

Unsupervised clustering individual growth

Marko Smiljanic

10/8/2021

Background

Short summary: we have two large sites selected to study pollution signal archived within tree growth signal. Overarching question is whether there is the significant difference in the pattern of the tree growth at two sites which could be potentially traced back to lowering pollution signal in the late 1980 to early 1990.

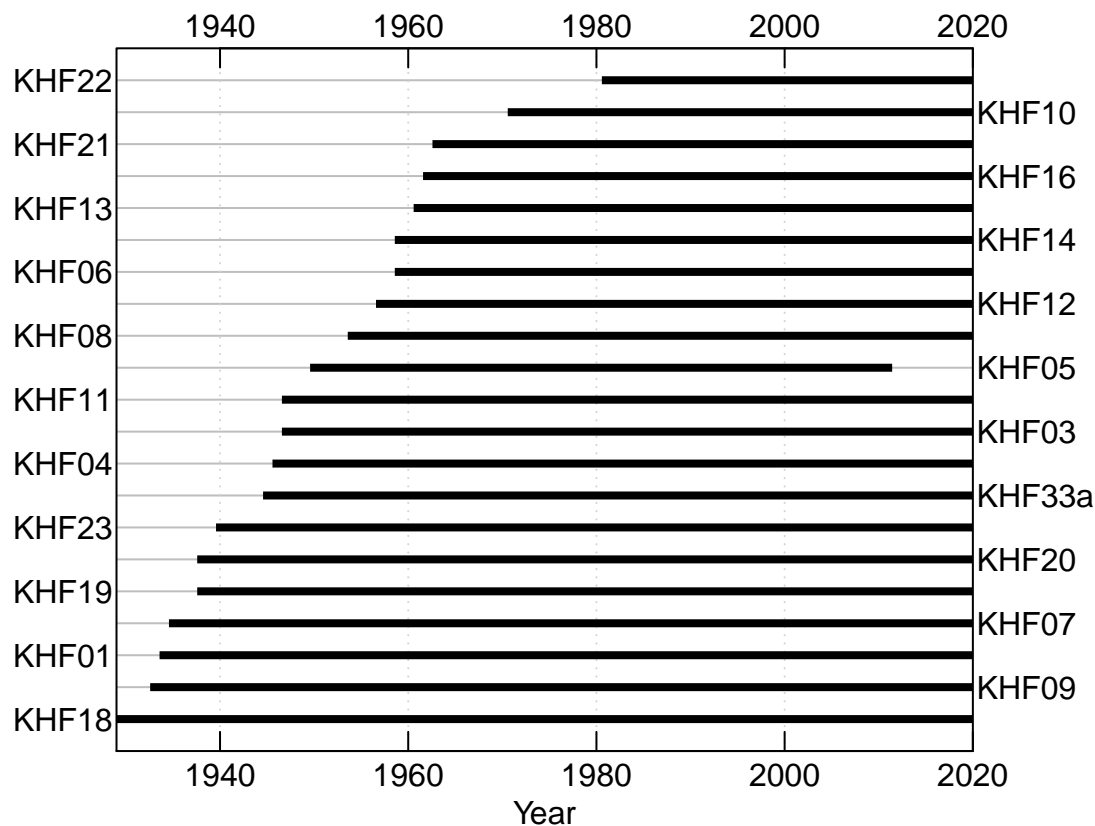
Introduction

To estimate the significance of the potential differences we decided to conduct hierarchical agglomerative clustering, HAC, of the time series. The analysis should help with the objective evaluation of the different growth patterns and indicate groups of the similar trees growth series. If those groupings are more or less clustered within two big areas, or smaller subsites, they might indicate dissimilar growth patterns within the sample.

Data import and transformation

Hierarchical clustering algorithm expects dissimilarity matrix between all tree series mostly generated by dist function. For this we first need to load all data, do individual transformations (i.e. detrending) of the tree series, and finally collect the data into one large data frame. First step is to import all the data into the list of the raw data frames.

```
library(dplR)
library(tidyverse)
files <- list.files("../Crossdated complete/", full.names = TRUE) #lists all files within the folder
raw_data <- lapply(files, function(x) { # returns the list of all raw files
  dplR::read.rawl(x) })
plot(raw_data[[1]]) # plot series from the first site using dplR/just for check
```



Second step is to detrend the series from our data frames. I have opted for the simple negative exponential here, but we can adjust later.

```
detrended_data <- lapply(raw_data, function(x) { # detrends all of the series individually with negative exponential
  x %>%
  as.data.frame() %>%
  dplR::detrend(method="ModNegExp") })
```

Finally, third step is to collect all of the detrended data frames into single data frame. We will use years for this. As the start and end years are not from the start I have to use the joining operation.

```
detrended_data_collection <- lapply(detrended_data, function(x) {
  x %>%
  rownames_to_column("Years") %>% # First make a column from rownames
  mutate(Years = as.numeric(Years)) %>% # Second transform it to numeric
  plyr::join_all() %>% # Use new column to stick all of the individual data frames together
  column_to_rownames("Years") # finally reverse the manipulation of the data frame in order to have rownames as years
```

With that we are done with data manipulation and have a data in the form we need and HAC algorithm expects.

Clustering analysis

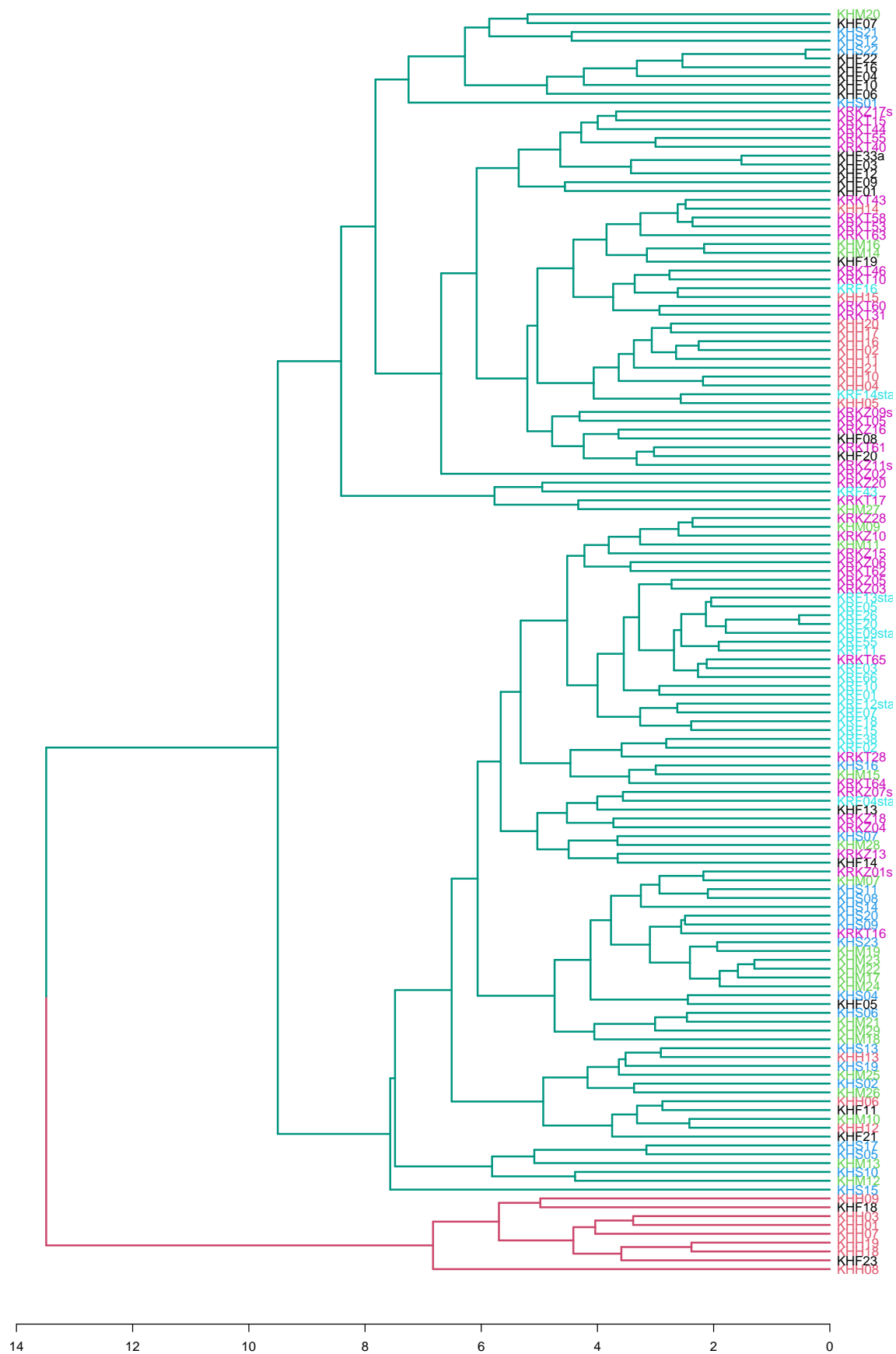
First clustering analysis. To start, we transpose the data collection so that the rows are tree series and columns years, then build distance matrix and pass it through an hclust algorithm. Finally, we save the plot as dendrogram for further colouring and annotation. Distance matrix evaluate common Euclidean distance for now and Ward process. Colour of the labels indicate the site of origin, while colour of the branches indicate two most distinct sites, i.e. two most pronounced groups...

```
library(dendextend) # for fancy dendrogram colouring

dd <- detrended_data_collection %>%
  t() %>% #transpose
  dist() %>% #create distance matrix
  hclust() %>% #algorithm tool
  as.dendrogram() #save as dendrogram

#
cols <- as.factor(substr(colnames(detrended_data_collection),1,3))
labels_colors(dd) <- as.numeric(cols)[order.dendrogram(dd)]
dd <- dd %>% colour_branches(k=2) %>% set("branches_lwd", 2)

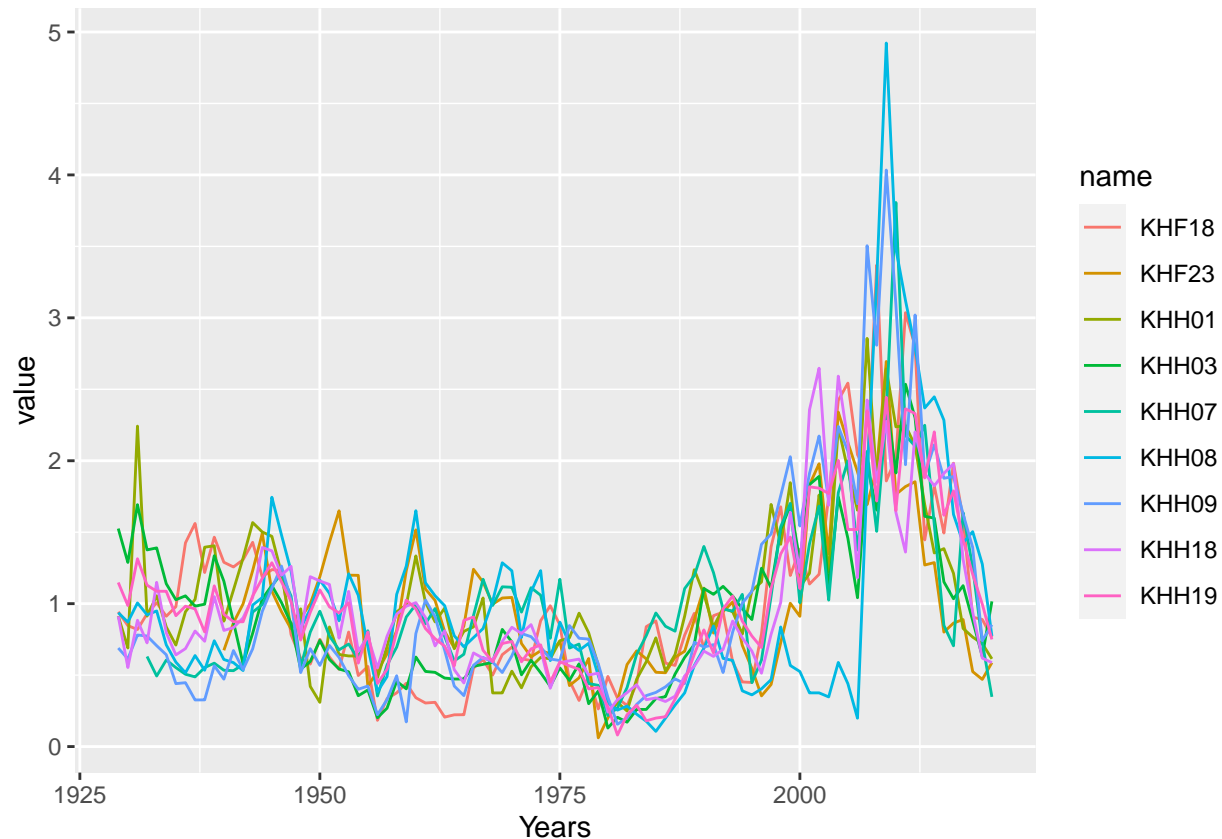
par(cex=0.8)
plot(dd, horiz=TRUE)
```



Seems that there are two very distinct clusters, however one much smaller than the other. Now what could be the reason of this? Most of the series from the small cluster (marked red in the figure above) are from 'KHH' site (additionally KHF23 and KHF18). Let's first all of the tree series from the smaller cluster:

```
detrended_data_collection %>%
  select(KHH01, KHH03, KHH07, KHH08, KHH09, KHH18, KHH19, KHF18, KHF23) %>%
  rownames_to_column("Years") %>% # First make a column from rownames
  mutate(Years = as.numeric(Years)) %>%
  pivot_longer(-Years) %>%
  ggplot(aes(x=Years, y=value, colour = name)) + geom_line()
```

Warning: Removed 14 row(s) containing missing values (geom_path).



Right. They seem to follow the similar pattern, i.e. high release from late 1970 being more pronounced after late 1990s early 2000s. To test this hypothesis let's do a second stage detrending using spline.

```
detrended_spline_data <- lapply(raw_data, function(x) { # detrends all of the series individually with
  x %>%
  as.data.frame() %>%
  dplR::detrend(method="Spline") })
```

Now repeat the manipulation of the two-stage detrending similar to before.

```
detrended_spline_data_collection <- lapply(detrended_spline_data, function(x) {
  x %>%
  rownames_to_column("Years") %>%
  mutate(Years = as.numeric(Years)) }) %>% plyr::join_all() %>% column_to_rownames("Years")
```

Clustering...

```

library(dendextend)

dd_spline <- as.dendrogram(hclust(dist(t(detrended_spline_data_collection))))

cols <- as.factor(substr(colnames(detrended_spline_data_collection),1,3))
labels_colors(dd_spline) <- as.numeric(cols)[order.dendrogram(dd_spline)]

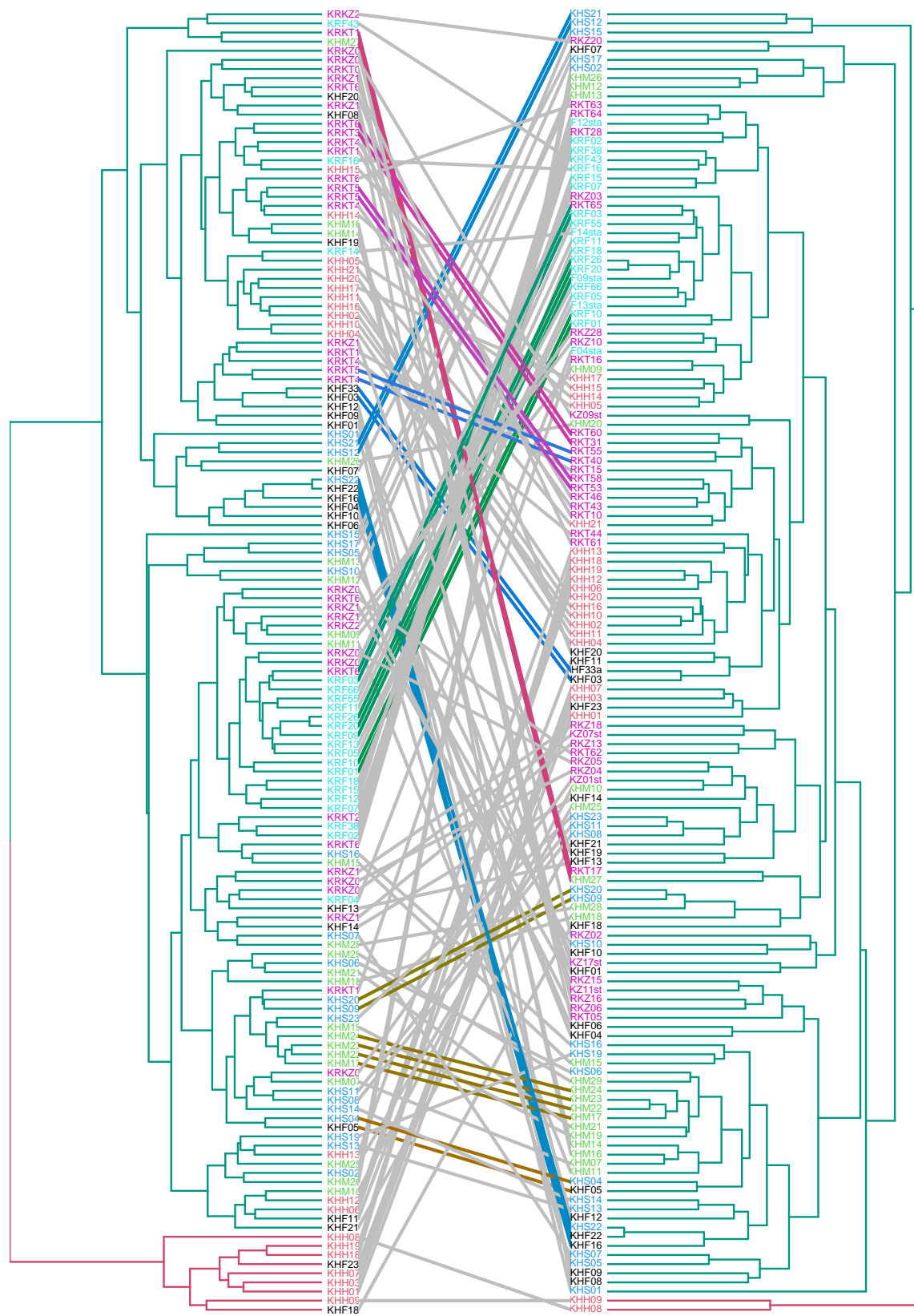
dd_spline <- dd_spline %>% colour_branches(k=2) %>% set("branches_lwd", 2)

par(cex=0.8)
plot(dd_spline, horiz=TRUE)

```


Finally let's compare two dendrograms.

```
d1 <- dendlist(dd, dd_spline)
tanglegram(d1, sort = TRUE, common_subtrees_color_lines = TRUE,
            highlight_distinct_edges = FALSE, highlight_branches_lwd = FALSE)
```

12 10 8 6 4 2 0

0 2 4 6

So it seems that there is a joint release of some, but not all trees at KHH site. Otherwise, tree series seem to be quite randomly distributed.

Discussion & some first conclusions

Without going too deep into the matter for the first analysis, it seems that tree series cluster more or less randomly together, which in turn probably indicates that there are no obviously different clusters. KHH and release event is one possible outlier, which should probably be investigated as soon as possible. Specifically, what caused the sudden release in those 7-9 trees. First step would be to investigate if there were one of large management pressure/fire/dieback events at KHH site in the late 1970s or early 1990s?

Second point is to investigate whether the pollution acted as a modifying parameter, interacting with others. On the individual scale that can probably be achieved with (non)linear mixed effect models. However, for that our sample size might actually not be enough...

Phase 2 - moving window clustering comparison

In phase 2 of the clustering analysis we wanted to do a moving window comparison of the clustering dendrograms. Basically we compare clustering dendrogram from the reference period (1930-1960 in this case) with one from periods 1931-1961, 1932-1962, ... , 1989-2019, 1990-2020. We expect the difference to constantly increase with time, as pollution depositions are accumulated. Structural breaks (sudden changes) in the series might indicate a change of the growing conditions. As a measurement of the difference between two dendrograms I am using the entanglement, which is estimating the difference in the orders of the labels of two dendrograms after searching for optimal rotation.

Important dates to remember

1960 probable pollution start at KRK

1968 environmental monitoring

1992 or 1993 sulfur filters

1995 or 1996 dust filters

```
two_period_compare <- function(period = 1931:1961, dat = detrended_spline_data_collection, ref_period =  
  #extract reference period  
  dat_ref <- dat %>%  
    rownames_to_column(var="Years") %>%  
    filter(Years %in% ref_period) %>%  
    column_to_rownames("Years")  
  
  #select the series which do not have near zero variance - removes trees that are too young for the re  
  select_vec <- as.vector(caret::nzv(dat_ref))  
  dat_ref <- dat_ref %>% select(-all_of(select_vec))  
  
  #reference period clustering and dendrogram save  
  dd_ref <- dat_ref %>%  
    t() %>%  
    dist() %>%  
    hclust() %>%  
    as.dendrogram()  
  
  #extract comparison period and select the same trees as in ref  
  dat_period <- dat %>%  
    rownames_to_column(var="Years") %>%  
    filter(Years %in% period) %>%
```

```

column_to_rownames("Years") %>%
select(-all_of(select_vec))

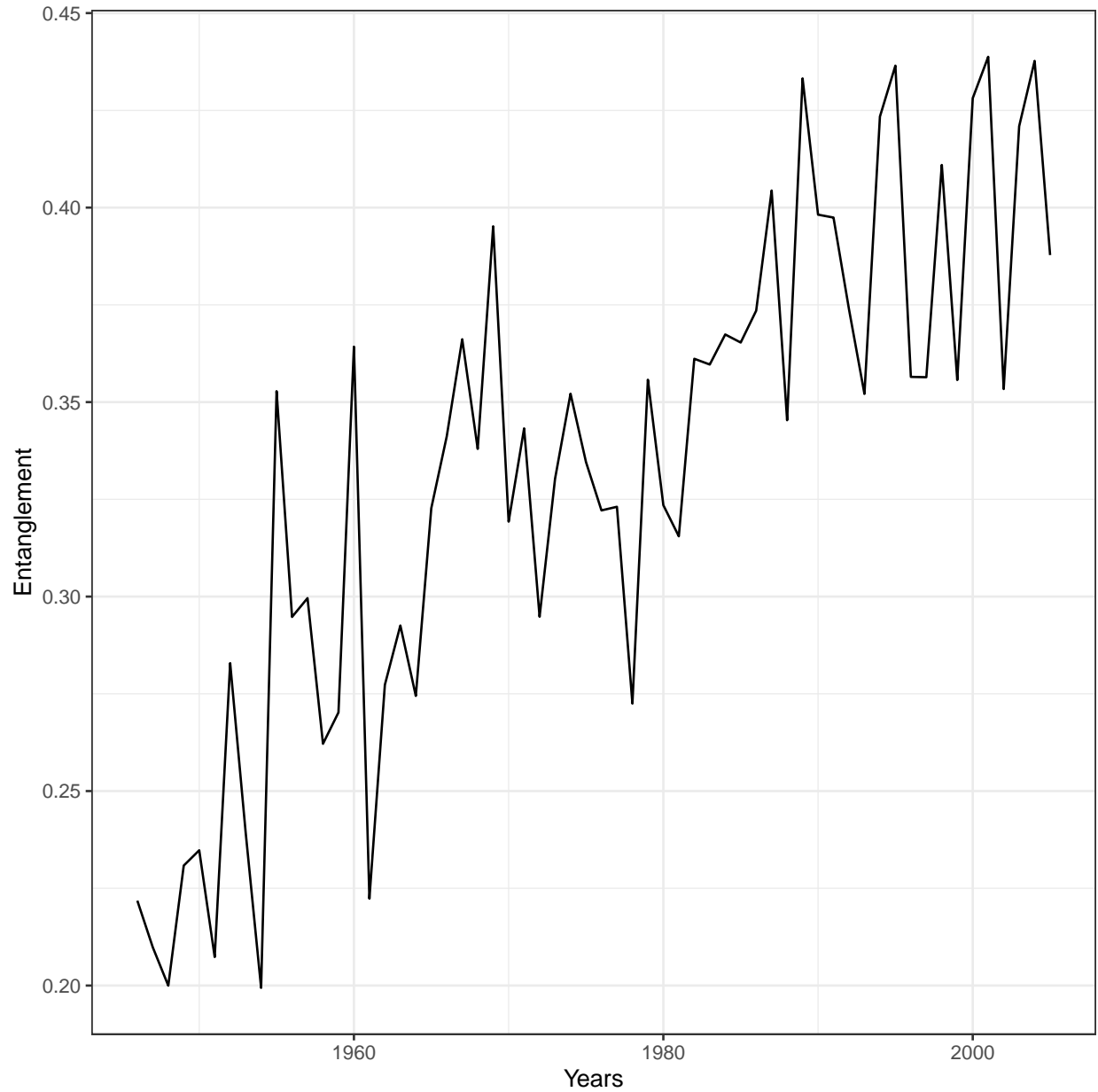
# comparison period clustering and dendrogram
dd_period <- dat_period %>%
  t() %>%
  dist() %>%
  hclust() %>%
  as.dendrogram

# entanglement estimation
dendlist(dd_ref, dd_period) %>% untangle(method="random", R=Replicates) %>% entanglement()
};

### generate periods
start <- 1931:1990
end <- start + 30
periods <- map2(start, end, seq)

### generate entanglements
set.seed(1234)
#entangles <- map_dbl(periods, two_period_compare, Replicates = 10000)
### Plot the figure
# pdf("Entanglement.pdf")
# data.frame(Years = floor((start+end)/2), Entanglement = entangles) %>%
#   ggplot(aes(x=Years, y=Entanglement)) + geom_line() +
#   theme_bw()
# dev.off()
knitr::include_graphics("Entanglement.pdf")

```



Discuasion Phase 2

The difference between cluster dendrograms has been monotonally increasing until 1990s (period 1975-2015), after which differences remained more or less constant. Therefore, we can not rule out the possibility that the increasing pollution was contributing to the increasing variability of the tree growth.

ToDo other dissimilarity measurements...