

Clustering real data with five customers and choosing optimal number of clusters

1 Clustering approach

1. *Compute quantiles of distributions across each hour of day*
2. *Compute JS distance between households for each hour of day*
3. *Total distance between households computed as sum of JS distances for all hours*
4. *Cluster using this distance with hierarchical clustering algorithm (method “complete”)*

2 Plots of raw data

Five households are considered for the following analysis from the SGSC data set. The data sets contain three columns `customer_id`, `reading_datetime` and `general_supply_kwh`. They consist of half-hourly data from 2012-2014. Figure 1 shows how the raw time plots for these five households look like. Since all the dataset looks squeezed on this linear time scale, Figure 2 shows how the raw plot looks for 2 months (Sept 2013 and Oct 2013).

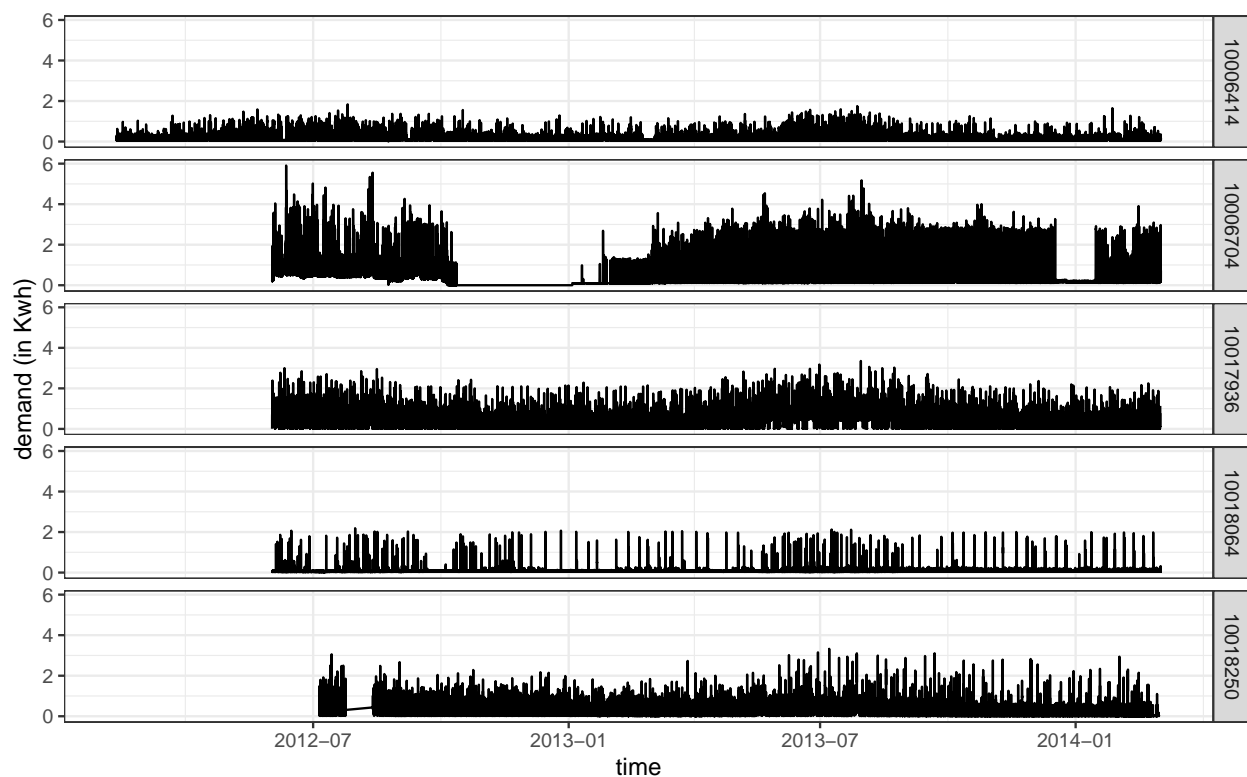


Figure 1: The raw time plots for demand shown for the entire observation period between 2012-2014 for 5 households (facets). The data is squeezed stopping us from seeing any behavioral patterns.

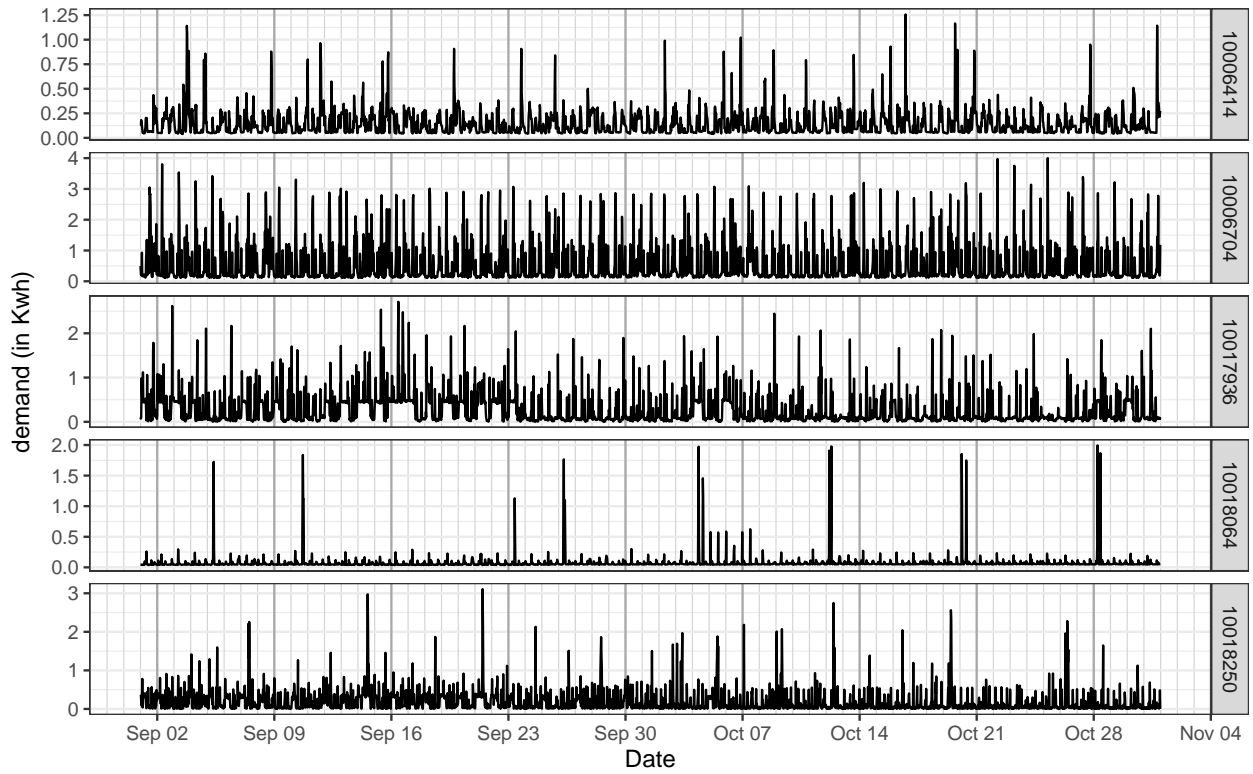


Figure 2: The raw plots for 5 a shown for Sept 2013-Oct 2013. The data is zoomed in and the y-scales made free to emphasize weekly, hourly or any behavior/patterns that they might have. There is some daily and weekly pattern in all households except the 4th one, that has spikes which seem to occur at irregular intervals.

3 Plots of distribution of data across categories of granularities

The distribution of demand across different hours of the day observed for these five households through heatmaps (Figure 4) and quantile plots (Figure 3). The following characterization is done to validate results of the clustering approach later on.

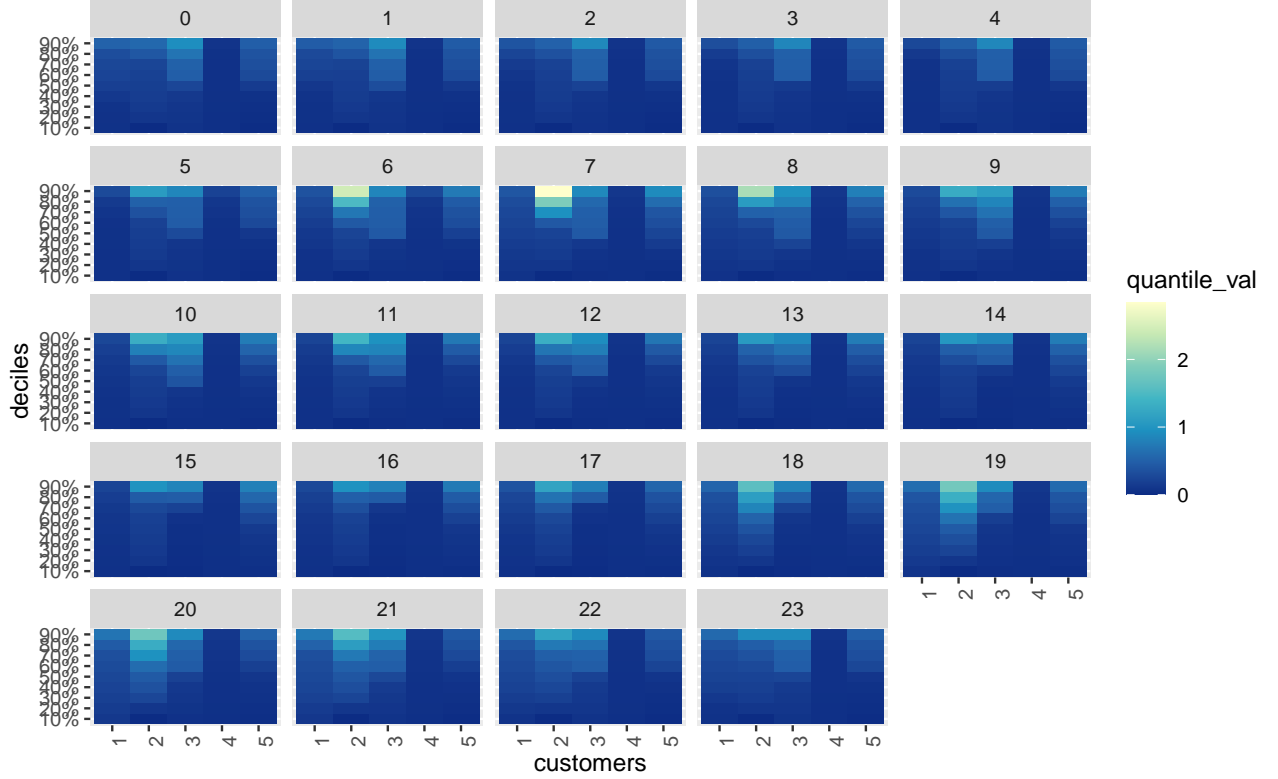


Figure 3: Heatmaps show deciles on the y-axis with customers on the x-axis with colors filled by the value of deciles and faceted by hours of the day. It again seems like all the households have gradual change in colors (almost for all hours) as we move to higher deciles except for the 4th household.

4 Iterations of each customer (adding some random noise for each iteration)

Iterations of each customer are considered by adding random noise ($N(0, \sigma^2)$, where σ^2 is very small relative to the variance of the data). Figure 5 shows that raw plots of the simulated iterations to show that their structure is same as the parent data set.

5 Dendrogram

The following figures show dendrograms when we are using optimal number of clusters (k) from $k = \text{fpc::nselectboot}()$, $k = 2, 3, 4$ and 5 . We observe that the 4th household tend to get split as different clusters as increase the number of clusters.

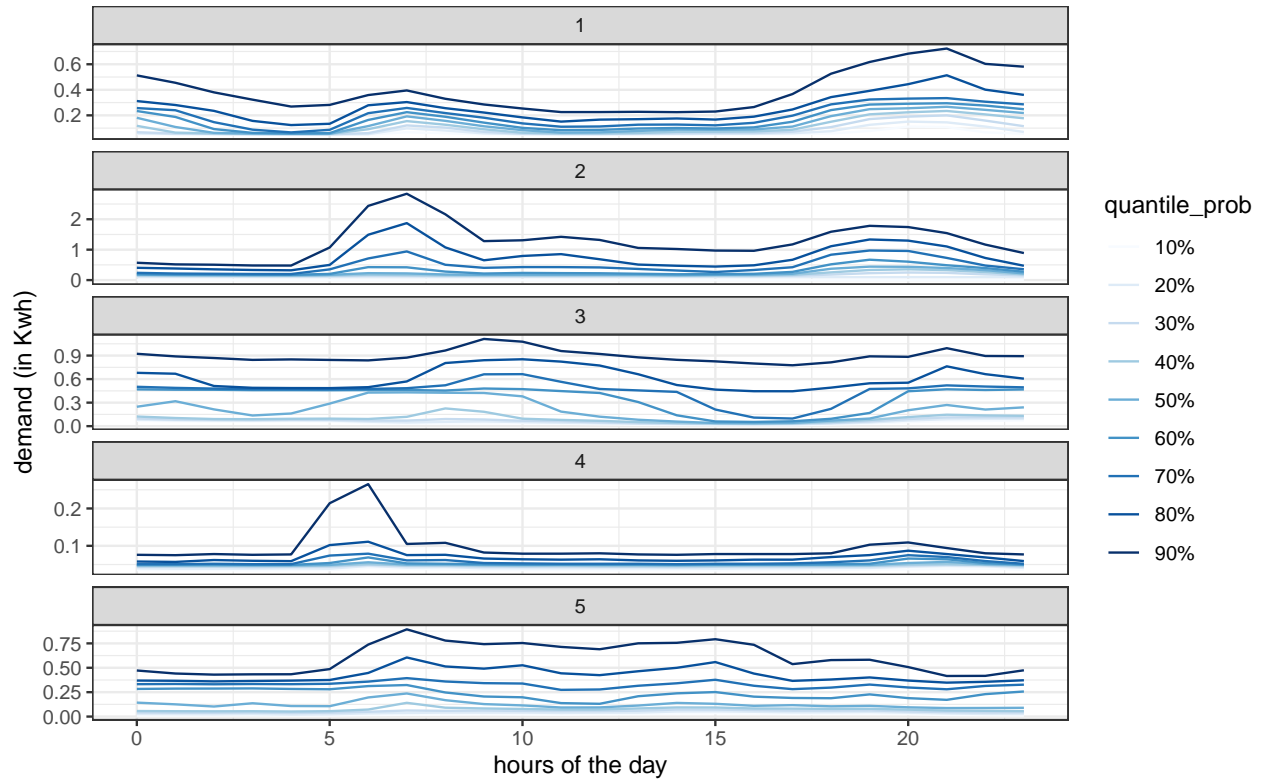


Figure 4: The deciles plots for 5 households are shown for entire observation period across hours of the day. The y-scales are made free to allow us to see daily patterns for each households. The deciles for all households tend to show an increase in the morning and evening hours, although the hours differ. Except for the 90th decile, the deciles for the 4th household look pretty flat.

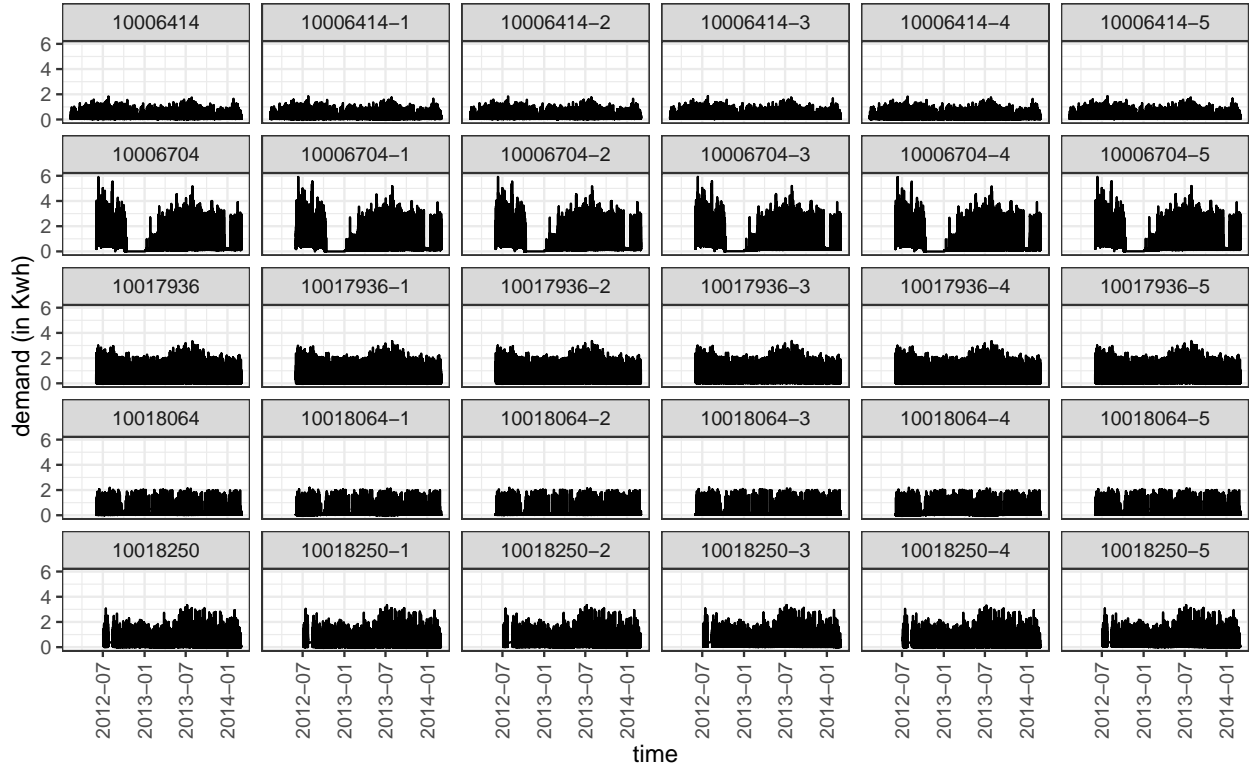
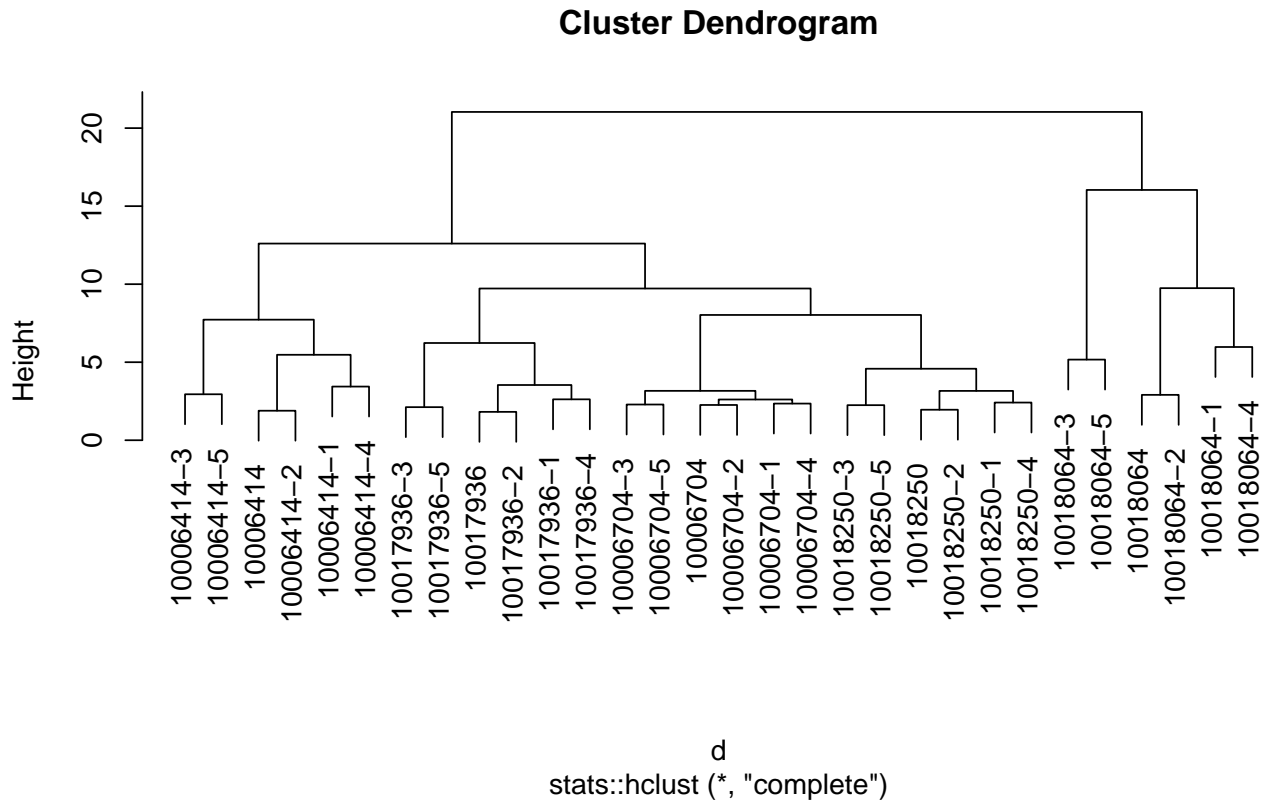
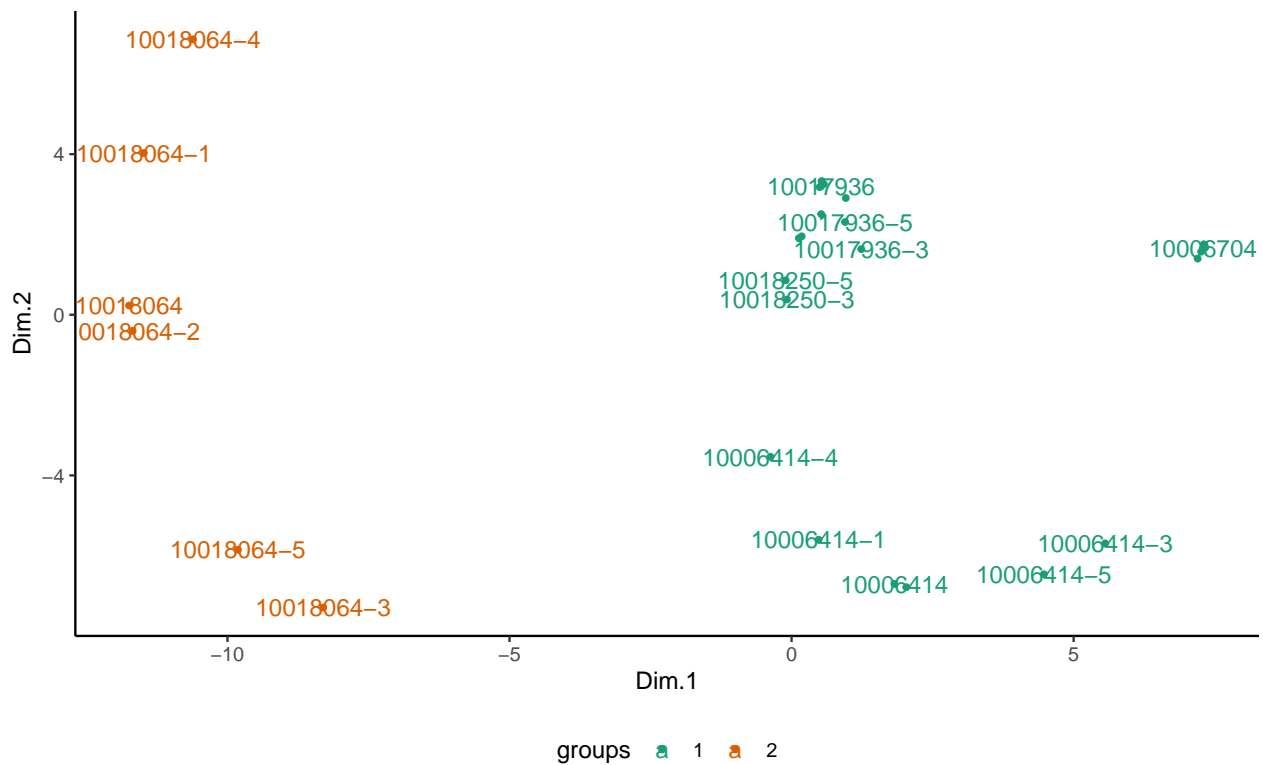


Figure 5: Five iterations for each customer are considered by adding random noise ($N(0, \sigma^2)$), where σ^2 is very small relative to the variance of the data. The raw plots for all simulated dataset is shown to make sure they look similar to the structure.



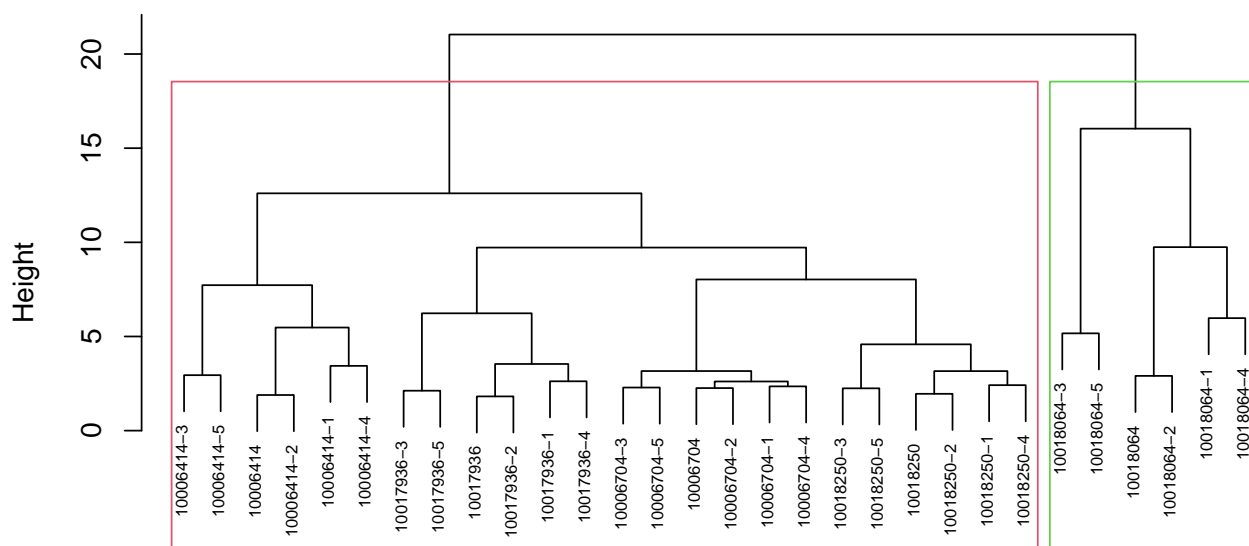
6 Multi-dimensional scaling



Optimal number of clusters as defined by the *fpc::nselectboot* is 2.

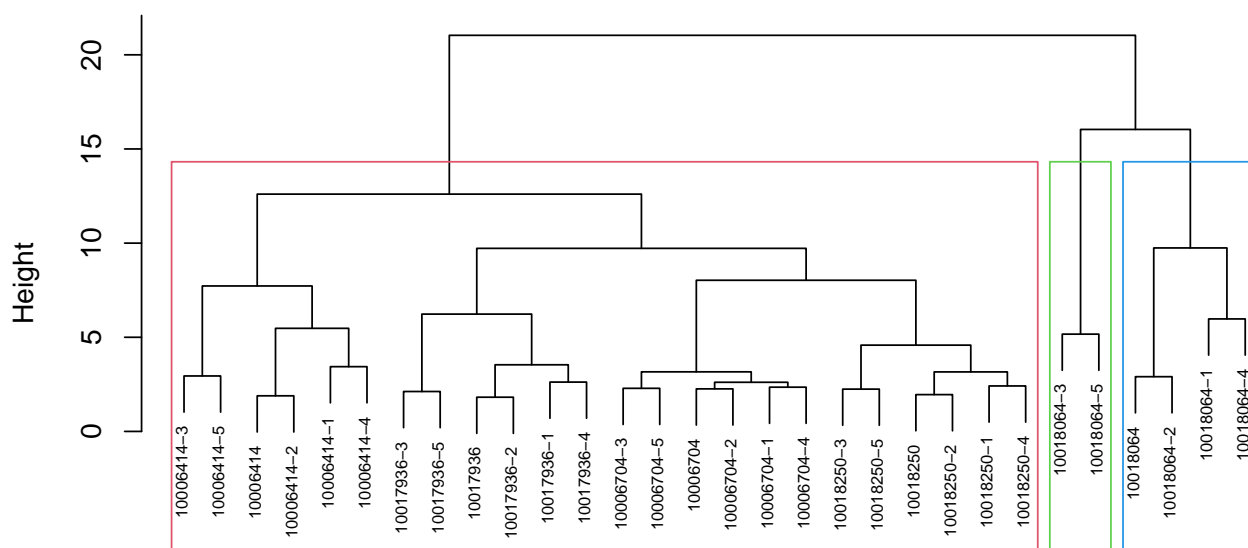
$k = 2$

Cluster Dendrogram



d
stats::hclust (*, "complete")

Cluster Dendrogram

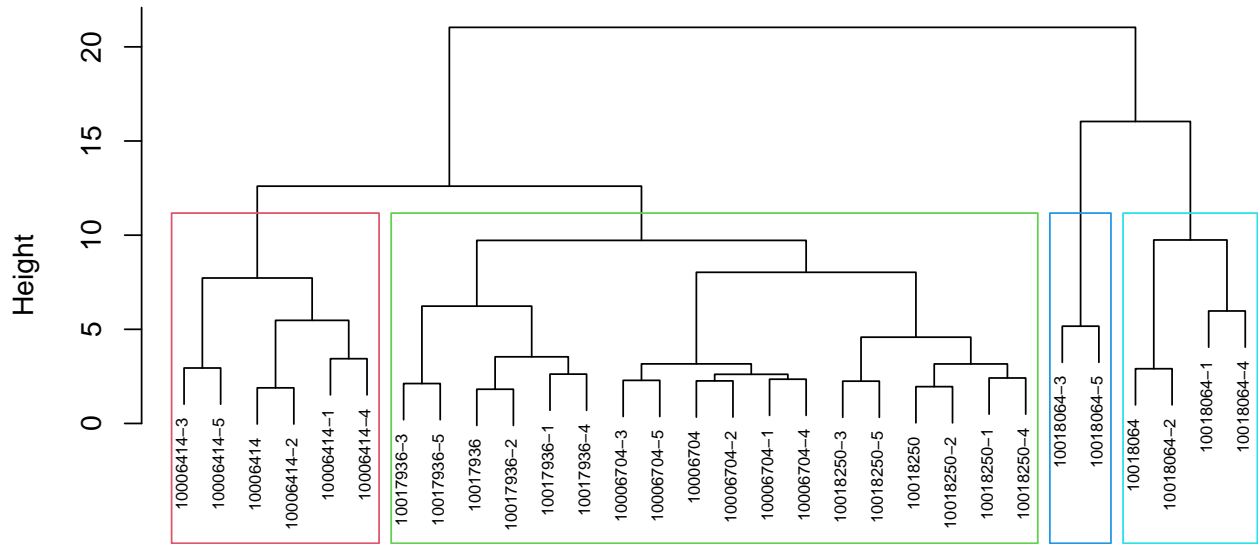


d
stats::hclust (*, "complete")

$k = 3$

$k = 4$

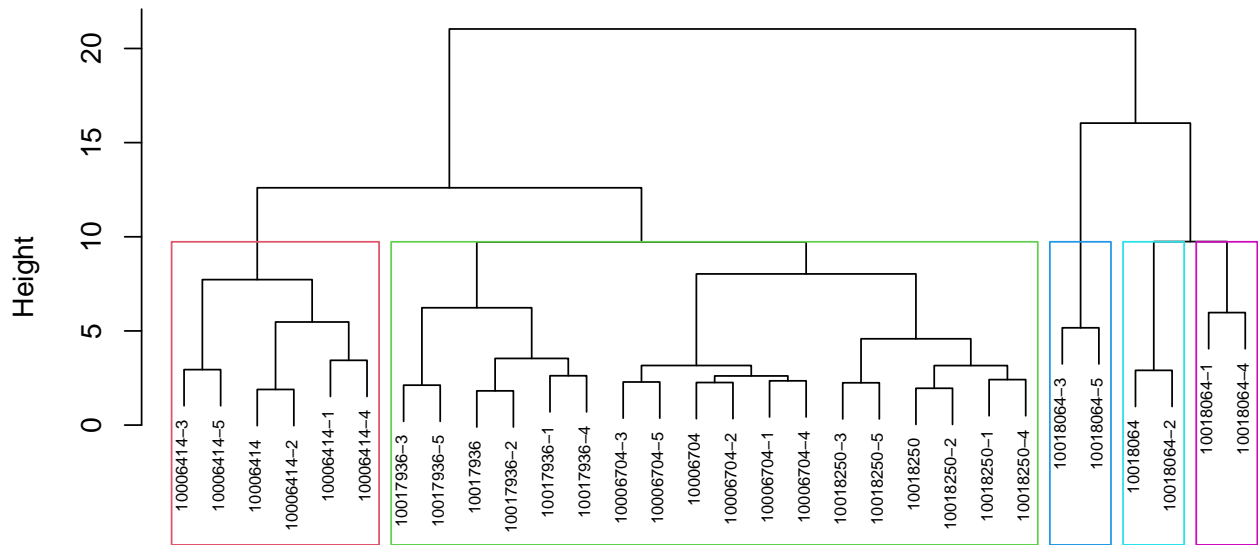
Cluster Dendrogram



d
stats::hclust (*, "complete")

$k = 5$

Cluster Dendrogram



d
stats::hclust (*, "complete")