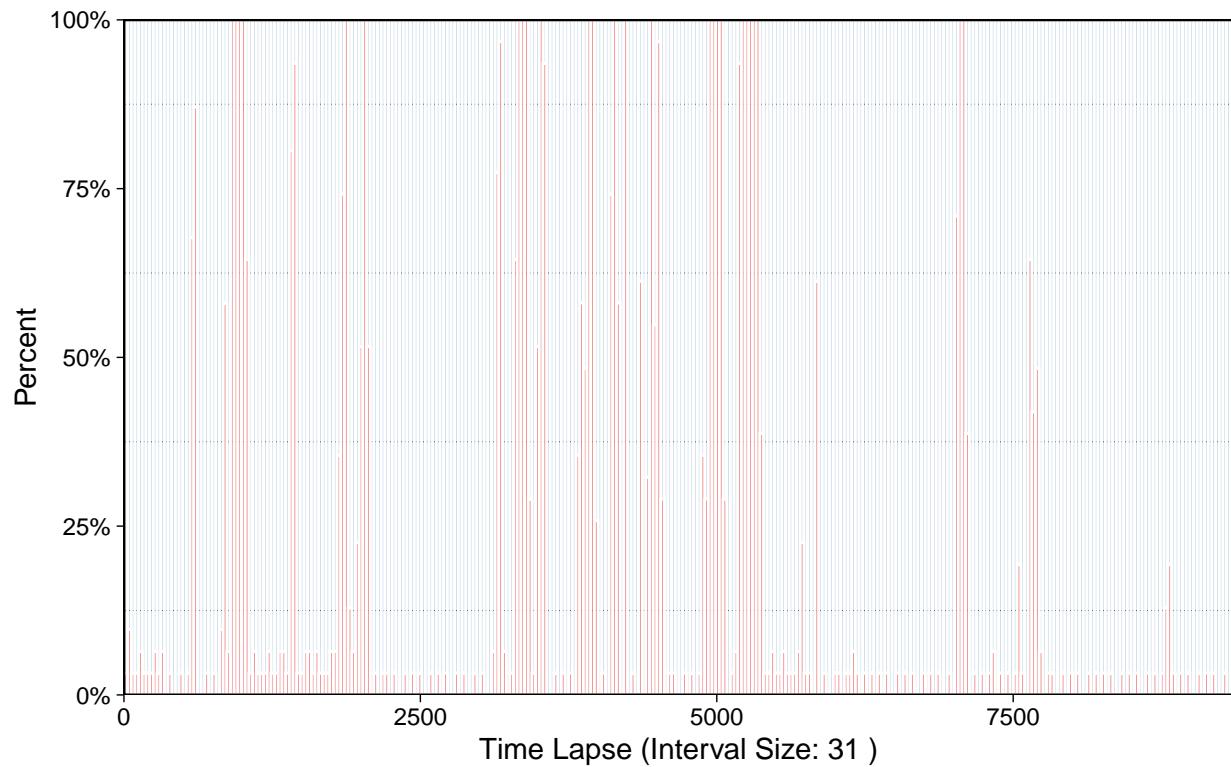
Amount of **NA** and **non-NA** for successive intervals



```
ggplot_na_distribution2(airqual_NA$CO.GT., interval_size = 31)
```

## Missing Values per Interval

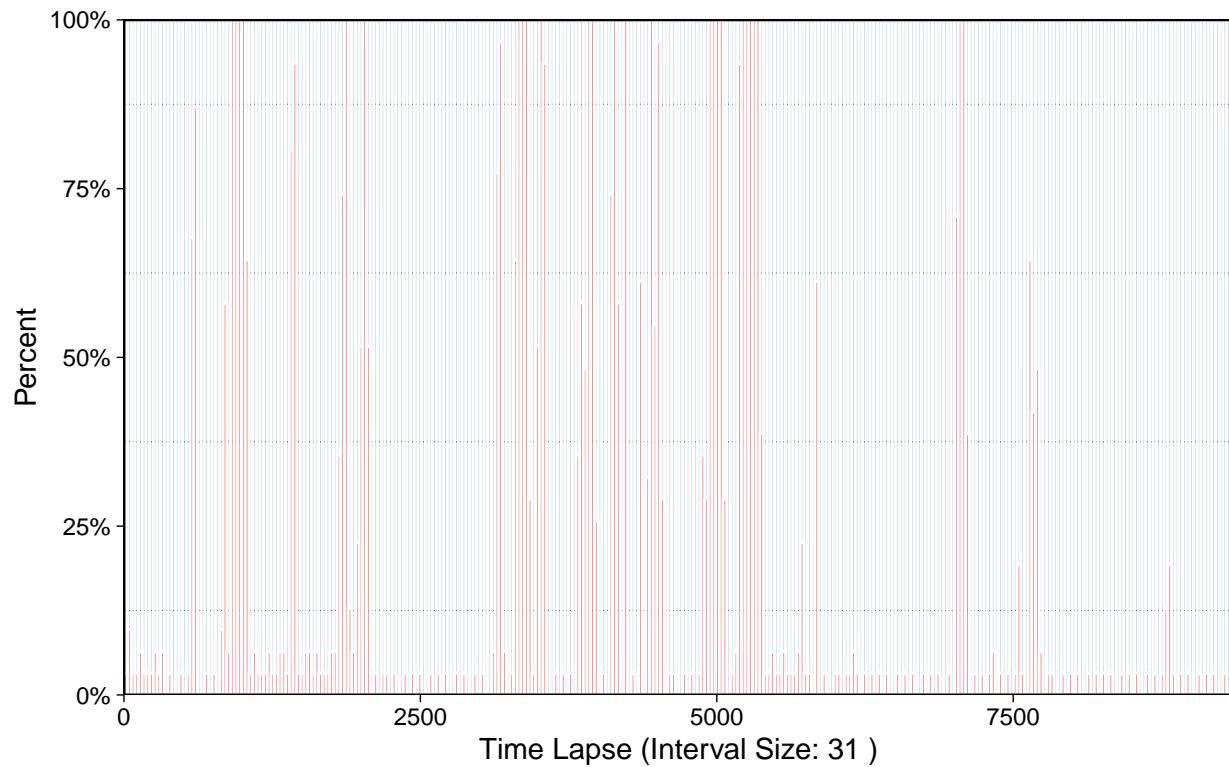
Amount of **NA** and **non-NA** for successive intervals



```
ggplot_na_distribution2(airqual_NA$CO.GT., interval_size = 31)
```

## Missing Values per Interval

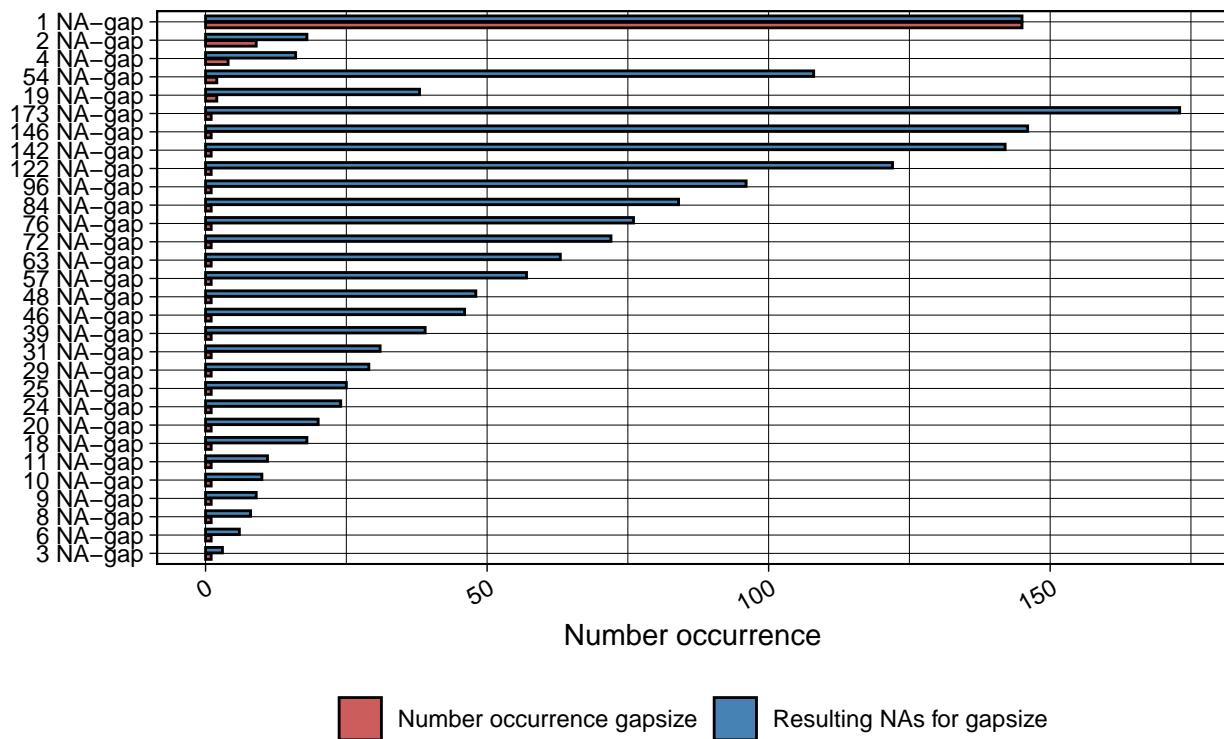
Amount of **NA** and **non-NA** for successive intervals



```
ggplot_na_gapsize(airqual_NA$CO.GT., limit = 31)
```

## Occurrence of gap sizes

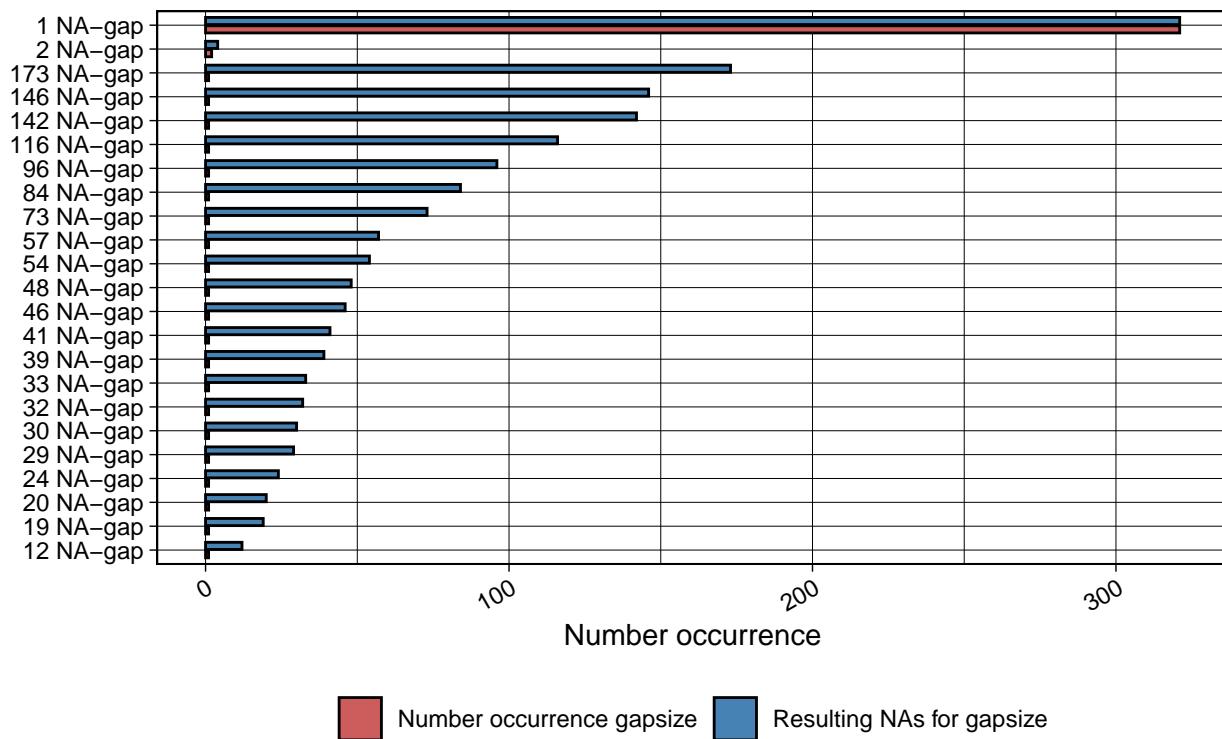
### Gap sizes (NAs in a row) ordered by most common



```
ggplot_na_gapsize(airqual_NA$NOx.GT., limit = 31)
```

## Occurrence of gap sizes

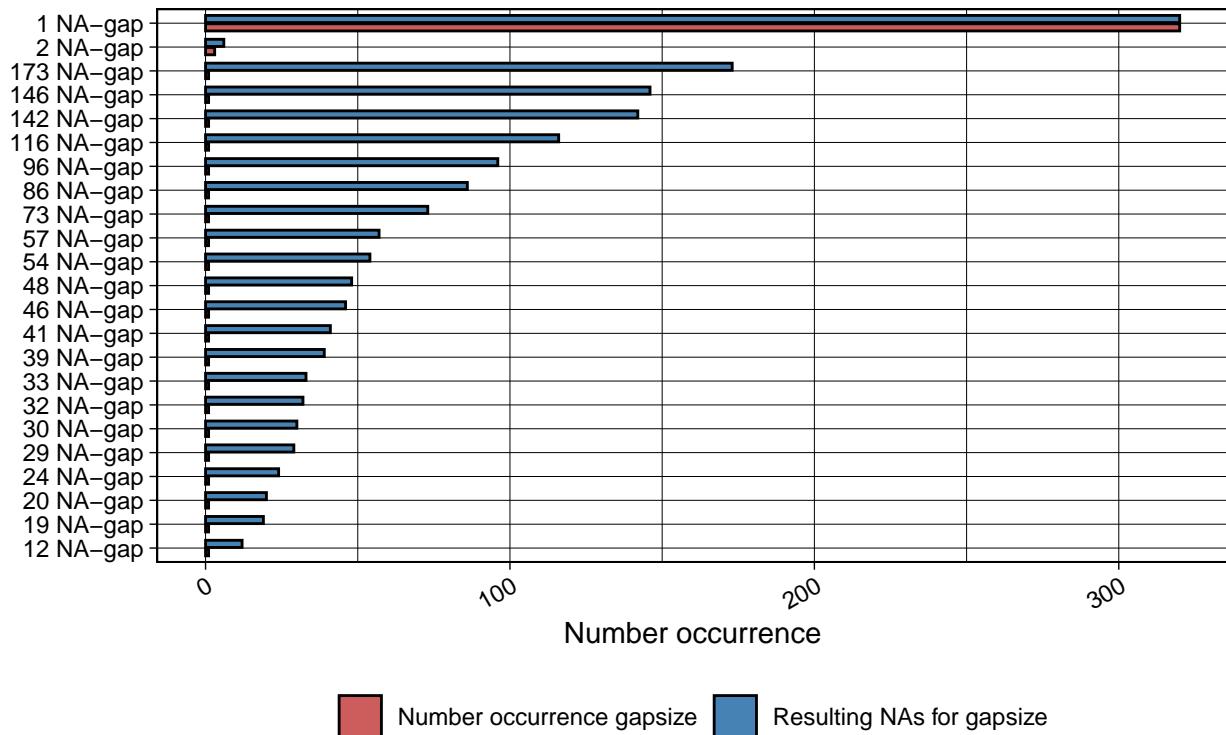
### Gap sizes (NAs in a row) ordered by most common



```
ggplot_na_gapsize(airqual_NA$NO2.GT., limit = 31)
```

## Occurrence of gap sizes

### Gap sizes (NAs in a row) ordered by most common



Can clearly see that there are a few wide gaps, but most of the missing values are distributed evenly, therefore impute values instead of removing the column or cutting of the date.

```

airqual_NA <- airqual
airqual_NA[airqual_NA == -200] = NA

missing_values <- colSums(is.na(airqual_NA))
print("Missing Values in Each Sensor:")

## [1] "Missing Values in Each Sensor:"
print(missing_values)

##      timedate      CO.GT.    PT08.S1.CO.      NMHC.GT.      C6H6.GT.
##          0        1683        366        8443        366
## PT08.S2.NMHC.      NOx.GT.    PT08.S3.NOx.      NO2.GT.    PT08.S4.NO2.
##         366        1639        366        1642        366
## PT08.S5.03.            T          RH          AH
##         366        366        366        366

# 2. Investigate the degree of missing values

# Visualizing the number of missing values over time

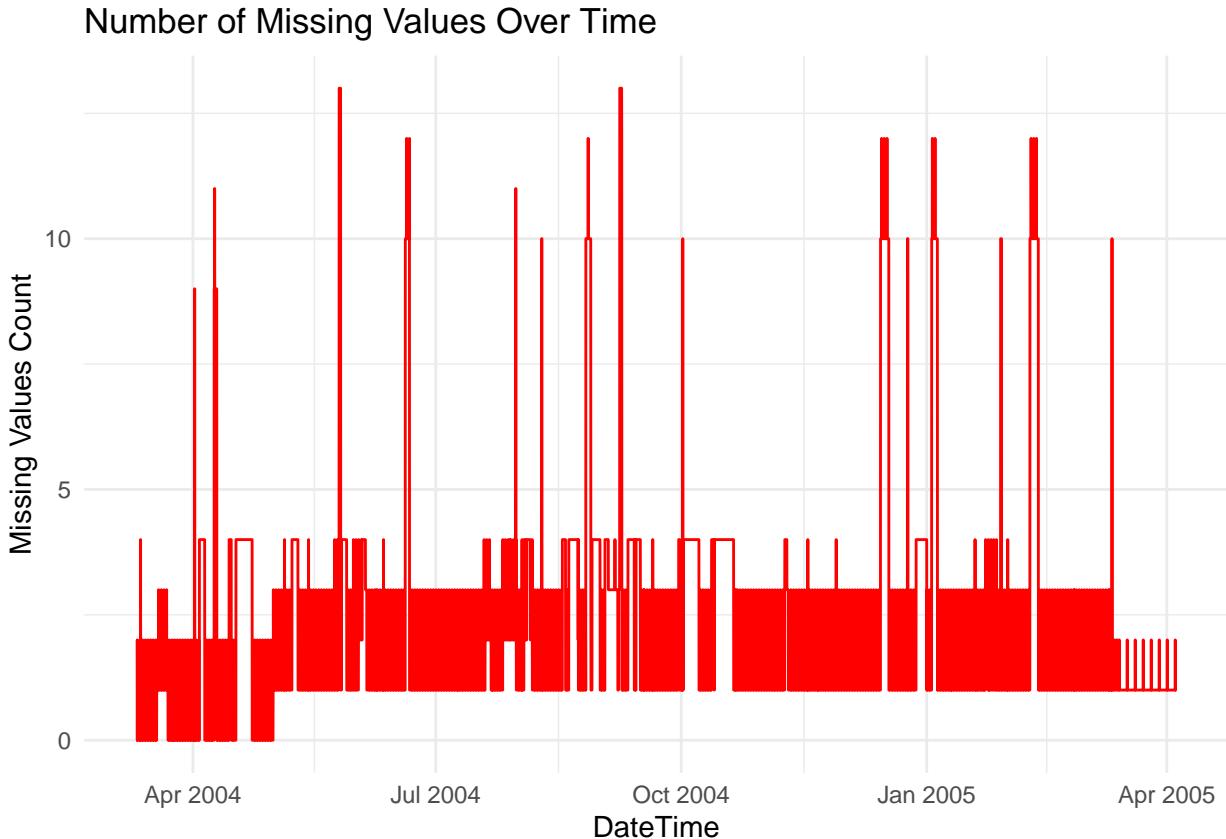
# Create a data frame to visualize missing values
missing_values <- airqual_NA %>%
  mutate(Missing = rowSums(is.na(airqual_NA))) %>%
  dplyr::select(timedate, Missing)

```

```

ggplot(missing_values, aes(x = timedate, y = Missing)) +
  geom_line(color = "red") +
  labs(title = "Number of Missing Values Over Time", x = "DateTime", y = "Missing Values Count") +
  theme_minimal()

```



- Is one or are multiple sensors behaving peculiarly? How would you handle this?
- NHMC.GT has almost all missing values after a certain time, as seen in the visualization, therefore remove to not skew the result away from real data.
- For the rest use imputations as there seems to be no big gaps with NA
- Discuss options for handling missing values: (a) drop all time points containing any missing value, (b) impute values for missing values. In case of (b) choose a method for imputation. Justify your decisions.
- **(a)** By dropping all NA we only get real data to do analysis on, but we miss out on number of data points so our overall accuracy becomes lower. In addition it is not possible as we need a regular time axis.
- **(b)** By imputing we can still use rows that are mostly NA free while having the model be skewed as little as possible by these values.
  - Choose to set the NA to the mean value of the overall data set, thereby making the data point the most probable if we were to get it randomly. By using rolling average we also get the most probable next step from the previous k steps, therefore making it more locally accurate than just using mean.
- At the end of this step, you should have a version of the data containing only valid values. Plot these data as in Part B.

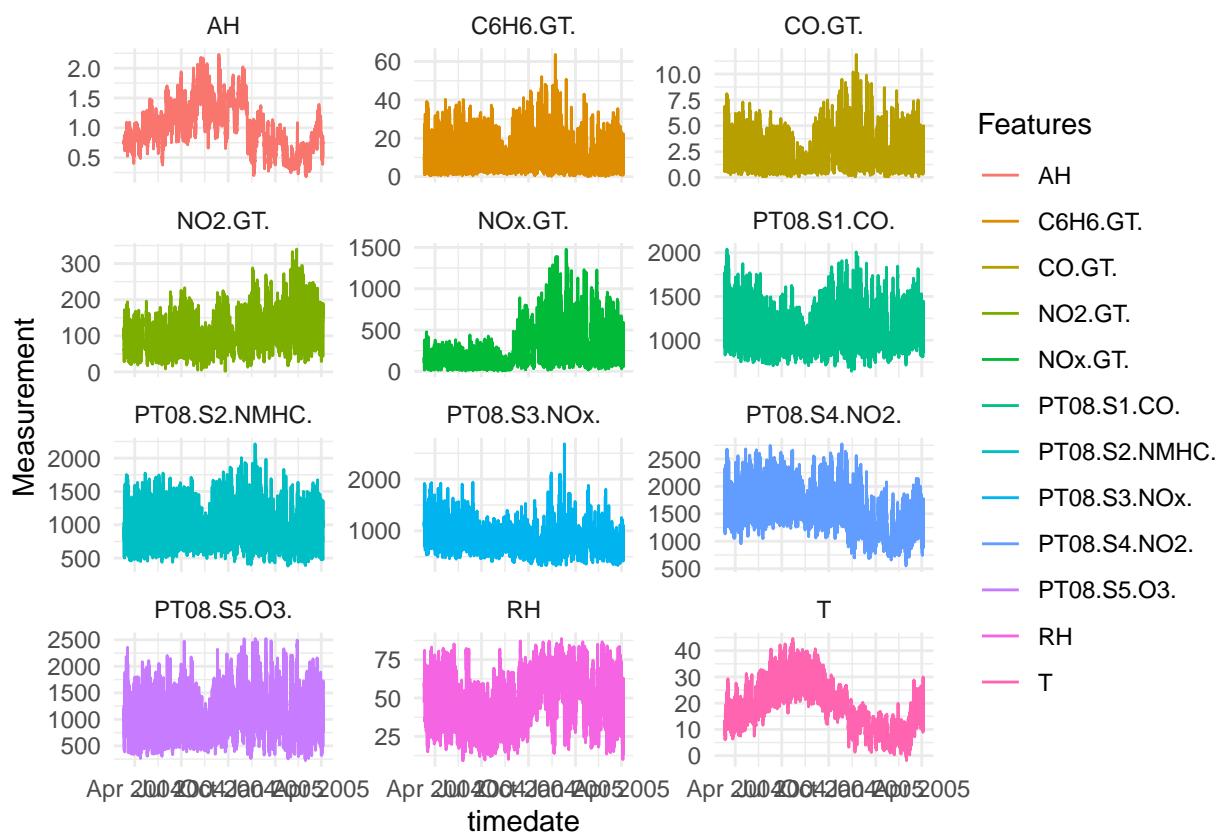
```

airqual_c <- airqual %>%
  dplyr::select(-c("NMHC.GT."))
  
airqual_c[airqual_c == -200] = NA

airqual_c <- imputeTS::na_ma(airqual_c, k = 4, weighting = "simple")

airqual_c %>%
  pivot_longer(!timedate, names_to="Features", values_to="Measurement") %>%
  ggplot(aes(x = timedate, y = Measurement, col = Features)) +
  geom_line() +
  theme_minimal() +
  facet_wrap(~ Features, scales = "free_y", ncol = 3)

```



Moving average: assumes that the current value ( $y_t$ ) is dependent on the error terms including the current error, non linear regression model, order ACF Auto regressive: assumes that the current value ( $y_t$ ) is dependent on previous values, linear regression model, order PACF

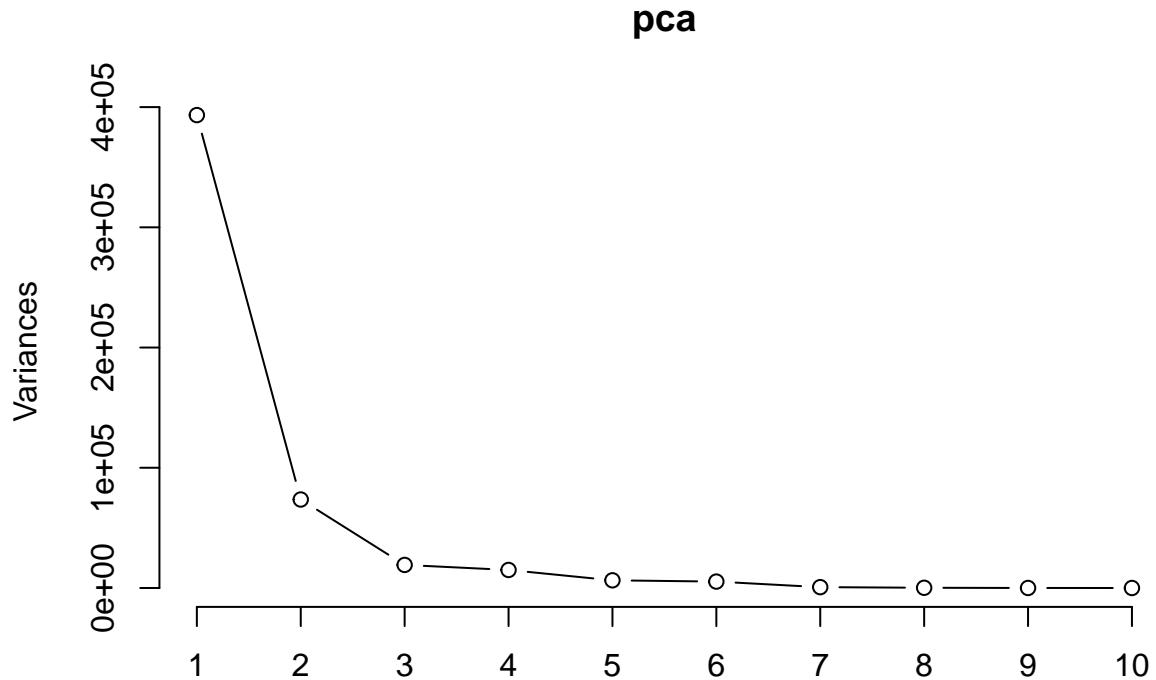
## Part E: PCA of cleaned data

- Perform PCA on the data as prepared in D

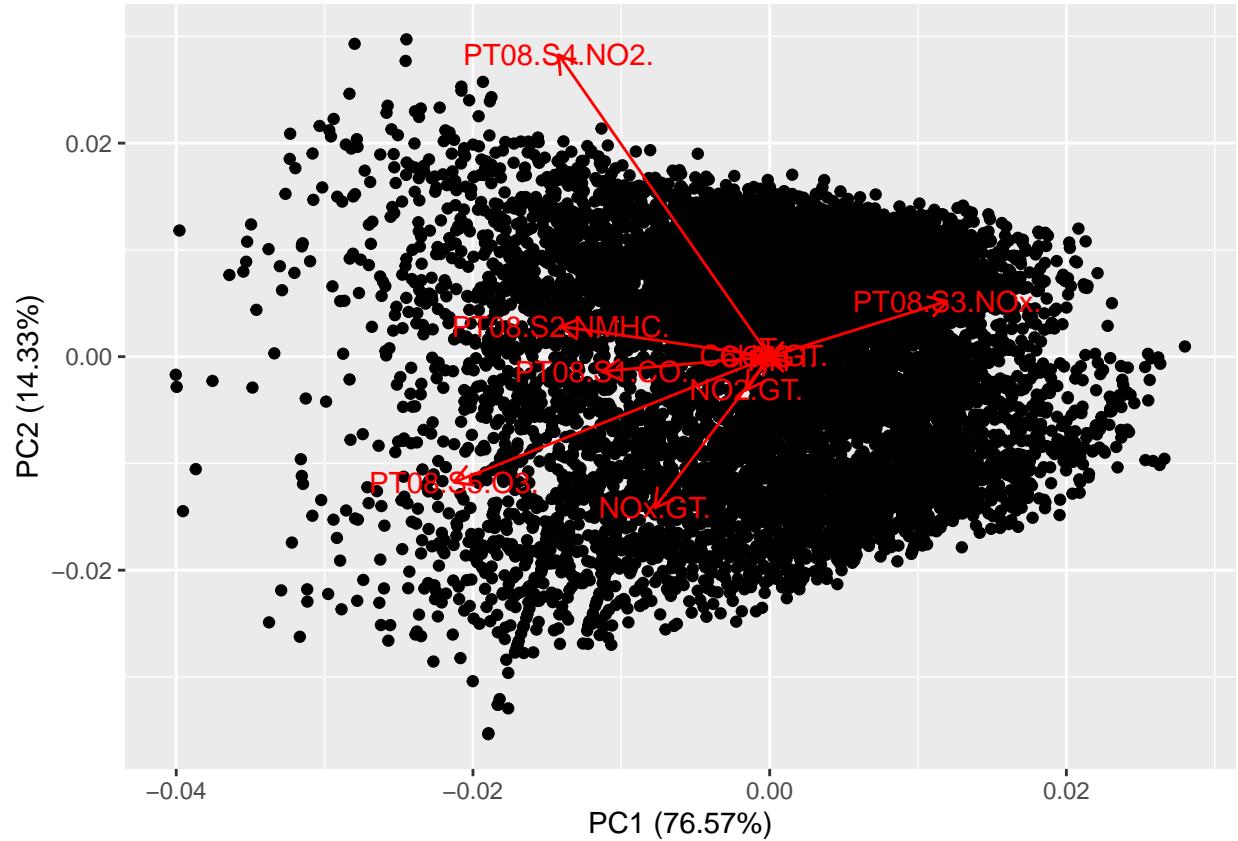
```
pca <- prcomp(airqual_c[,-1]) #, center = TRUE, scale. = TRUE
```

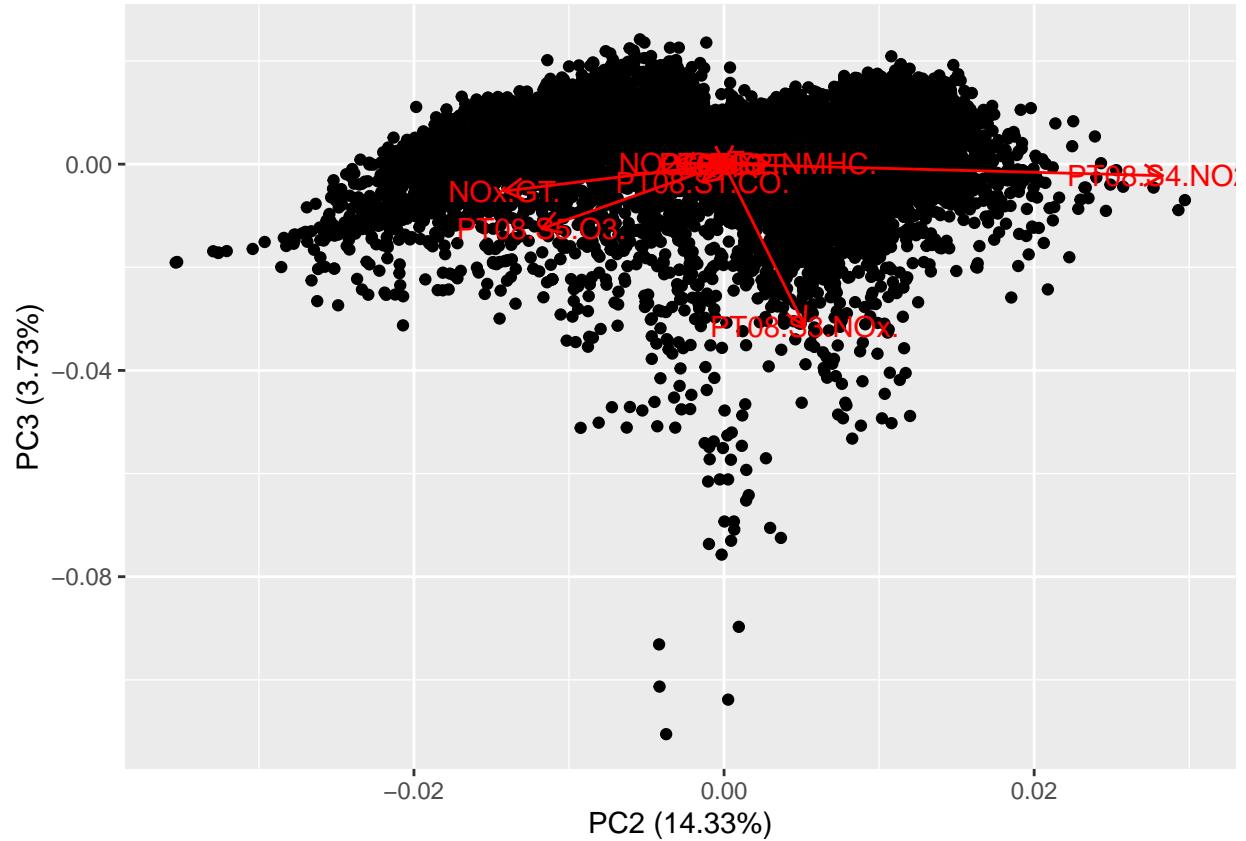
- Create a screeplot and biplots for 1st/2nd, 2nd/3rd, 3rd/4th PC

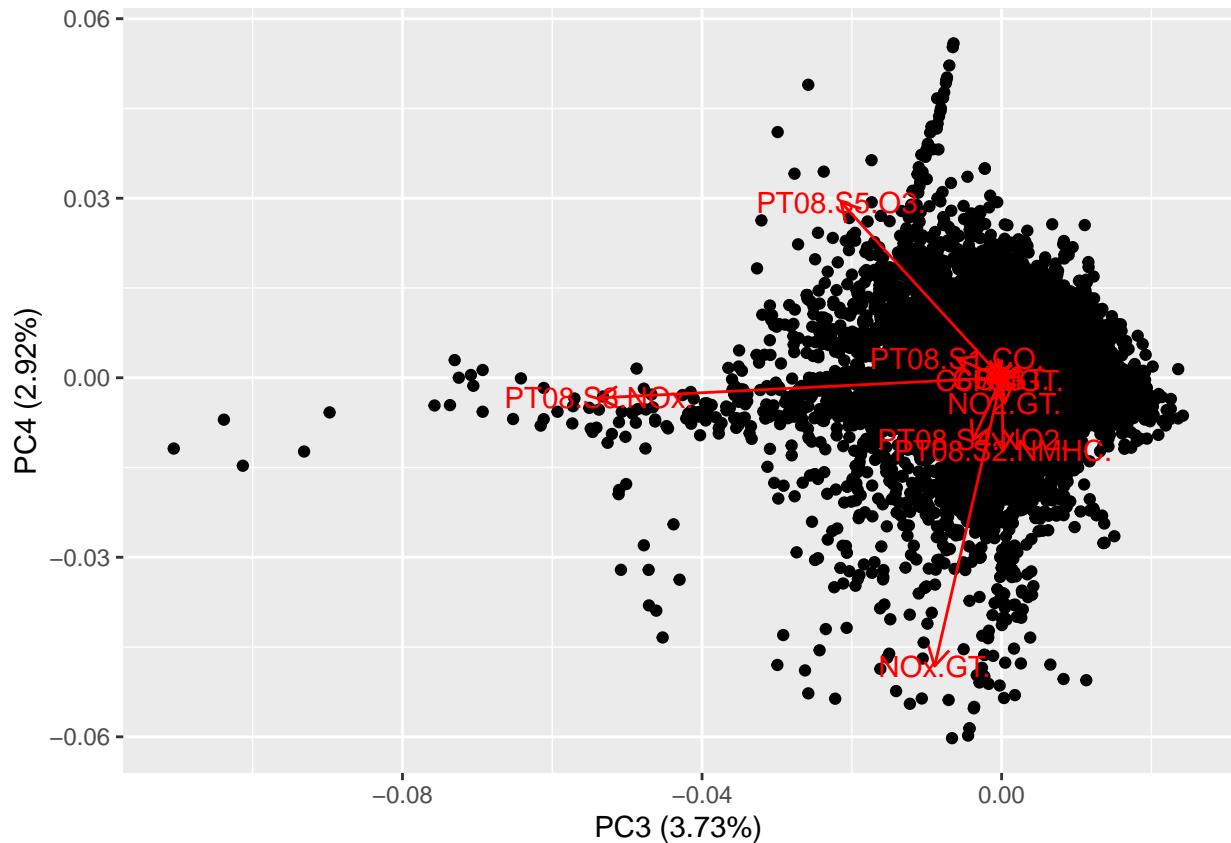
```
# screeplot
pca <- prcomp(airqual_c[,-1])
plot(pca, type = "l")
```



```
# biplots
autoplot(pca, x=1,y=2, loadings=TRUE, loadings.label=TRUE)
```







Biplot red arrow original axis of data

- Compute total variance explained by 1st, 1st and 2nd, 1st to 3rd, ... PCs

```
#Variance explained
summary(pca)$importance[3,]
```

```
##      PC1      PC2      PC3      PC4      PC5      PC6      PC7      PC8      PC9      PC10
## 0.76570 0.90897 0.94631 0.97548 0.98791 0.99826 0.99964 0.99997 1.00000 1.00000
##      PC11     PC12
## 1.00000 1.00000
```

- Choose how many PCs to keep and transform data back to original sample space

Keep to 0.99 mark, so PC1 to PC6

```
t <- airqual_c$timedate
x_1 <- pca$x[, 1:6] %*% t(pca$rotation[, 1:6])
x_2 <- t(pca$center + pca$scale * t(x_1))
```

$$S = XL \Leftrightarrow SL^{-1} = XLL^{-1} \Leftrightarrow X = SL^{-1} = SL^T$$

S : Scores L : Loadings X : Original matrix

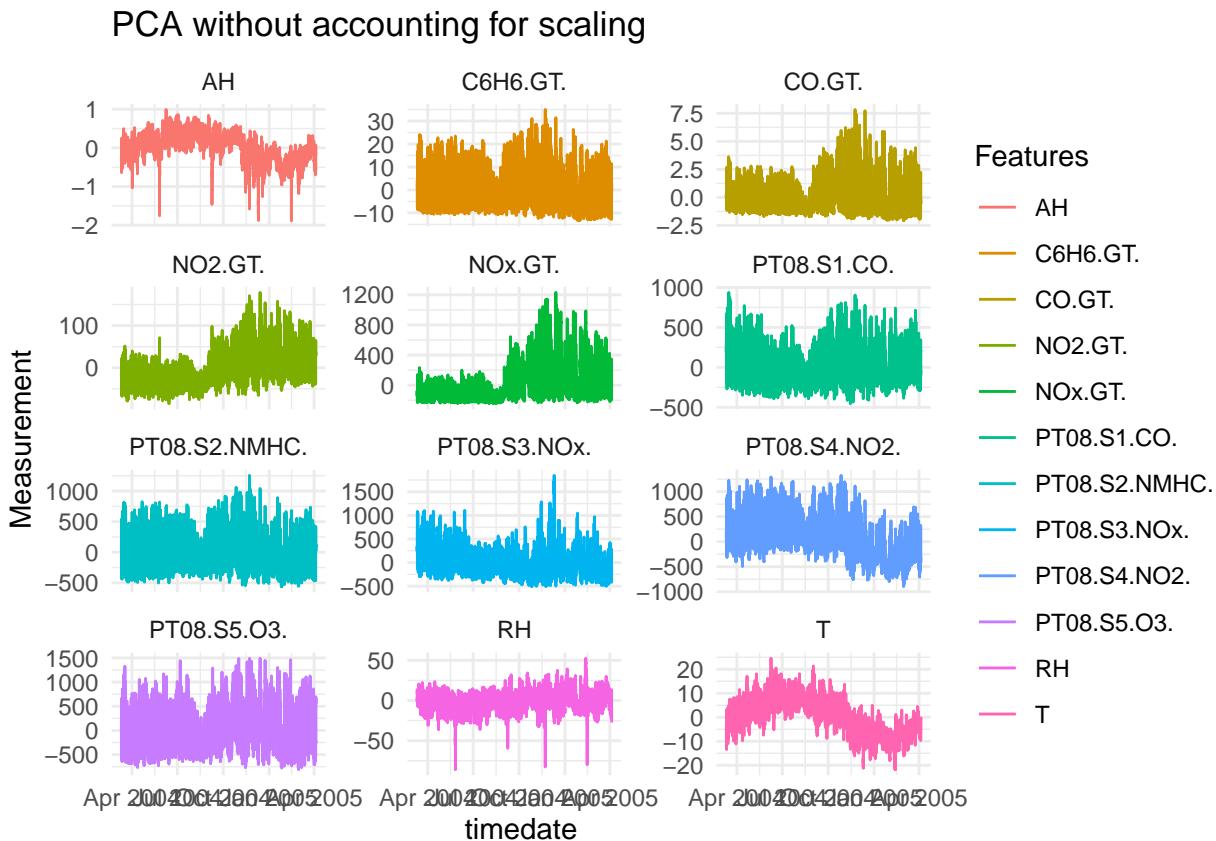
- Plot the result against the cleaned data, compare and discuss

```
data.frame(timedate = airqual_c$timedate, x_1) %>%
  pivot_longer(!timedate, names_to="Features", values_to="Measurement") %>%
  ggplot(aes(x = timedate, y = Measurement, col = Features)) +
  geom_line() +
```

```

theme_minimal() +
facet_wrap(~ Features, scales = "free_y", ncol = 3) +
ggtitle("PCA without accounting for scaling")

```



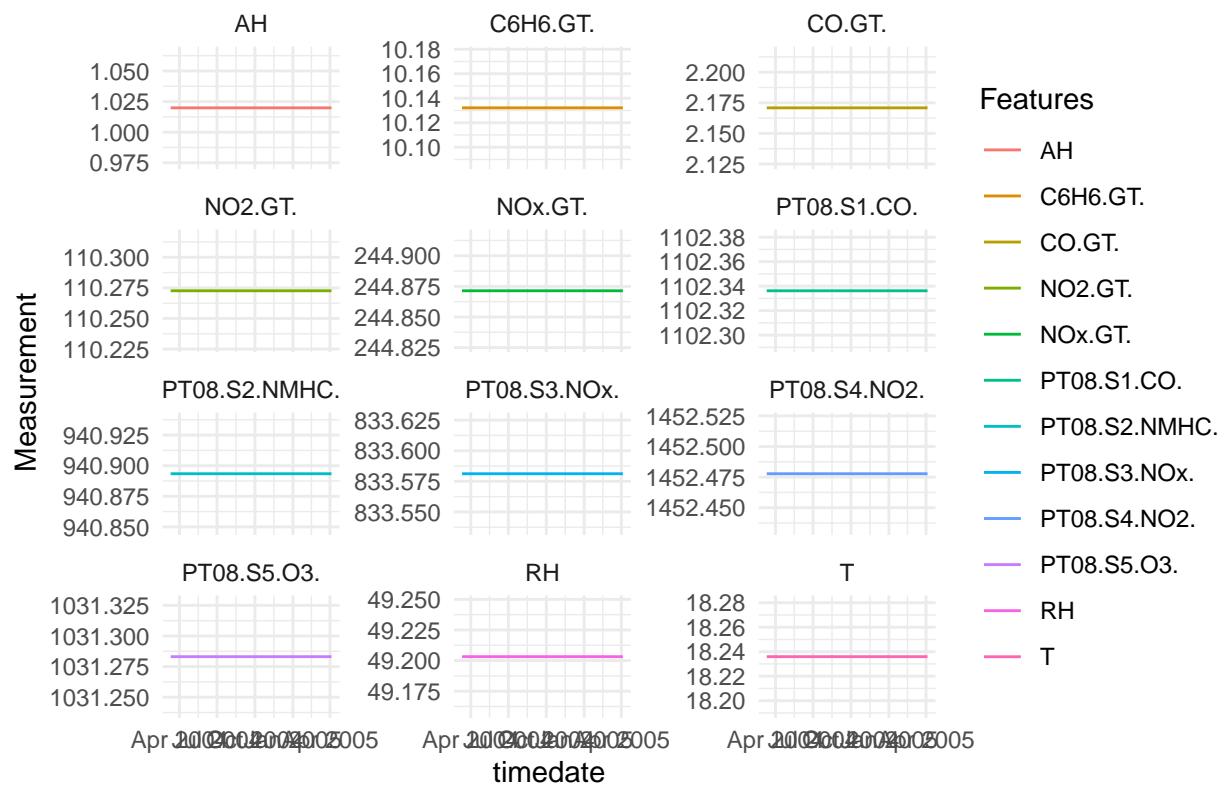
Can see certain dips into negative value, especially in RH and AH. Might be caused by the missing PC that accounted for that variability.

```

data.frame(timedate = airqual_c$timedate, x_2) %>%
pivot_longer(!timedate, names_to="Features", values_to="Measurement") %>%
ggplot(aes(x = timedate, y = Measurement, col = Features)) +
geom_line() +
theme_minimal() +
facet_wrap(~ Features, scales = "free_y", ncol = 3) +
ggtitle("PCA accounting for scaling")

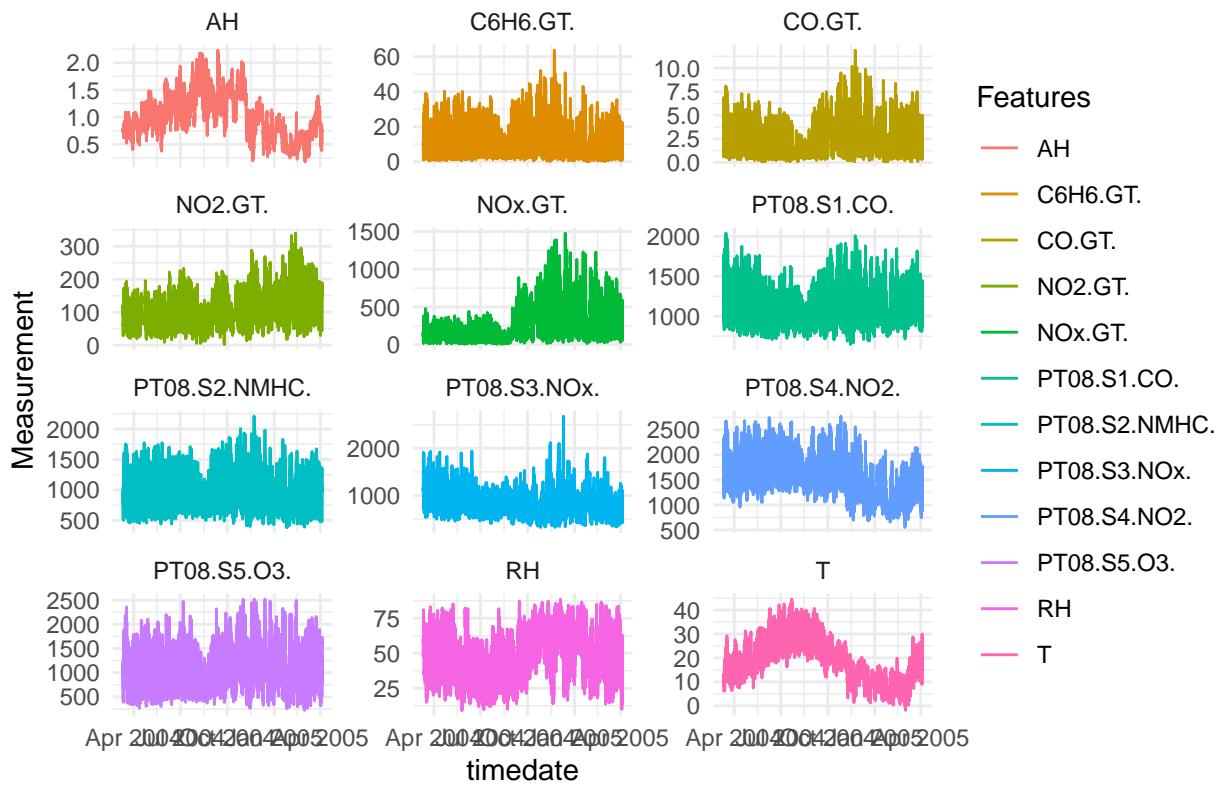
```

## PCA accounting for scaling



```
airqual_c %>%
  pivot_longer(!timedate, names_to="Features", values_to="Measurement") %>%
  ggplot(aes(x = timedate, y = Measurement, col = Features)) +
  geom_line() +
  theme_minimal() +
  facet_wrap(~ Features, scales = "free_y", ncol = 3) +
  ggtitle("Original data")
```

## Original data

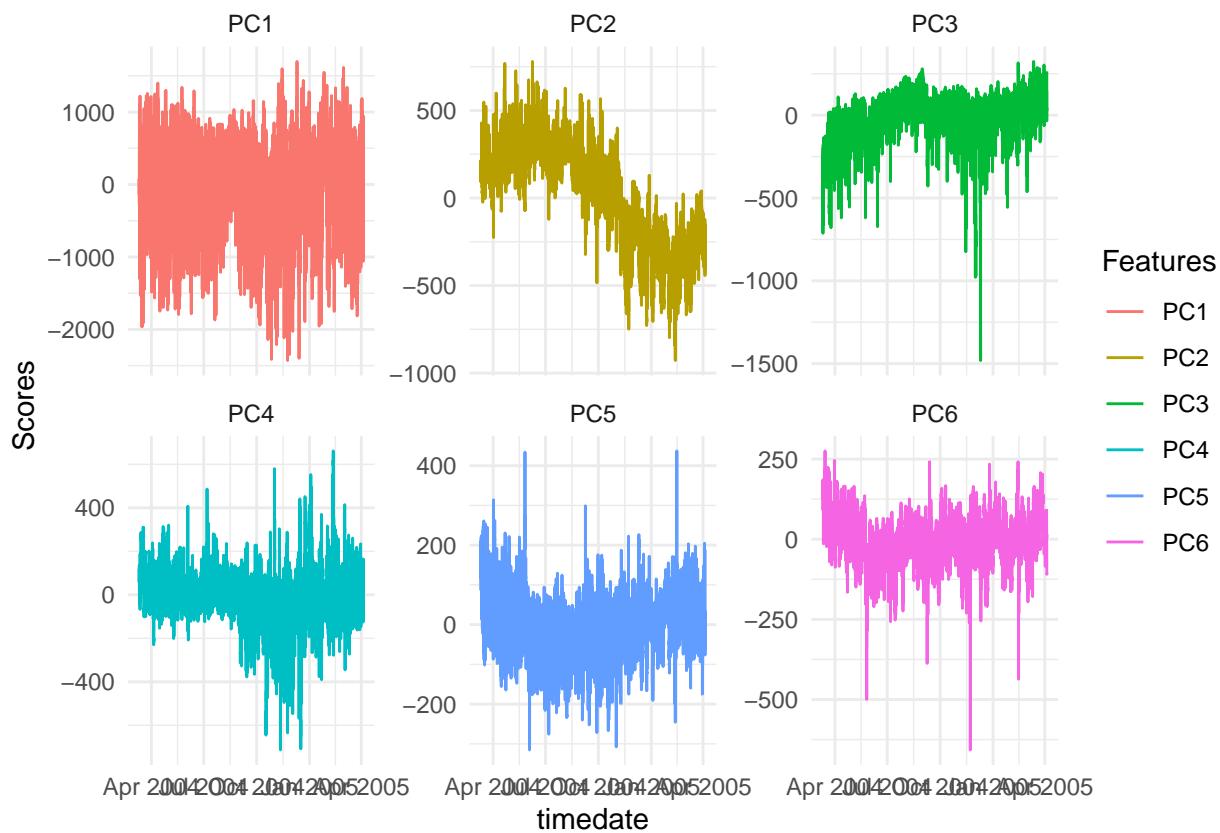


Looks mostly the same, but deviates slightly as we removed all after PC7.

- Also plot the scores, zoom in to short time intervals and look at periodicity

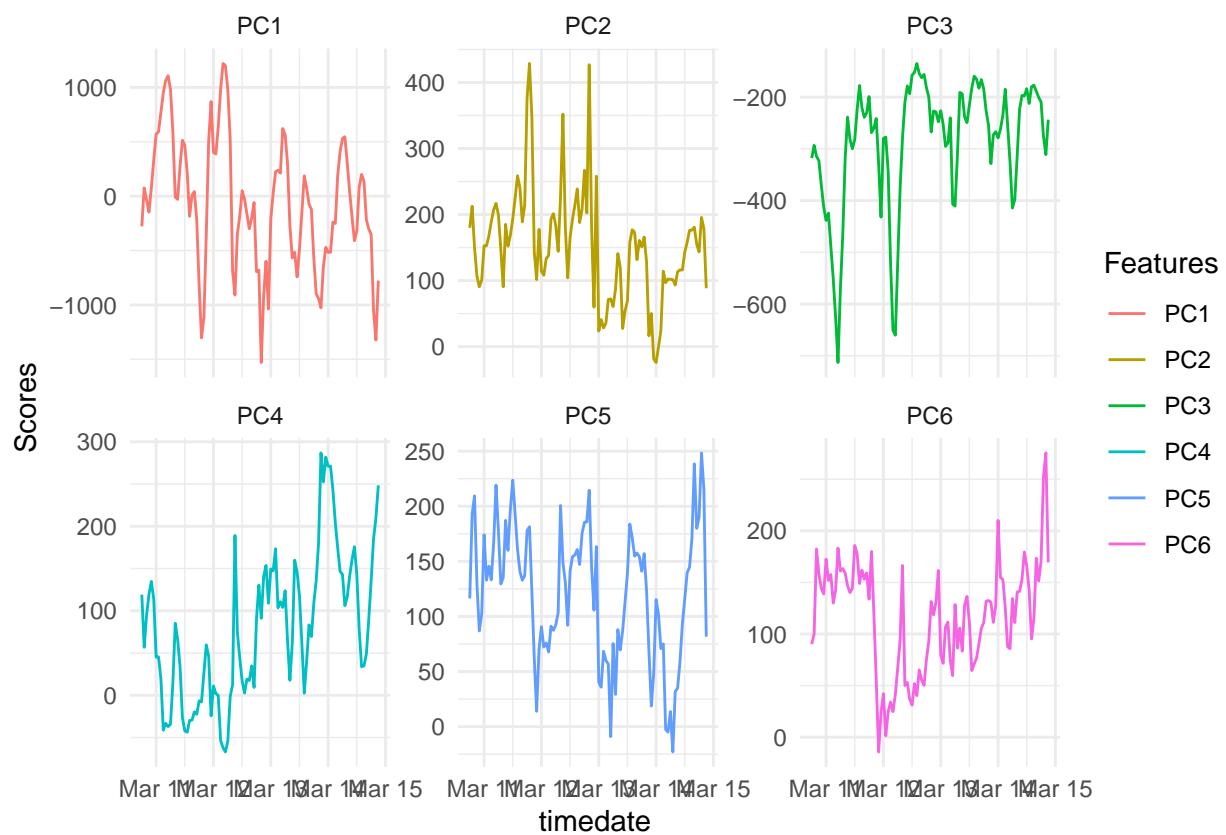
## Whole

```
data.frame(timedate = airqual_c$timedate, pca$x[,1:6]) %>%
  pivot_longer(!timedate, names_to="Features", values_to="Scores") %>%
  ggplot(aes(x = timedate, y = Scores, col = Features)) +
  geom_line() +
  theme_minimal() +
  facet_wrap(~ Features, scales = "free_y", ncol = 3)
```



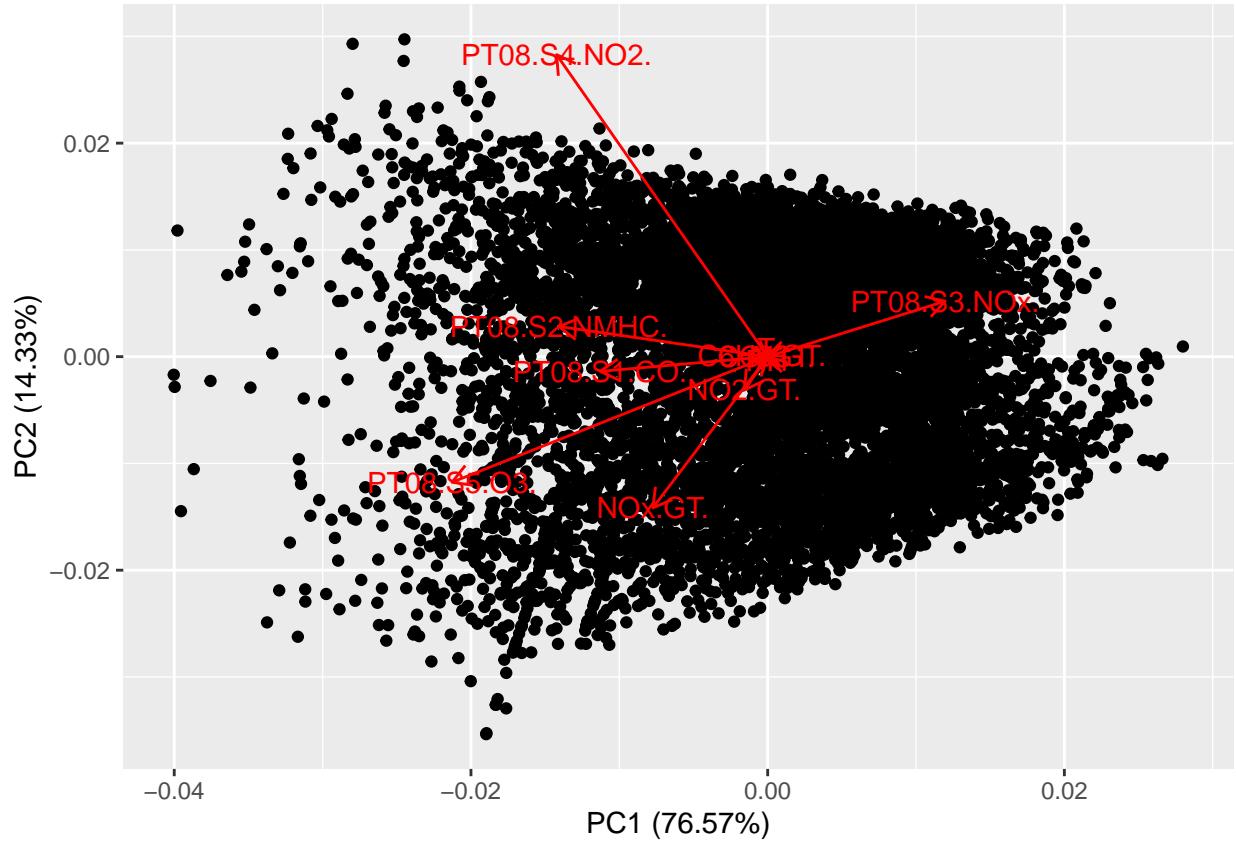
### Zoomed

```
data.frame(timedate = airqual_c$timedate[1:100], pca$x[1:100,1:6]) %>%
  pivot_longer(!timedate, names_to="Features", values_to="Scores") %>%
  ggplot(aes(x = timedate, y = Scores, col = Features)) +
  geom_line() +
  theme_minimal() +
  facet_wrap(~ Features, scales = "free_y", ncol = 3)
```



- Can you interpret certain PCs?

```
autoplot(pca, x=1,y=2, loadings=TRUE, loadings.label=TRUE)
```



```
data.frame(pca$rotation)
```

	PC1	PC2	PC3	PC4
## CO.GT.	-0.0018481988	-0.0007049452	-0.0014077536	-0.0043626668
## PT08.S1.CO.	-0.3276959483	-0.0400440912	-0.1009593770	0.0537630030
## C6H6.GT.	-0.0111497193	0.0023373054	-0.0045012384	-0.0078980451
## PT08.S2.NMHC.	-0.4085995243	0.0821947866	0.0037025890	-0.2073677556
## NOx.GT.	-0.2247579859	-0.4110090314	-0.1516711696	-0.8168812736
## PT08.S3.NOx.	0.3475027920	0.1508453517	-0.9111126863	-0.0570822819
## NO2.GT.	-0.0452313753	-0.0903548818	0.0065398011	-0.0728416677
## PT08.S4.NO2.	-0.4131515654	0.8212890426	-0.0647945566	-0.1760872468
## PT08.S5.O3.	-0.6176661288	-0.3415204901	-0.3633580682	0.4970953230
## T	-0.0021215778	0.0223286223	0.0146078597	-0.0068587449
## RH	-0.0019124088	-0.0116763960	-0.0125220967	-0.0031264697
## AH	-0.0001449796	0.0009679888	0.0006321034	-0.0003673498
	PC5	PC6	PC7	PC8
## CO.GT.	0.0002686689	0.000757493	-0.005358532	6.260908e-06
## PT08.S1.CO.	0.6259430252	0.695331376	0.014987171	-4.858268e-02
## C6H6.GT.	0.0163460121	-0.017403082	0.013341054	7.606726e-03
## PT08.S2.NMHC.	0.5577650907	-0.669313381	0.114878545	9.781349e-02
## NOx.GT.	-0.2518500677	0.132534035	0.073466372	-6.084502e-02
## PT08.S3.NOx.	0.1312507403	-0.072666427	-0.017993617	9.409716e-03
## NO2.GT.	0.0633409652	-0.041565772	-0.965138888	2.139533e-01
## PT08.S4.NO2.	-0.3082878182	0.135577313	-0.073581584	-8.415517e-03
## PT08.S5.O3.	-0.3340438218	-0.102865230	-0.001629221	-2.331928e-02
## T	-0.0070602392	-0.022954286	-0.040469919	-3.383485e-01

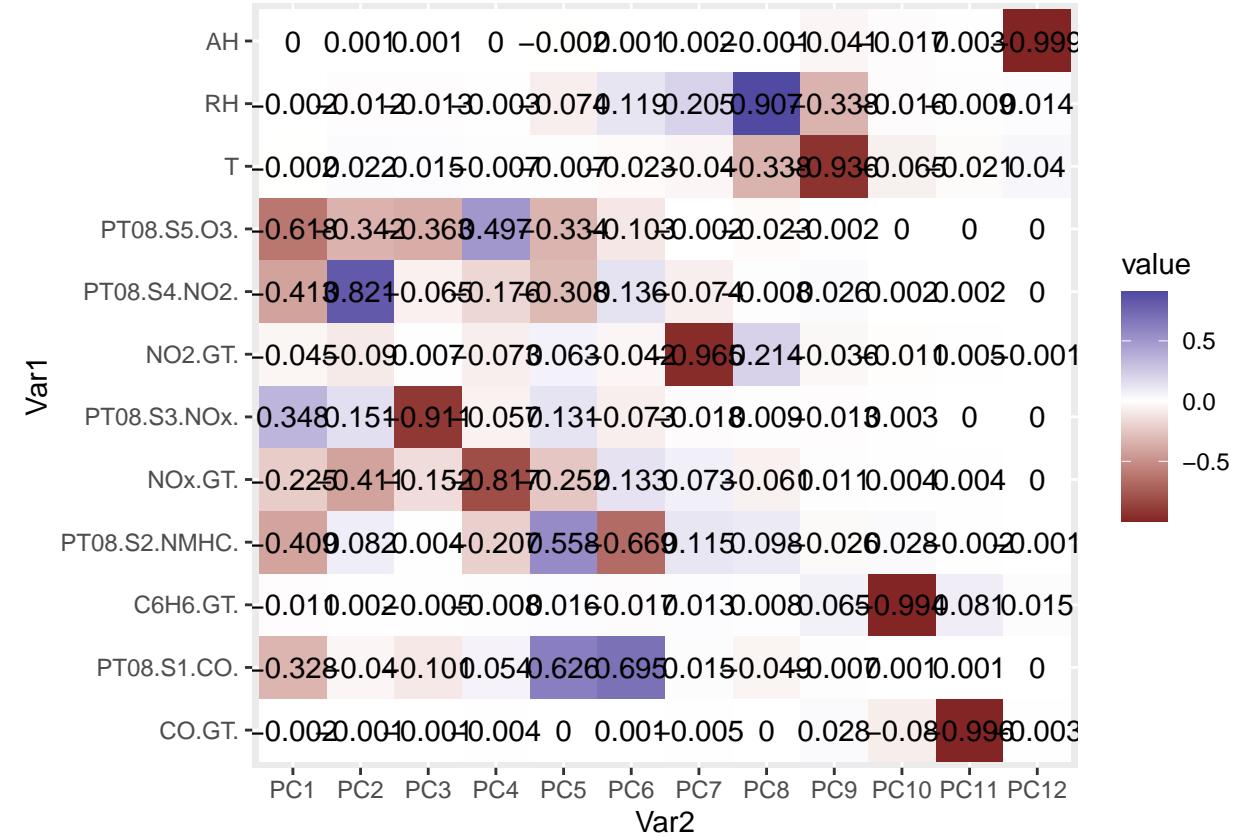
```

## RH          -0.0740118957  0.119022190  0.205178331  9.073875e-01
## AH          -0.0021484402  0.001182079  0.002142285 -1.130502e-03
##             PC9          PC10          PC11          PC12
## CO.GT.      0.028421694 -0.0798924563 -0.9963659489 -3.164062e-03
## PT08.S1.CO. -0.007152100  0.0011586314  0.0008637521 -2.323526e-04
## C6H6.GT.    0.065223128 -0.9939169941  0.0814893649  1.470071e-02
## PT08.S2.NMHC. -0.025579860  0.0277063380 -0.0023210316 -1.066060e-03
## NOx.GT.     0.011434225  0.0038849340  0.0041500425  2.374803e-04
## PT08.S3.NOx. -0.013083805  0.0034508345  0.0002142262 -3.947865e-04
## NO2.GT.     -0.035903855 -0.0106244017  0.0054654120 -8.600059e-04
## PT08.S4.NO2.  0.025888303  0.0016380067  0.0020691370  4.622803e-04
## PT08.S5.O3.  -0.002484763 -0.0002733566 -0.0004846717  7.147926e-05
## T            -0.936057264 -0.0653693331 -0.0213934536  4.010257e-02
## RH          -0.337718440 -0.0162504714 -0.0093587511  1.391709e-02
## AH          -0.041331962 -0.0172440281  0.0033651508 -9.989843e-01

melt(pca$rotation) %>%
  ggplot(aes(Var2, Var1)) +
  geom_tile(aes(fill = value)) +
  geom_text(aes(fill = value, label = round(value, 3))) +
  scale_fill_gradient2()

```

## Warning in geom\_text(aes(fill = value, label = round(value, 3))): Ignoring  
## unknown aesthetics: fill



By looking at the loadings of PC1 one can see that PT08.S5.O3 is the most important feature with respect to explainable variance. Scores determine which weight each PC enter into each row of X

## Task 2: STL and correlation on weather data

### Part A: Data collection for a single station

Based on material from the lectures, write an R function that can obtain a daily average temperature series for a meteorological station from the Norwegian Met Institute's Frost service. The function shall return a tibble.

```
.client_id <- str_trim(read_file("client_id.txt"))

# Server to collect data from and resource we want from server
server <- "frost.met.no"
resource <- "observations/v0.jsonld"

# Station(s) we want data for. SN17850 is the station ID for Ås (Blindern is SN18700)
sources <- 'SN17850'

# Type of data we want, P1D means daily data
elements <- 'mean(air_temperature P1D)'

# Time range we want data for
reference_time <- '1874-01-01/2023-12-31'

# Specify that we want mean temperature calculated from midnight to midnight
timeoffsets <- 'PTOH'

.query_url <- str_glue("https://.{client_id}@{server}/{resource}?sources={sources}&referencetime={reference_time}&elements={elements}&timeoffsets={timeoffsets}&get_data_from_frost={get_data_from_frost}")

# Set this to TRUE to generate a json file
get_data_from_frost = TRUE
weather_file = "weather_data_json.rds.bz2"
if ( get_data_from_frost ) {
  raw_data <- try(fromJSON(URLencode(.query_url), flatten=TRUE))

  if ( class(raw_data) != 'try-error' ) {
    print("Data retrieved from frost.met.no!")
    write_rds(raw_data, weather_file, compress="bz2", text=TRUE) # JSON represents data as text
    print(str_glue("Raw data (JSON) written to '{weather_file}'"))
  } else {
    print("Error: the data retrieval was not successful!")
  }
} else {
  raw_data <- read_rds(weather_file)
  print(str_glue("Raw data (JSON) read from '{weather_file}'"))
}

## [1] "Data retrieved from frost.met.no!"
## Raw data (JSON) written to 'weather_data_json.rds.bz2'
df <- unnest(raw_data$data, cols = c(observations))

head(df)

## # A tibble: 6 x 14
##   sourceId referenceTime      elementId value unit  timeOffset timeResolution
##   <chr>     <chr>          <chr>     <dbl> <chr> <chr>       <chr>
```

```

## 1 SN17850:0 1874-01-01T00:00:00~ mean(air~ 3.2 degC PTOH P1D
## 2 SN17850:0 1874-01-02T00:00:00~ mean(air~ 4.1 degC PTOH P1D
## 3 SN17850:0 1874-01-03T00:00:00~ mean(air~ 2.5 degC PTOH P1D
## 4 SN17850:0 1874-01-04T00:00:00~ mean(air~ 0.3 degC PTOH P1D
## 5 SN17850:0 1874-01-05T00:00:00~ mean(air~ -3.9 degC PTOH P1D
## 6 SN17850:0 1874-01-06T00:00:00~ mean(air~ 1.9 degC PTOH P1D
## # i 7 more variables: timeSeriesId <int>, performanceCategory <chr>,
## # exposureCategory <chr>, qualityCode <int>, level.levelType <chr>,
## # level.unit <chr>, level.value <int>
df |> dplyr::select(referenceTime, value) |>
  mutate(referenceTime=as.Date(referenceTime)) |>
  rename(Date=referenceTime, Temp=value) |>
  as_tsibble(index = Date) -> dc
head(dc)

## # A tsibble: 6 x 2 [1D]
##   Date      Temp
##   <date>    <dbl>
## 1 1874-01-01  3.2
## 2 1874-01-02  4.1
## 3 1874-01-03  2.5
## 4 1874-01-04  0.3
## 5 1874-01-05 -3.9
## 6 1874-01-06  1.9

```

## Part B: Data preparation for a single station

- Identify gaps in the time series.

```
has_gaps(dc)
```

```

## # A tibble: 1 x 1
##   .gaps
##   <lg1>
## 1 TRUE

gaps <- count_gaps(dc)
head(gaps)

## # A tibble: 6 x 3
##   .from      .to       .n
##   <date>    <date>   <int>
## 1 1874-08-01 1874-08-31     31
## 2 1876-01-01 1877-07-31    578
## 3 1879-01-01 1879-01-31     31
## 4 1880-10-31 1880-10-31      1
## 5 1881-05-01 1881-05-31     31
## 6 1900-12-31 1900-12-31      1

```

- Assume that gaps up to 31 days are acceptable. Find the earliest date in the time series such that all following data have no gaps longer than 31 days. Limit the time series to this.

```

# cutoff_date <- as.character(tail(gaps[gaps$.n >= 31,], n=1)$to)
# cutoff_date
if (nrow(gaps) != 0) {
  cutoff_date <- as.character(tail(gaps[gaps$.n >= 31,], n=1)$to)
}

```

```

dc_cut <- dc %>% tsibble::filter_index(cutoff_date ~ .)
} else {
  dc_cut <- dc
}

```

identify that all dates after, even if 31 is included, is 1988-06-17

```

# dc_cut <- dc %>% tsibble::filter_index(cutoff_date ~ .)
head(dc_cut)

```

```

## # A tsibble: 6 x 2 [1D]
##   Date      Temp
##   <date>    <dbl>
## 1 1988-06-18  18.7
## 2 1988-06-21  19.7
## 3 1988-06-22  18
## 4 1988-06-23  21.2
## 5 1988-06-25  21.3
## 6 1988-06-26  23

```

- Create a regular time series by filling gaps in the tsibble with n/a-s.

```

dc_cut_filled <- tsibble::fill_gaps(dc_cut, .full = TRUE)

head(dc_cut_filled)

```

```

## # A tsibble: 6 x 2 [1D]
##   Date      Temp
##   <date>    <dbl>
## 1 1988-06-18  18.7
## 2 1988-06-19  NA
## 3 1988-06-20  NA
## 4 1988-06-21  19.7
## 5 1988-06-22  18
## 6 1988-06-23  21.2

```

- Impute values for the n/a-s. Justify your choice of imputation method.

Choose to impute using simple rolling average, as you get the advantage of getting the most probable next step as well as being locally probable as opposed to using mean or median for the whole dataset.

```

dc_cut_filled_imp <- imputeTS::na_ma(dc_cut_filled, k = 4, weighting = "simple")

```

- You should now have a regular time series with only numeric values.

```

tsibble::has_gaps(dc_cut_filled_imp)

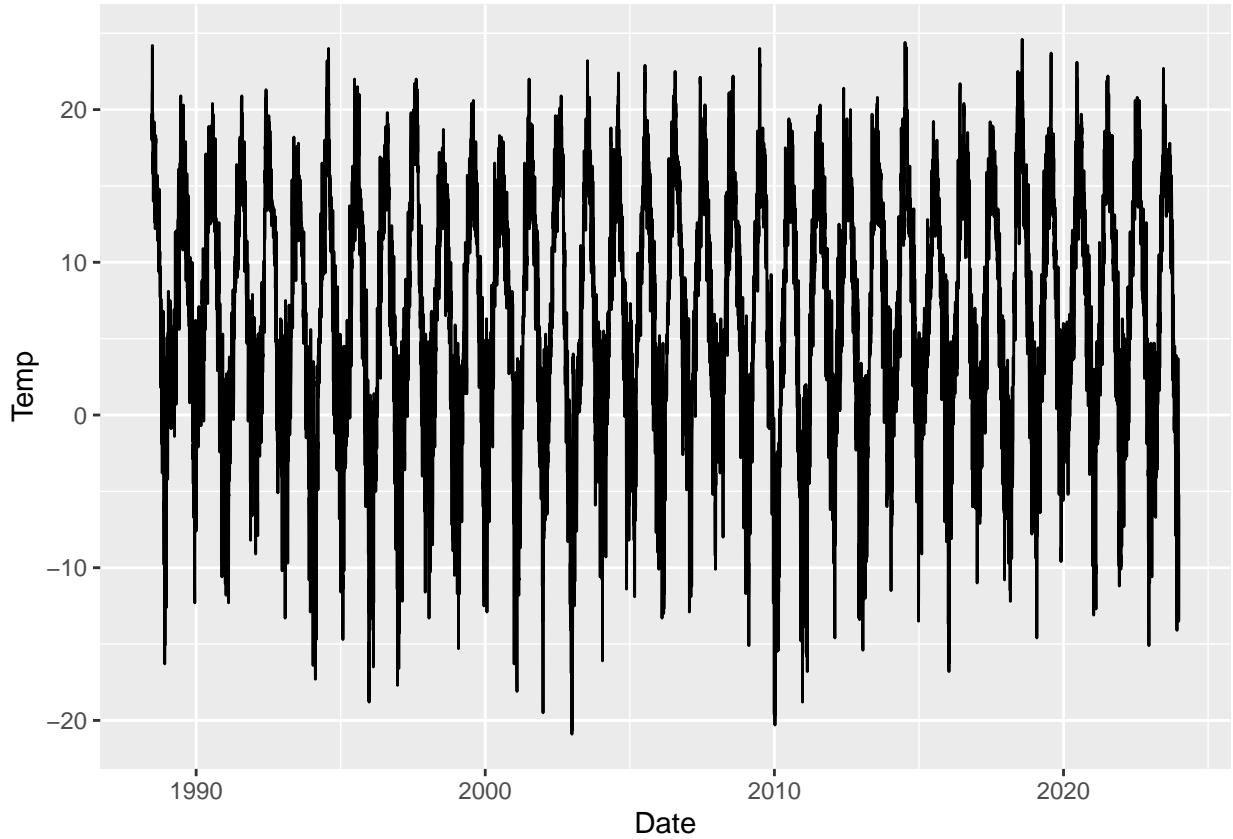
```

```

## # A tibble: 1 x 1
##   .gaps
##   <lgl>
## 1 FALSE
is.numeric(dc_cut_filled_imp$Temp)

## [1] TRUE
dc_cut_filled_imp %>%
  ggplot(aes(Date, Temp)) +
  geom_line()

```



- Remove all data for 29 February so all years have data for exactly 365 days.

```
dc_cut_filled_imp_noLeap <- dc_cut_filled_imp %>%
  dplyr::filter(!(month(Date) == 2 & day(Date) == 29))
head(dc_cut_filled_imp_noLeap)
```

```
## # A tsibble: 6 x 2 [1D]
##   Date      Temp
##   <date>    <dbl>
## 1 1988-06-18 18.7
## 2 1988-06-19 19.4
## 3 1988-06-20 19.4
## 4 1988-06-21 19.7
## 5 1988-06-22 18
## 6 1988-06-23 21.2
```

- Combine all this code into a function for re-use later. The function should receive the original tsibble from part A as input and return a new tsibble.

```
timeseries_cleaner <- function(table) {
  gaps <- count_gaps(table)
  gaps_31 <- gaps[gaps$n >= 31,]
  if (nrow(gaps_31) != 0) {
    cutoff_date <- as.character(tail(gaps_31, n=1)$to)
    table_cut <- table %>% tsibble::filter_index(cutoff_date ~ .)
  } else {
    table_cut <- table
  }
}
```

```

table_cut_filled <- tsibble::fill_gaps(table_cut, .full = TRUE)
table_cut_filled_imp <- imputeTS::na_ma(table_cut_filled, k = 4, weighting = "simple")
table_cut_filled_imp_noLeap <- table_cut_filled_imp %>%
  dplyr::filter(!(month(Date) == 2 & day(Date) == 29))

return(table_cut_filled_imp_noLeap)
}

identical(timeseries_cleaner(dc), dc_cut_filled_imp_noLeap)

## [1] TRUE

```

### Hints

- tidyverse provides functions such as has\_gaps() and count\_gaps()

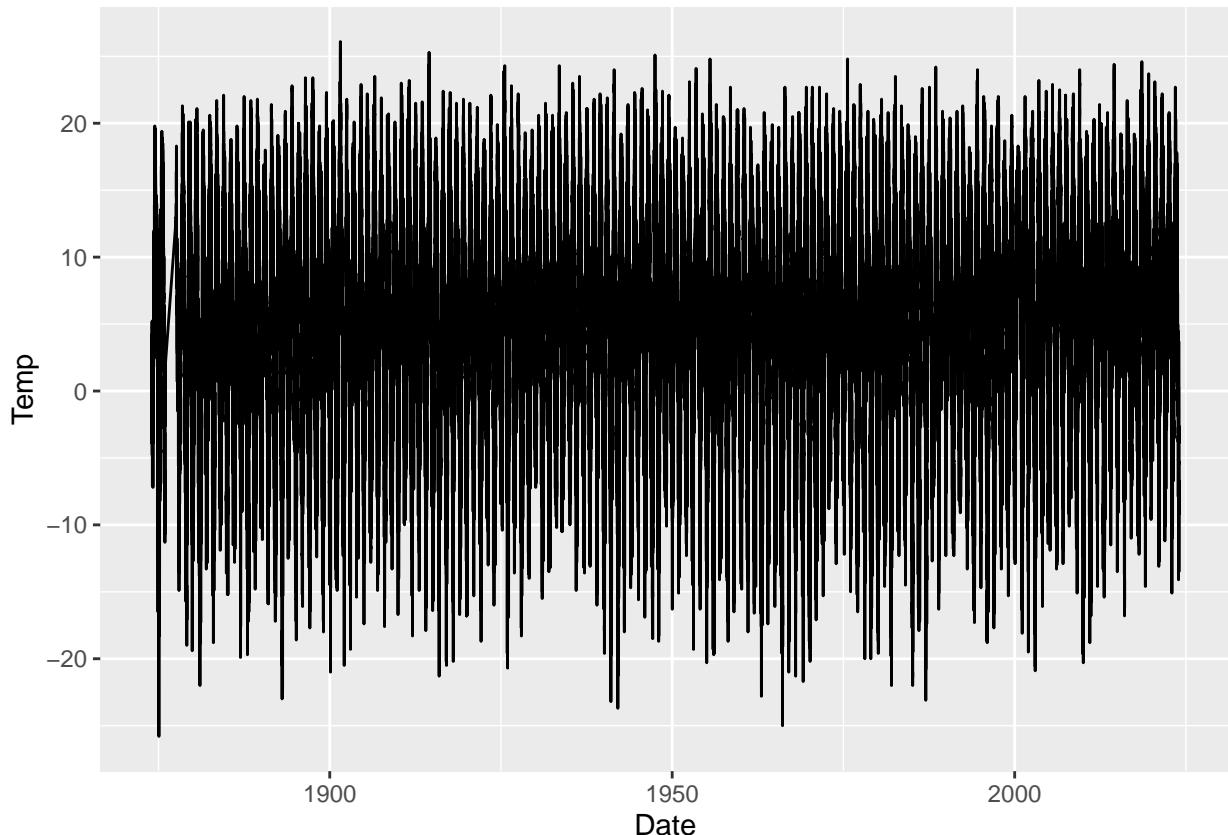
## Part C: Exploratory analysis for a single station

- Plot the temperature data as function of time

```

dc %>%
  ggplot(aes(Date, Temp)) +
  geom_line()

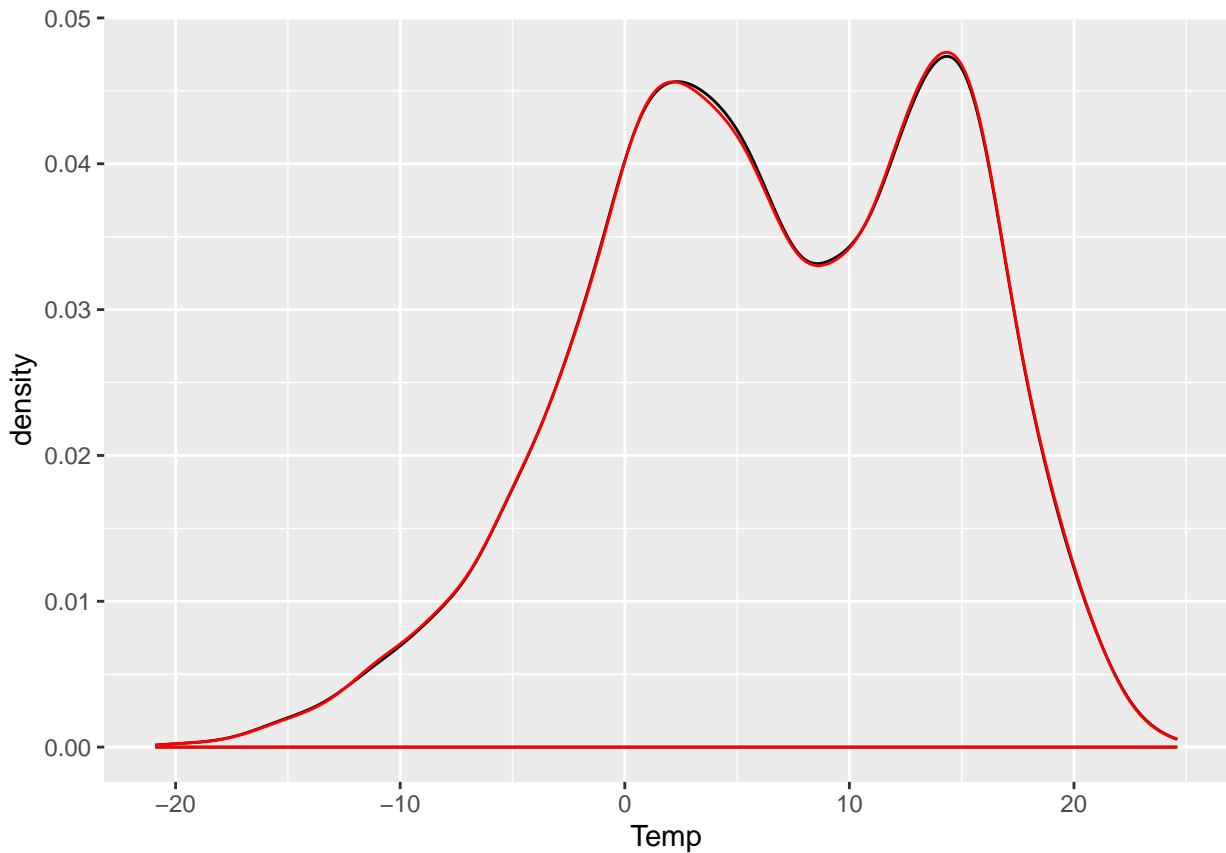
```



- Create density plots of original data and data with imputed values

```
#dc_stuff
```

```
ggplot(dc %>% tsibble::filter_index(cutoff_date ~ .), aes(Temp)) +
  stat_density(kernel = "gaussian", fill = NA, colour = "black") +
  stat_density(data = dc_cut_filled_imp, kernel = "gaussian", fill = NA, color = "red") +
  labs("t", "a")
```



Can see that even with imputed values, the density shape of the Temp is still almost the same. So we get more data points with very minimal skew.

Temperature seems to be bimodal

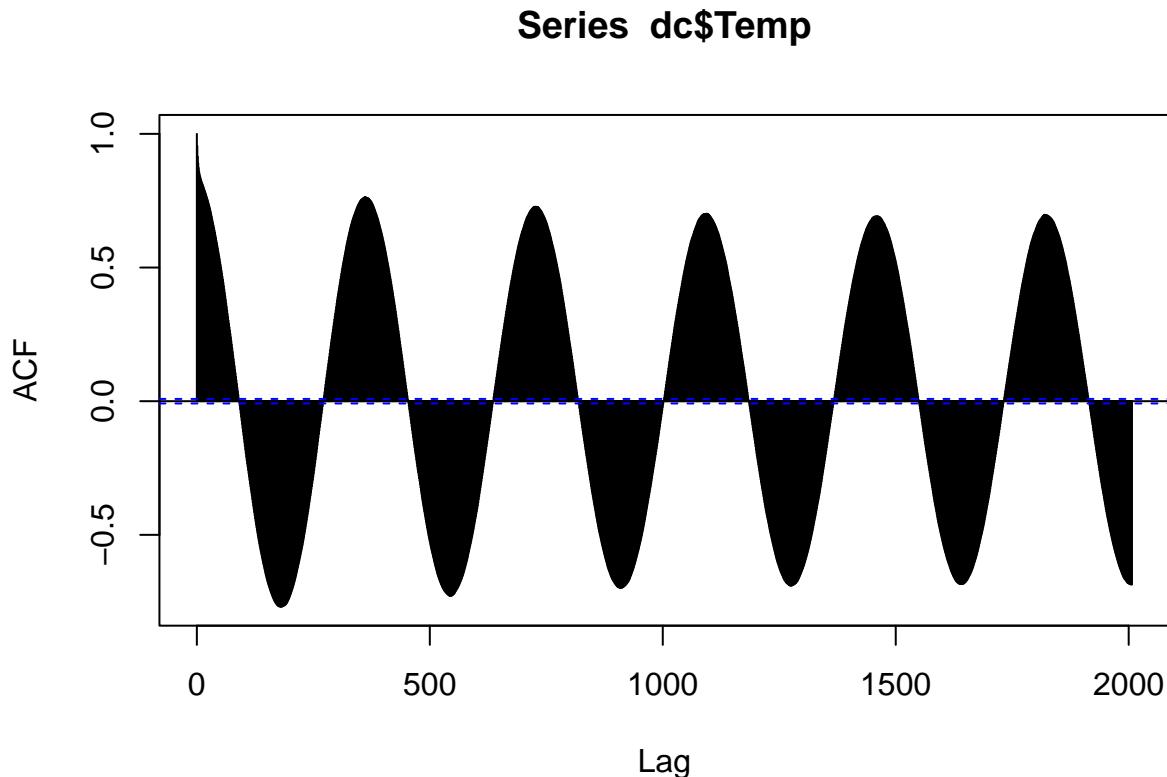
- Turn the temperature data into a timeseries (ts) object

```
ts_dc <- ts(dc_cut_filled_imp_noLeap, frequency = 365)
head(ts_dc)
```

```
## Time Series:
## Start = c(1, 1)
## End = c(1, 6)
## Frequency = 365
##           Date Temp
## 1.000000 6743 18.7
## 1.002740 6744 19.4
## 1.005479 6745 19.4
## 1.008219 6746 19.7
## 1.010959 6747 18.0
## 1.013699 6748 21.2
```

- Plot the autocorrelation function for lags up to 5.5 years; describe and discuss your observations

```
acf(dc$Temp, lag.max = 365*5.5)
```



- How correlated each Temperature is to the temperature 5.5 years ago
- Can see clear seasonality in the plot with a time interval of about a year (365 days)
- From eyeballing it seems that there is a positive correlation from Jan to Jul that decreases over time and a negative correlation that increases over time from Jul to Dec
- The correlation seems to be decreasing slightly over time, with the first years being perfectly near perfectly correlated (1.0) and the last ones being ca 0.7

The autocorrelation analysis helps in detecting hidden patterns and seasonality and in checking for randomness correlation between a time series and its own lagged version. Clear pattern in regression error, can be fixed. Breaks the assumption of normally distributed errors.

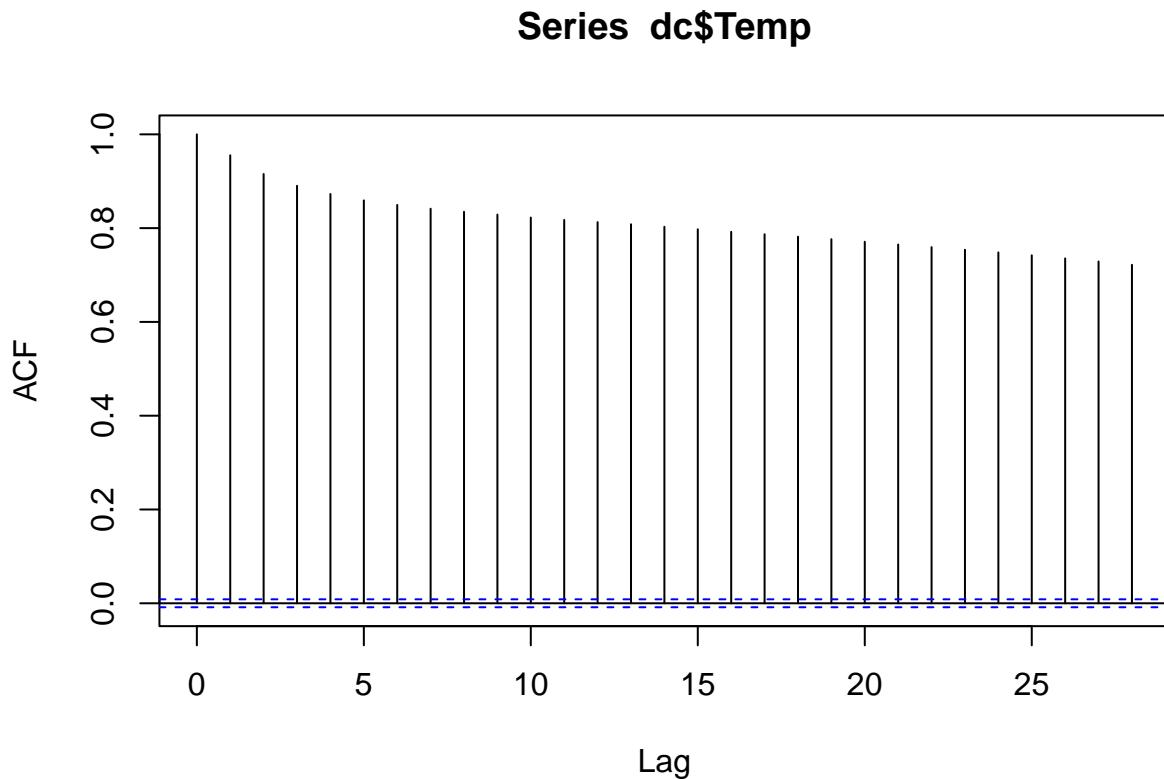
Autocorrelation represents the degree of similarity between a given time series and a lagged version of itself over successive time intervals.

Autocorrelation measures the relationship between a variable's current value and its past values.

An autocorrelation of +1 represents a perfect positive correlation, while an autocorrelation of -1 represents a perfect negative correlation.

- Also plot the ACF only for short lags, up to four weeks

```
acf(dc$Temp, lag.max = 7*4)
```



Hard to see pattern, decreasing correlation over time? (As 1 is perfect positive correlation)

- Select some days distributed throughout the year and plot temperature as function of year for, e.g., 1 October, as a scatter plot. This plot can be useful to choose the seasonality window later (see Figs 7 and 8 in Cleveland et al, 1990)

```

months <- c("Janaury", "February")
month_plot <- function(data, date) {
  out_plot <- data %>%
    dplyr::filter(month(Date) == date & day(Date) == 1)) %>%
    ggplot(aes(Date, Temp)) +
    geom_point() +
    ylim(-25, 25) +
    ggtitle(month.name[i])

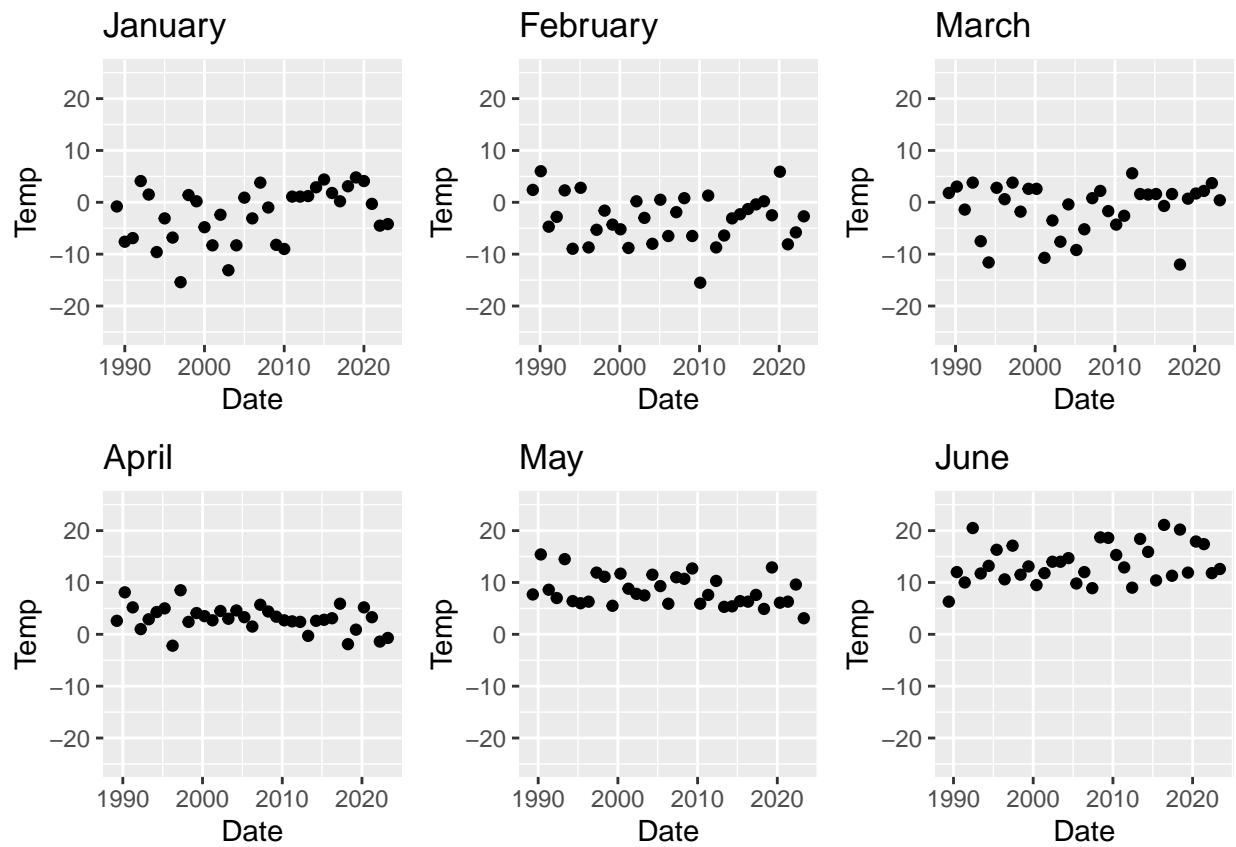
  return(out_plot)
}

month_plots_list <- list()

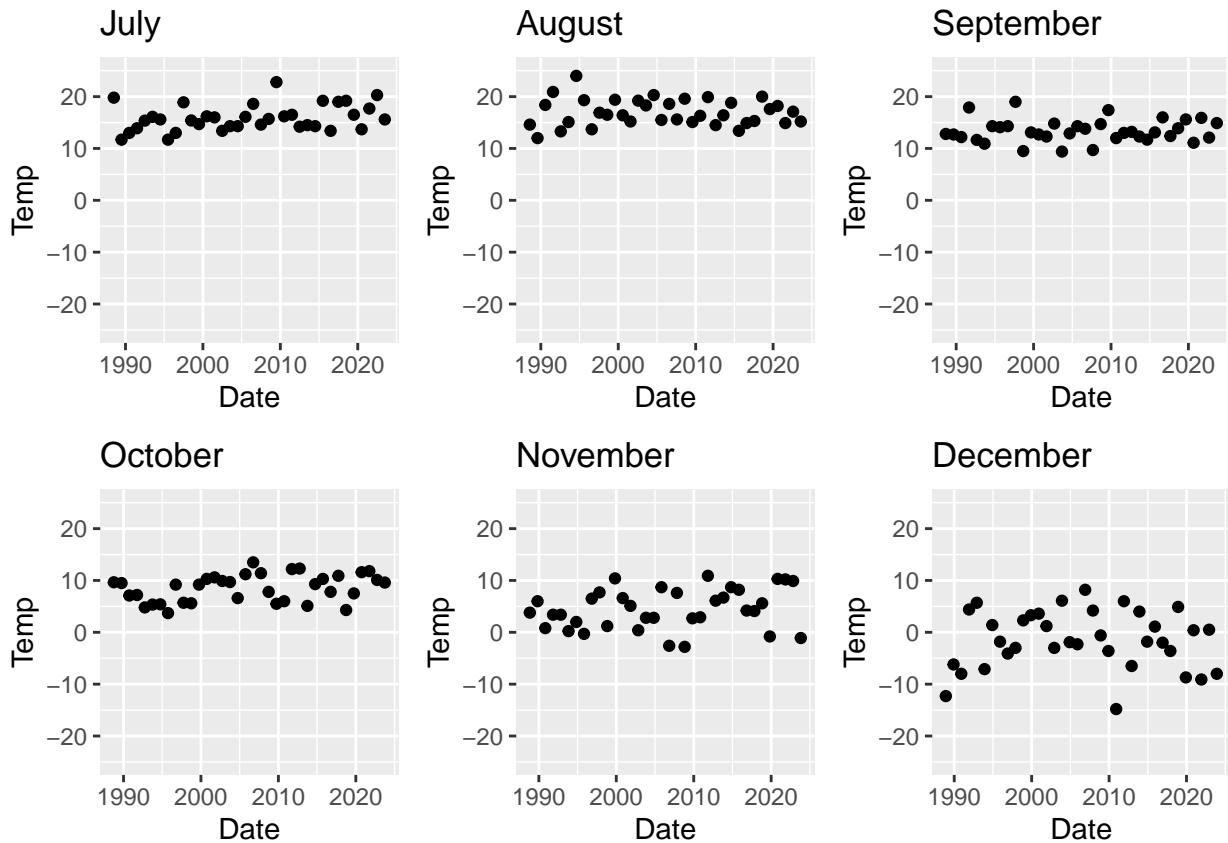
for (i in 1:12) {
  month_plots_list[[i]] <- month_plot(dc_cut_filled_imp_noLeap, i)
}

gridExtra::grid.arrange(grobs = month_plots_list[1:6], ncol = 3)

```



```
gridExtra::grid.arrange(grobs = month_plots_list[7:12], ncol = 3)
```



Can clearly see that the temperature remains high in the summer seasons, and low in winter seasons. But it seems some months are more stable than others. The more stable ones seems to fluctuate between -5 5 room while the unstable ones seems to fluctuate between -10 and 10. Over time it seems that the heat remains stable, as in the temperature doesn't seems to have a trend. Would recommend using a wide smoothing window because of this `s.window = "periodic"` or something similar

In (Cleveland et al, fig .7 and .8 , 1990) it is shown that by plotting data points for specific days in a month every year you can visually inspect how good of a fit the loess model is to the different months, very useful for determining `s.window`

## Part D: STL analysis

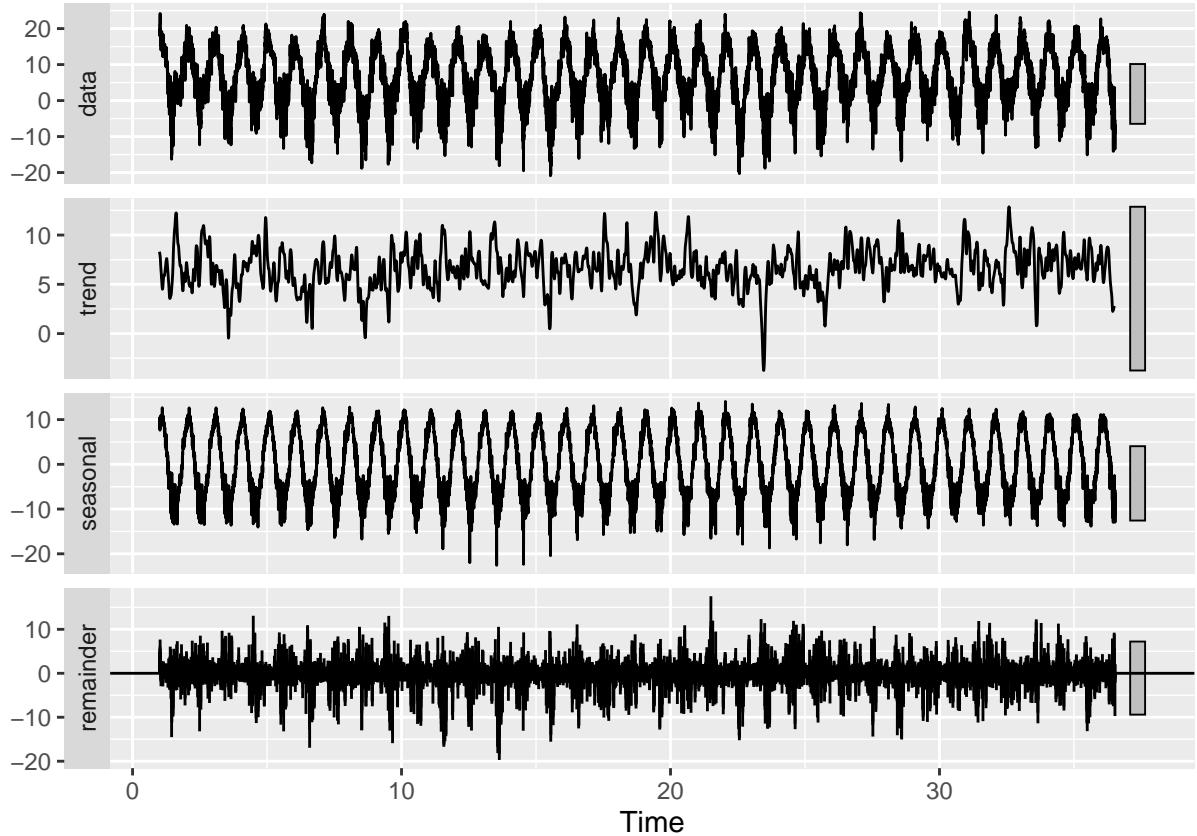
Seasonal Trend Residual, has inner and outer loop from loess using smoothing inner loop: updates trend and seasonality, use hyperparams to separate trend and seasonality, trend superposition seasonality increase in curve outer loop: compute robustness weights, assessing the residuals. Loess : local regression, idea: for every data point, fit a polynomial to make a linear fit, quadratic fit Iteratively improve fit assigning reduced weight to outliers. Use loess to estimate seasonal and trend components, and residuals

- Perform STL on the data. Explore different values for the seasonality and trend windows (remember that we want to look at trends over many years!), the choice between robust STL or not, and possibly the low pass filter window. Describe your observations. It might be interesting to look at the ACF of the remainder in the STL result.

If the autocorrelation are small and close to zero, it indicates that the remainders are uncorrelated (resemble white noise)

```
# Trend over years
dc_c_ts <- ts(dc_cut_filled_imp_noLeap$Temp, frequency = 365) #365 Leap years removed, no need for 365..
```

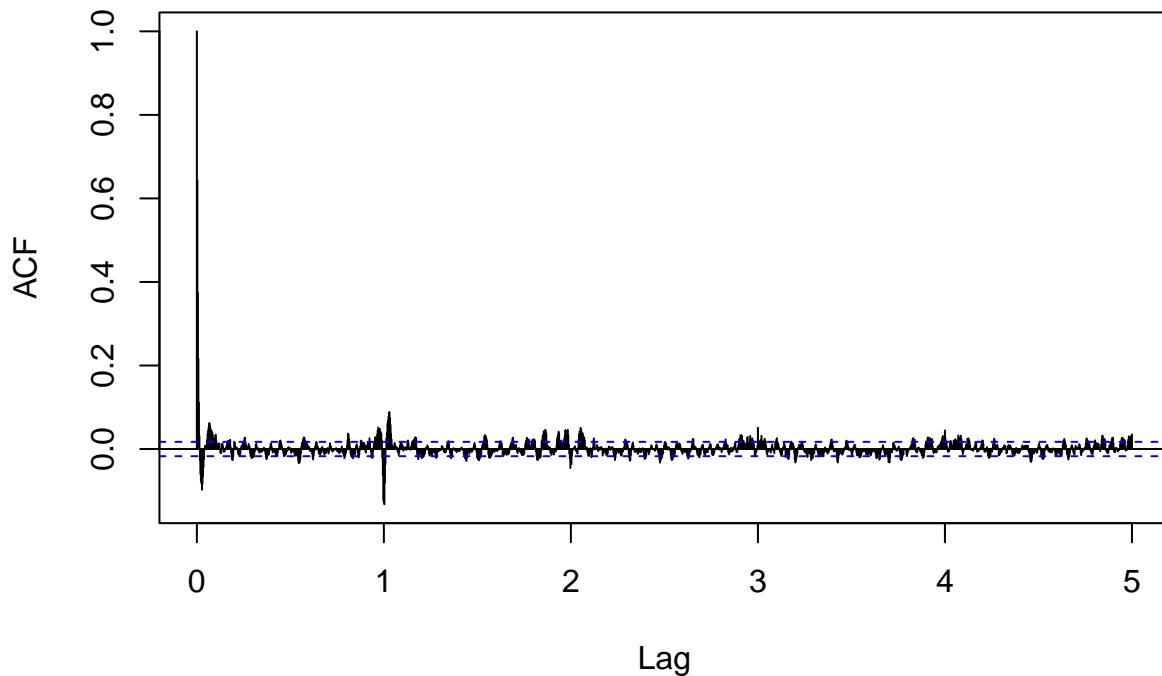
```
dc_stl <- stl(dc_c_ts, s.window=7, t.window=35, robust=TRUE)
autoplot(dc_stl)
```



Can clearly see the seasonal component in the data and its own plot, but the trend line seems more uncertain. Remainders seems to spike on every seasonal top and bottom.

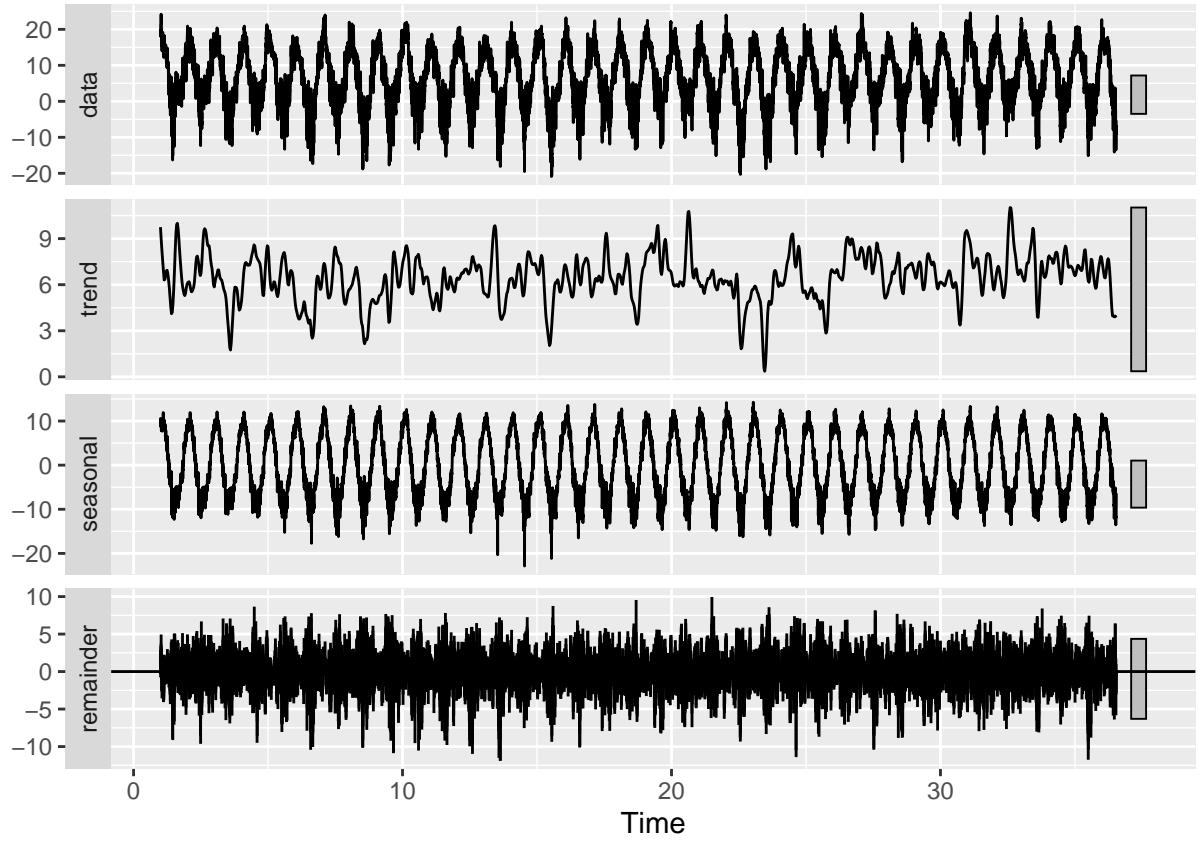
```
# autocorrelation on remainders
remainders <- dc_stl$time.series[, "remainder"]
acf(remainders, lag.max = 365*5)
```

## Series remainders



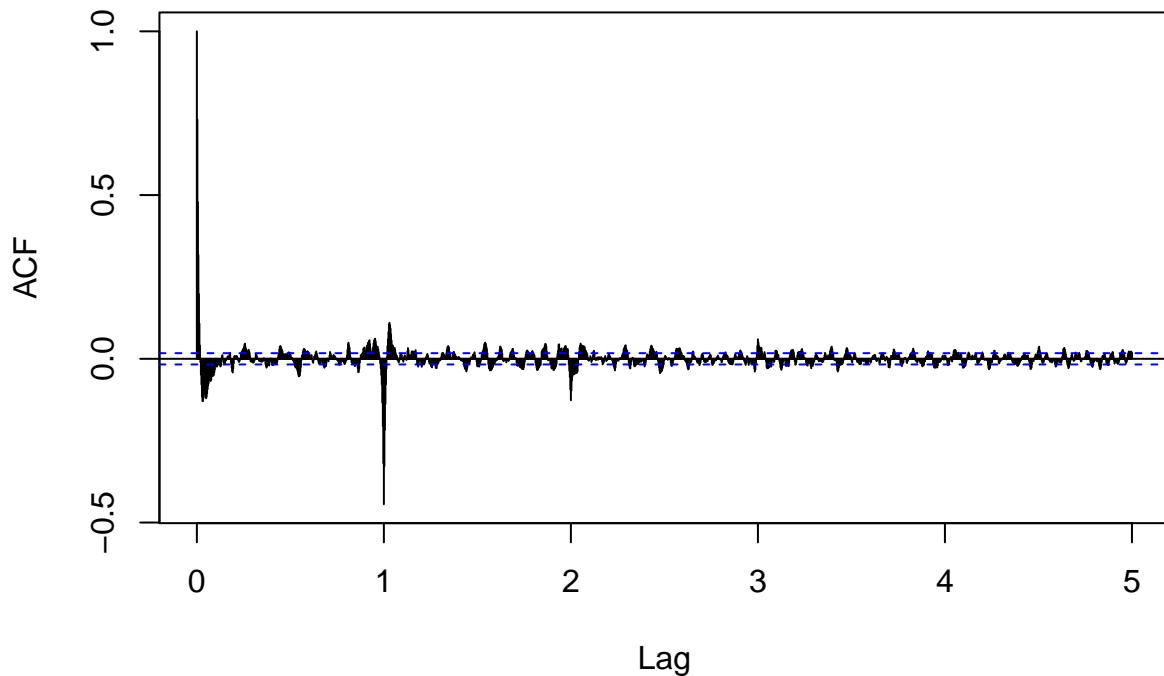
Autocorrelation seems to indicate that the remainders are mostly noise, except at 0 and 1 (0 and 5.5 years).

```
dc_stl <- stl(dc_c_ts, s.window=7, t.window=100, robust=FALSE)
autoplot(dc_stl)
```



```
# autocorrelation on remainders
remainders <- dc_stl$time.series[, "remainder"]
acf(remainders, lag.max = 365*5)
```

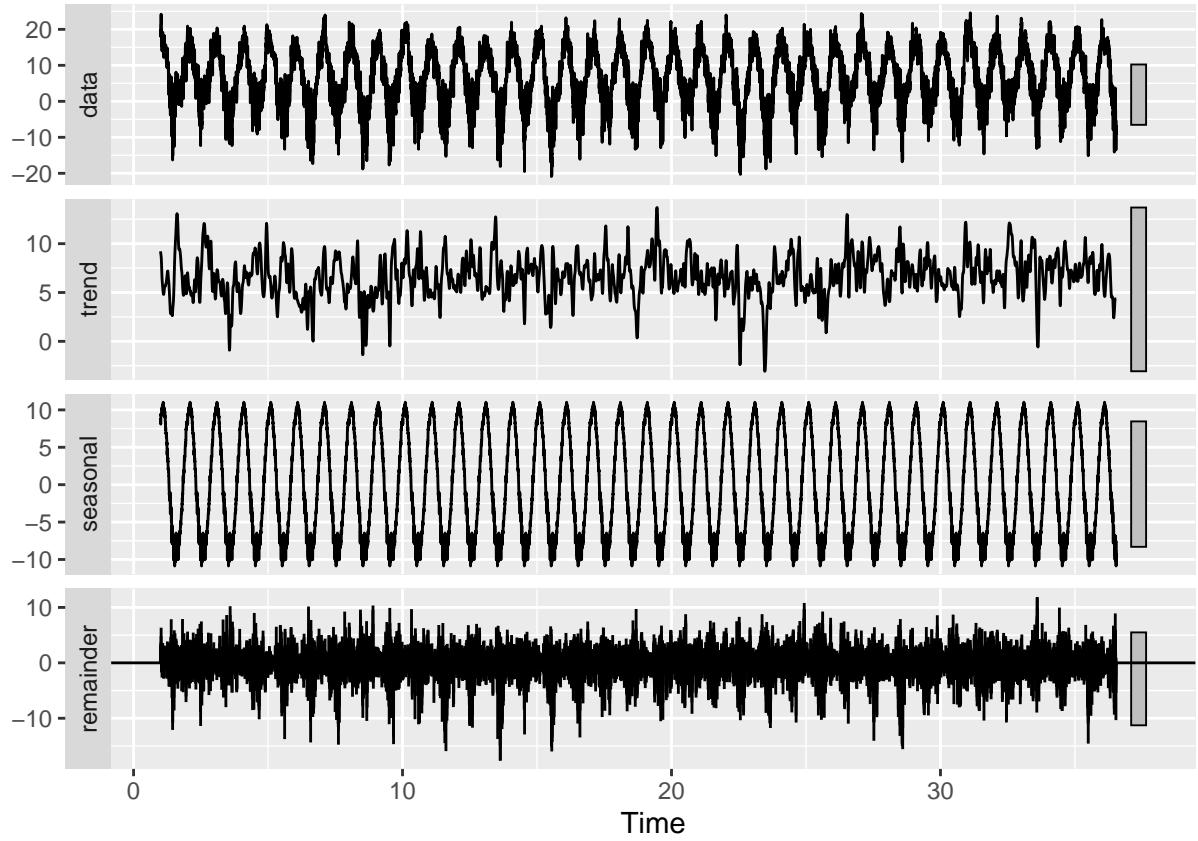
## Series remainders



Autocorrelation seems to indicate that the remainders are mostly noise, except at lags 0 and 1 (0 and 5.5 years).

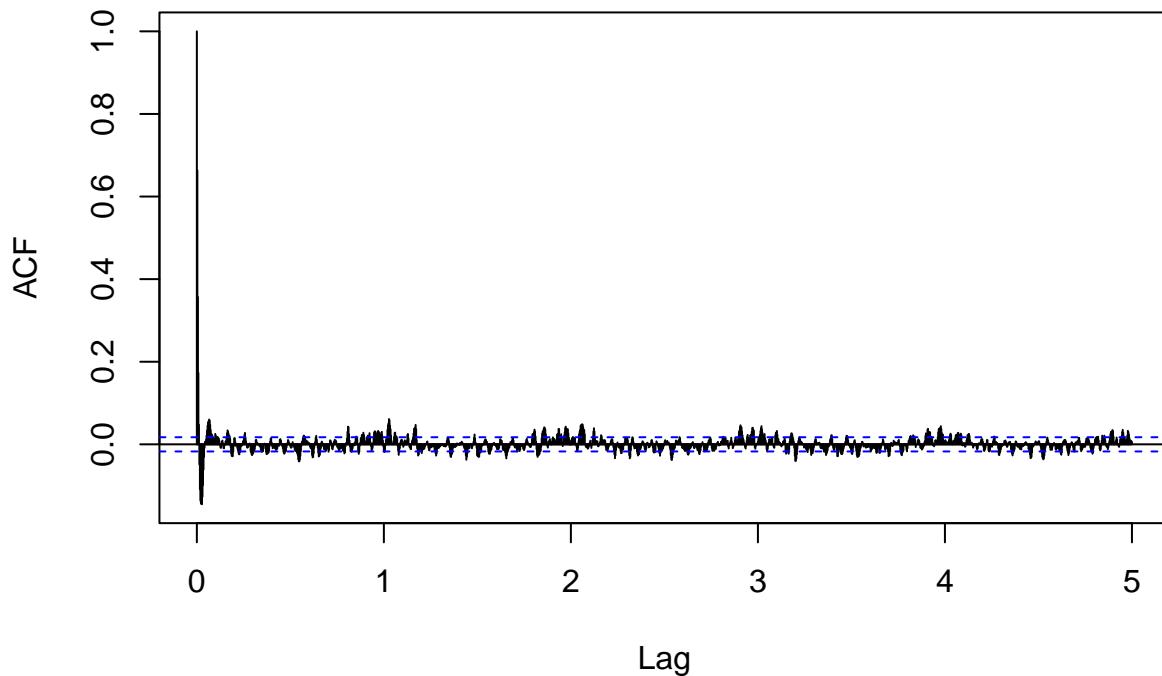
### Wide seasonality window

```
# STL wide seasonality window
dc_stl <- stl(dc_c_ts, s.window="periodic", t.window=35, robust=TRUE)
autoplot(dc_stl)
```



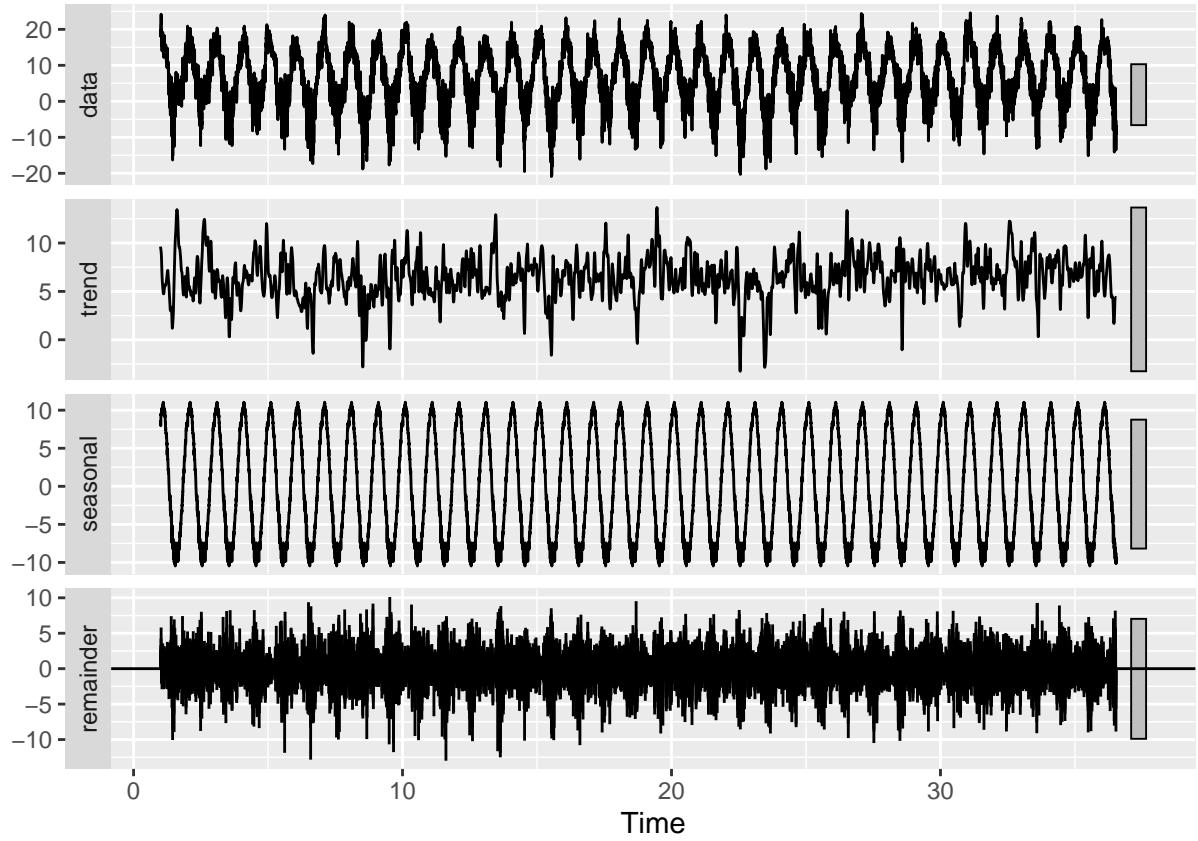
```
# autocorrelation on remainders
remainders <- dc_stl$time.series[, "remainder"]
acf(remainders, lag.max = 365*5)
```

## Series remainders



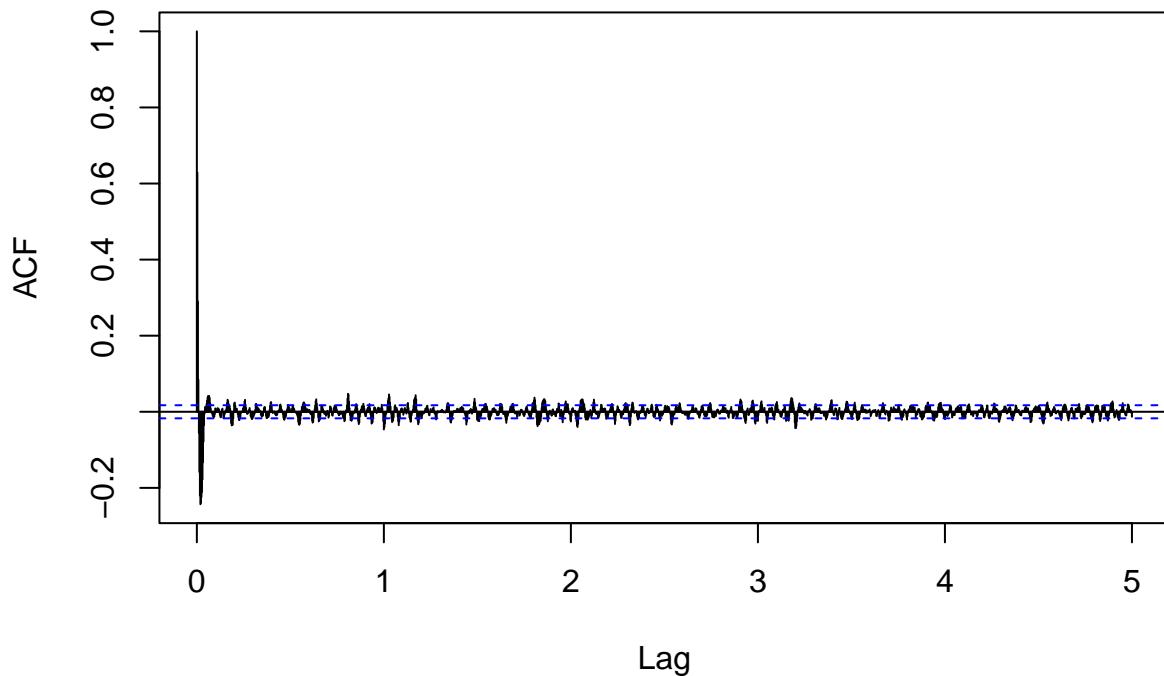
Autocorrelation seems to indicate that the remainders are mostly noise, except at 0 and 1 (0 and 5.5 years).

```
# STL wide seasonality window
dc_stl <- stl(dc_c_ts, s.window="periodic", t.window=35, robust=FALSE)
autoplot(dc_stl)
```



```
# autocorrelation on remainders
remainders <- dc_stl$time.series[, "remainder"]
acf(remainders, lag.max = 365*5)
```

## Series remainders

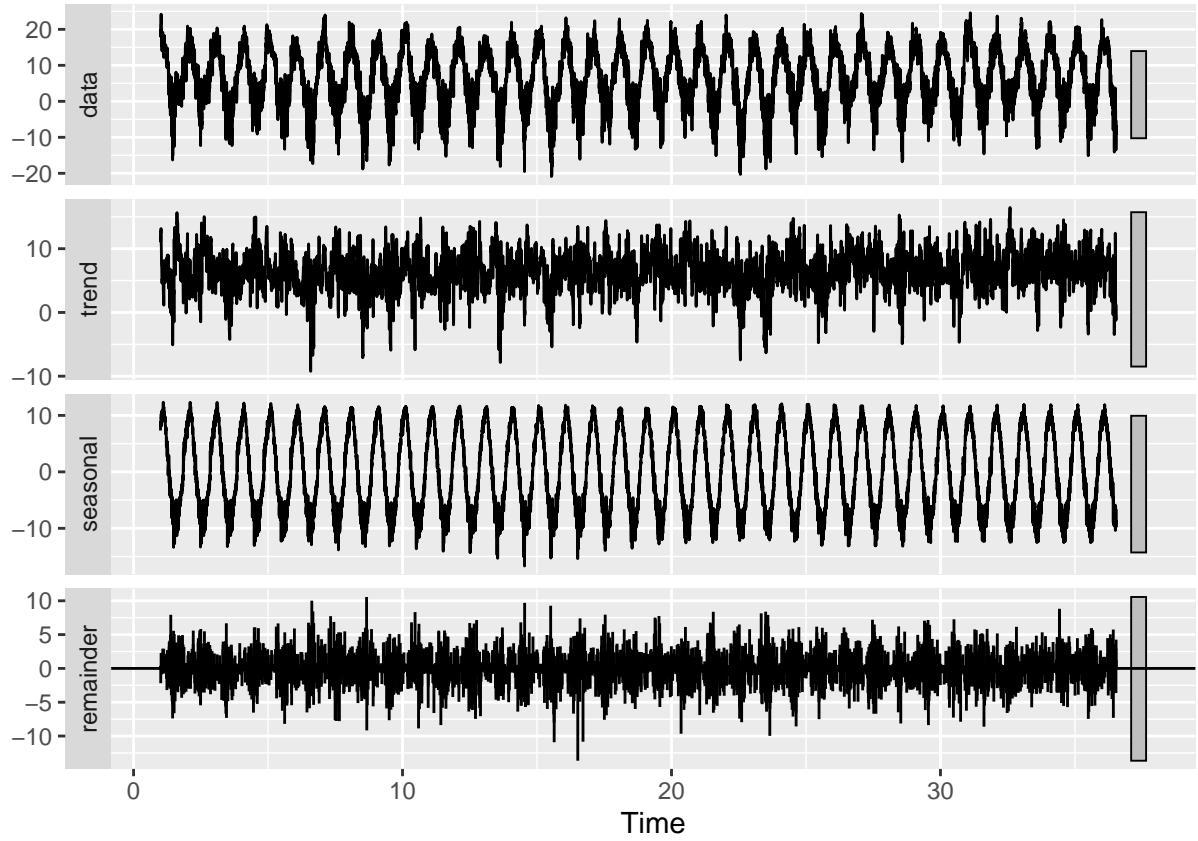


Autocorrelation seems to indicate that the remainders are mostly noise, except at 0 and 1 (0 and 5.5 years).

From eyeballing, seems robust makes little difference

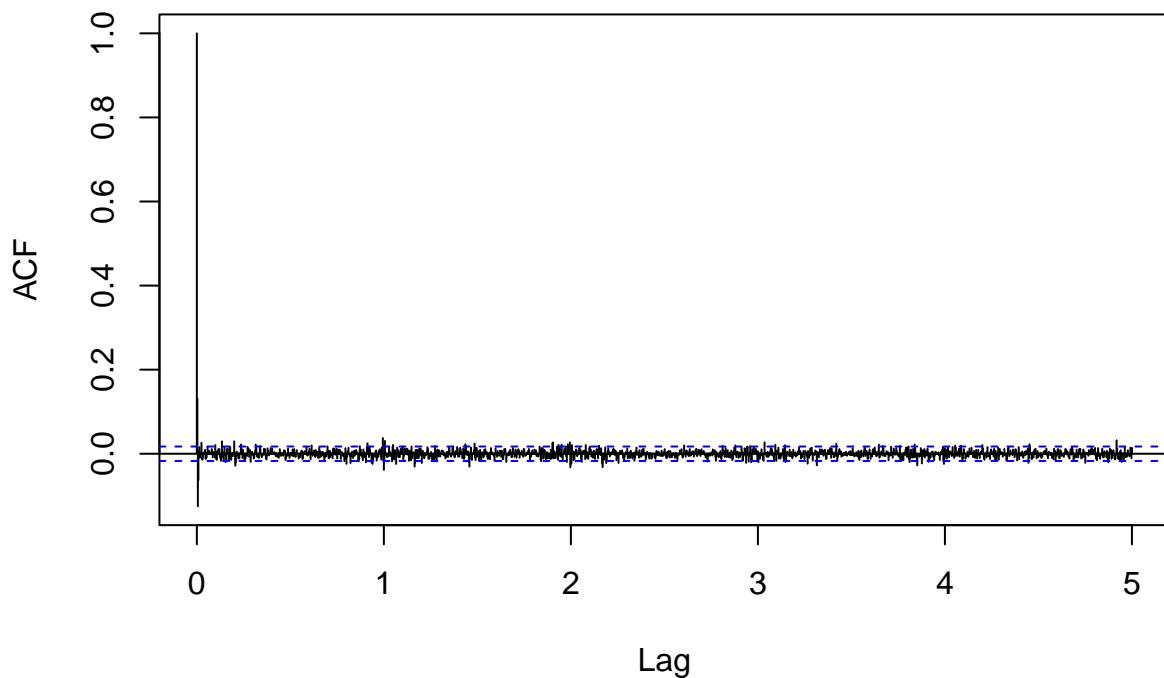
### STL short seasonality window

```
dc_stl <- stl(dc_c_ts, s.window=7, t.window=5, robust=TRUE)
autoplot(dc_stl)
```

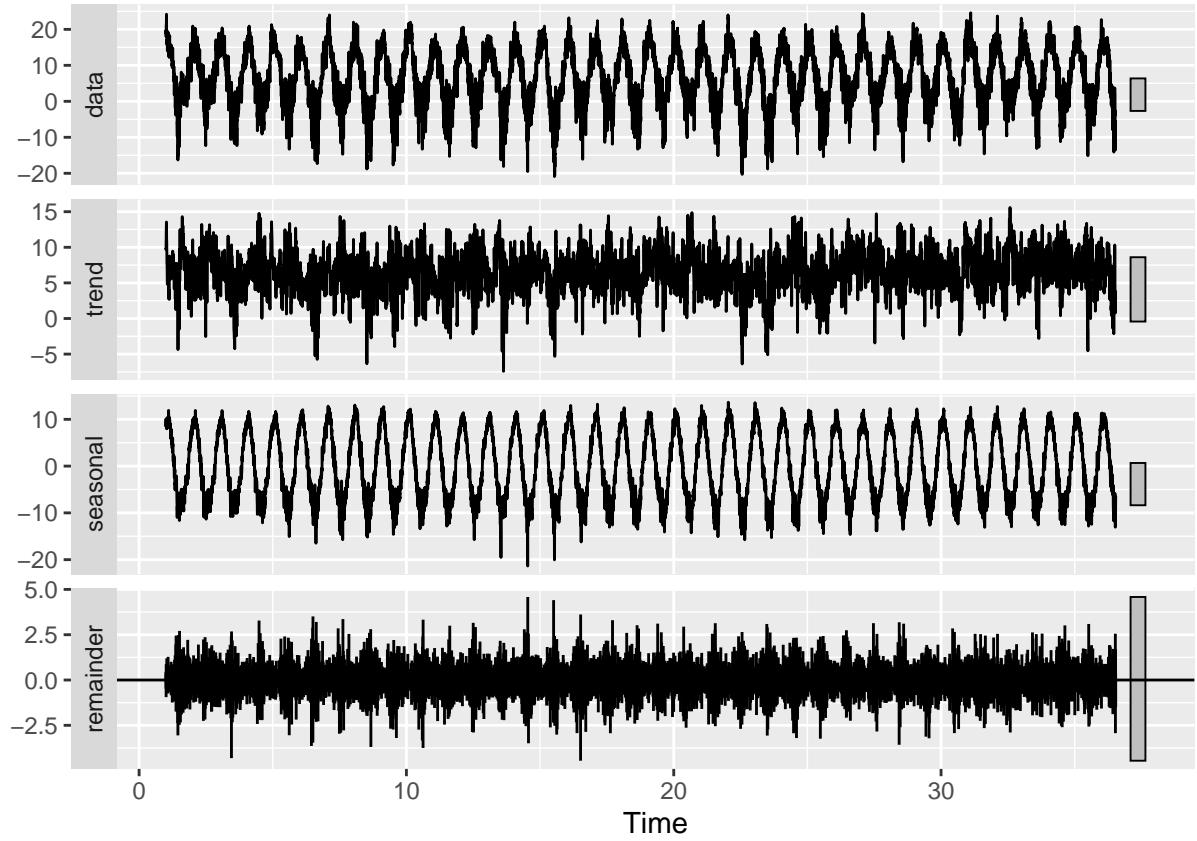


```
# autocorrelation on remainders
remainders <- dc_stl$time.series[, "remainder"]
acf(remainders, lag.max = 365*5)
```

## Series remainders

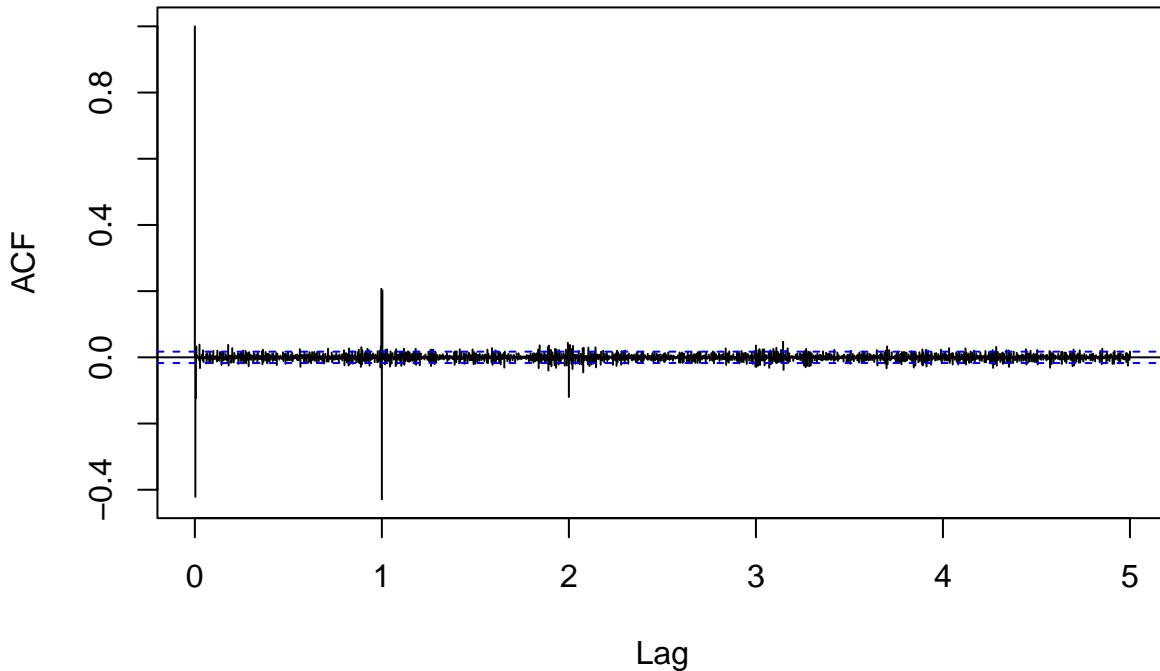


```
dc_stl <- stl(dc_c_ts, s.window=7, t.window=5, robust=FALSE)
autoplot(dc_stl)
```



```
# autocorrelation on remainders
remainders <- dc_stl$time.series[, "remainder"]
acf(remainders, lag.max = 365*5)
```

## Series remainders



- The remainder acf in all the plots seems to imply that the residuals are normally distributed
- The trend seems relatively flat in all the plots, therefore longer window seems better
- Seems that using a small window is better for getting smaller residuals, as `s.window = 7` with ‘robusts = FALSE’ seems to have lower range of residuals (-2.5, 5.0)
- Consult the original STL paper by Cleveland et al. (1990) for suggestions on how to choose STL parameters.

`s.window = "periodic"` means take all the data into account and wrap around , will make the seasonal component periodic, aka amplitude will not scale over time. Do not use if the seasonality changes significantly over time.

If you assume that the seasonal pattern is constant through time, you should set `s.window` to a be a big odd number, so that you use your entire data to perform your analysis. If seasonal patterns evolve quickly, you should use a small number but bigger than 7 and odd, thereby only using the most recent data such that the analysis is not affected by old seasonal pattern that are not relevant anymore.

$n_{(s)}$  : the smoothing parameter for the seasonal component “must be carefully tailored to each application”(Cleveland et al. p.9, (1990)). In the `stl()` function those parameters are `s.window` and `s.degree`.

Assume yearly periodicity with daily data, therefore `n_(p) = frequency = 365`

`s.window` is the seasonal smoothing parameter

- “We will always take  $n_{(s)}$  to be odd.”(Cleveland et al. p.15, (1990))
- “we also want  $n_{(s)}$  to be at least 7.”(Cleveland et al. p.15, (1990))
- “The choice of appropriate variation depend critically on the characteristics of the series.”

diagnostic method:

- seasonal-diagnostic plot : “helps us decide how much of the variation in the data other than trend should go into the seasonal component and how much in the remainder”(Cleveland et al. p.15, (1990))

“Without robustness, the seasonal component is distorted. This is illustrated by seasonal-diagnostic display”(Cleveland et al. p.15, (1990))

- s.window : Should be odd and at least 7, according to Cleveland et al.
- s.degree : Should be zero or none
- t.window : span in lags of the loess window for trend extraction, should be odd.
- 
- Based on your analysis, can you suggest a set of STL parameters to use for further work?

Since the data seems to be seasonally consistent over the years, I would recommend using a big odd number for `s.window` or just `periodic`

### Trend smoothing parameter $n_{(t)}$

## Part E: Multiple station analysis

- Obtain data from eight more stations. Two should be in the same part of Norway as the station from part A; then choose three stations each from two other parts of Norway. Data should cover several decades at least, so look for stations with long series.

```
.stations_url = str_glue("https://{.client_id}@frost.met.no/sources/v0.jsonld")
raw_stations <- fromJSON(URLencode(.stations_url), flatten=TRUE)
df_stations <- data.frame(raw_stations)

# 2 from Ås SN17850

unique(df_stations[, "data.county"])

## [1] "ROGALAND"           "INNLANDET"          "VESTLAND"
## [4] "BUSKERUD"           "TRØNDELAG"          "OSLO"
## [7] "VESTFOLD"            NA                  "SVALBARD"
## [10] "MØRE OG ROMSDAL"   "NORDLAND"          "AGDER"
## [13] "AKERSHUS"           "NORSKEHAVET"       "TROMS"
## [16] "ØSTFOLD"             "NORDSJØEN"         "TELEMARK"
## [19] "KONTINENTALSOKKELEN" "FINNMARK"          "BARENTSHAVET"
## [22] "JAN MAYEN"

COUNTYS = c("AKERSHUS", "FINNMARK", "TROMS")    # replace with names of "fylker" you are interested in

stations <- unnest(raw_stations$data, cols='id') |>
  dplyr::select(id, validFrom, country, county, municipality, name, masl, `@type`) |>
  mutate(validFrom=as.Date(validFrom)) |>
  filter(`@type` == "SensorSystem" & validFrom <= "1950-01-01" & country == "Norge" & county %in% COUNTYS)

cut_stations <- rbind(
  stations[stations$county == "AKERSHUS",][1:2,],
  stations[stations$county == "TROMS",][c(2,3,5),],
  stations[stations$county == "FINNMARK",][c(1,2,4),]
)
cut_stations

## # A tibble: 8 x 8
##   id      validFrom  country county  municipality name      masl `@type`
##   <chr>   <date>     <chr>   <chr>     <chr>      <chr>    <int> <chr>
## 1 SN19710 1913-01-01 Norge   AKERSHUS ASKER      ASKER      163 Sensor~
```

```

## 2 SN17850 1874-01-01 Norge    AKERSHUS ÅS           ÅS           92 Sensor-
## 3 SN90800 1933-09-11 Norge    TROMS    KARLSØY      TORSVÅG FYR      21 Sensor-
## 4 SN90450 1895-08-01 Norge    TROMS    TROMSØ       TROMSØ        100 Sensor-
## 5 SN89350 1940-07-01 Norge    TROMS    MÅLSELV      BARDUFOSS     76 Sensor-
## 6 SN93700 1868-01-01 Norge    FINNMARK KAUTOKEINO KAUTOKEINO 307 Sensor-
## 7 SN99370 1940-11-01 Norge    FINNMARK SØR-VARANGER KIRKENES LUFTH~ 89 Sensor-
## 8 SN98550 1829-06-01 Norge    FINNMARK VARDØ      VARDØ RADIO     10 Sensor-

```

- Preprocess the data as described in Part B. Find the latest starting date of any series and create a multivariate time series with data from all nine stations starting at this date.

```

get_dc_from_frost <- function(src_id, cutoff_date_glob = NA) {
  # Set this to TRUE to generate a json file
  weather_file = paste(src_id, "_weather_data_json.rds.bz2", sep="")
  get_data_from_frost = !file.exists(file = weather_file)

  .query_url <- str_glue("https://{{client_id}}@{{server}}/{{resource}}?sources={{src_id}}&referencetime={{referencetime}}")
  raw_data <- list()

  if (get_data_from_frost) {
    raw_data <- try(fromJSON(URLencode(.query_url), flatten=TRUE))

    if (class(raw_data) != 'try-error') {
      print("Data retrieved from frost.met.no!")
      write_rds(raw_data, weather_file, compress="bz2", text=TRUE) # JSON represents data as text
      print(str_glue("Raw data (JSON) written to '{weather_file}'"))
    } else {
      print("Error: the data retrieval was not successful!")
      stop()
    }
  } else {
    raw_data <- read_rds(weather_file)
    print(str_glue("Raw data (JSON) read from '{weather_file}'"))
  }

  df <- unnest(raw_data$data, cols = c(observations))

  df |> dplyr::select(referenceTime, value) |>
    mutate(referenceTime=as.Date(referenceTime)) |>
    rename(Date=referenceTime, Temp=value) |>
    as_tsibble(index = Date) -> dc

  if (!is.na(cutoff_date_glob)) {
    dc <- dc %>% tsibble::filter_index(cutoff_date_glob ~ .)
  }

  return(dc)
}

station_list <- c()
for (i in 1:8) {
  station_list <- c(station_list, as.character(timeseries_cleaner(get_dc_from_frost(as.character(cut_st
} 

## Raw data (JSON) read from 'SN19710_weather_data_json.rds.bz2'
## Raw data (JSON) read from 'SN17850_weather_data_json.rds.bz2'

```

```

## Raw data (JSON) read from 'SN90800_weather_data_json.rds.bz2'
## Raw data (JSON) read from 'SN90450_weather_data_json.rds.bz2'
## Raw data (JSON) read from 'SN89350_weather_data_json.rds.bz2'
## Raw data (JSON) read from 'SN93700_weather_data_json.rds.bz2'
## Raw data (JSON) read from 'SN99370_weather_data_json.rds.bz2'
## Raw data (JSON) read from 'SN98550_weather_data_json.rds.bz2'
station_list

## [1] "1983-01-01" "1988-06-18" "1984-11-01" "1920-07-01" "1946-01-01"
## [6] "1996-07-08" "1984-06-01" "1945-05-06"

start_date_dc <- max(station_list)
start_date_dc

## [1] "1996-07-08"

dc_stations <- timeseries_cleaner(get_dc_from_frost(as.character(cut_stations$id[1]), start_date_dc)) %>%
  rename(!!paste("Temp", cut_stations$name[1], sep=".") := "Temp")

## Raw data (JSON) read from 'SN19710_weather_data_json.rds.bz2'
for (i in 2:8) {
  dc_stations[paste("Temp", cut_stations$name[i], sep=".")] <- timeseries_cleaner(get_dc_from_frost(as.}

## Raw data (JSON) read from 'SN17850_weather_data_json.rds.bz2'
## Raw data (JSON) read from 'SN90800_weather_data_json.rds.bz2'
## Raw data (JSON) read from 'SN90450_weather_data_json.rds.bz2'
## Raw data (JSON) read from 'SN89350_weather_data_json.rds.bz2'
## Raw data (JSON) read from 'SN93700_weather_data_json.rds.bz2'
## Raw data (JSON) read from 'SN99370_weather_data_json.rds.bz2'
## Raw data (JSON) read from 'SN98550_weather_data_json.rds.bz2'

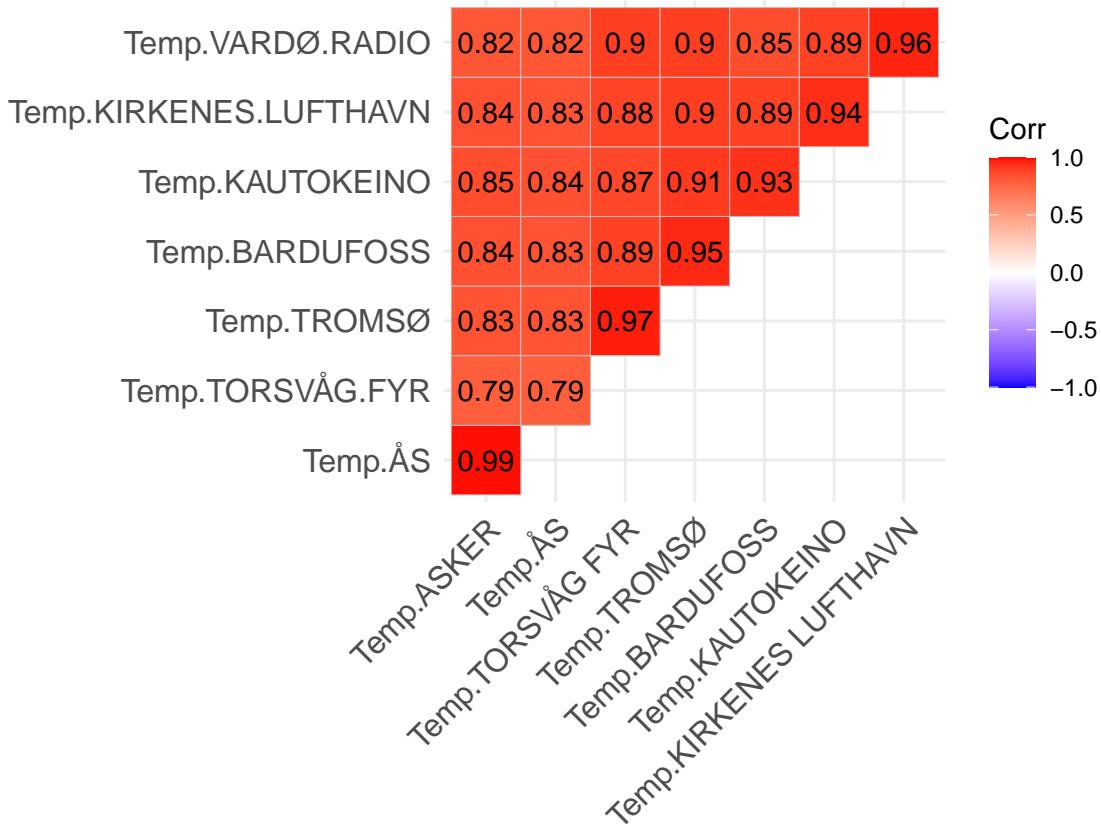
head(dc_stations)

## # A tsibble: 6 x 9 [1D]
##   Date      Temp.ASKER Temp.ÅS `Temp.TORSVÅG FYR` Temp.TROMSØ Temp.BARDUFOSS
##   <date>     <dbl>    <dbl>           <dbl>     <dbl>        <dbl>
## 1 1996-07-08     12.9     12.7          6.6       6.4        9.4
## 2 1996-07-09     16.5     15.9          7.5       8.3        9.9
## 3 1996-07-10     18.1     17.6          10        10.9       14.1
## 4 1996-07-11     16.9     16.7          10.9      10.8       13.3
## 5 1996-07-12     18.3     17.7          11.3      11         13.9
## 6 1996-07-13     15.2     14.2          12.4      13.1       13.9
## # i 3 more variables: Temp.KAUTOKEINO <dbl>, `Temp.KIRKENES LUFTHAVN` <dbl>,
## #   `Temp.VARDØ RADIO` <dbl>

  • Obtain the cross-correlation matrix between the nine stations. Is there any structure in this 9x9 matrix?

dc_stations %>%
  subset(select = -c(Date)) %>%
  cor() %>%
  round(2) %>%
  data.frame() %>%
  ggcormplot(type = "upper", lab = TRUE)

```



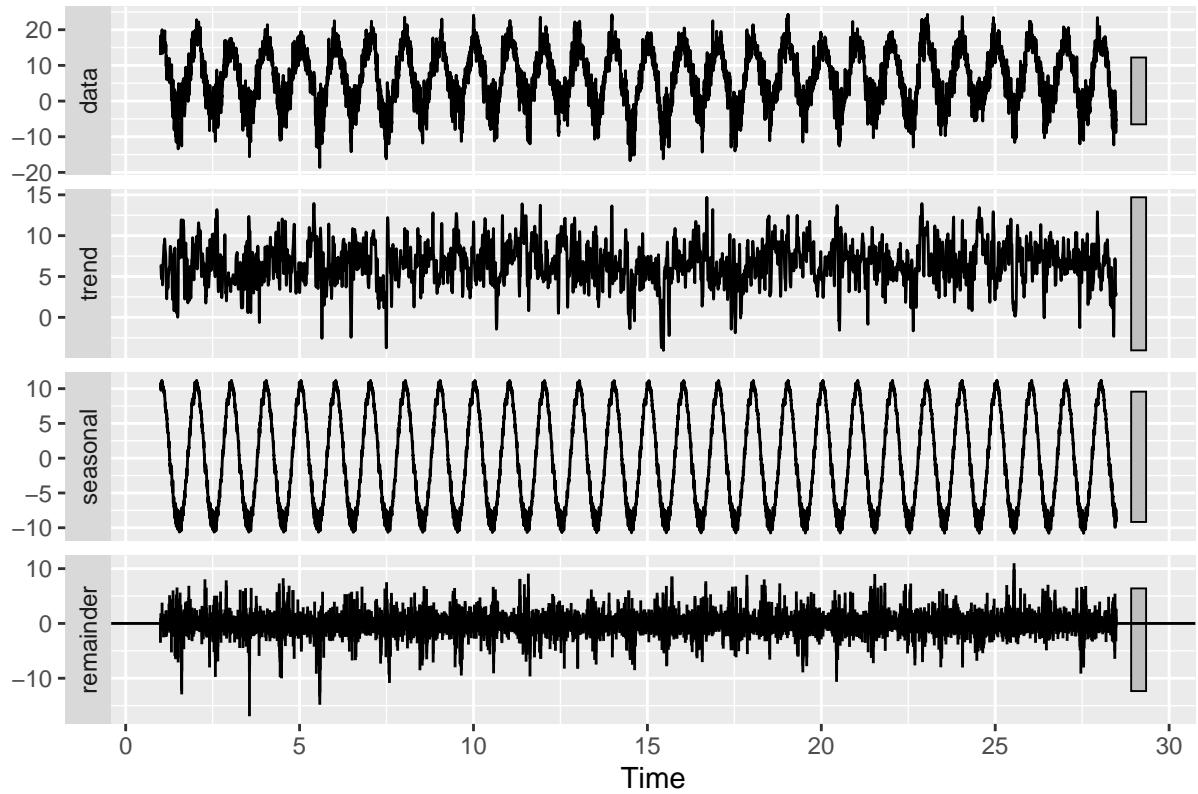
- The temperatures that comes from the same fylke seems to be more correlated, and the closer the fylke is to each other the more correlated they are.
- Everything positively correlated
- Perform STL individually on each of the nine stations using the parameters from part D. Compare the resulting trends. Are all STL results of equal quality?

```
p_list <- list()
for (colname in colnames(dc_stations)) {
  dc_stl <- stl(ts(dc_stations[[colname]]), frequency = 365), s.window=41, t.window=11, robust=TRUE)
  pl <- autoplot(dc_stl) +
    ggtitle(colname)
  p_list[[colname]] <- pl
}
```

As the trend seems to be relatively flat in all examples, I can't see any significant difference in quality of the trend lines. Remainders seems to be between 10 and -10 in all the plots so the quality seems to be about the same in all plots. Except Kirekenes Lufthavn and Kautokeino where the remainder range seems to be between [10, -20], a significant deviant from the other plots.

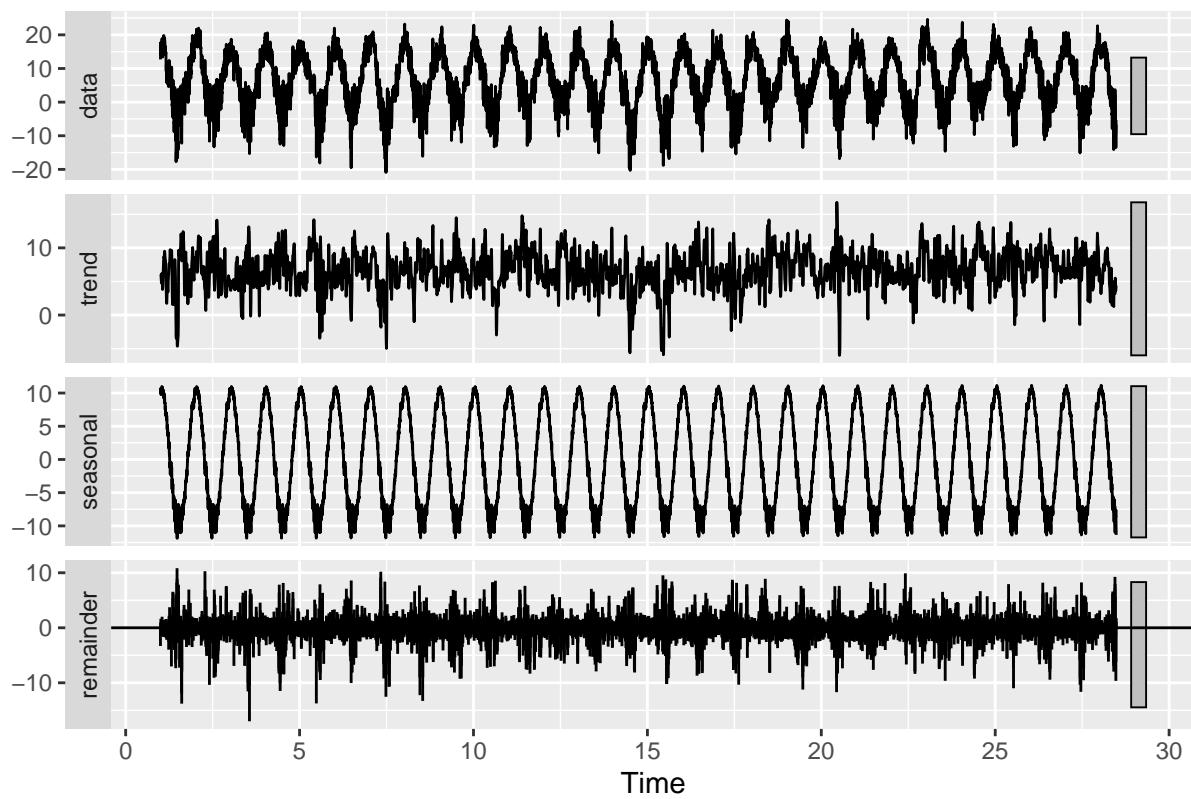
```
p_list[[2]]
```

Temp.ASKER



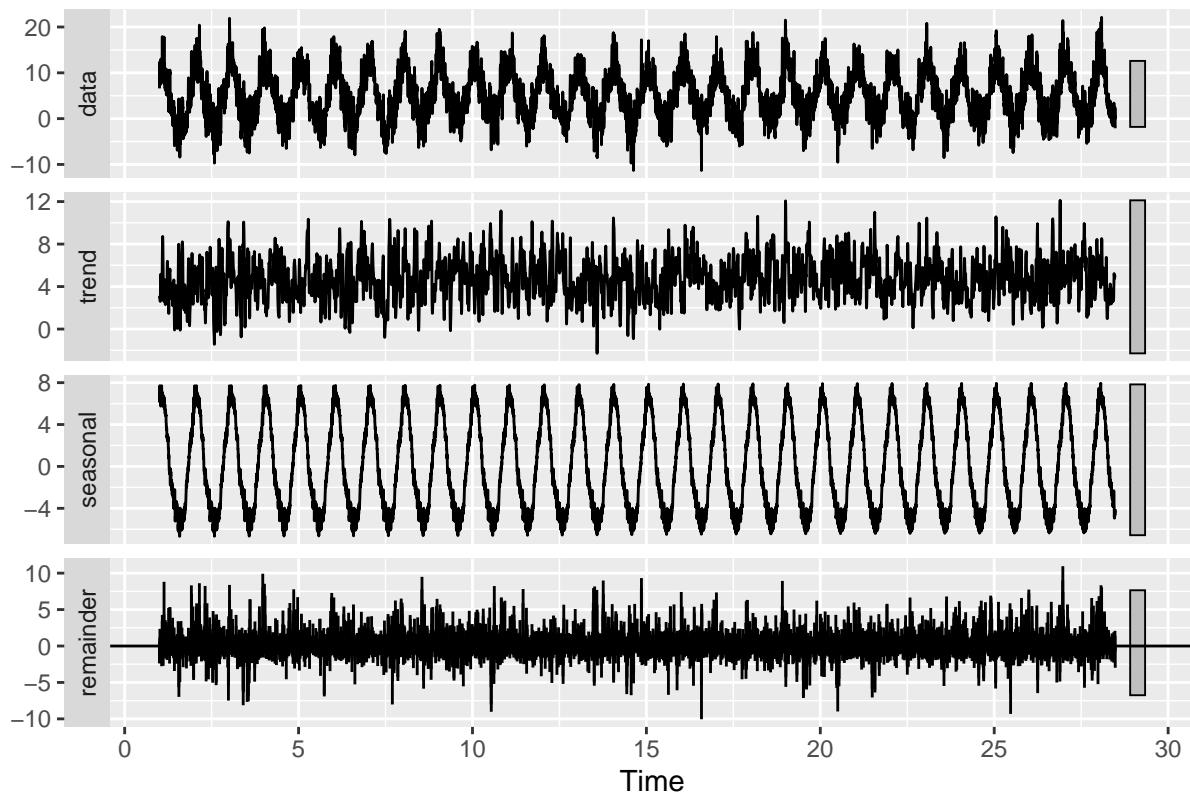
p\_list[[3]]

Temp.ÅS



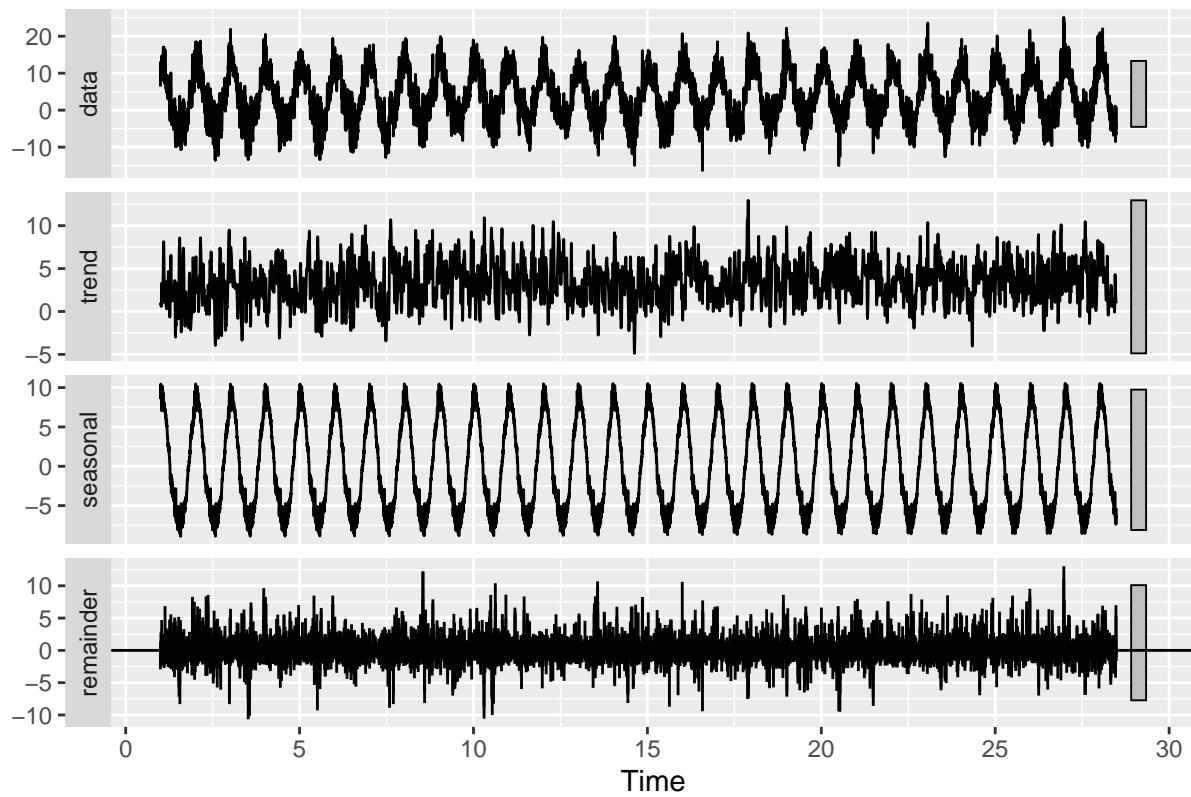
p\_list[[4]]

### Temp.TORSVÅG FYR



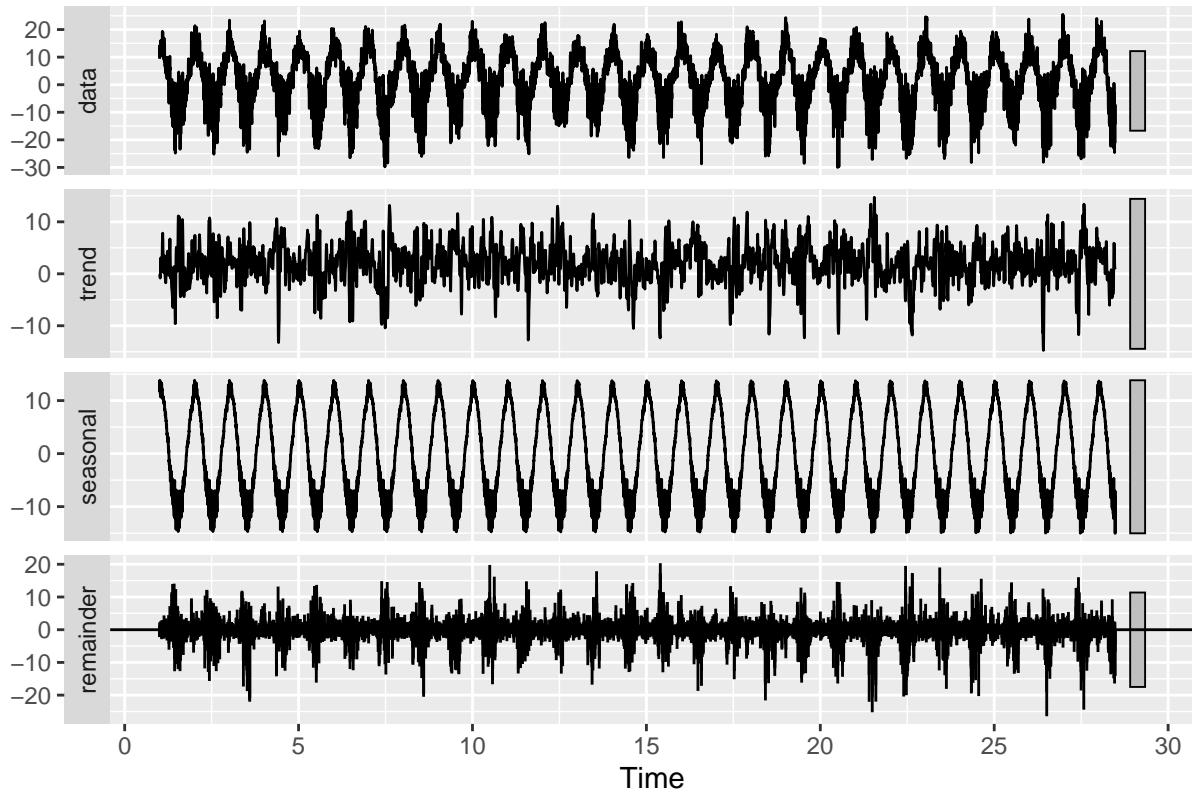
```
p_list[[5]]
```

### Temp.TROMSØ



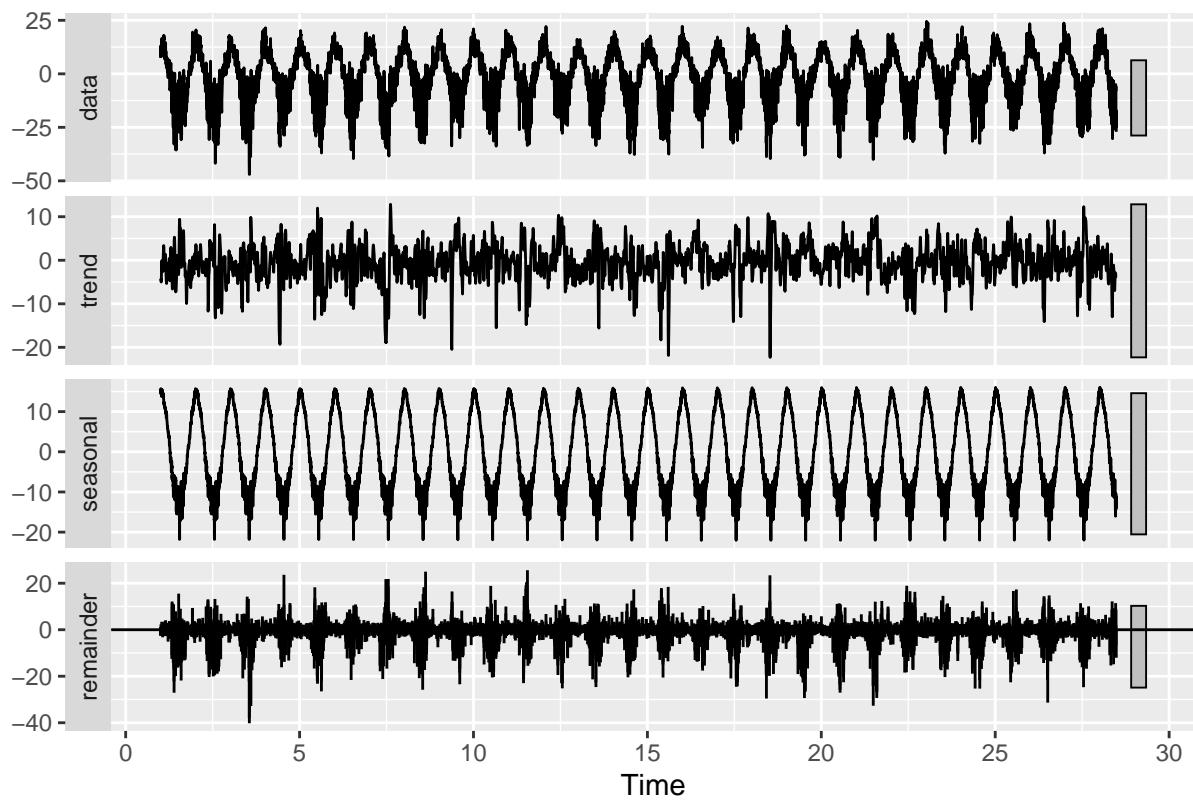
```
p_list[[6]]
```

### Temp.BARDUFOSS



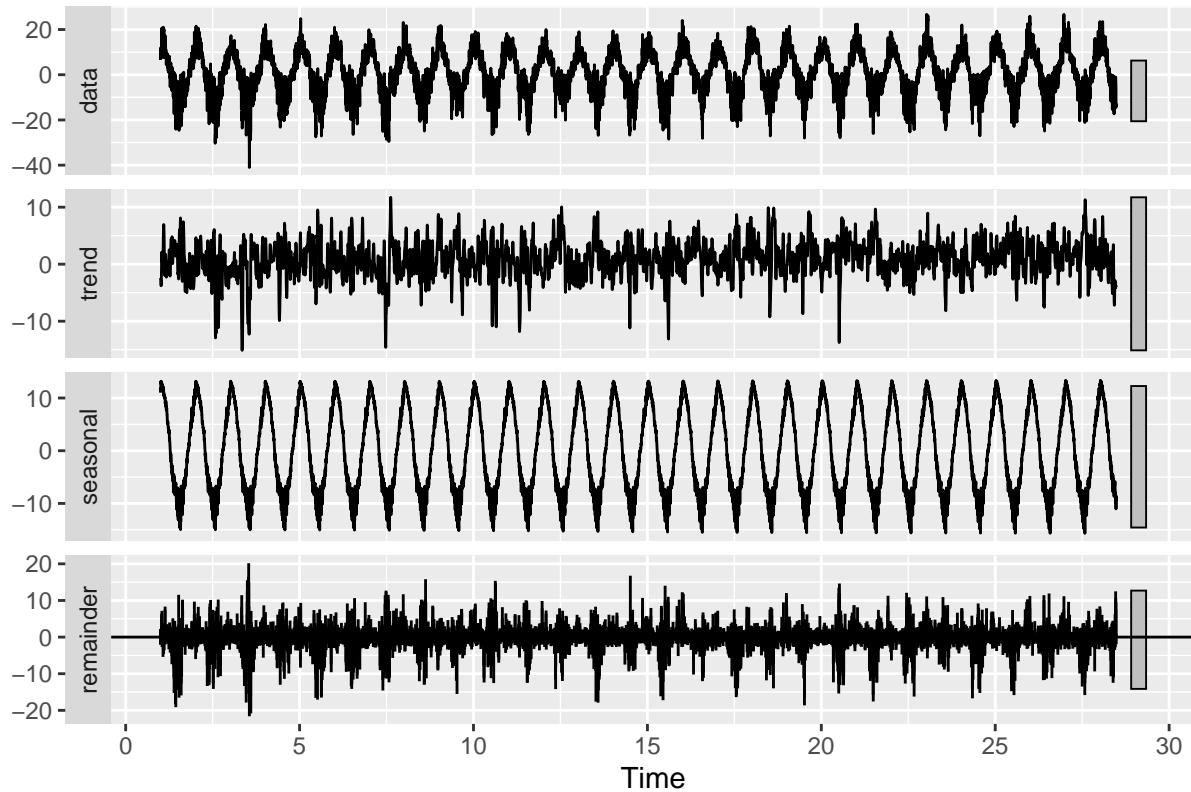
p\_list[[7]]

### Temp.KAUTOKEINO



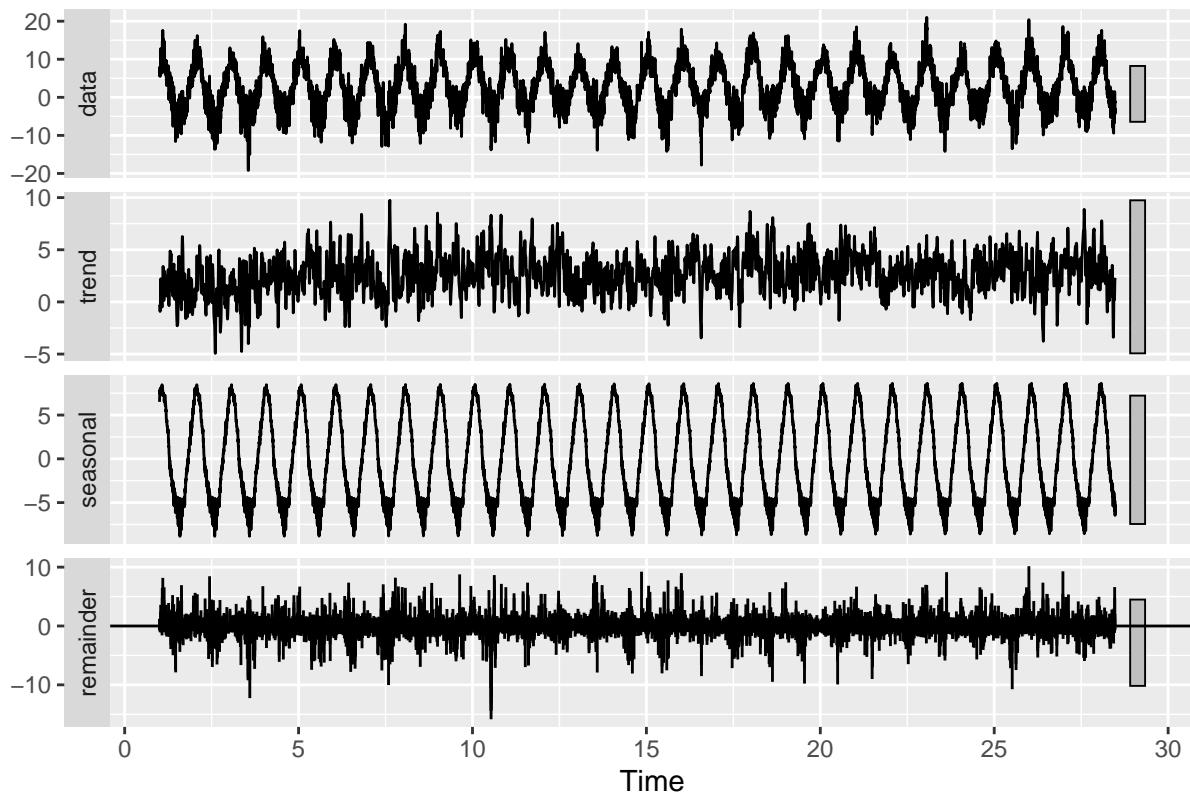
```
p_list[[8]]
```

### Temp.KIRKENES LUFTHAVN



```
p_list[[9]]
```

## Temp.VARDØ RADIO



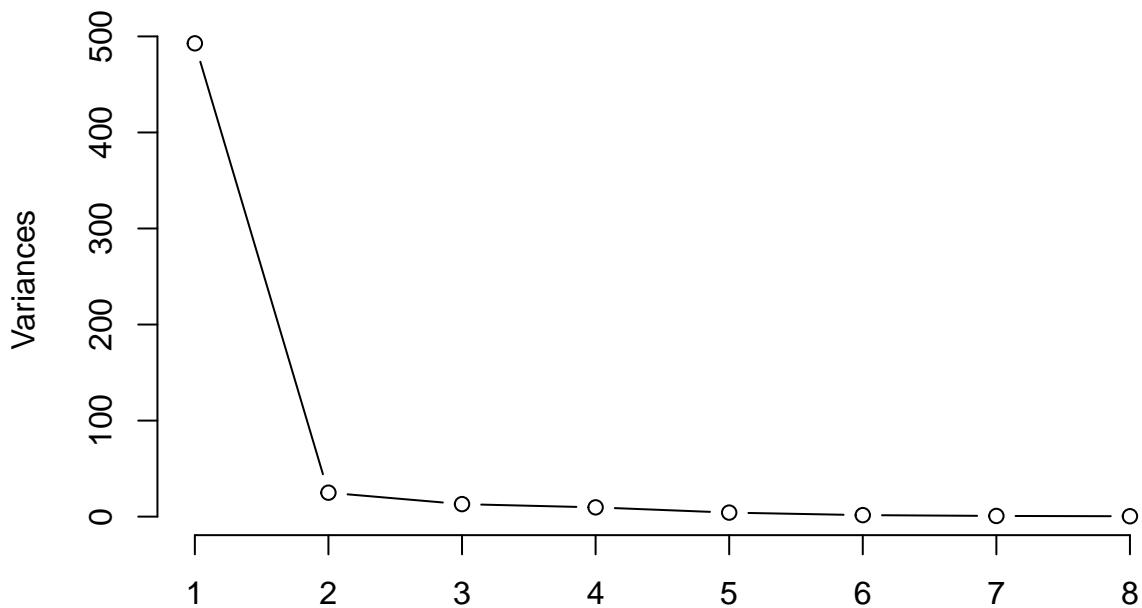
### Part F (bonus): PCA

- Perform PCA on the multivariate time series.

```
mts_pca <- dc_stations %>%
  subset(select = -c(Date)) %>%
  prcomp(center = TRUE)

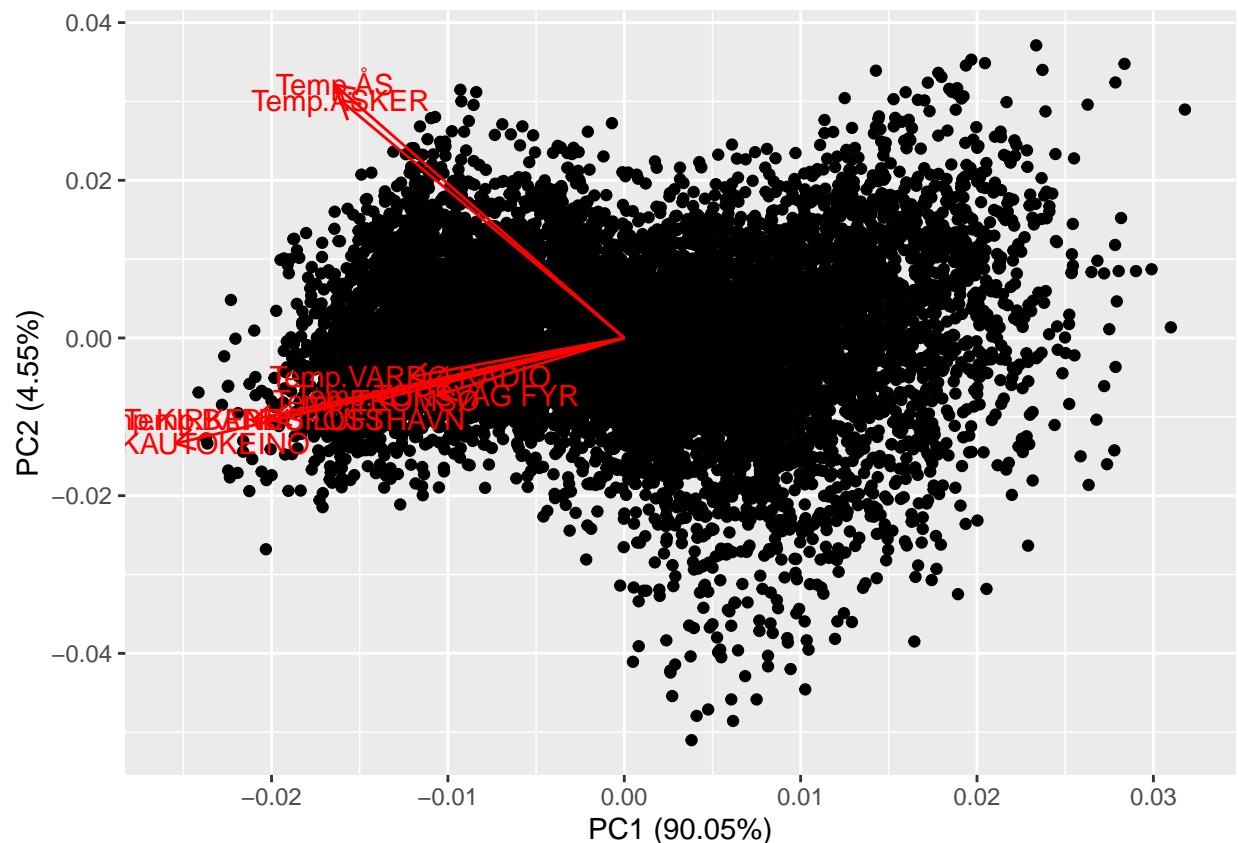
plot(mts_pca, type = "l")
```

## mts\_pca

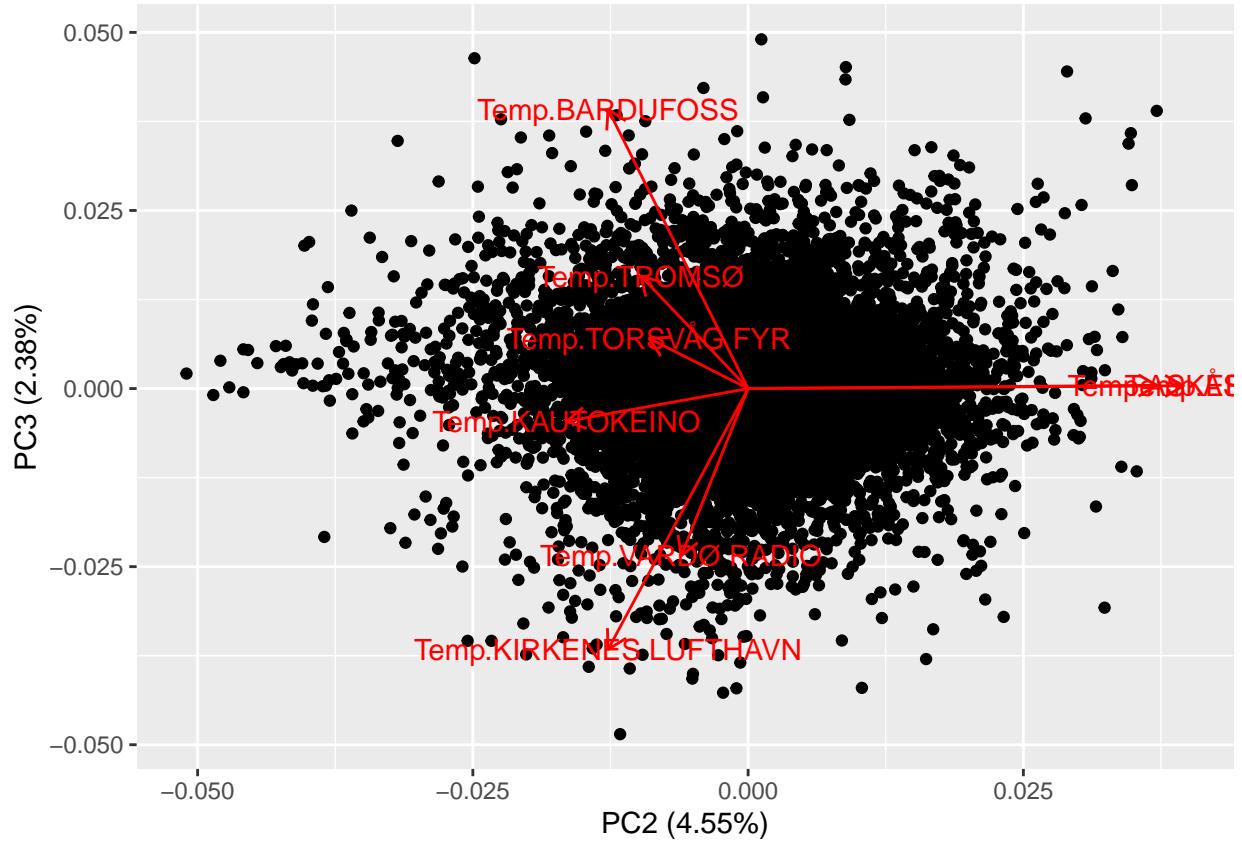


Can clearly see that most the variance is captured in the first two PC

```
autoplot(mts_pca, x=1,y=2, loadings=TRUE, loadings.label=TRUE)
```



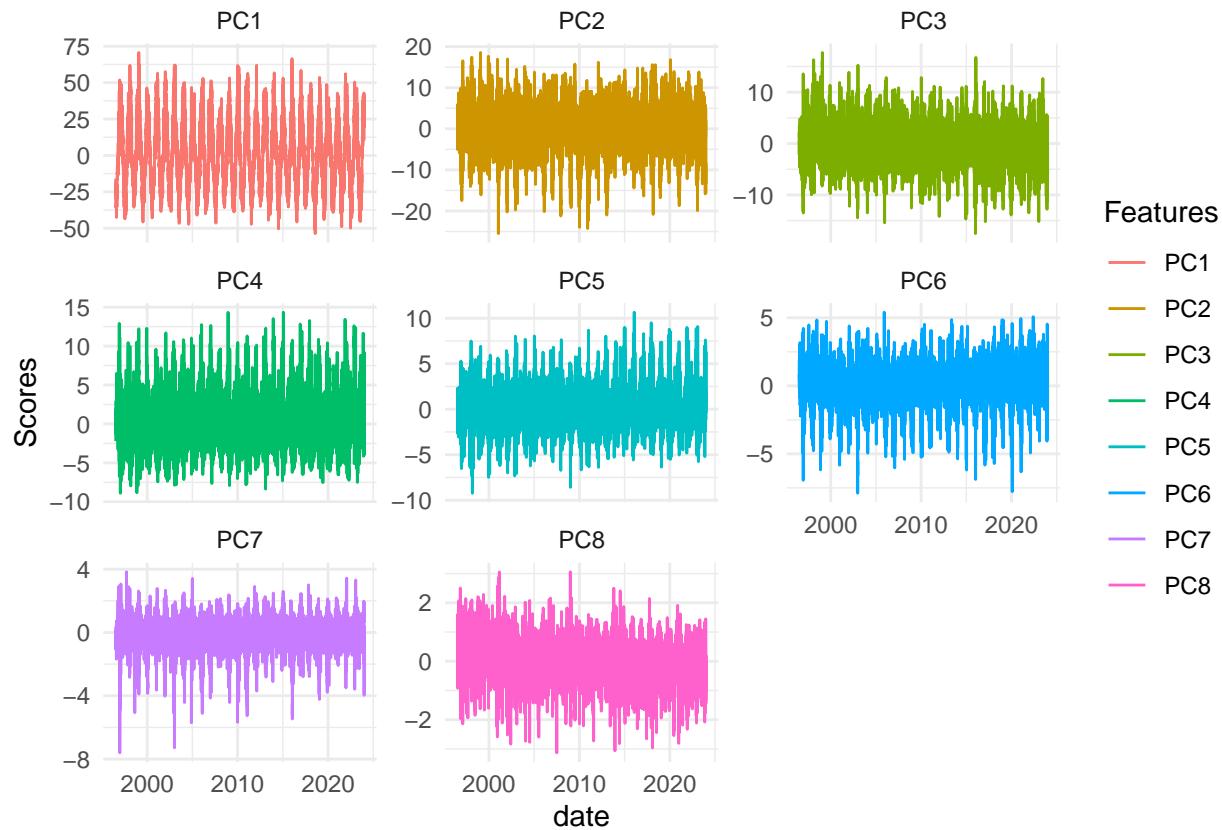
```
autoplot(mts_pca, x=2,y=3, loadings=TRUE, loadings.label=TRUE)
```



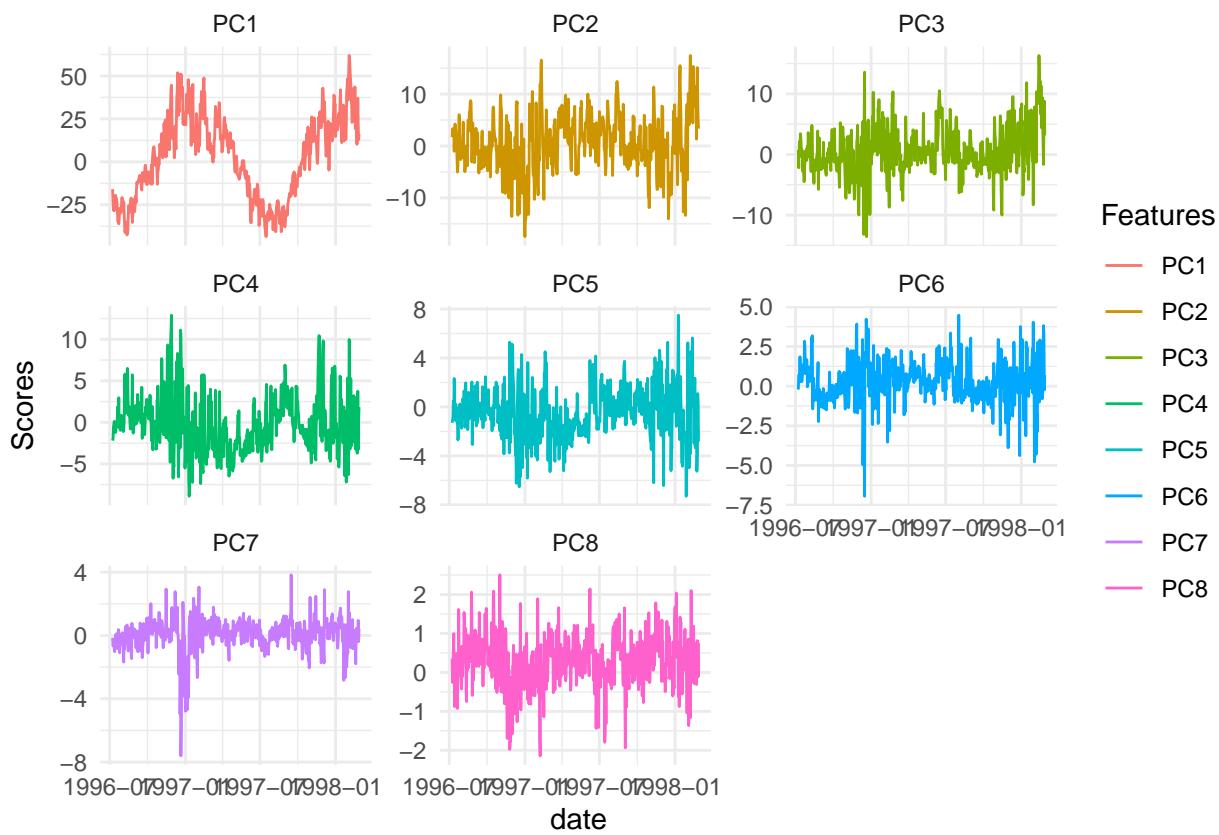
```

data.frame(date = dc_stations$Date, mts_pca$x) %>%
  pivot_longer(!date, names_to="Features", values_to="Scores") %>%
  ggplot(aes(x = date, y = Scores, col = Features)) +
  geom_line() +
  theme_minimal() +
  facet_wrap(~ Features, scales = "free_y", ncol = 3)

```



```
data.frame(date = dc_stations$Date[1:600], mts_pca$x[1:600,]) %>%
  pivot_longer(!date, names_to="Features", values_to="Scores") %>%
  ggplot(aes(x = date, y = Scores, col = Features)) +
  geom_line() +
  theme_minimal() +
  facet_wrap(~ Features, scales = "free_y", ncol = 3)
```

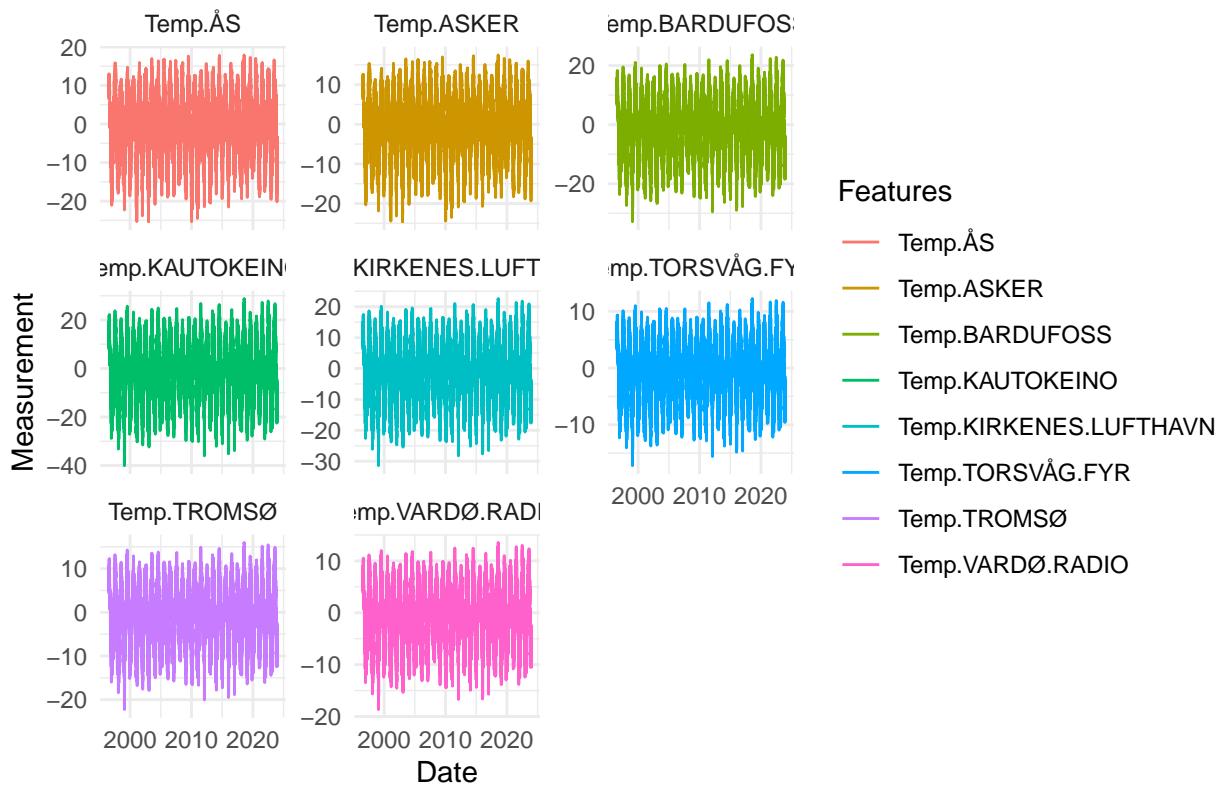


It seems that PC1 is seasonality and PC2 holds most of the trend, the rest of the PC seems to be the residuals.

```
t <- dc_stations$Date
mts_x_1 <- mts_pca$x[, 1:2] %*% t(mts_pca$rotation[, 1:2])
mts_x_2 <- t(mts_pca$center + mts_pca$scale * t(mts_x_1))

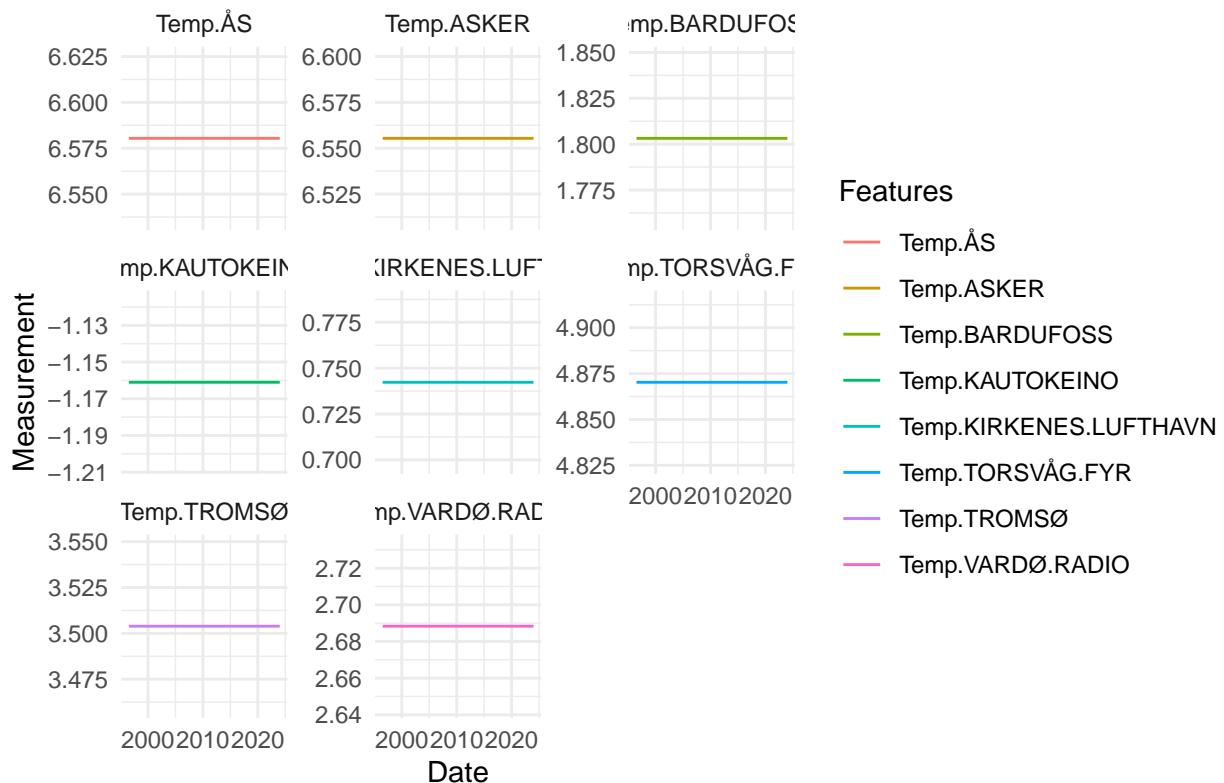
data.frame(Date = dc_stations$Date, mts_x_1) %>%
  pivot_longer(!Date, names_to="Features", values_to="Measurement") %>%
  ggplot(aes(x = Date, y = Measurement, col = Features)) +
  geom_line() +
  theme_minimal() +
  facet_wrap(~ Features, scales = "free_y", ncol = 3) +
  ggtitle("mts_pca without accounting for scaling")
```

## mts\_pca without accounting for scaling



```
data.frame(Date = dc_stations$Date, mts_x_2) %>%
  pivot_longer(!Date, names_to="Features", values_to="Measurement") %>%
  ggplot(aes(x = Date, y = Measurement, col = Features)) +
  geom_line() +
  theme_minimal() +
  facet_wrap(~ Features, scales = "free_y", ncol = 3) +
  ggtitle("mts_pca accounting for scaling")
```

## mts\_pca accounting for scaling



```
dc_stations %>%
  pivot_longer(!Date, names_to="Features", values_to="Measurement") %>%
  ggplot(aes(x = Date, y = Measurement, col = Features)) +
  geom_line() +
  theme_minimal() +
  facet_wrap(~ Features, scales = "free_y", ncol = 3) +
  ggtitle("Original data")
```

## Original data

