

# Vertiefende statistische Verfahren

## 5. Übungsblatt SS 2024

Stefan Kolb, Joachim Walzl

### Allgemeine Information

Alle Aufgaben sind mit R zu lösen, wenn nicht explizit anders angegeben. Die Berechnungen sollen nachvollziehbar und dokumentiert sein. Um die vollständige Punktezahl zu erreichen, müssen alle Ergebnisse und Fragen entsprechend interpretiert bzw. beantwortet werden. Code alleine ist nicht ausreichend! Die Abgabe erfolgt über Moodle entsprechend der Abgaberrichtlinien als pdf und Rmd File. Bitte inkludieren Sie namentlich alle beteiligten Gruppenmitglieder sowohl im Bericht als auch im Source Code. Die jeweiligen Datensätze die für diese Übung relevant sind finden Sie ebenfalls in Moodle.

### 1 Clustering - Diabetes[3P]

Verwenden Sie den Datensatz `diabetes_RM.csv`. Der Datensatz enthält fünf Messungen, die an 145 nicht adipösen erwachsenen Patienten durchgeführt wurden (Beschreibung siehe UE4). Reaven und Miller [ref] wendeten in Anlehnung an Friedman und Rubin (1967) eine Clusteranalyse auf die drei primären Variablen (`insulin`, `glucose` und `sspg`) an und identifizierten drei Cluster: “normal”, “chemical” und “overt” diabetische Probanden. Die Variable `group` enthält die Klassifizierungen der Probanden in diese drei Gruppen und dient hier als Ground Truth.

- Führen Sie eine Clusteranalyse durch. Verwenden Sie eine Clusteranzahl von 3 und vergleichen Sie die Genauigkeit (gegenüber Ground Truth) folgender Cluster-Algorithmen:
  - k-means, k-medoids
  - hierarchisches Clustering
  - hierarchischer k-means
  - Modell-basiertes Clustering

```
load_source()

# Load the data
diabetes <- read.csv("diabetes_RM.csv", header = TRUE, sep = ",")
groups <- c("normal", "chemical", "overt")

# Scale the selected columns
diabetes_scaled_cols <- scale(diabetes[, c("rw", "fpg", "glucose", "insulin", "sspg")])

# Create a new data frame from the scaled columns
diabetes_scaled_df <- as.data.frame(diabetes_scaled_cols)

# Set the row names of the scaled data frame to match the row names of the original data frame
rownames(diabetes_scaled_df) <- rownames(diabetes)
```

```

# Bind the scaled columns with the unscaled columns
diabetes_scaled <- cbind(diabetes_scaled_df, group = diabetes$group)

# Perform k-means clustering
kmeans_result <- kmeans(diabetes_scaled_cols, centers = 3)

# Perform k-medoids clustering
kmedoids_result <- pam(diabetes_scaled_cols, k = 3)

# Perform hierarchical clustering
hclust_result <- cutree(hclust(dist(diabetes_scaled_cols)), k = 3)

# Perform hierarchical k-means clustering
hkmeans_result <- hkmeans(diabetes_scaled_cols, 3)

# Perform model-based clustering
mclust_result <- Mclust(diabetes_scaled_cols, G = 3)

kmeans_cluster <- kmeans_result$cluster
kmedoids_cluster <- kmedoids_result$cluster
hclust_cluster <- hclust_result
hkmeans_cluster <- hkmeans_result$cluster
mclust_cluster <- mclust_result$classification

# Create a named vector for recoding
recode_vector <- setNames(1:length(groups), groups)

# Recode the 'group' variable to the matching cluster
ground_truth <- recode(diabetes$group, !!!recode_vector)

#ground_truth <- diabetes$group

# Create a list of the cluster vectors
list_of_cluster_vectors <- list(ground_truth = ground_truth,
                                kmeans = kmeans_cluster,
                                kmedoids = kmedoids_cluster,
                                hclust = hclust_cluster,
                                hkmeans = hkmeans_cluster,
                                mclust = mclust_cluster)

kmeans_accuracy <- calculate_accuracy(diabetes_scaled, kmeans_cluster, groups)

##           [,1]      [,2]      [,3]
## [1,] 0.371128 1.7749877 4.1700617
## [2,] 1.977689 0.2475693 3.7832967
## [3,] 3.637292 3.2219678 0.6135777

```

```
## [1] 1 2 3
```

```
kmedoids_accuracy <- calculate_accuracy(diabetes_scaled, kmedoids_cluster, groups)
```

```
##           [,1]      [,2]      [,3]
## [1,] 0.2750971 1.9263863 4.731498
## [2,] 1.8825475 0.4080579 4.383376
## [3,] 3.5838033 2.9547495 1.223532
## [1] 1 2 3
```

```
hclust_accuracy <- calculate_accuracy(diabetes_scaled, hclust_cluster, groups)
```

```
##           [,1]      [,2]      [,3]
## [1,] 0.4881984 4.048269 3.9939866
## [2,] 1.1937538 2.880662 3.6054722
## [3,] 3.3203439 5.118442 0.4326986
## [1] 1 2 3
```

```
hkmeans_accuracy <- calculate_accuracy(diabetes_scaled, hkmeans_cluster, groups)
```

```
##           [,1]      [,2]      [,3]
## [1,] 0.371128 1.7749877 4.1700617
## [2,] 1.977689 0.2475693 3.7832967
## [3,] 3.637292 3.2219678 0.6135777
## [1] 1 2 3
```

```
mclust_accuracy <- calculate_accuracy(diabetes_scaled, mclust_cluster, groups)
```

```
##           [,1]      [,2]      [,3]
## [1,] 0.09247276 1.8287537 4.0804257
## [2,] 1.72876397 0.2333928 3.6906429
## [3,] 3.59758685 3.0974983 0.5223698
## [1] 1 2 3
```

```
# Put the results into a list
```

```
accuracy_results <- list(
  ground_truth = 1,
  kmeans = kmeans_accuracy,
  kmedoids = kmedoids_accuracy,
  hclust = hclust_accuracy,
  hkmeans = hkmeans_accuracy,
  mclust = mclust_accuracy
)
```

```
# Print the accuracy results
```

```
cat("k-means accuracy: ", kmeans_accuracy, "\n")
```

```
## k-means accuracy: 0.7862069
```

```
cat("k-medoids accuracy: ", kmedoids_accuracy, "\n")
```

```
## k-medoids accuracy: 0.7517241
```

```
cat("Hierarchical clustering accuracy: ", hclust_accuracy, "\n")
```

```
## Hierarchical clustering accuracy: 0.7448276
```

```
cat("Hierarchical k-means accuracy: ", hkmeans_accuracy, "\n")
```

```
## Hierarchical k-means accuracy: 0.7862069
```

```
cat("Model-based clustering accuracy: ", mclust_accuracy, "\n")
```

```
## Model-based clustering accuracy: 0.862069
```

- Welches Verfahren ist am besten geeignet?

Model-based Clustering erreicht mit etwa 86% die höchste Genauigkeit. Auf dem zweiten Platz liegen gleich auf das k-means Clustering und das hierarchische k-means Clustering mit etwa 78% Genauigkeit. Es folgen mit etwa 75% Genauigkeit das k-medoids Clustering und das hierarchische Clustering.

- Stellen Sie die Ergebnisse grafisch dar (Scatter Plot).

```
# Load the data
diabetes <- read.csv("diabetes_RM.csv", header = TRUE, sep = ",")
groups <- c("normal", "chemical", "overt")

# Scale the selected columns
diabetes_scaled_cols <- scale(diabetes[, c("rw", "fpg", "glucose", "insulin", "sspg")])

# Create a new data frame from the scaled columns
diabetes_scaled_df <- as.data.frame(diabetes_scaled_cols)

# Set the row names of the scaled data frame to match the row names of the original data frame
rownames(diabetes_scaled_df) <- rownames(diabetes)

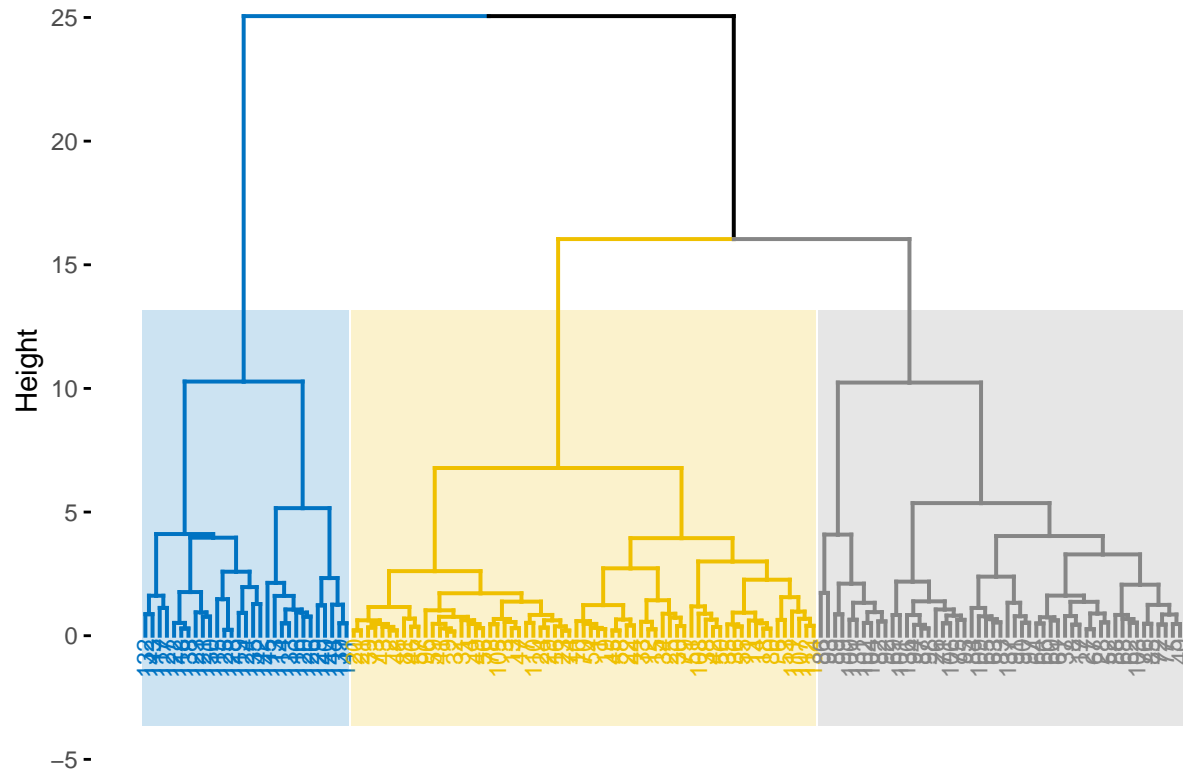
# Bind the scaled columns with the unscaled columns
diabetes_scaled <- cbind(diabetes_scaled_df, group = diabetes$group)

#-----
# Dendrogram

# Visualize the tree
fviz_dend(hkmeans_result, cex = 0.6, palette = "jco",
          rect = TRUE, rect_border = "jco", rect_fill = TRUE)
```

```
## Warning: The 'scale' argument of 'guides()' cannot be 'FALSE'. Use "none" instead as
## of ggplot2 3.3.4.
## i The deprecated feature was likely used in the factoextra package.
## Please report the issue at <https://github.com/kassambara/factoextra/issues>.
## This warning is displayed once every 8 hours.
## Call 'lifecycle::last_lifecycle_warnings()' to see where this warning was
## generated.
```

## Cluster Dendrogram

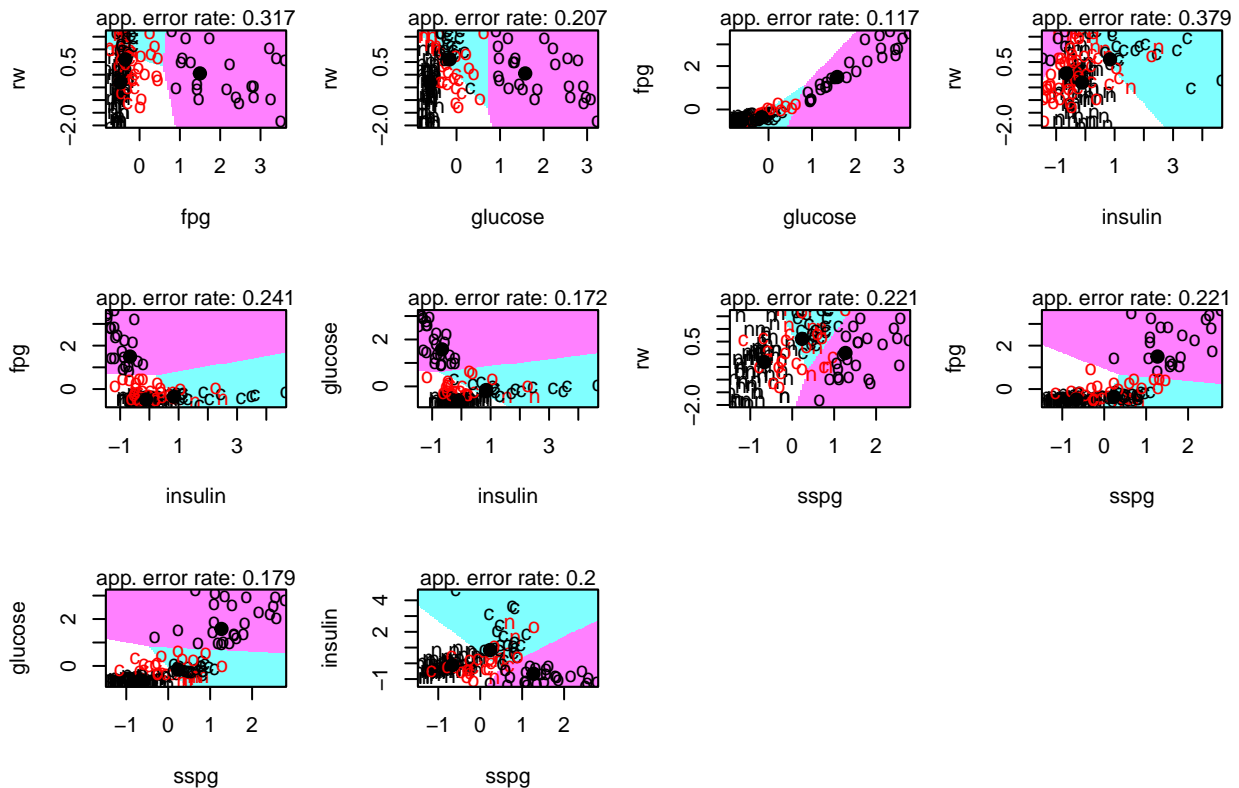


```
#-----
# Perform LDA

diabetes_scaled$group <- as.factor(diabetes_scaled$group)

# Partition plot
partimat(group ~ ., data = diabetes_scaled, method = "lda")
```

## Partition Plot



```
ml <- lda(group ~ ., data = diabetes_scaled)
print(ml)
```

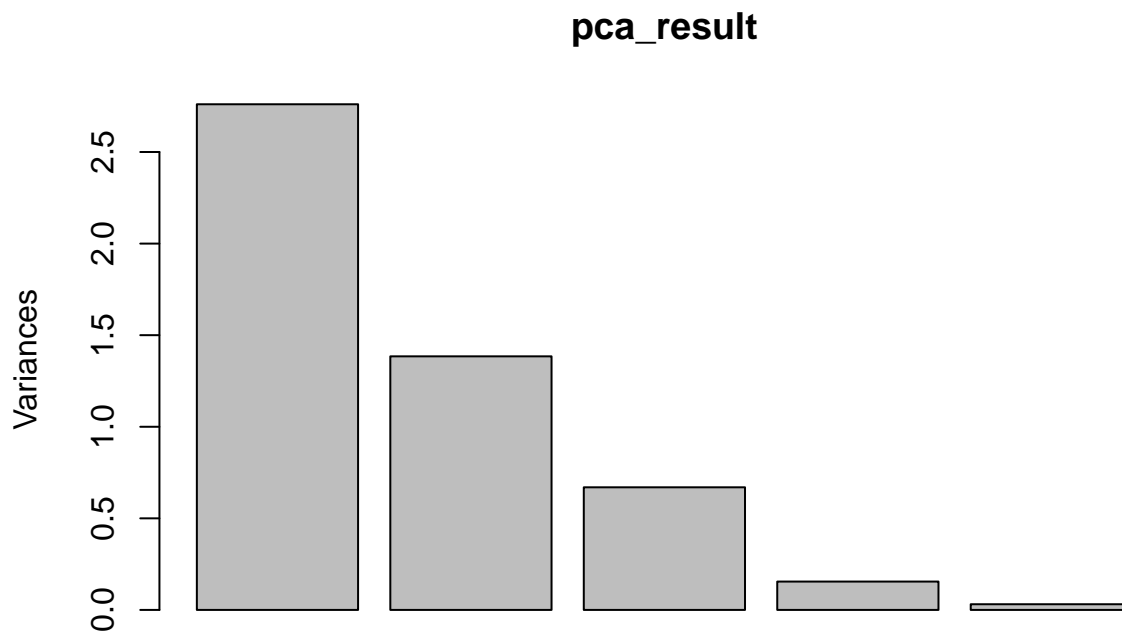
```
## Call:
## lda(group ~ ., data = diabetes_scaled)
##
## Prior probabilities of groups:
##   chemical   normal   overt
## 0.2482759 0.5241379 0.2275862
##
## Group means:
##           rw      fpg    glucose    insulin    sspg
## chemical 0.60759742 -0.3547709 -0.1567099 0.8424577 0.2335693
## normal  -0.31008189 -0.4818051 -0.6109468 -0.1114027 -0.6621427
## overt    0.05129444 1.4966346 1.5779852 -0.6624810 1.2701317
##
## Coefficients of linear discriminants:
##           LD1      LD2
## rw      0.17607468 0.4890445
## fpg     -2.15118075 -2.3419829
## glucose 3.98609901 2.2478209
## insulin -0.01236254 0.7465840
## sspg     0.44990449 -0.1202453
##
```

```
## Proportion of trace:  
##   LD1   LD2  
## 0.8812 0.1188
```

```
#-----  
# Perform PCA  
# Using subset  
diabetes_scaled <- subset(diabetes_scaled, select = -group)  
pca_result <- prcomp(diabetes_scaled, center = TRUE, scale. = TRUE)  
  
# Print summary of the PCA result  
summary(pca_result)
```

```
## Importance of components:  
##              PC1      PC2      PC3      PC4      PC5  
## Standard deviation  1.6615 1.1766 0.8181 0.39341 0.17589  
## Proportion of Variance 0.5521 0.2769 0.1339 0.03095 0.00619  
## Cumulative Proportion 0.5521 0.8290 0.9629 0.99381 1.00000
```

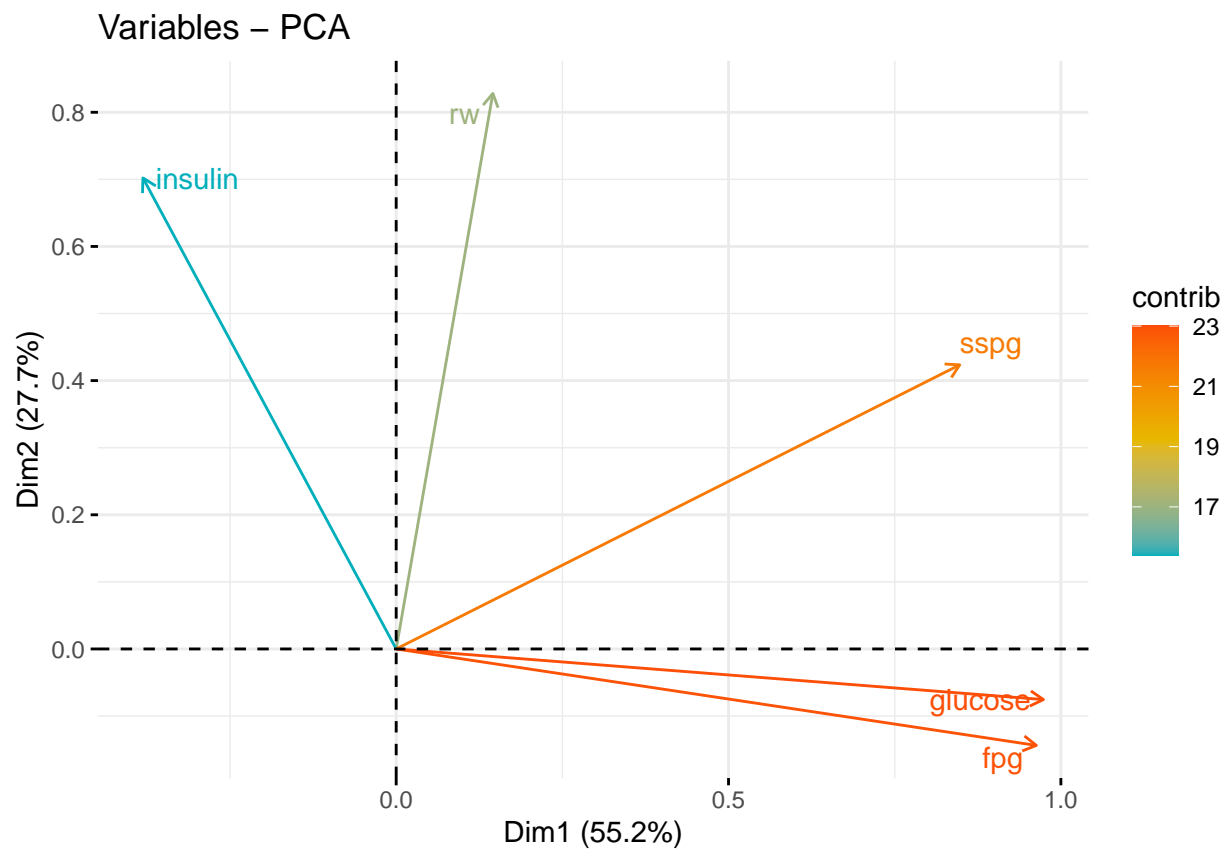
```
# Plot the variance explained by each principal component  
plot(pca_result)
```



```
# Perform PCA  
pca_result <- prcomp(diabetes_scaled[sapply(diabetes_scaled, is.numeric)])
```

```
# Plot the correlation circle
```

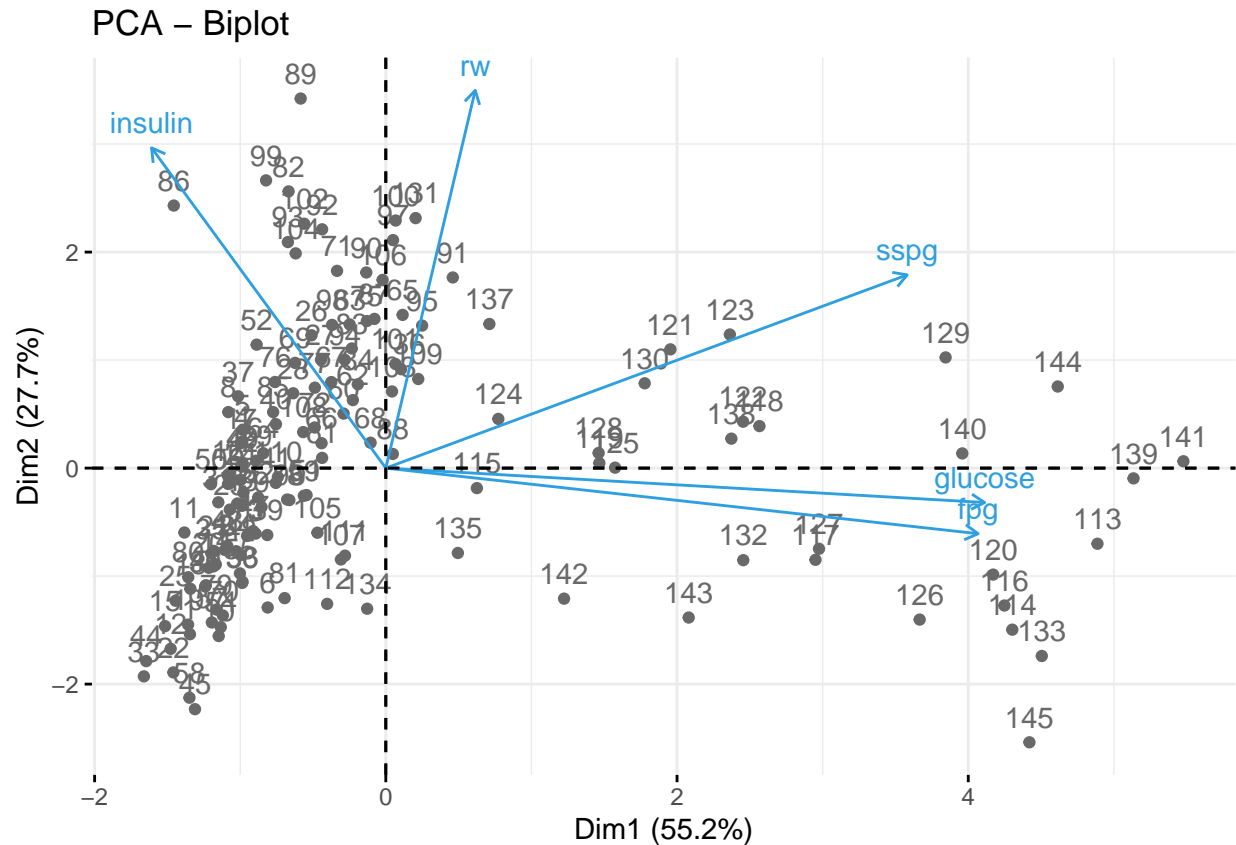
```
fviz_pca_var(pca_result, col.var="contrib", gradient.cols = c("#00AFBB", "#E7B800", "#FC4E07"), repel =
```



```
# Plot the biplot
```

```
fviz_pca_biplot(pca_result, col.var="#2E9FDF", col.ind="#696969")
```





Die LDA zeigt, dass fpg und glucose sowohl zu LD1 und LD2 am stärksten beitragen (d.h. die höchsten Absolutbeträge aufweisen). Die PCA bestätigt dies, da fpg und glucose die stärksten Beiträge zu PC1 und PC2 liefern. Deshalb werden die 2-D Scatter plots *direkt* für fpg und glucose erstellt.

```
load_source()

# Load the data
diabetes <- read.csv("diabetes_RM.csv", header = TRUE, sep = ",")
groups <- c("normal", "chemical", "overt")

# Scale the selected columns
diabetes_scaled_cols <- scale(diabetes[, c("rw", "fpg", "glucose", "insulin", "sspg")], center = TRUE)

# Create a new data frame from the scaled columns
diabetes_scaled_df <- as.data.frame(diabetes_scaled_cols)

# Set the row names of the scaled data frame to match the row names of the original data frame
rownames(diabetes_scaled_df) <- rownames(diabetes)

# Bind the scaled columns with the unscaled columns
diabetes_scaled <- cbind(diabetes_scaled_df, group = diabetes$group)

ignore <- TRUE
# Loop over the list and create a scatter plot for each set of cluster labels
```

```

for (name in names(list_of_cluster_vectors)) {
  if(name == "") {
    ignore <- FALSE
  }
  plot(create_scatter_plot_with_accuracy(diabetes_scaled,
    list_of_cluster_vectors[[name]], name, accuracy_results[[name]],
    groups,
    ignore))

  ignore <- TRUE
}

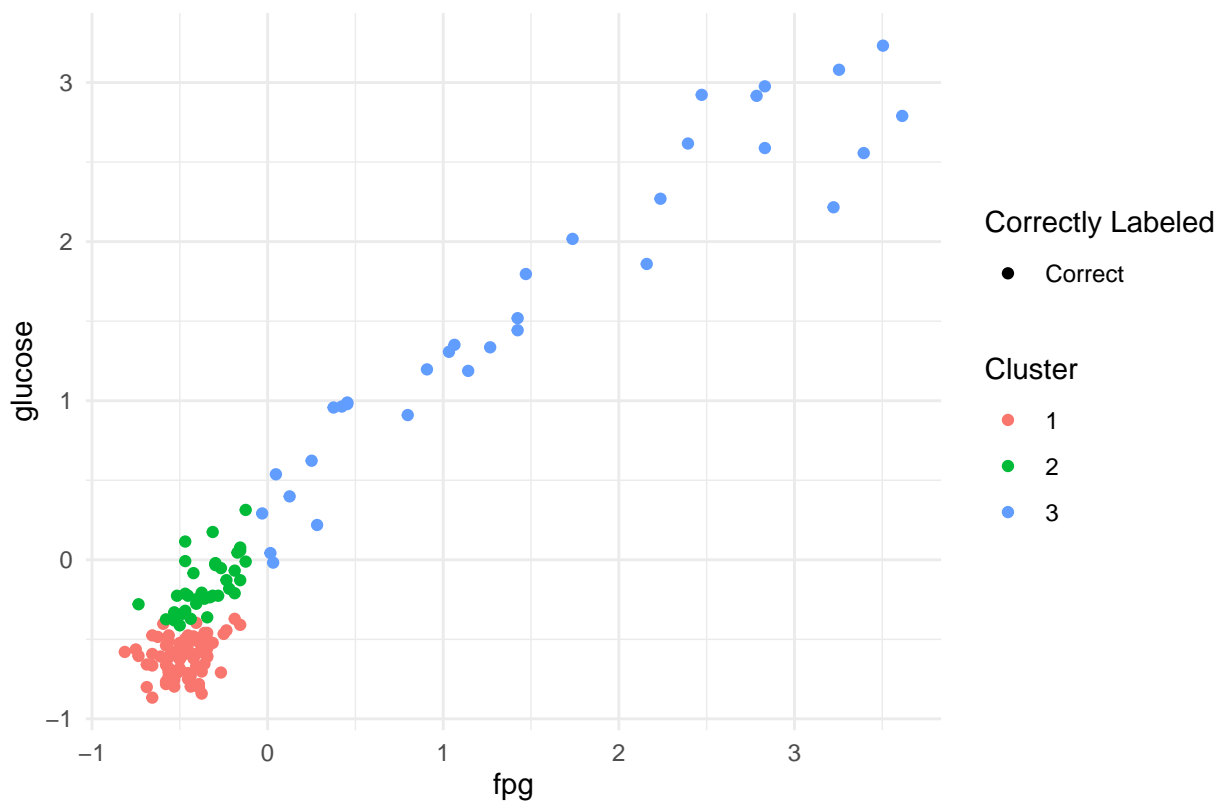
```

```

##           [,1]      [,2]      [,3]
## [1,] 0.000000 1.666359 3.587982
## [2,] 1.666359 0.000000 3.175790
## [3,] 3.587982 3.175790 0.000000
## [1] 1 2 3

```

Scatter plot of fpg vs glucose: ground\_truth (Accuracy: 100%)

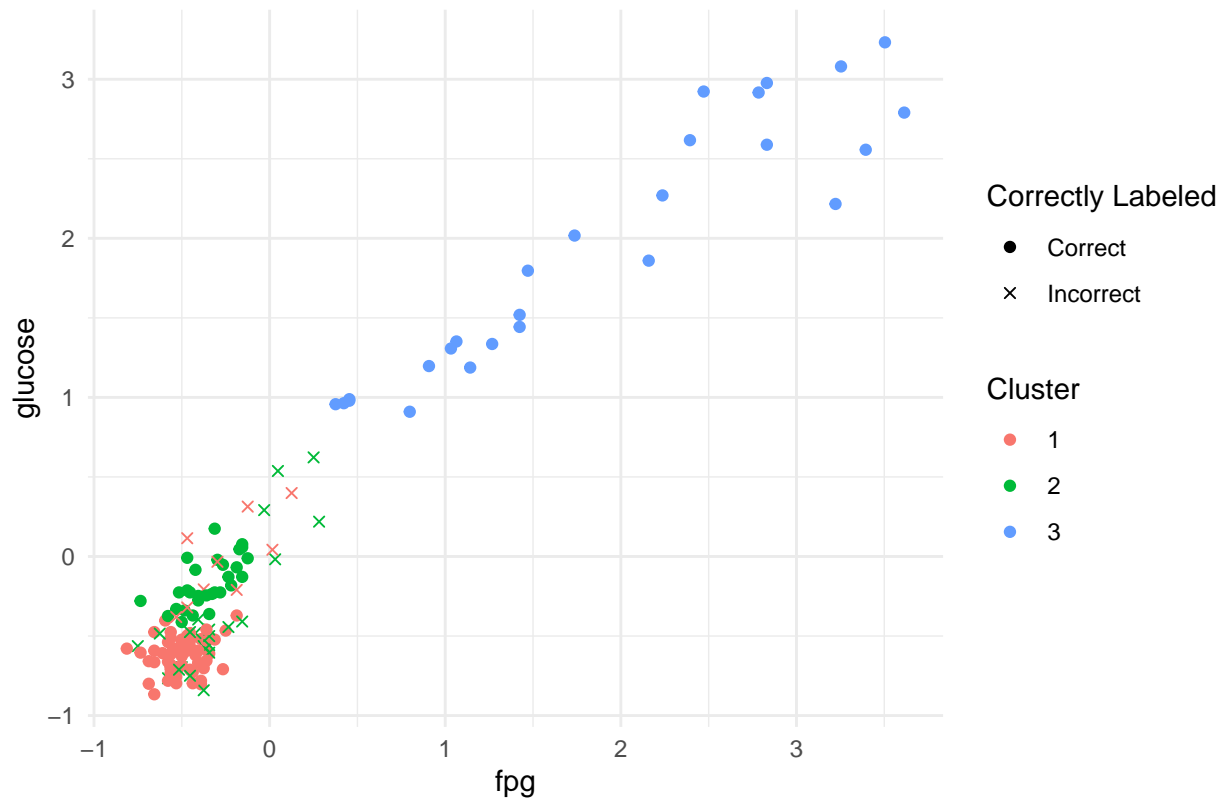


```

##           [,1]      [,2]      [,3]
## [1,] 0.371128 1.7749877 4.1700617
## [2,] 1.977689 0.2475693 3.7832967
## [3,] 3.637292 3.2219678 0.6135777
## [1] 1 2 3

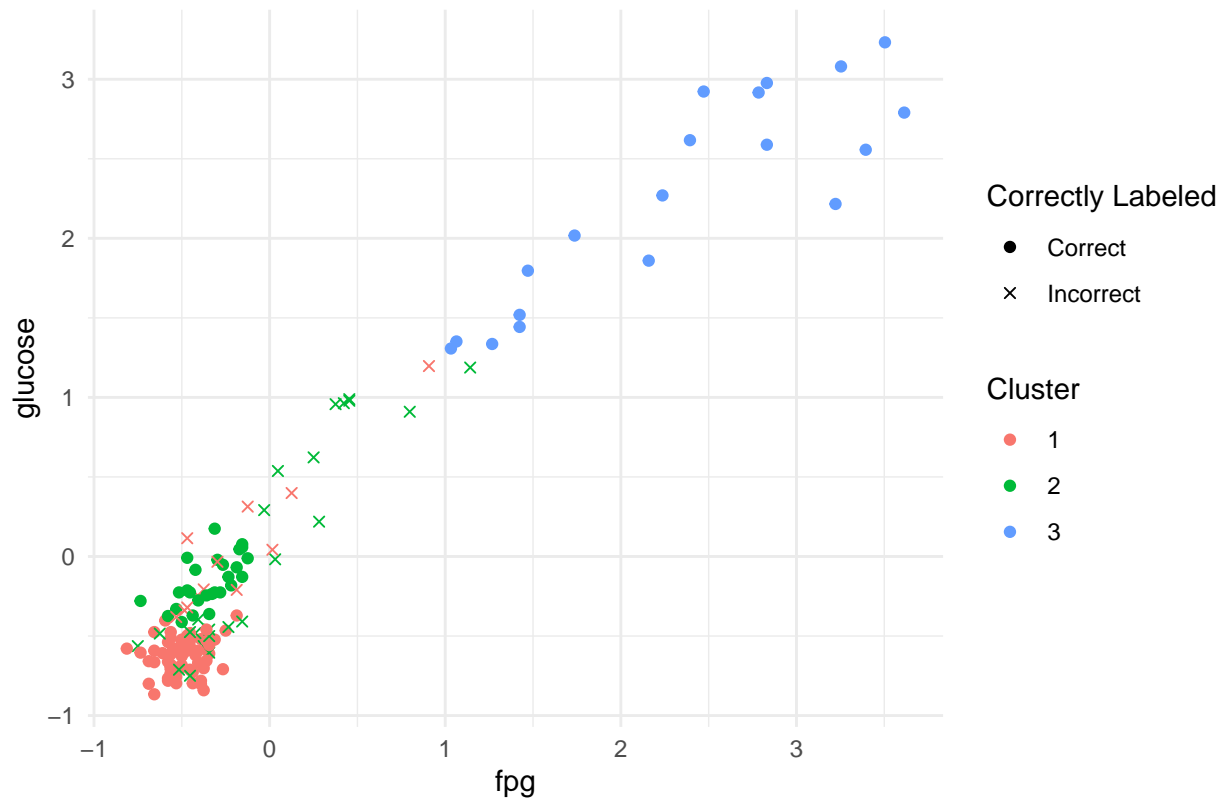
```

Scatter plot of fpg vs glucose: kmeans (Accuracy: 79%)



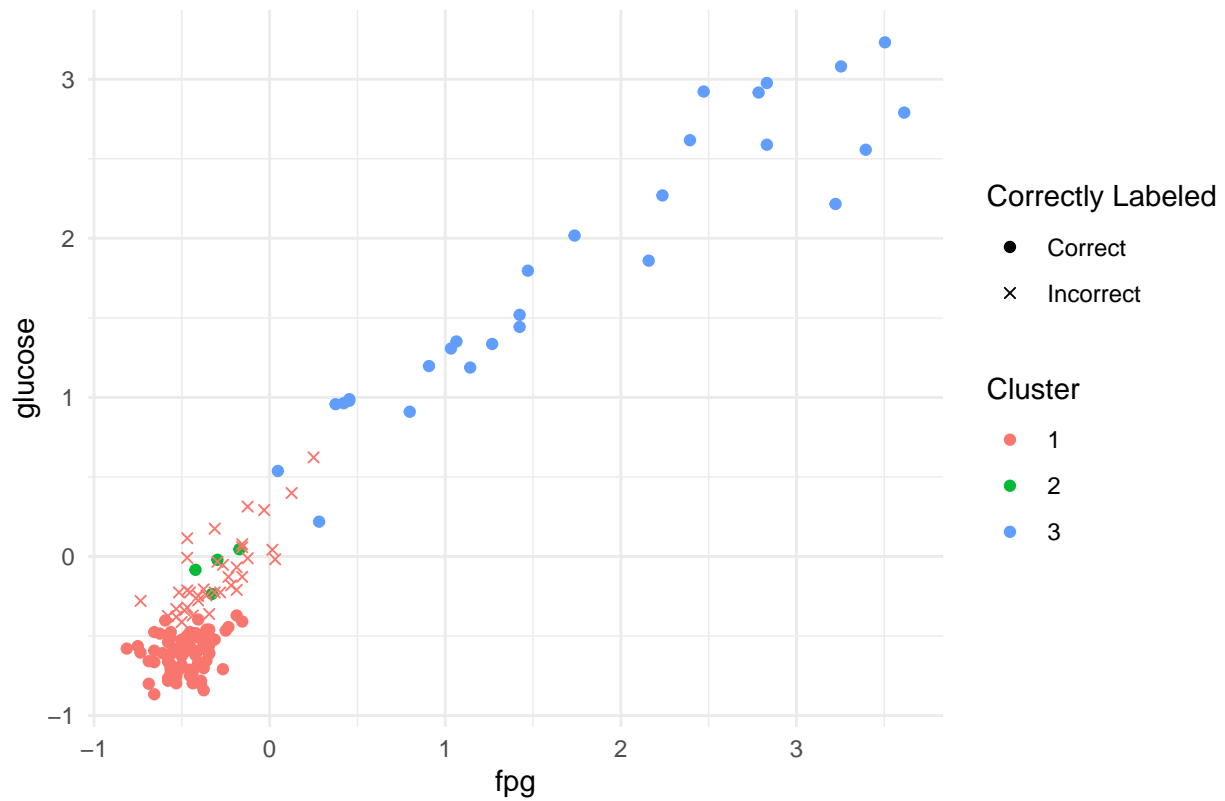
```
##          [,1]      [,2]      [,3]
## [1,] 0.2750971 1.9263863 4.731498
## [2,] 1.8825475 0.4080579 4.383376
## [3,] 3.5838033 2.9547495 1.223532
## [1] 1 2 3
```

Scatter plot of fpg vs glucose: kmedoids (Accuracy: 75%)



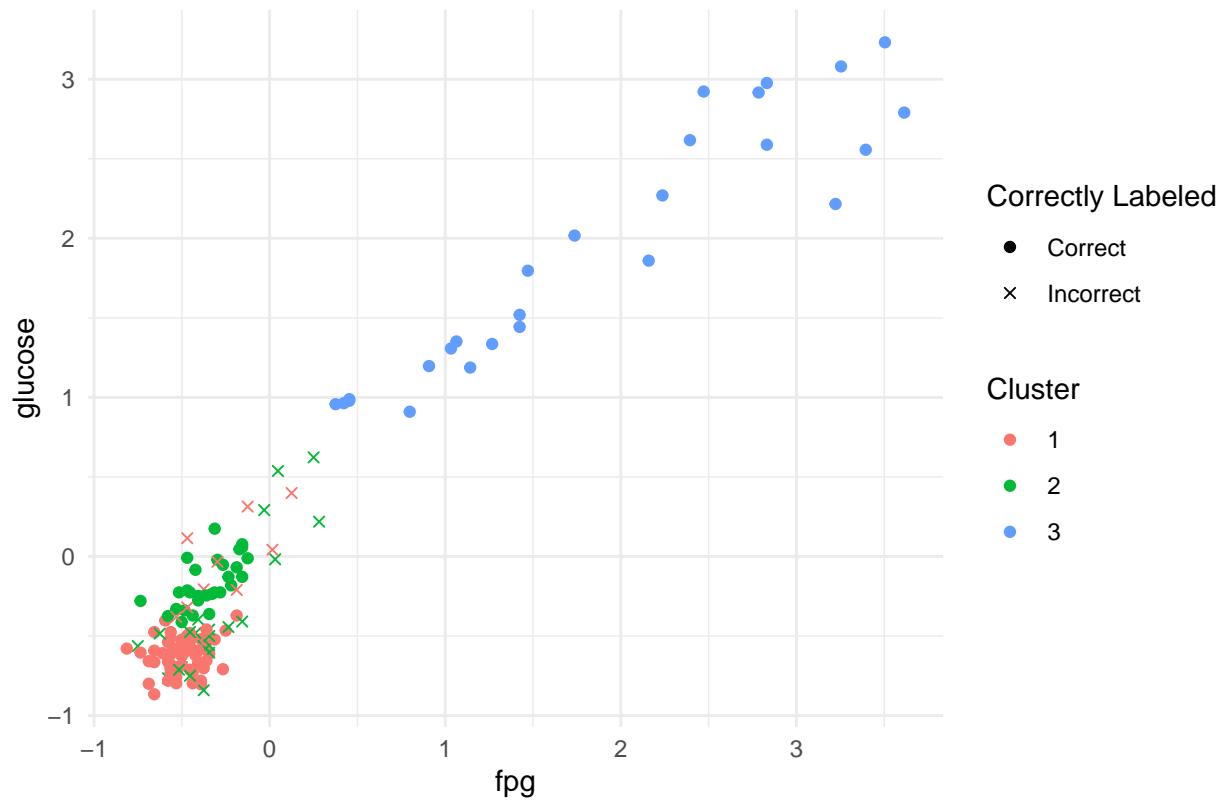
```
##          [,1]      [,2]      [,3]
## [1,] 0.4881984 4.048269 3.9939866
## [2,] 1.1937538 2.880662 3.6054722
## [3,] 3.3203439 5.118442 0.4326986
## [1] 1 2 3
```

Scatter plot of fpg vs glucose: hclust (Accuracy: 74%)



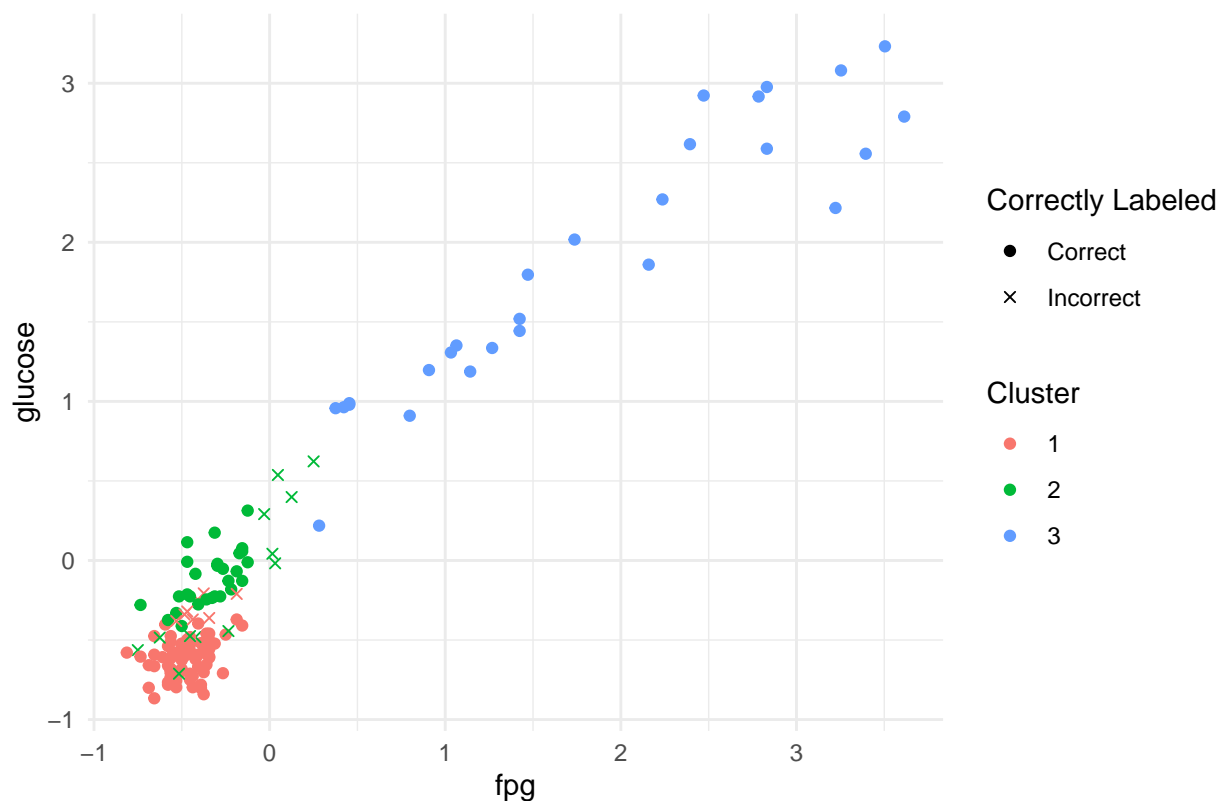
```
##          [,1]      [,2]      [,3]
## [1,] 0.371128 1.7749877 4.1700617
## [2,] 1.977689 0.2475693 3.7832967
## [3,] 3.637292 3.2219678 0.6135777
## [1] 1 2 3
```

Scatter plot of fpg vs glucose: hkmeans (Accuracy: 79%)



```
##           [,1]      [,2]      [,3]
## [1,] 0.09247276 1.8287537 4.0804257
## [2,] 1.72876397 0.2333928 3.6906429
## [3,] 3.59758685 3.0974983 0.5223698
## [1] 1 2 3
```

Scatter plot of fpg vs glucose: mclust (Accuracy: 86%)



- Finden Sie für k-means die optimale Anzahl an Cluster, und beurteilen Sie ob sich die Ground Truth Clusterstruktur reproduzieren lässt.

```
load_source()

# Load the data
diabetes <- read.csv("diabetes_RM.csv", header = TRUE, sep = ",")
groups <- c("normal", "chemical", "overt")

# Scale the selected columns
diabetes_scaled_cols <- scale(diabetes[, c("rw", "fpg", "glucose", "insulin", "sspg")])

# Create a new data frame from the scaled columns
diabetes_scaled_df <- as.data.frame(diabetes_scaled_cols)

# Set the row names of the scaled data frame to match the row names of the original data frame
rownames(diabetes_scaled_df) <- rownames(diabetes)

# Bind the scaled columns with the unscaled columns
diabetes_scaled <- cbind(diabetes_scaled_df, group = diabetes$group)

#-----
```

```

n = 6
set.seed(123)
sse <- numeric(n)
list_of_cluster_vectors <- list()
list_of_cluster_results <- list()

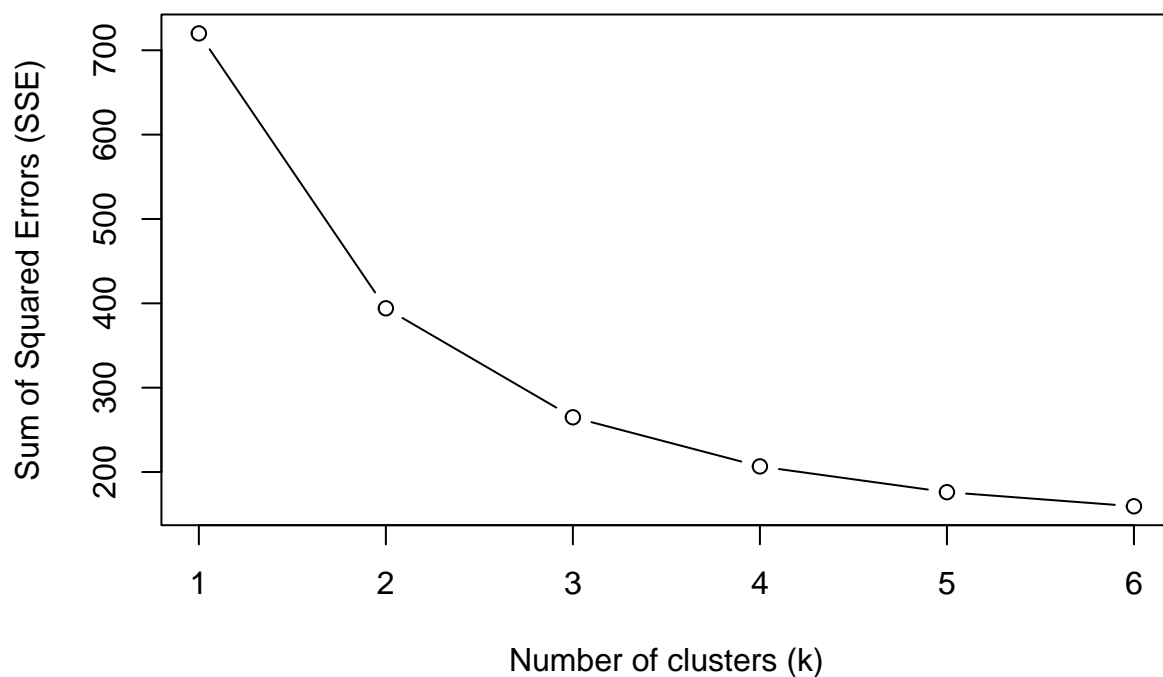
for(k in 1:n) {
  kmeans_result <- kmeans(diabetes_scaled_cols, centers = k)
  sse[k] <- kmeans_result$tot.withinss

  list_of_cluster_vectors[[k]] <- kmeans_result$cluster
  list_of_cluster_results[[k]] <- kmeans_result
}

#-----
# Plot the SSE/WSS

# Determine the optimal number of clusters using SSE
plot(1:n, sse, type = "b", xlab = "Number of clusters (k)", ylab = "Sum of Squared Errors (SSE)")

```

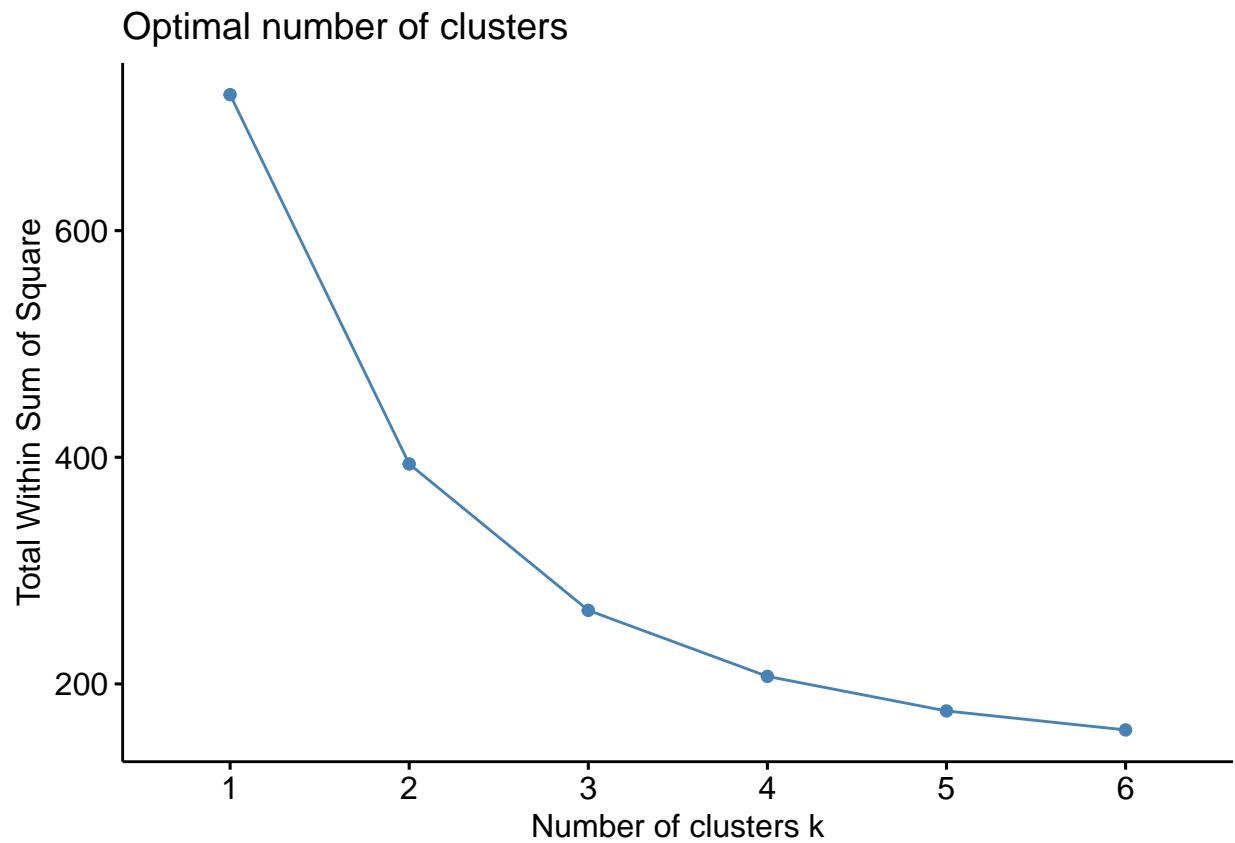


```

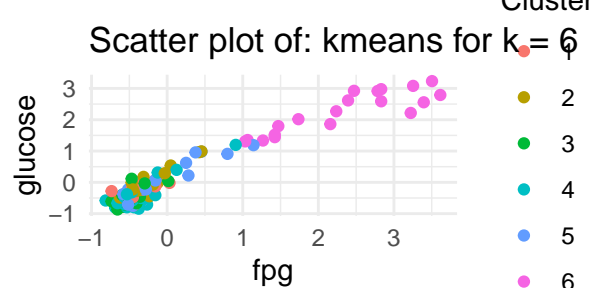
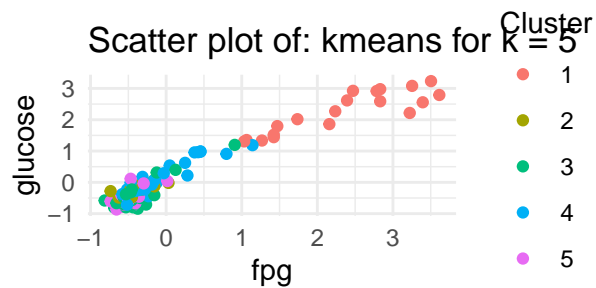
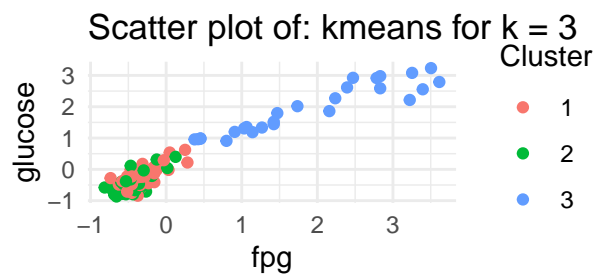
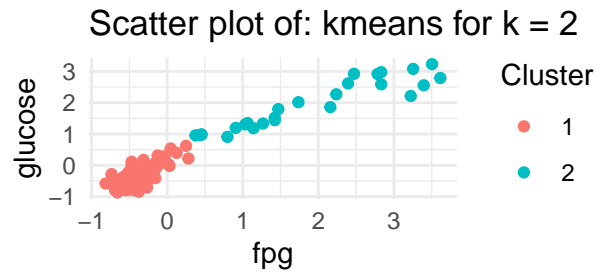
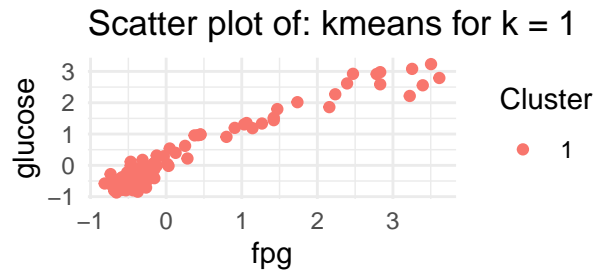
# Alternaitvely, determine the optimal number of clusters using the within-cluster sum of squares (WSS)
# Both methods should give the same result
fviz_nbclust(x = diabetes_scaled_cols, FUNcluster = kmeans, method = "wss", k.max = n)

```





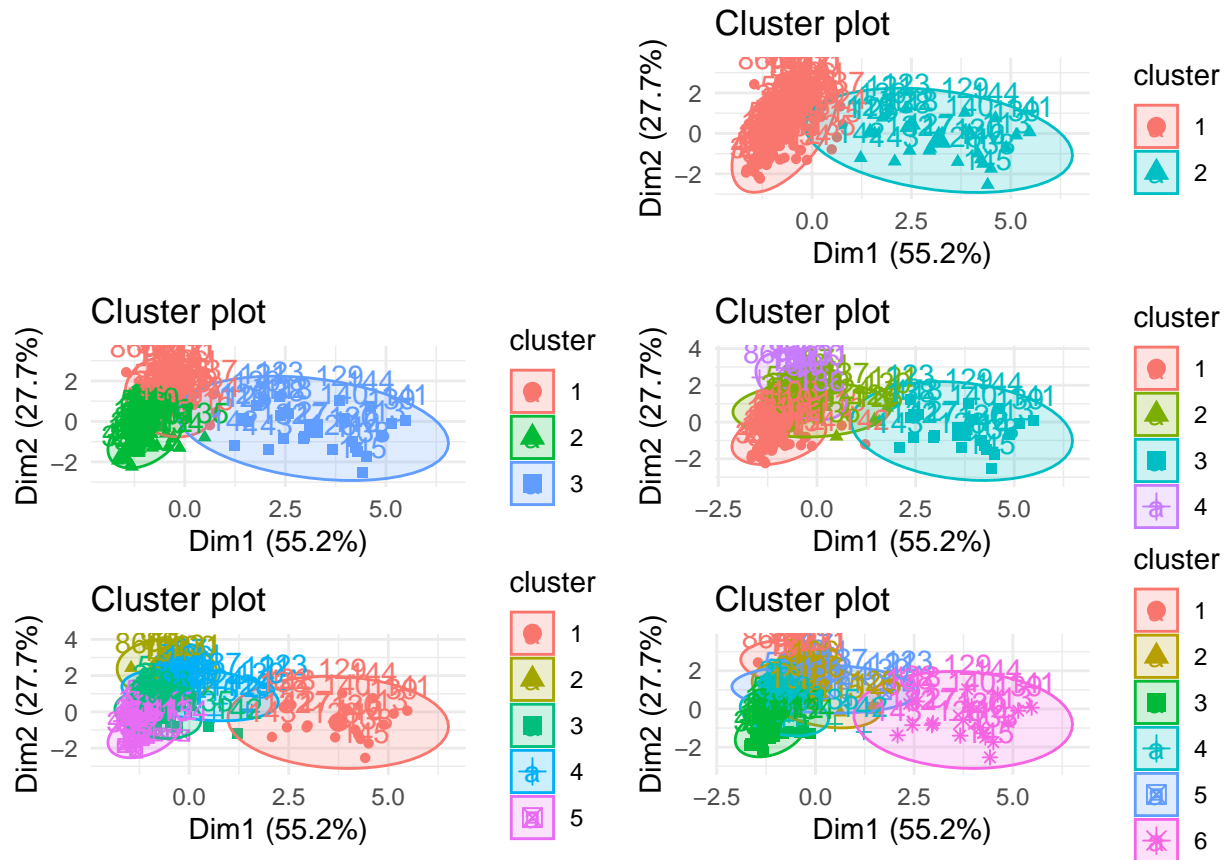
```
#-----  
# Visualize the k-means clustering results  
  
# Depiction 1  
plot_list <- list()  
  
for(k in 1:n) {  
  # Generate the plot and store it in the list  
  plot_list[[k]] <- create_scatter_plot(diabetes_scaled, list_of_cluster_vectors[[k]])  
}  
  
# Combine all the plots into a single plot  
grid.arrange(grobs = plot_list, ncol = 2)
```



```
# Depiction 2
plot_list <- list()

for(k in 2:n) {
  # Visualize kmeans clustering
  plot_list[[k]] <- fviz_cluster(list_of_cluster_results[[k]], diabetes_scaled_cols, ellipse.type = "none",
    theme_minimal()
  )
}

# Combine all the plots into a single plot
grid.arrange(grobs = plot_list, ncol = 2)
```



Für  $k = 3$  lässt sich die Ground Truth Clusterstruktur reproduzieren. Die Cluster sind klar voneinander getrennt und entsprechen visuell den Ground Truth Klassen “normal”, “chemical” und “overt”.

```
set.seed(123)

n = 6
list_of_cluster_results <- list()
list_of_cluster_vectors <- list()
silhouette_list <- list()
silhouette_score <- numeric(n)

for(k in 2:n) { # silhouette score is undefined for k = 1
  kmeans_result <- kmeans(diabetes_scaled_cols, centers = k)
  cluster_stats <- cluster.stats(dist(diabetes_scaled_cols), kmeans_result$cluster)
  silhouette_values <- silhouette(kmeans_result$cluster, dist(diabetes_scaled_cols))
  silhouette_score[k] <- mean(cluster_stats$avg.silwidth)
  list_of_cluster_results[[k]] <- kmeans_result
  list_of_cluster_vectors[[k]] <- kmeans_result$cluster
  silhouette_list[[k]] <- fviz_silhouette(silhouette_values)
}
```

```
## cluster size ave.sil.width
## 1 1 119 0.53
## 2 2 26 0.42
## cluster size ave.sil.width
## 1 1 26 0.38
```

```
## 2      2  51      0.23
## 3      3  68      0.48
##   cluster size ave.sil.width
## 1      1  26      0.37
## 2      2  52      0.30
## 3      3  23      0.31
## 4      4  44      0.36
##   cluster size ave.sil.width
## 1      1  31      0.16
## 2      2  19      0.44
## 3      3  11      0.38
## 4      4  38      0.34
## 5      5  46      0.31
##   cluster size ave.sil.width
## 1      1  10      0.35
## 2      2  19      0.46
## 3      3  42      0.34
## 4      4  34      0.33
## 5      5  22      0.17
## 6      6  18      0.20
```

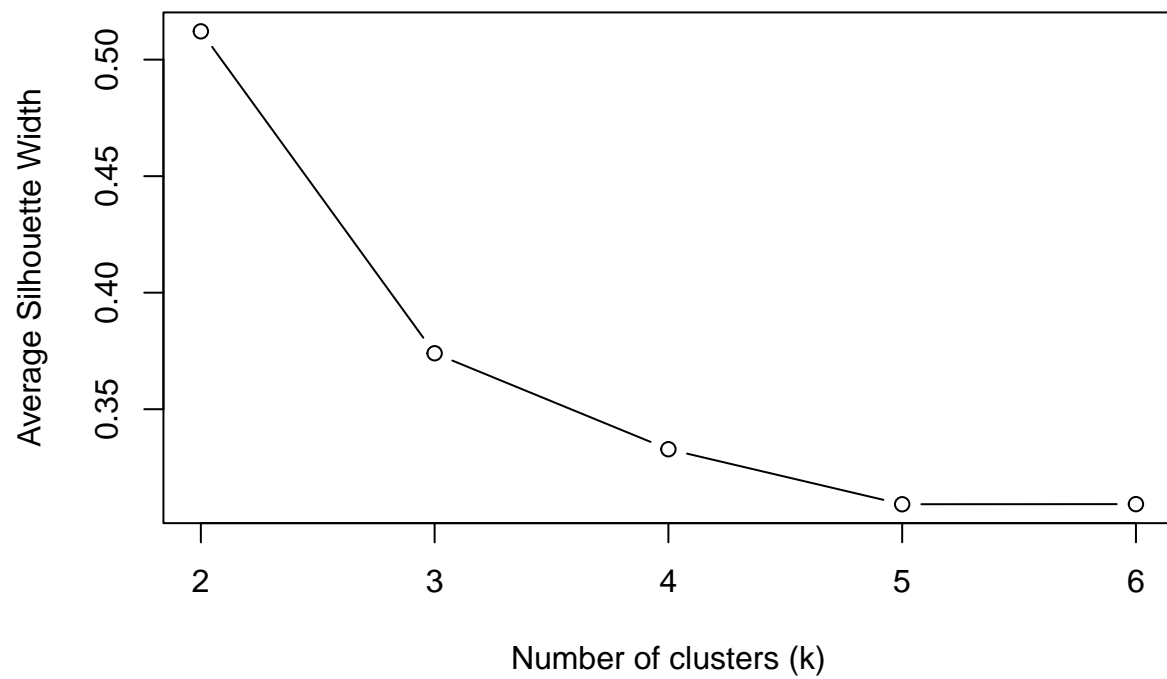
```
# Find the number of clusters that gives the highest silhouette score
best_k <- which.max(silhouette_score)
```

```
# Print the best number of clusters
print(paste("Best number of clusters (k):", best_k))
```

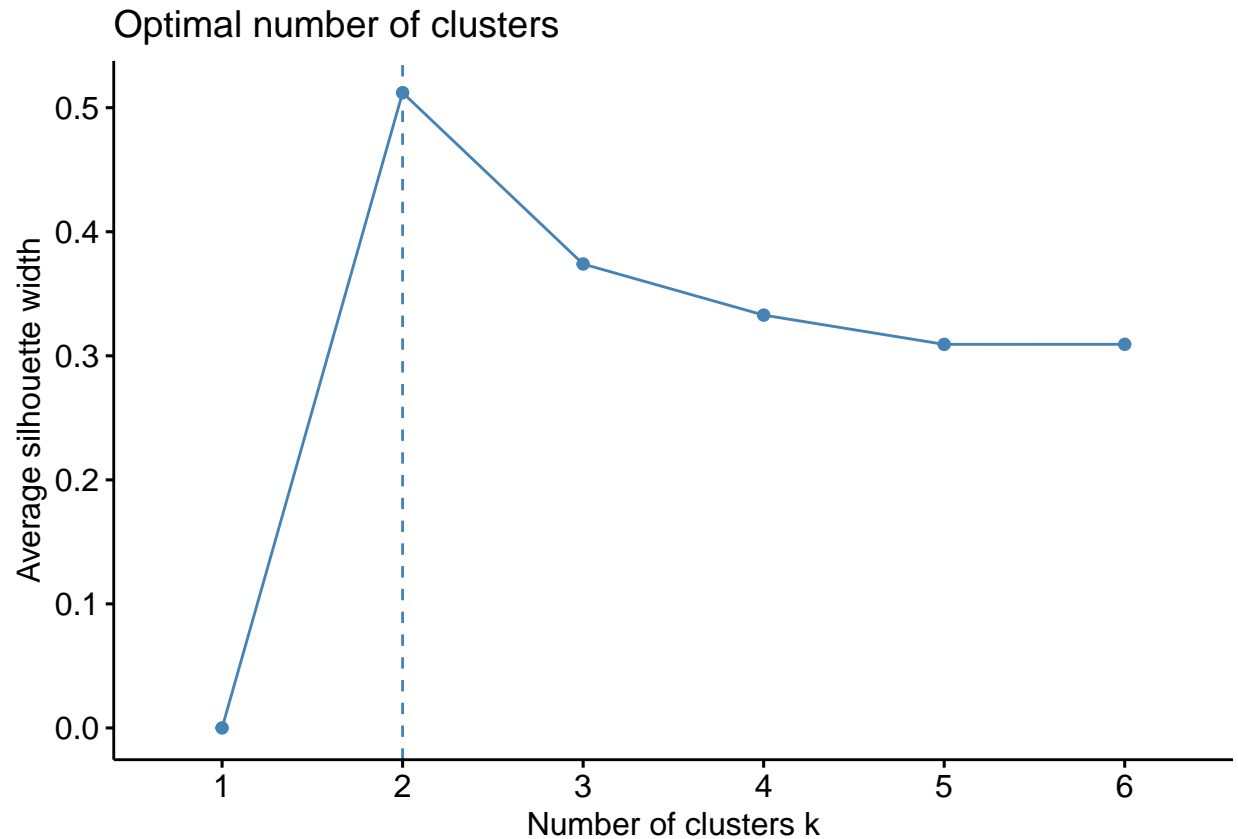
```
## [1] "Best number of clusters (k): 2"
```

```
#-----
# Plot the silhouette scores

# Determine the optimal number of clusters using the silhouette score
plot(2:n, silhouette_score[2:n], type = "b", xlab = "Number of clusters (k)", ylab = "Average Silhouette Score")
```



```
# Alternatively, determine the optimal number of clusters using the silhouette method  
# Both methods should give the same result  
fviz_nbclust(x = diabetes_scaled_cols, FUNcluster = kmeans, method = "silhouette", k.max = n)
```



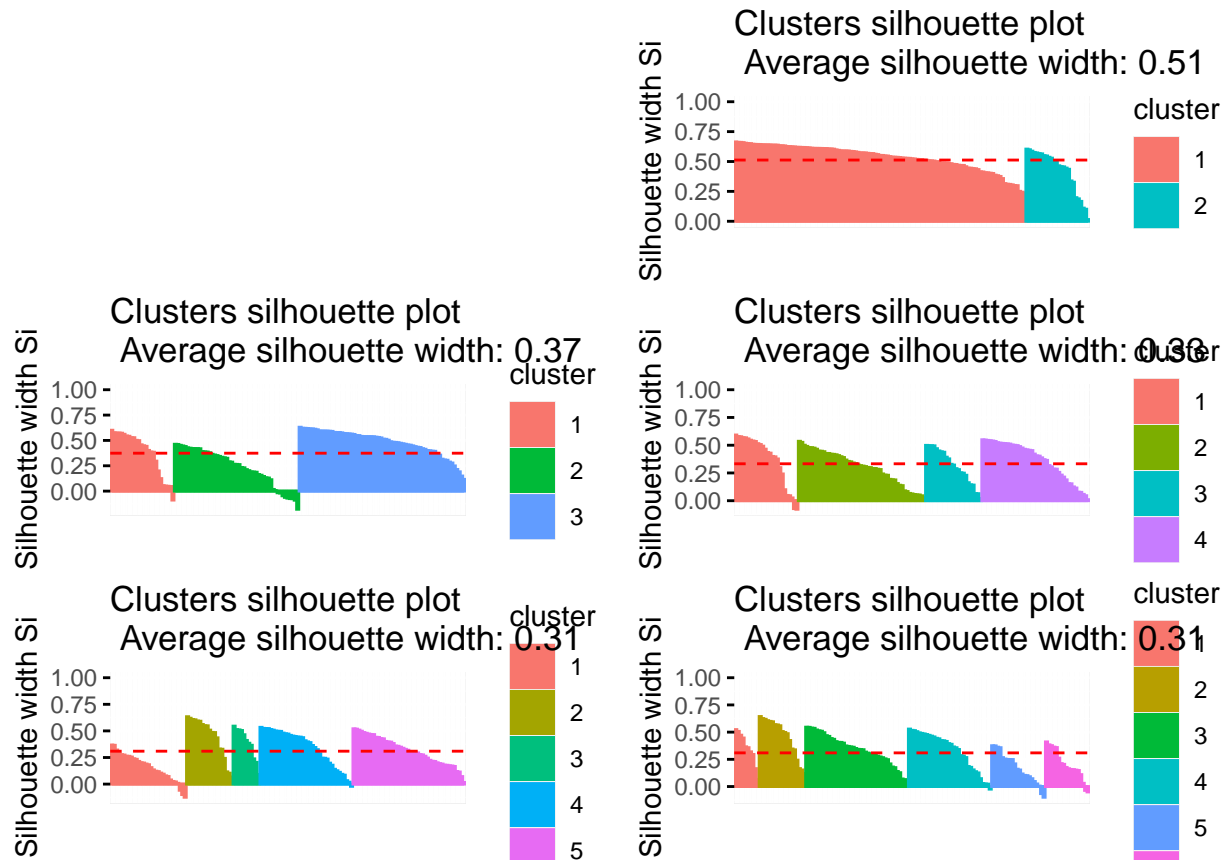
```
#-----
# Silhouette plots

# Initialize an empty list to store the plots
plot_list <- list()

for(k in 2:n) {
  # Visualize kmeans clustering
  # Print the silhouette plots for each number of clusters
  plot_list[[k]] <- silhouette_list[[k]]
  score <- round(cluster.stats( dist(diabetes_scaled_cols), list_of_cluster_vectors[[k]])$avg.silwidth, 2)
  print(paste("For k = ",k," Silhouette Score: ", score))
}
```

```
## [1] "For k = 2 Silhouette Score: 0.51"
## [1] "For k = 3 Silhouette Score: 0.37"
## [1] "For k = 4 Silhouette Score: 0.33"
## [1] "For k = 5 Silhouette Score: 0.31"
## [1] "For k = 6 Silhouette Score: 0.31"
```

```
# Combine all the plots into a single plot
grid.arrange(grobs = plot_list, ncol = 2)
```



Bestimmung des optimalen k-Werts für k-Means-Clustering mit Hilfe des des Silhouettenkoeffizienten. Laut dem Silhouettenkoeffizienten ist der optimale k-Wert 2 (0.51). Die Elbow-Methode spricht ebenso, dass der optimale k-Wert eher 2 ist. Hätten wir keine Labelinformation zur Verfügung, würde wir den k-Wert 2 wählen.

## 2 Clustering - Breast Cancer [4P]

Brustkrebs ist weltweit die häufigste bösartige Erkrankung bei Frauen und eine der Hauptursachen für krebsbedingte Todesfälle sowohl in Entwicklungs- als auch in Industrieländern. Verwenden Sie den Datensatz `breast_cancer.csv`. Die Merkmale werden aus einem digitalisierten Bild eines Feinnadelaspirats einer Brustmasse berechnet. Sie beschreiben Merkmale der im Bild vorhandenen Zellkerne. Das Zielmerkmal erfasst die Prognose gutartig (B) oder bösartig (M) und dient hier als Ground Truth. Achten Sie auf uninformative Features (z.B. ID) und fehlende Daten (Missing Values).

- Führen Sie eine Clusteranalyse durch, um eine etwaige Clusterstruktur zwischen gutartigen und bösartigen Zellen zu identifizieren.
- Vergleichen Sie die Genauigkeit folgender Cluster-Algorithmen:
  - k-means
  - hierarchisches Clustering
  - Modell-basiertes Clustering
  - DBSCAN
- Welches Verfahren ist am besten geeignet?
- Stellen Sie die Ergebnisse grafisch dar (Scatter Plot z.B. `radius_mean` vs. `texture_mean`).

- Lässt sich das Ergebnis verbessern, wenn vor dem Clustering der Merkmalsraum mittels PCA reduziert wird? Vergleichen Sie die Ergebnisse mit den vorherigen Resultaten.

```
# Load the data
breast_cancer <- read.csv("breast_cancer.csv", header = TRUE, sep = ",")
str(breast_cancer)

## 'data.frame': 569 obs. of 31 variables:
## $ diagnosis : chr "M" "M" "M" "M" ...
## $ radius_mean : num 18 20.6 19.7 11.4 20.3 ...
## $ texture_mean : num 10.4 17.8 21.2 20.4 14.3 ...
## $ perimeter_mean : num 122.8 132.9 130 77.6 135.1 ...
## $ area_mean : num 1001 1326 1203 386 1297 ...
## $ smoothness_mean : num 0.1184 0.0847 0.1096 0.1425 0.1003 ...
## $ compactness_mean : num 0.2776 0.0786 0.1599 0.2839 0.1328 ...
## $ concavity_mean : num 0.3001 0.0869 0.1974 0.2414 0.198 ...
## $ concave.points_mean : num 0.1471 0.0702 0.1279 0.1052 0.1043 ...
## $ symmetry_mean : num 0.242 0.181 0.207 0.26 0.181 ...
## $ fractal_dimension_mean : num 0.0787 0.0567 0.06 0.0974 0.0588 ...
## $ radius_se : num 1.095 0.543 0.746 0.496 0.757 ...
## $ texture_se : num 0.905 0.734 0.787 1.156 0.781 ...
## $ perimeter_se : num 8.59 3.4 4.58 3.44 5.44 ...
## $ area_se : num 153.4 74.1 94 27.2 94.4 ...
## $ smoothness_se : num 0.0064 0.00522 0.00615 0.00911 0.01149 ...
## $ compactness_se : num 0.049 0.0131 0.0401 0.0746 0.0246 ...
## $ concavity_se : num 0.0537 0.0186 0.0383 0.0566 0.0569 ...
## $ concave.points_se : num 0.0159 0.0134 0.0206 0.0187 0.0188 ...
## $ symmetry_se : num 0.03 0.0139 0.0225 0.0596 0.0176 ...
## $ fractal_dimension_se : num 0.00619 0.00353 0.00457 0.00921 0.00511 ...
## $ radius_worst : num 25.4 25 23.6 14.9 22.5 ...
## $ texture_worst : num 17.3 23.4 25.5 26.5 16.7 ...
## $ perimeter_worst : num 184.6 158.8 152.5 98.9 152.2 ...
## $ area_worst : num 2019 1956 1709 568 1575 ...
## $ smoothness_worst : num 0.162 0.124 0.144 0.21 0.137 ...
## $ compactness_worst : num 0.666 0.187 0.424 0.866 0.205 ...
## $ concavity_worst : num 0.712 0.242 0.45 0.687 0.4 ...
## $ concave.points_worst : num 0.265 0.186 0.243 0.258 0.163 ...
## $ symmetry_worst : num 0.46 0.275 0.361 0.664 0.236 ...
## $ fractal_dimension_worst : num 0.1189 0.089 0.0876 0.173 0.0768 ...

#summary(breast_cancer)

# check for missing values
sum(is.na(breast_cancer))

## [1] 0

# Preprocess the data
# Remove the diagnosis column (ground truth)
breast_cancer_scaled <- scale(breast_cancer[, -c(1)])
breast_cancer_scaled <- as.data.frame(breast_cancer_scaled)
```



```
# Perform k-means clustering
kmeans_result <- kmeans(breast_cancer_scaled, centers = 2, nstart = 12)
groups <- kmeans_result$cluster
groups <- car::recode(groups, "1='M'; 2='B'", as.factor = TRUE)

table(breast_cancer$diagnosis, groups)
```

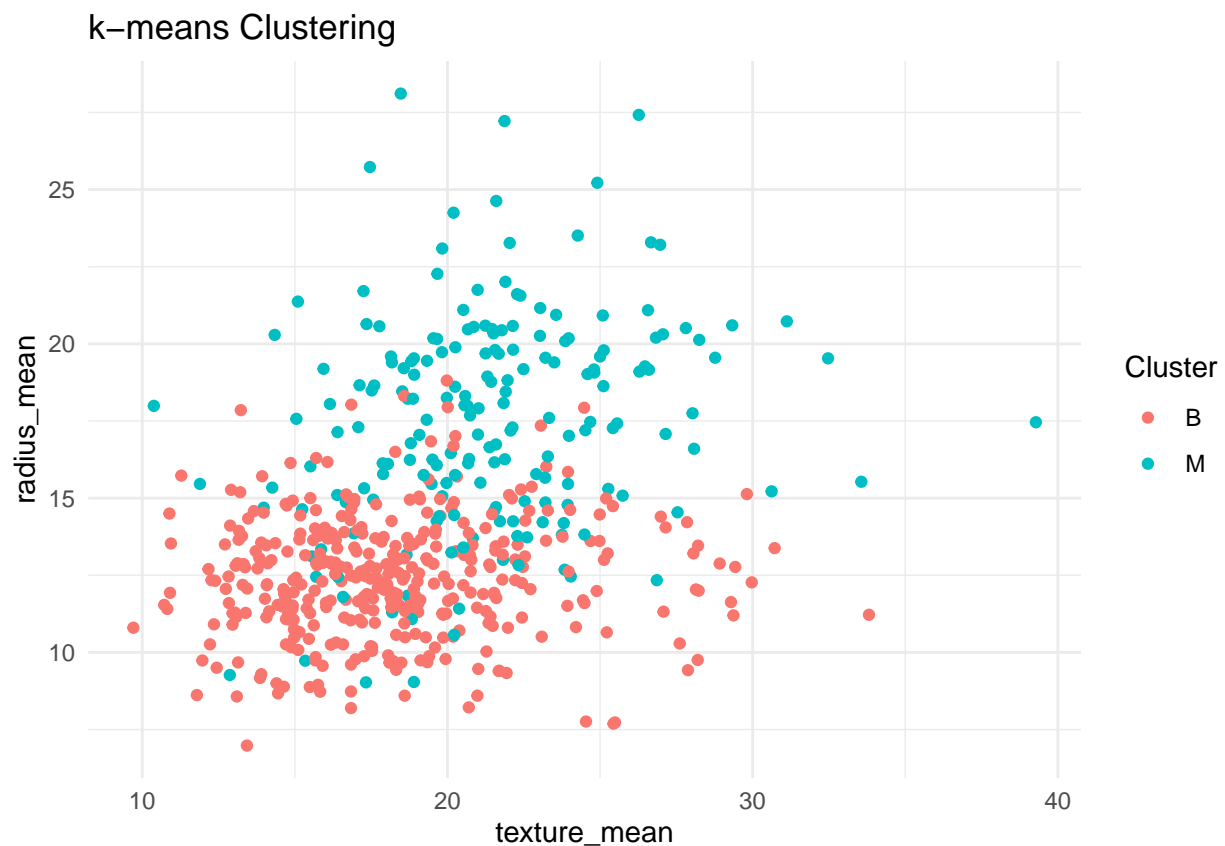
```
##      groups
##      B    M
## B 343  14
## M  37 175
```

```
acc_kmeans <- mean(breast_cancer$diagnosis==groups)

cat("Accuracy of k-means clustering: ", acc_kmeans, "\n")
```

```
## Accuracy of k-means clustering:  0.9103691
```

```
# Visualize the clusters (radius_mean vs. texture_mean)
ggplot(breast_cancer, aes(x = texture_mean, y = radius_mean, color = groups)) +
  geom_point() +
  labs(title = "k-means Clustering", color = "Cluster") +
  theme_minimal()
```



```

# Perform k-medoids clustering (PAM)
ncl.pam<-cluster::pam(breast_cancer_scaled,k = 2)

groups<-ncl.pam$clustering
groups

##      [1] 1 1 1 1 1 1 1 1 1 1 2 1 1 2 1 1 2 1 1 2 2 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##     [38] 2 2 2 2 2 1 2 2 1 2 1 2 2 2 2 2 1 2 2 1 1 2 2 2 2 1 2 1 1 2 2 1 2 1 2 1 2
##     [75] 2 1 2 1 1 2 2 1 1 1 2 1 2 1 2 2 2 2 2 1 1 2 2 2 2 2 2 2 2 2 2 1 2 2 1 2 2
##    [112] 2 1 2 2 2 2 1 1 2 2 1 1 2 2 2 2 2 1 1 2 1 2 2 1 2 2 2 1 2 2 2 2 2 2 1 2
##    [149] 2 2 2 1 1 2 2 2 1 2 2 2 2 1 1 2 1 2 2 2 1 2 2 2 1 2 2 2 1 1 2 2 1 1 2 2 2
##    [186] 2 2 2 2 2 1 2 2 1 1 2 1 2 1 1 2 1 1 1 2 2 2 2 2 1 2 1 1 1 1 2 2 1 1 2 2
##    [223] 2 1 2 2 2 2 2 1 1 2 2 1 2 2 1 1 2 1 2 2 1 2 1 2 2 2 2 1 2 1 2 1 2 1 1 1
##    [260] 1 1 2 1 2 1 1 2 2 2 2 2 2 1 2 2 2 2 2 2 2 1 2 1 1 2 2 2 2 2 2 1 2 2 2 2
##    [297] 2 2 2 2 1 2 1 2 2 2 2 2 2 2 2 2 2 2 2 2 2 1 1 2 2 2 2 1 2 2 2 2 1 1 1 2 2
##    [334] 2 2 1 2 1 2 1 2 2 2 1 2 2 2 2 2 2 2 1 1 1 2 2 2 2 2 2 2 2 2 2 2 1 1 2 1 1
##    [371] 1 2 1 1 2 2 1 2 2 1 2 2 2 2 2 2 2 2 2 1 2 2 1 1 2 2 2 2 2 2 1 2 2 2 2 2
##    [408] 2 1 2 2 2 2 2 2 2 2 1 2 2 2 1 2 2 2 2 2 2 2 2 1 2 1 1 2 2 2 2 2 2 1 2 2
##    [445] 2 2 1 2 2 1 2 1 2 2 2 2 2 2 2 2 2 1 1 2 2 2 1 2 2 1 1 2 2 2 2 2 2 2 1 2
##    [482] 2 2 2 2 1 2 1 2 2 2 2 1 2 2 2 2 2 1 1 2 1 2 1 1 1 2 2 2 1 2 2 1 2 2 1 1
##    [519] 2 2 2 1 2 2 2 2 2 2 2 2 2 2 2 2 1 2 1 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
##    [556] 2 2 2 2 2 2 2 1 1 1 1 2 1 2

```

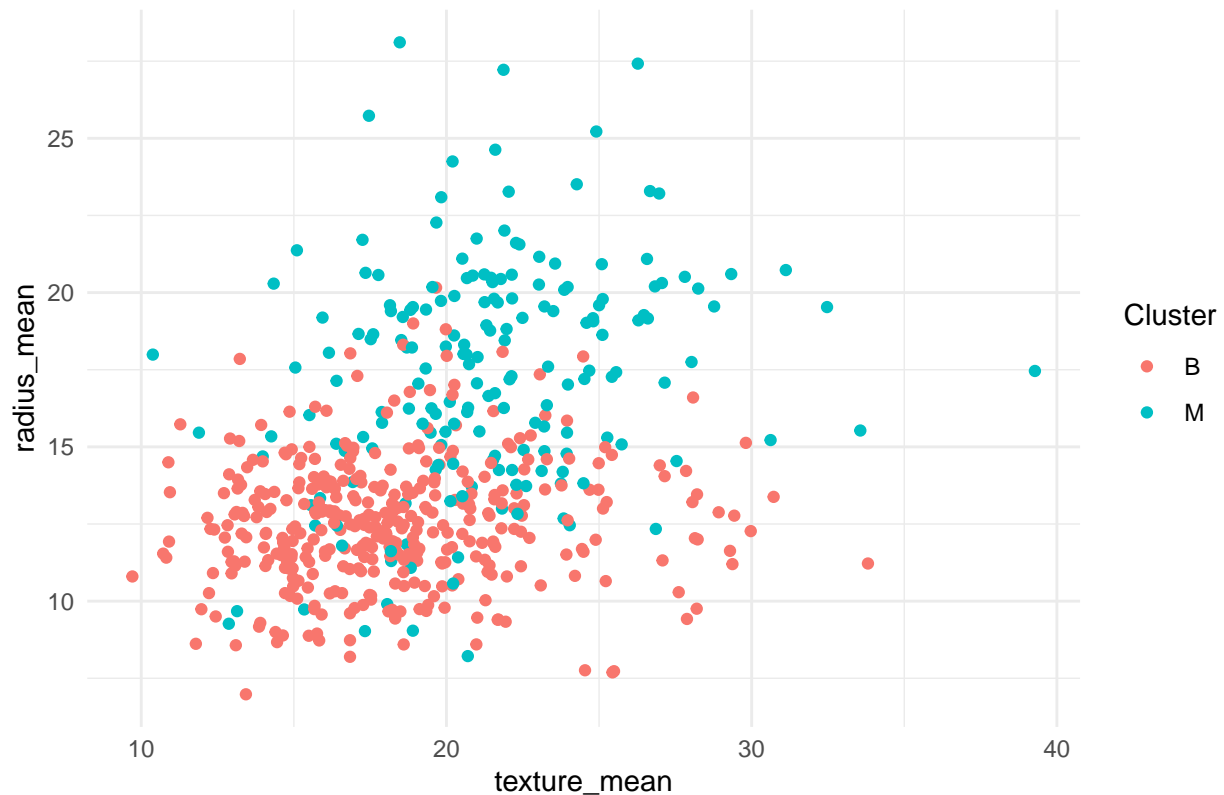
```

groups<-car::recode(groups,recodes="1='M';2='B'",as.factor = T)
acc_kmediod <- mean(breast_cancer$diagnosis==groups)

# Visualize the clusters (area_mean vs. smoothness_mean)
ggplot(breast_cancer, aes(x = texture_mean, y = radius_mean, color = groups)) +
  geom_point() +
  labs(title = "k-mediods Clustering", color = "Cluster") +
  theme_minimal()

```

## k-medoids Clustering

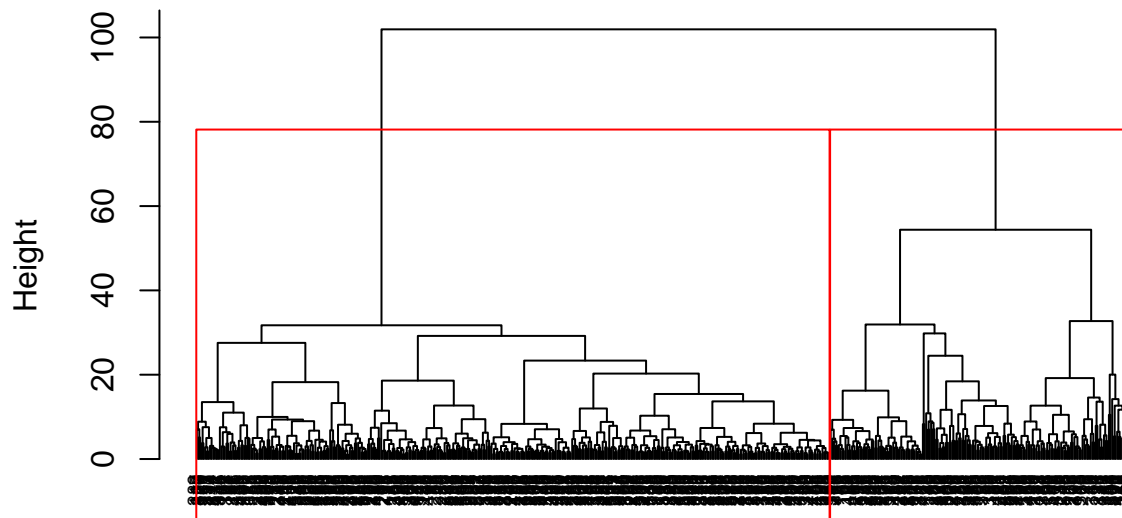


Interpretation: Die k-means-Clustering-Methode hat eine Genauigkeit von 0,9103, was bedeutet, dass 91% der Daten dem korrekten Cluster zugeordnet wurden. Die Genauigkeit kann in unserem Fall ermittelt werden, da wir die Ground-Truth-Informationen haben. Die Visualisierung anhand der Merkmale `texture_mean` und `radius_mean` zeigt, dass die Cluster stark überlappen und diese zufällig gewählten Merkmale eine geringe Trennschärfe aufweisen.

Beim k-medoids-Clustering beträgt die Genauigkeit 0,8910. Im Vergleich zu k-means ist die Genauigkeit etwas niedriger, aber immer noch relativ hoch. Die Visualisierung anhand der Merkmale `area_mean` und `smoothness_mean` zeigt ebenfalls eine starke Überlappung der Cluster.

```
# Perform hierarchical clustering
hclust <- hclust(dist(breast_cancer_scaled), method = "ward.D2")
plot(hclust, cex=0.6, hang = -1)
# Cluster einzeichnen
rect.hclust(hclust, k = 2, border = "red")
```

## Cluster Dendrogram



```
dist(breast_cancer_scaled)
hclust(*, "ward.D2")
```

```
# Cluster zuweisen
groups <- cutree(hclust, k = 2)
groups <- car::recode(groups, "1='M'; 2='B'", as.factor = TRUE)

acc_hclust <- mean(breast_cancer$diagnosis==groups)
```

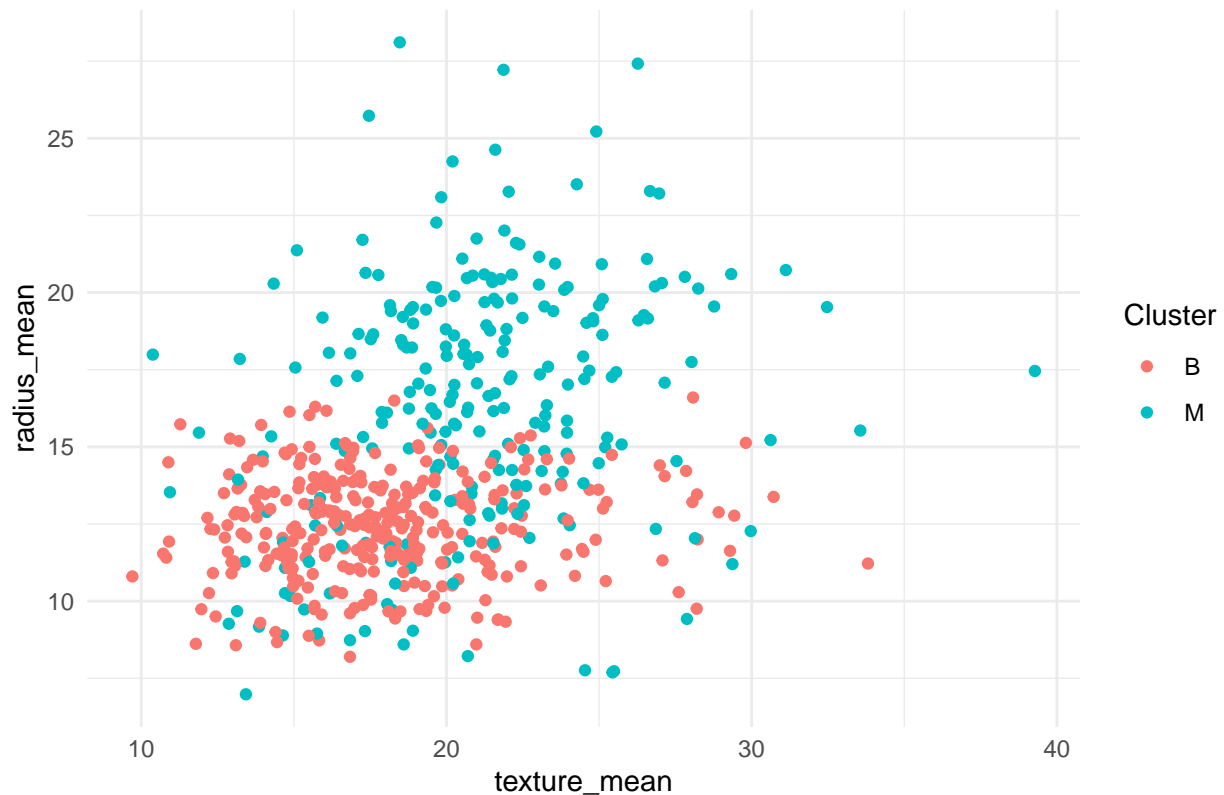
Das hierarchische Clustering hat erfolgreich zwei Hauptcluster identifiziert, die mit einer Genauigkeit von 88% den tatsächlichen Diagnosen (benign vs. malignant) entsprechen. Auch durch hierarchisches Clustering ist also eine klare Trennung der Datenpunkte möglich. Die Genauigkeit ist jedoch etwas niedriger als bei k-means und k-medoids.

```
# Perform model-based clustering
mclust_result <- Mclust(breast_cancer_scaled, G = 2)
groups <- mclust_result$classification
groups <- car::recode(groups, "1='M'; 2='B'", as.factor = TRUE)

acc_mclust <- mean(breast_cancer$diagnosis==groups)

# Visualize the clusters (texture_mean vs. radius_mean)
ggplot(breast_cancer, aes(x = texture_mean, y = radius_mean, color = groups)) +
  geom_point() +
  labs(title = "Model-based Clustering", color = "Cluster") +
  theme_minimal()
```

## Model-based Clustering



Das modellbasierte Clustering hat eine Genauigkeit von 0,8910, was der Genauigkeit des k-mediod-Clustering entspricht. Dies bedeutet, dass das modellbasierte Clustering in diesem Fall ähnlich effektiv ist wie k-medoids. Die Visualisierung zeigt eine klare Trennung der Cluster anhand der Merkmale `texture_mean` und `radius_mean`.

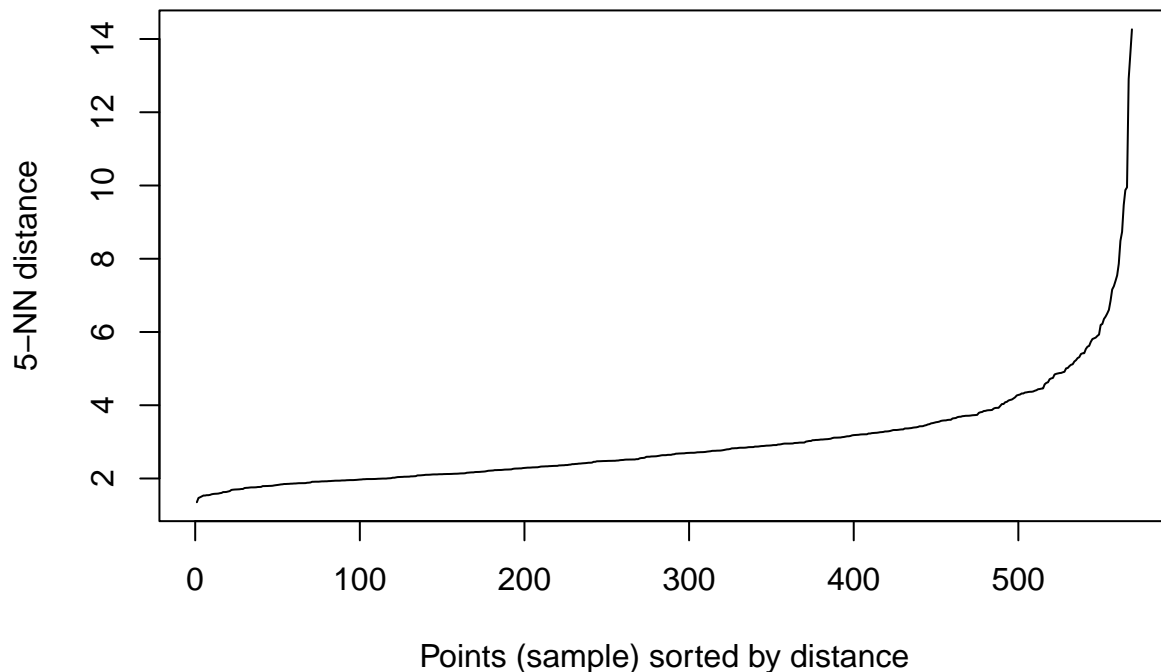
```
# Perform DBSCAN clustering
set.seed(123)
dbscan_result <- dbscan(breast_cancer_scaled, eps = 5, MinPts = 50)

## Warning in dbscan(breast_cancer_scaled, eps = 5, MinPts = 50): converting
## argument MinPts (fpc) to minPts (dbscan)!

groups <- dbscan_result$cluster
groups <- car::recode(groups, "0='M'; 1='B'", as.factor = TRUE)

acc_dbscan <- mean(breast_cancer$diagnosis==groups)

# Berechnung der k-nearest Neighbour-Distanz für k = MinPts
kNNdistplot(breast_cancer_scaled, k = 5) # Hier setzen wir k = MinPts (Startwert)
abline(h = 0.5, col = "red", lty = 2)
```



DBSCAN scheint für diesen Datensatz nicht optimal zu sein, um zwei Cluster (benign vs. malignant) zu identifizieren. Unabhängig von den Parametern (eps, MinPts) konnte immer nur ein Cluster und Rauschen identifiziert werden. Die ermittelte Genauigkeit beträgt demnach nur 0.6432. Alternative Methoden wie k-Means und hierarchisches Clustering könnten somit besser geeignet sein, um die Daten in zwei Hauptcluster zu unterteilen.

```
# Accuracy comparison
accuracy <- c(acc_kmeans, acc_kmediod, acc_hclust, acc_mclust, acc_dbscan)
method <- c("k-means", "k-medoids", "hierarchichal", "model-based", "DBSCAN")
accuracy_df <- data.frame(method, accuracy)
accuracy_df
```

```
##           method accuracy
## 1      k-means 0.9103691
## 2    k-medoids 0.8910369
## 3 hierarchical 0.8804921
## 4 model-based 0.8629174
## 5      DBSCAN 0.6432337
```

Die Untersuchung unterschiedlicher Clustering-Methoden auf den Brustkrebsdatensatz zeigt, dass das k-Means-Verfahren mit einer Genauigkeit von 91% am besten abschneidet. Das k-Medoids-Verfahren und das modellbasierte Clustering folgen dicht dahinter mit einer Genauigkeit von jeweils etwa 89%. Hierarchisches Clustering zeigt ebenfalls gute Ergebnisse mit einer Genauigkeit von 88%. Im Vergleich dazu zeigt das dichte-basierte DBSCAN-Verfahren eine deutlich geringere Genauigkeit von 64%, was darauf hinweist, dass dieses Verfahren für den gegebenen Datensatz nicht geeignet ist.

Jetzt wird überprüft, ob durch eine vorhergehende PCA eine bessere Trennung der Cluster erreicht werden kann.

```
# Dimensionality Reduction with PCA
```

```
pca_result <- prcomp(breast_cancer_scaled, scale. = TRUE)
summary(pca_result)
```

```
## Importance of components:
```

```
##          PC1      PC2      PC3      PC4      PC5      PC6      PC7
## Standard deviation  3.6444 2.3857 1.67867 1.40735 1.28403 1.09880 0.82172
## Proportion of Variance 0.4427 0.1897 0.09393 0.06602 0.05496 0.04025 0.02251
## Cumulative Proportion 0.4427 0.6324 0.72636 0.79239 0.84734 0.88759 0.91010
##          PC8      PC9      PC10     PC11     PC12     PC13     PC14
## Standard deviation  0.69037 0.6457 0.59219 0.5421 0.51104 0.49128 0.39624
## Proportion of Variance 0.01589 0.0139 0.01169 0.0098 0.00871 0.00805 0.00523
## Cumulative Proportion 0.92598 0.9399 0.95157 0.9614 0.97007 0.97812 0.98335
##          PC15     PC16     PC17     PC18     PC19     PC20     PC21
## Standard deviation  0.30681 0.28260 0.24372 0.22939 0.22244 0.17652 0.1731
## Proportion of Variance 0.00314 0.00266 0.00198 0.00175 0.00165 0.00104 0.0010
## Cumulative Proportion 0.98649 0.98915 0.99113 0.99288 0.99453 0.99557 0.9966
##          PC22     PC23     PC24     PC25     PC26     PC27     PC28
## Standard deviation  0.16565 0.15602 0.1344 0.12442 0.09043 0.08307 0.03987
## Proportion of Variance 0.00091 0.00081 0.0006 0.00052 0.00027 0.00023 0.00005
## Cumulative Proportion 0.99749 0.99830 0.9989 0.99942 0.99969 0.99992 0.99997
##          PC29     PC30
## Standard deviation  0.02736 0.01153
## Proportion of Variance 0.00002 0.00000
## Cumulative Proportion 1.00000 1.00000
```

```
# create a data frame with the first 10 principal components
```

```
pca_data <- data.frame(pca_result$x[, 1:10])
```

```
# k-Means Clustering
```

```
set.seed(123)
```

```
kmeans_pca <- kmeans(pca_data, centers = 2, nstart = 20)
```

```
groups.kmeans_pca <- kmeans_pca$cluster
```

```
groups.kmeans_pca <- car::recode(groups.kmeans_pca, "1='B'; 2='M'", as.factor = TRUE)
```

```
# k-Medoids Clustering
```

```
kmedoids_pca <- pam(pca_data, k = 2)
```

```
groups.kmedoids_pca <- kmedoids_pca$clustering
```

```
groups.kmedoids_pca <- car::recode(groups.kmedoids_pca, "1='M'; 2='B'", as.factor = TRUE)
```

```
# Hierarchisches Clustering
```

```
dist_pca <- dist(pca_data)
```

```
hc_pca <- hclust(dist_pca, method = "ward.D2")
```

```
groups.hc_pca <- cutree(hc_pca, k = 2)
```

```
groups.hc_pca <- car::recode(groups.hc_pca, "1='M'; 2='B'", as.factor = TRUE)
```

```
# Modell-basiertes Clustering
```

```
mbc_pca <- Mclust(pca_data, G = 2)
```

```
groups.mbc_pca <- mbc_pca$classification
```

```

groups.mbc_pca <- car::recode(groups.mbc_pca, "1='M'; 2='B'", as.factor = TRUE)

# DBSCAN Clustering
dbscan_pca <- dbscan(pca_data, eps = 1, MinPts = 5)

## Warning in dbscan(pca_data, eps = 1, MinPts = 5): converting argument MinPts
## (fpc) to minPts (dbscan)!

groups.dbscan_pca <- dbscan_pca$cluster
groups.dbscan_pca <- car::recode(groups.dbscan_pca, "0='B'; 1='M'", as.factor = TRUE)

# Berechnung der Genauigkeiten
acc_kmeans_pca <- mean(breast_cancer$diagnosis == groups.kmeans_pca)
acc_kmedoids_pca <- mean(breast_cancer$diagnosis == groups.kmedoids_pca)
acc_hc_pca <- mean(breast_cancer$diagnosis == groups.hc_pca)
acc_mbc_pca <- mean(breast_cancer$diagnosis == groups.mbc_pca)
acc_dbscan_pca <- mean(breast_cancer$diagnosis == groups.dbscan_pca)

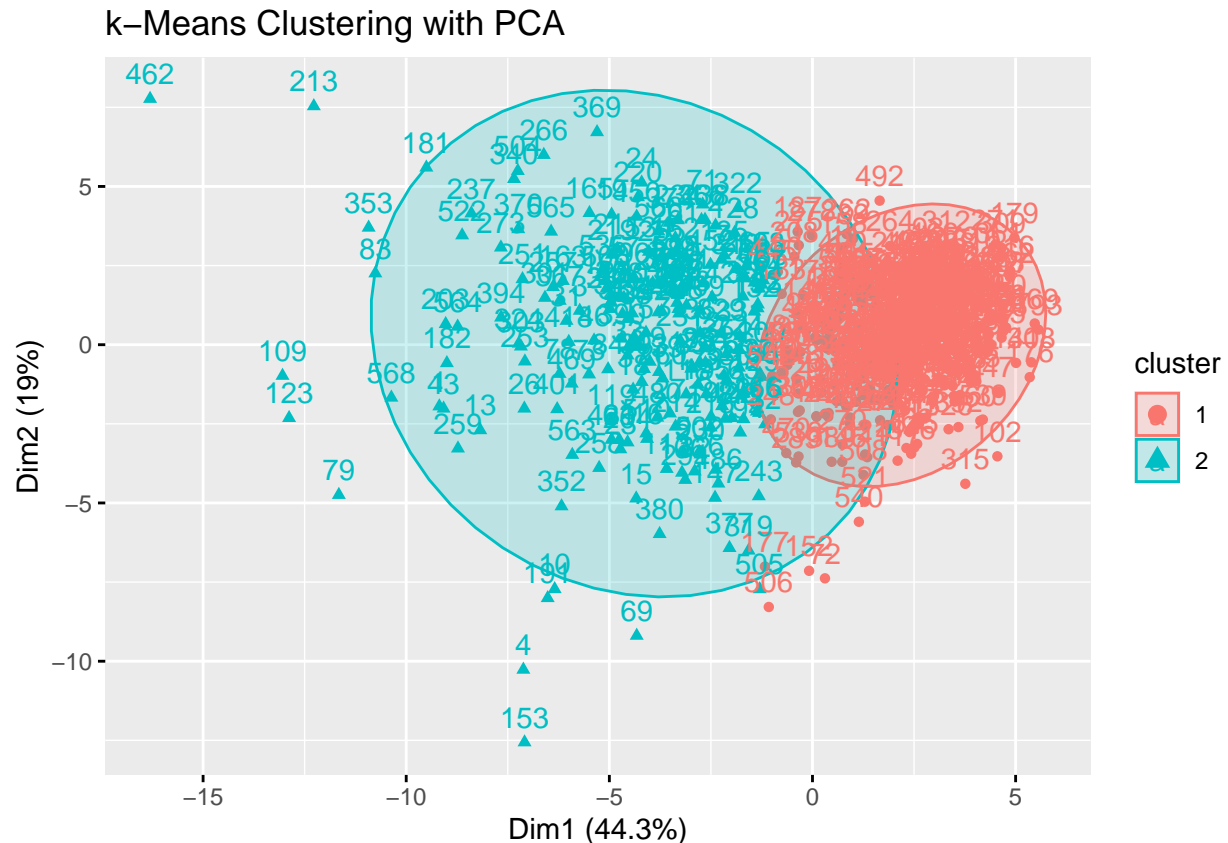
# Accuracy comparison
accuracy_pca <- c(acc_kmeans_pca, acc_kmedoids_pca, acc_hc_pca, acc_mbc_pca, acc_dbscan_pca)
method_pca <- c("k-means", "k-medoids", "hierarchical", "model-based", "DBSCAN")
accuracy_df_pca <- data.frame(method_pca, accuracy_pca)
accuracy_df_pca

##      method_pca accuracy_pca
## 1      k-means    0.9103691
## 2    k-medoids    0.9138840
## 3 hierarchical    0.9191564
## 4  model-based    0.6362039
## 5      DBSCAN    0.6274165

# visualize the clusters of kmeans with PCA - Dimensionality Reduction
fviz_cluster(object=kmeans_pca, data=breast_cancer_scaled,ellipse.type = "norm",repel = F) + ggtitle("k")

```





Die Auswirkung der PCA auf die Genauigkeit der Clusteranalyse zeigt deutliche Unterschiede bei den einzelnen Methoden. Während sich die Genauigkeit des k-Means-Verfahrens nicht verbessert hat, zeigt die k-Medoids eine Verbesserung von 89.1% auf 91.38%. Die hierarchische Clusteranalyse zeigt die stärkste Verbesserung, mit einem Sprung von 88% auf 91.92%. Im Gegensatz dazu zeigt das modellbasierte Clustering eine deutliche Verschlechterung der Genauigkeit von 89.1% auf 63.62%. Offenbar ist das modellbasierte Clustering weniger gut geeignet, um die Daten in zwei Hauptcluster zu unterteilen, nachdem die Dimensionalität reduziert wurde. Das DBSCAN-Verfahren zeigt ebenfalls eine Verschlechterung der Genauigkeit von 64% auf nur 36%, wobei hier wieder nur ein Cluster und Rauschen iden

### 3 Clustering - Heart Disease Patients [3P]

Verwenden Sie den Datensatz `heart_disease_patients.csv`. Der Datensatz enthält anonymisierte Daten von Patienten, bei denen eine Herzerkrankung diagnostiziert wurde. Patienten mit ähnlichen Merkmalen könnten auf die gleichen Behandlungen ansprechen, und Ärzte könnten davon profitieren, etwas über die Behandlungsergebnisse von Patienten zu erfahren, die denen ähneln, die sie behandeln. Zu diesem Zweck führen Sie bitte eine Clusteranalyse durch. Vergleichen Sie unterschiedliche Algorithmen und versuchen Sie ein bestmögliches Ergebnis zu erreichen. Begründen Sie ihre Entscheidungen. Verwenden Sie nur numerische Merkmale und achten Sie auf uninformative Features und fehlende Daten. Versuchen Sie die resultierenden Cluster zu interpretieren.

PCA and Summary Stats

```
load_source()
# Data Preparation
```

```

# Load the data
heart_disease_patients_tot <- read.csv("heart_disease_patients.csv")

str(heart_disease_patients_tot)

## 'data.frame': 303 obs. of 12 variables:
## $ id      : int  1 2 3 4 5 6 7 8 9 10 ...
## $ age     : int  63 67 67 37 41 56 62 57 63 53 ...
## $ sex     : int  1 1 1 1 0 1 0 0 1 1 ...
## $ cp      : int  1 4 4 3 2 2 4 4 4 4 ...
## $ trestbps: int  145 160 120 130 130 120 140 120 130 140 ...
## $ chol    : int  233 286 229 250 204 236 268 354 254 203 ...
## $ fbs     : int  1 0 0 0 0 0 0 0 0 1 ...
## $ restecg : int  2 2 2 0 2 0 2 0 2 2 ...
## $ thalach : int  150 108 129 187 172 178 160 163 147 155 ...
## $ exang   : int  0 1 1 0 0 0 0 1 0 1 ...
## $ oldpeak : num  2.3 1.5 2.6 3.5 1.4 0.8 3.6 0.6 1.4 3.1 ...
## $ slope   : int  3 2 2 3 1 1 3 1 2 3 ...

# Remove uninformative non-numeric features
heart_disease_patients <- subset(heart_disease_patients_tot, select= -c(id,sex,cp,fbs,restecg,exang,slope))
# print(heart_disease_patients)

# Perform exploratory data analysis
# summary(heart_disease_patients)

# Handle missing data
heart_disease_patients <- na.omit(heart_disease_patients)

# Normalize the numerical features (assuming all features are numerical)
heart_disease_patients_scaled <- scale(heart_disease_patients, center = TRUE)

#-----
# Perform PCA
# Using subset
pca_result <- prcomp(heart_disease_patients_scaled)

# Print summary of the PCA result
summary(pca_result)

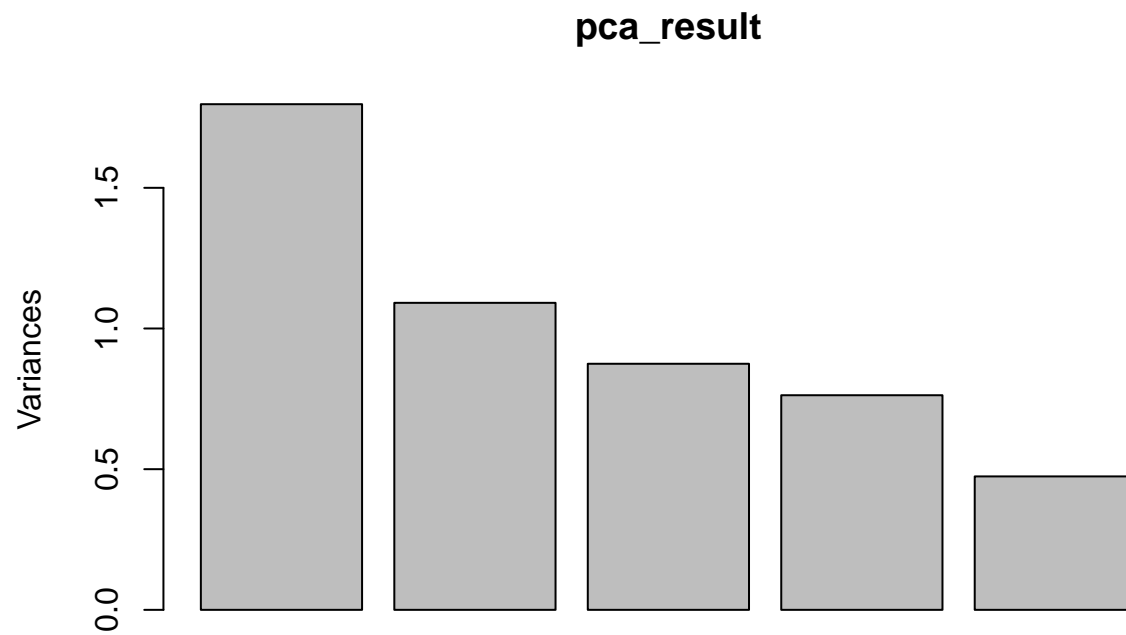
## Importance of components:
##              PC1      PC2      PC3      PC4      PC5
## Standard deviation  1.3406 1.0446 0.9352 0.8734 0.68865
## Proportion of Variance 0.3594 0.2182 0.1749 0.1525 0.09485
## Cumulative Proportion 0.3594 0.5777 0.7526 0.9052 1.00000

# Get the data in the PC space
data <- pca_result$x

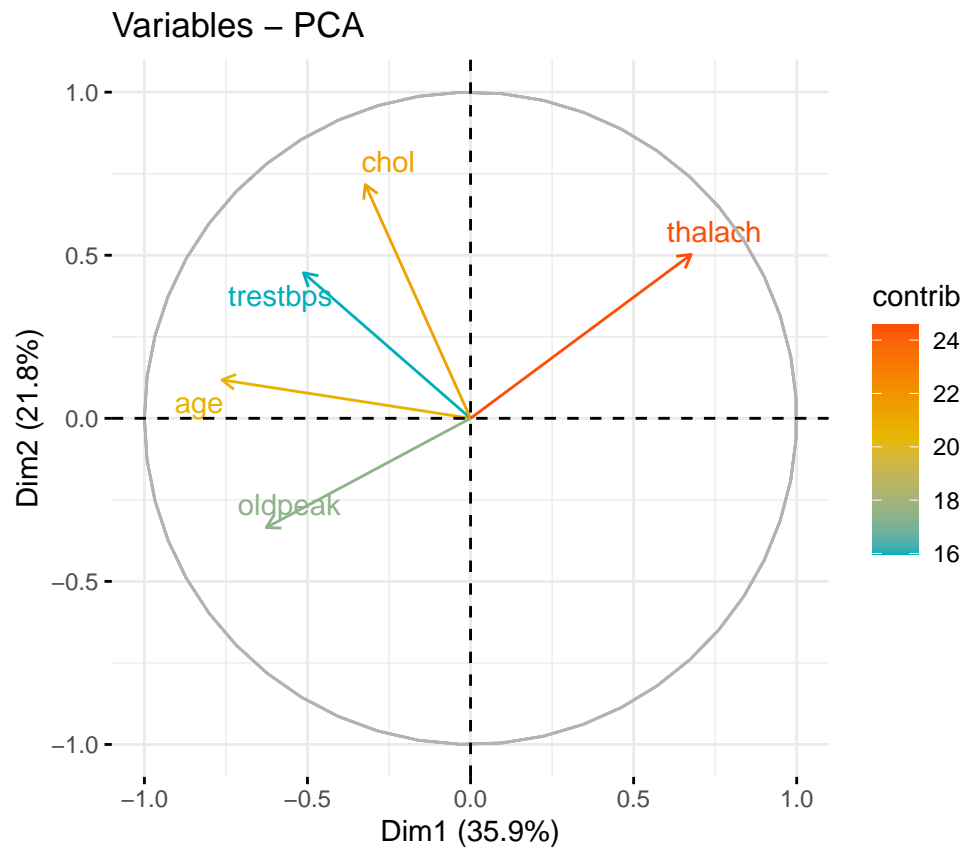
# Max number of clusters
n = 6

```

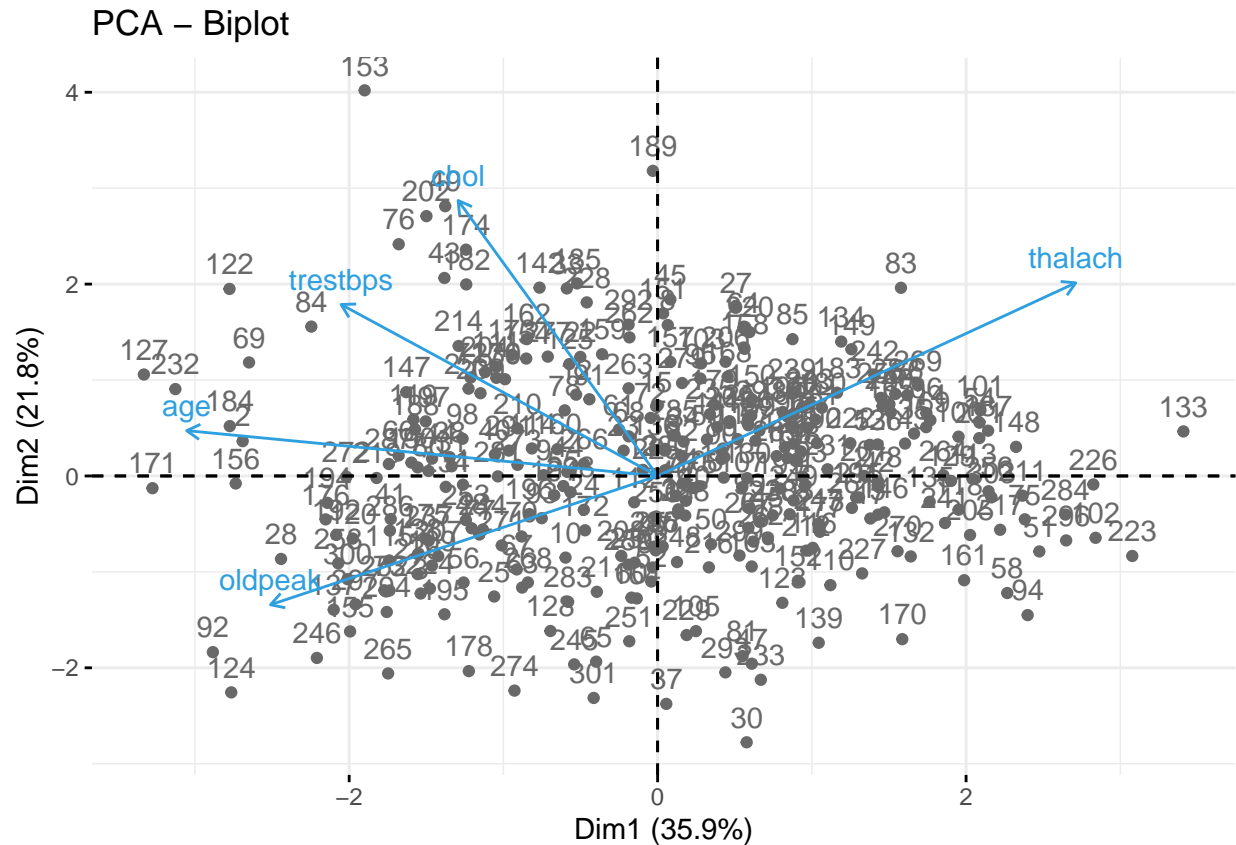
```
#-----  
# PCA Plots  
  
# Plot the variance explained by each principal component  
plot(pca_result)
```



```
# Plot the correlation circle  
fviz_pca_var(pca_result, col.var="contrib", gradient.cols = c("#00AFBB", "#E7B800", "#FC4E07"), repel =
```

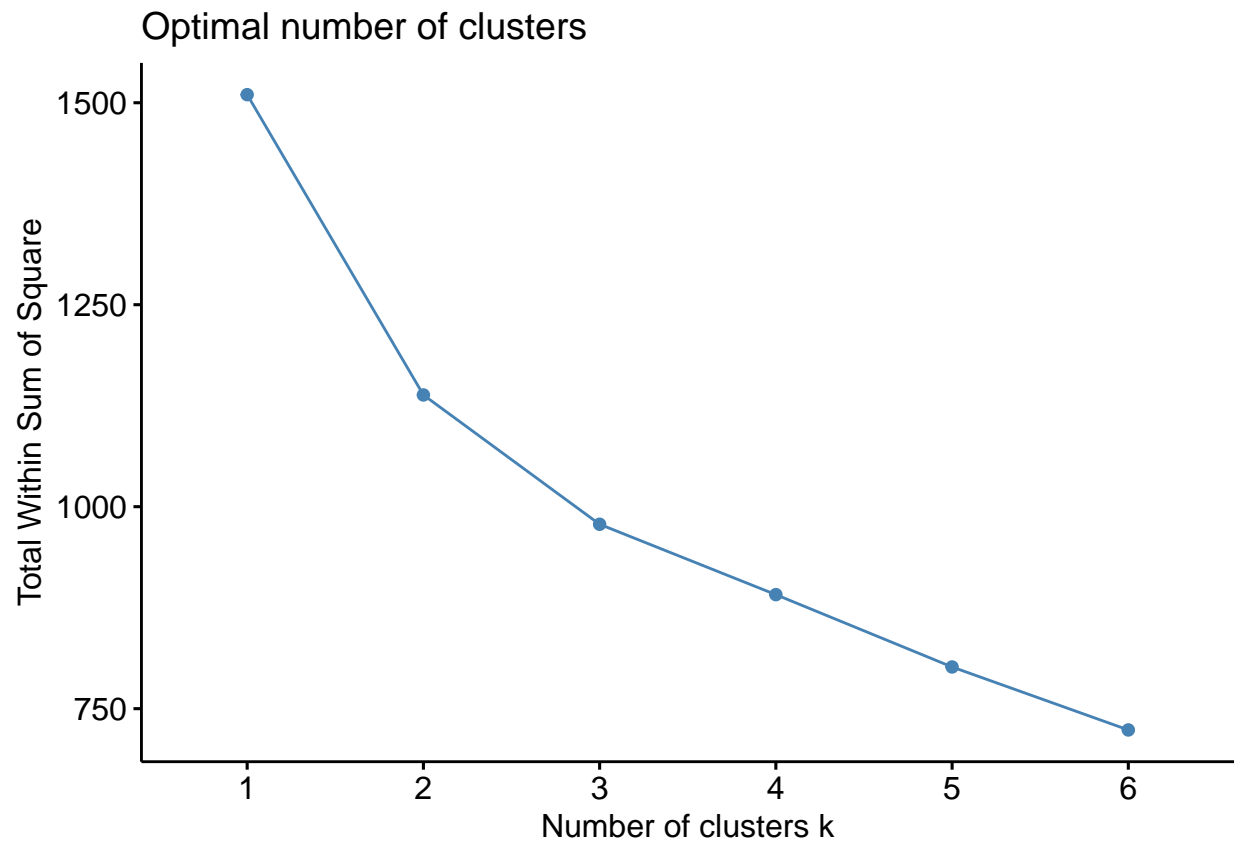


```
# Plot the biplot  
fviz_pca_biplot(pca_result, col.var="#2E9FDF", col.ind="#696969")
```

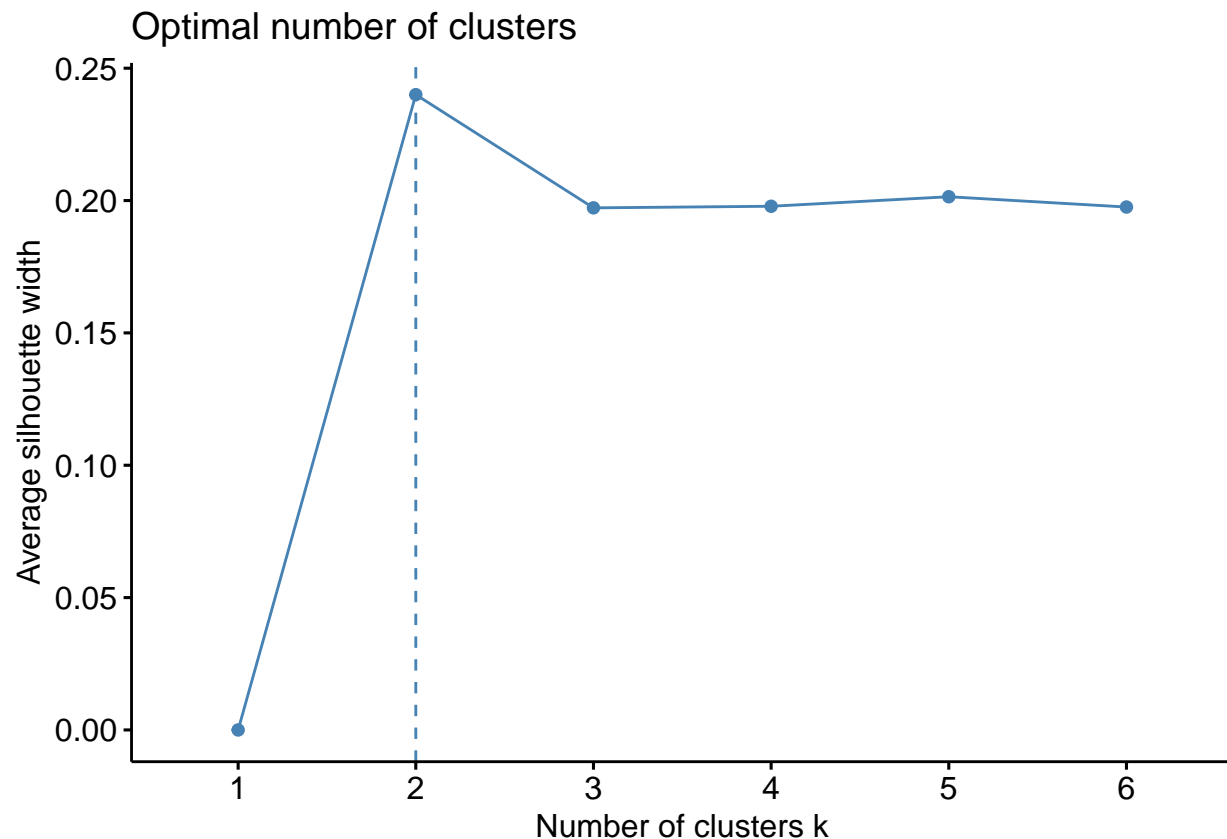


```
#-----
# Plots for Number of Clusters

# Determine the optimal number of clusters using the within-cluster sum of squares (WSS)
fviz_nbclust(x = data, FUNcluster = kmeans, method = "wss", k.max = n)
```



```
# Determine the optimal number of clusters using the silhouette method  
fviz_nbclust(x = data, FUNcluster = kmeans, method = "silhouette", k.max = n)
```



```
#-----
# Silhouette Scores

n = 6 # max number of clusters

# Cluster vectors are generated by different clustering methods for a given k
# They are stored in a list which is generated for each k
# These lists are stored inside another list
lists_of_cluster_vectors <- list()

# Create a list of lists of cluster vectors for different k
for (k in 1:n) {
  lists_of_cluster_vectors[[k]] <- get_list_of_cluster_vectors(data, k) # for m clustering methods and
  # Name each list of cluster vectors 'k'
  names(lists_of_cluster_vectors)[k] <- as.character(k)
}

m = length(lists_of_cluster_vectors[[1]])
# Loop over k lists of cluster vectors
# and pick the i-th cluster vector from that list
for( i in 1:m) {

  print(paste("Method Used: ", names(lists_of_cluster_vectors[[1]])[i]))
  for(k in 2:n) {
    score <- round(cluster.stats( dist(data), lists_of_cluster_vectors[[k]][[i]])$avg.silwidth,2)
    print(paste("For k = ",names(lists_of_cluster_vectors)[k]," Silhouette Score: ", score))
  }
}
```

```

}

}

## [1] "Method Used: kmeans"
## [1] "For k = 2 Silhouette Score: 0.24"
## [1] "For k = 3 Silhouette Score: 0.19"
## [1] "For k = 4 Silhouette Score: 0.2"
## [1] "For k = 5 Silhouette Score: 0.21"
## [1] "For k = 6 Silhouette Score: 0.19"
## [1] "Method Used: kmedoids"
## [1] "For k = 2 Silhouette Score: 0.19"
## [1] "For k = 3 Silhouette Score: 0.17"
## [1] "For k = 4 Silhouette Score: 0.16"
## [1] "For k = 5 Silhouette Score: 0.13"
## [1] "For k = 6 Silhouette Score: 0.14"
## [1] "Method Used: hkmeans"
## [1] "For k = 2 Silhouette Score: 0.23"
## [1] "For k = 3 Silhouette Score: 0.2"
## [1] "For k = 4 Silhouette Score: 0.2"
## [1] "For k = 5 Silhouette Score: 0.21"
## [1] "For k = 6 Silhouette Score: 0.17"
## [1] "Method Used: mclust"
## [1] "For k = 2 Silhouette Score: 0.17"
## [1] "For k = 3 Silhouette Score: 0.09"
## [1] "For k = 4 Silhouette Score: 0.05"
## [1] "For k = 5 Silhouette Score: 0.15"
## [1] "For k = 6 Silhouette Score: 0.11"
## [1] "Method Used: hclust"
## [1] "For k = 2 Silhouette Score: 0.36"
## [1] "For k = 3 Silhouette Score: 0.22"
## [1] "For k = 4 Silhouette Score: 0.19"
## [1] "For k = 5 Silhouette Score: 0.19"
## [1] "For k = 6 Silhouette Score: 0.14"

```

Die Methode `get_list_of_cluster_vectors` wird verwendet, um Clusterzuweisungen für verschiedene Clustering-Methoden zu berechnen und diese in einer Liste zu speichern. Diese Liste enthält die Clusterzuweisungen für eine gegebene Anzahl von Clustern  $k$ . Der Outcome wird verwendet, um die Silhouettenwerte für verschiedene Clustering-Methoden und Clusteranzahlen zu berechnen und zu vergleichen. Dies hilft dabei, die optimale Anzahl von Clustern für die gegebenen Daten zu identifizieren.

Die Ergebnisse deuten darauf hin, dass  $k=2$  die beste Clusteranzahl zu sein scheint, da die meisten Methoden die höchsten Silhouettenwerte. Auch der WSS-Plot sowie der Silhouette-Plot lassen auf 2 als optimale Clusteranzahl schließen.

Grouped Scatter/Cluster Plots for 1 to  $n$  Clusters for Each (Suitable) Clustering Method

```

load_source()
# Data Preparation

# Load the data
heart_disease_patients_tot <- read.csv("heart_disease_patients.csv")

# Remove uninformativ non-numeric features

```



```

heart_disease_patients <- subset(heart_disease_patients_tot, select= -c(id,sex,cp,fbs,restecg,exang,slop

# Handle missing data
heart_disease_patients <- na.omit(heart_disease_patients)

# Normalize the numerical features (assuming all features are numerical)
heart_disease_patients_scaled <- scale(heart_disease_patients, center = TRUE)

#-----
# Perform PCA

# Perform PCA on the scaled data
pca_result <- prcomp(heart_disease_patients_scaled)

# Get the data in the PC space
data <- pca_result$x

#-----
# Perform different clustering methods
n = 6 # max number of clusters

# Cluster vectors are generated by different clustering methods for a given k
# They are stored in a list which is generated for each k
# These lists are stored inside another list
lists_of_cluster_vectors <- list()

# Create a list of lists of cluster vectors for different k
for (k in 1:n) {
  lists_of_cluster_vectors[[k]] <- get_list_of_cluster_vectors(data, k) # for m methods of clustering a
  # Name each list of cluster vectors 'k'
  names(lists_of_cluster_vectors)[k] <- as.character(k)
}

# Cluster results are generated by different clustering methods for a given k
# They are stored in a list which is generated for each k
# These lists are stored inside another list
lists_of_cluster_results <- list()

# Create a list of lists of cluster results for different k
for (k in 1:n) {
  lists_of_cluster_results[[k]] <- get_list_of_cluster_results(data, k) # for m methods of clustering a
  # Name each list of cluster results 'k'
  names(lists_of_cluster_results)[k] <- as.character(k)
}

#-----
# Visualize the cluster vectors and results

#--- Methods ---#

create_grid_plot_vectors <- function(data, lists_of_cluster_vectors, i, n) {
  # Depiction 1
  plot_list <- list()

```

```

for(k in 1:n) {
  # Generate the plot and store it in the list
  plot_list[[k]] <- create_scatter_plot_pca( data, lists_of_cluster_vectors[[k]][[i]], names(lists_of_
}

# Combine all the plots into a single plot
grid.arrange(grobs = plot_list, ncol = 2)
}

create_grid_plot_results <- function(data, lists_of_cluster_results, i, n) {
  # Depiction 2
  plot_list <- list()

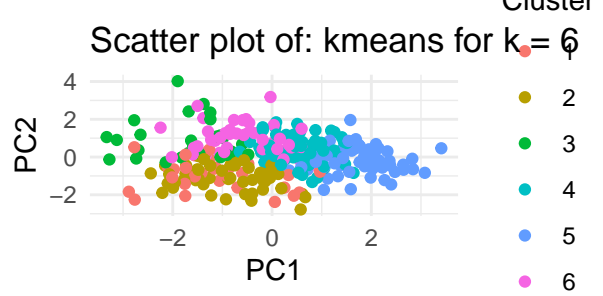
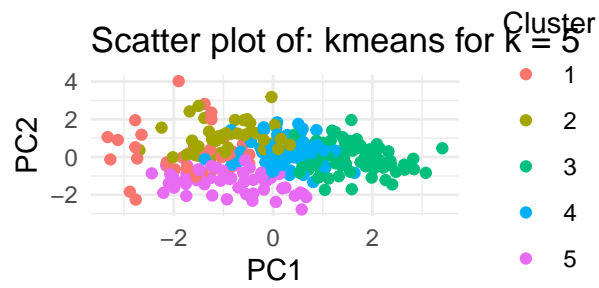
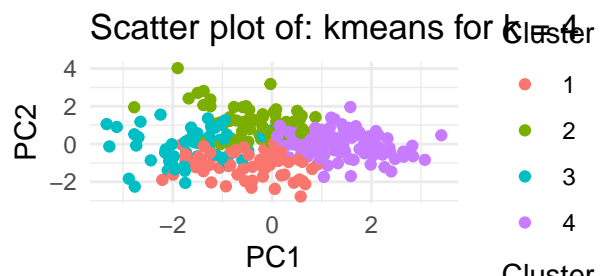
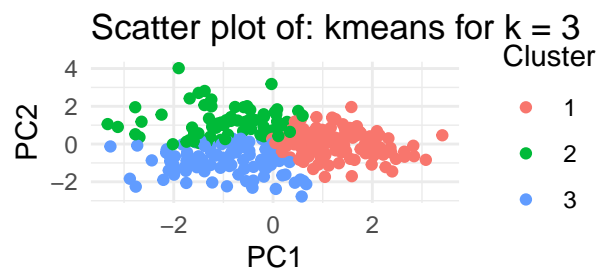
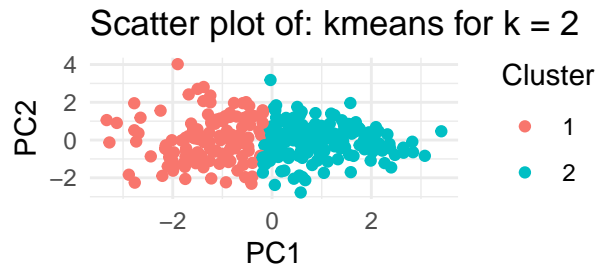
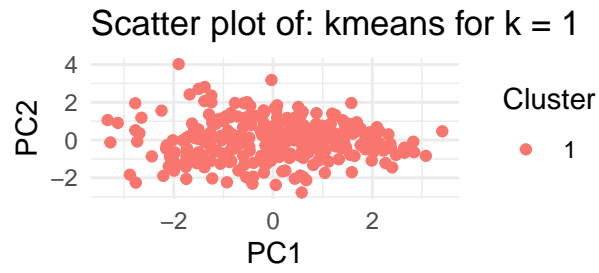
  for(k in 2:n) {
    # Visualize kmeans clustering
    plot_list[[k]] <- fviz_cluster(lists_of_cluster_results[[k]][[i]], data = data[, c(1,2)], ellipse.t
    theme_minimal() +
    ggtitle(paste("Cluster plot of:", names(lists_of_cluster_results[[k]])[i], "for k = ", names(lists_
    xlab("PC1") +
    ylab("PC2")
  }

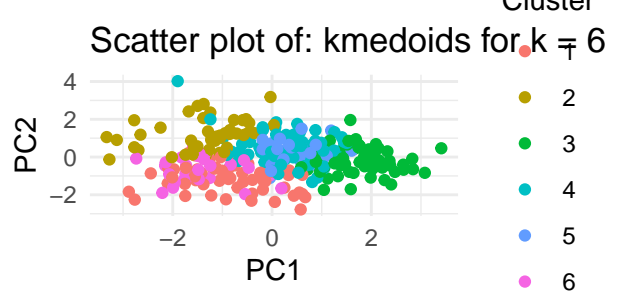
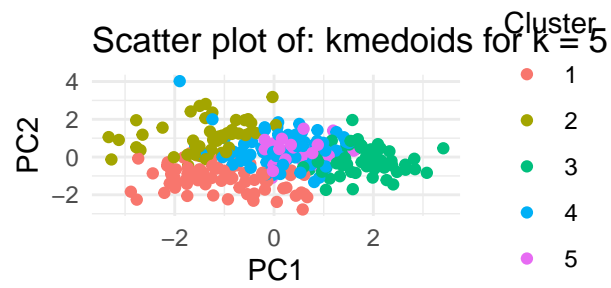
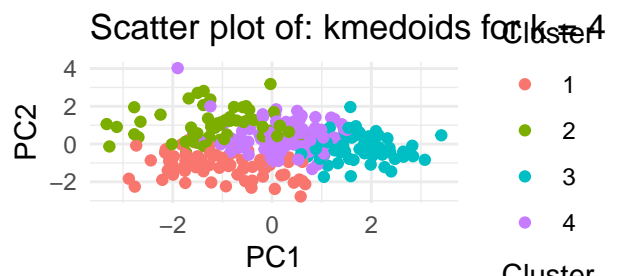
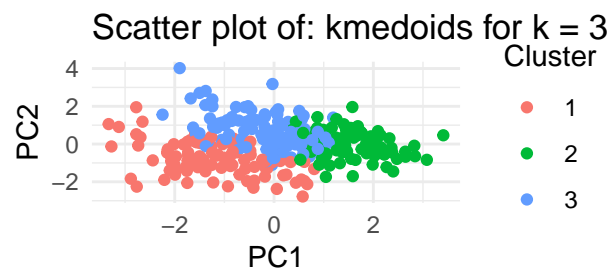
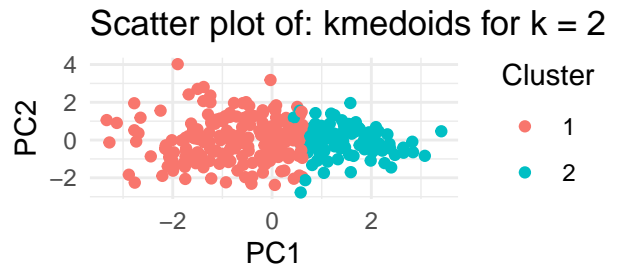
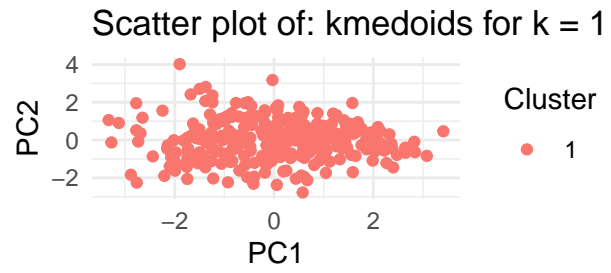
  # Combine all the plots into a single plot
  grid.arrange(grobs = plot_list, ncol = 2)
}

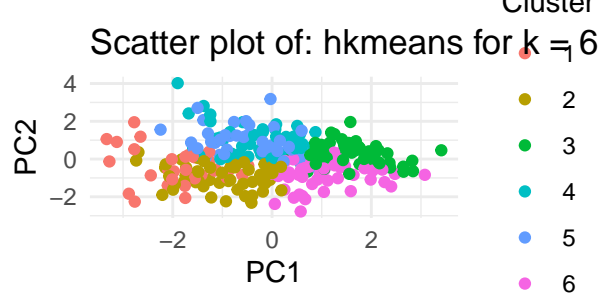
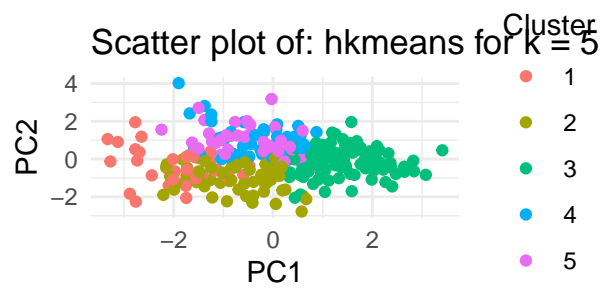
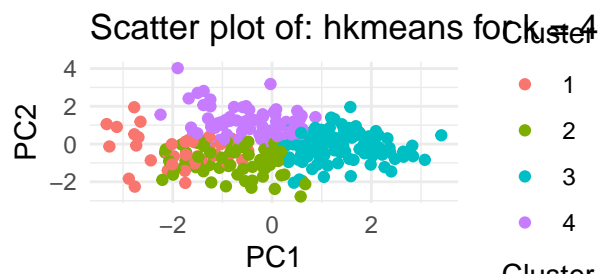
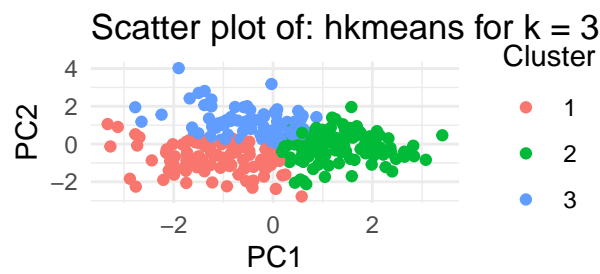
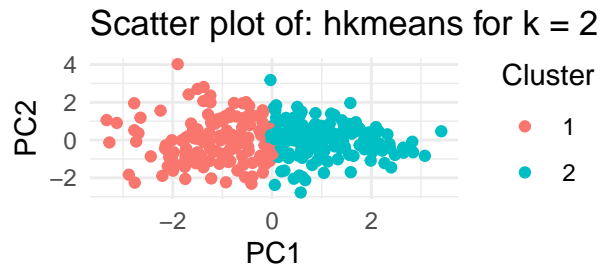
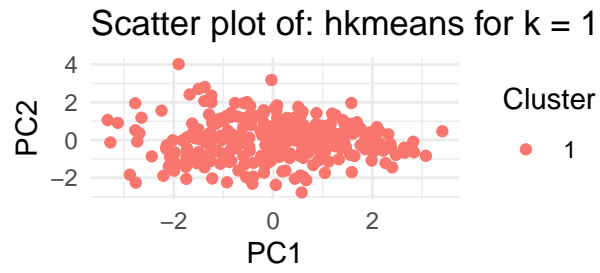
#--- End of Methods ---#

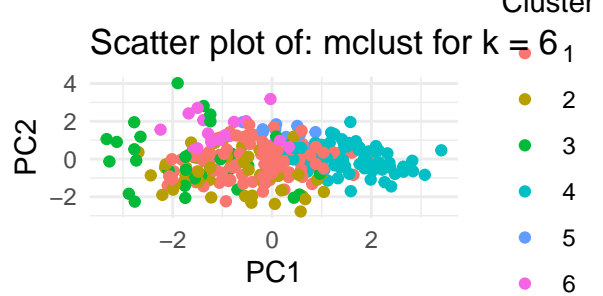
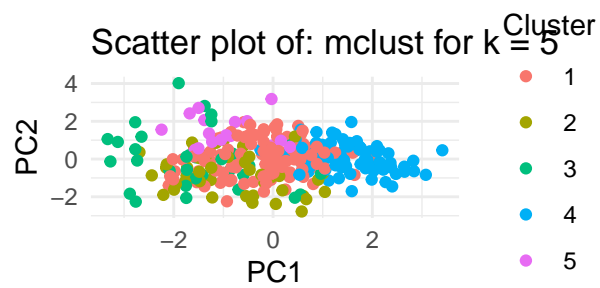
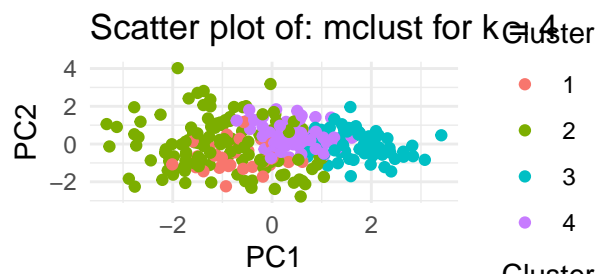
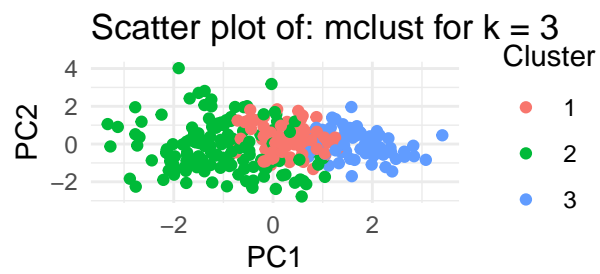
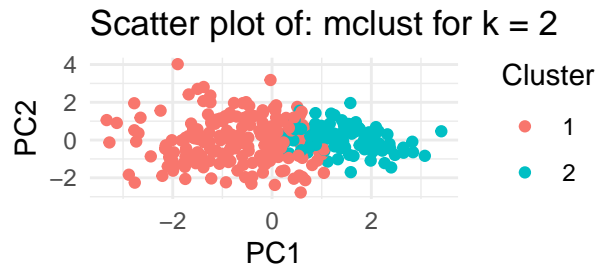
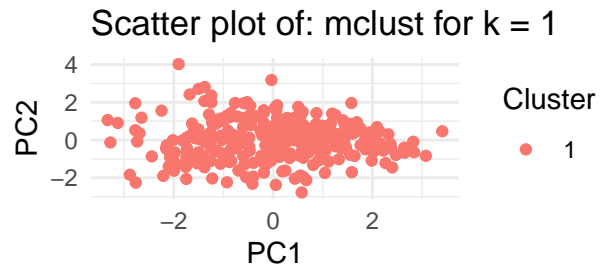
# All cluster vectors of different clustering methods have the same length independently of k
m = length(lists_of_cluster_vectors[[1]])
# Loop over k lists of cluster vectors
# and pick the i-th cluster vector from that list
for( i in 1:m) {
  create_grid_plot_vectors(data,lists_of_cluster_vectors,i,n)
}

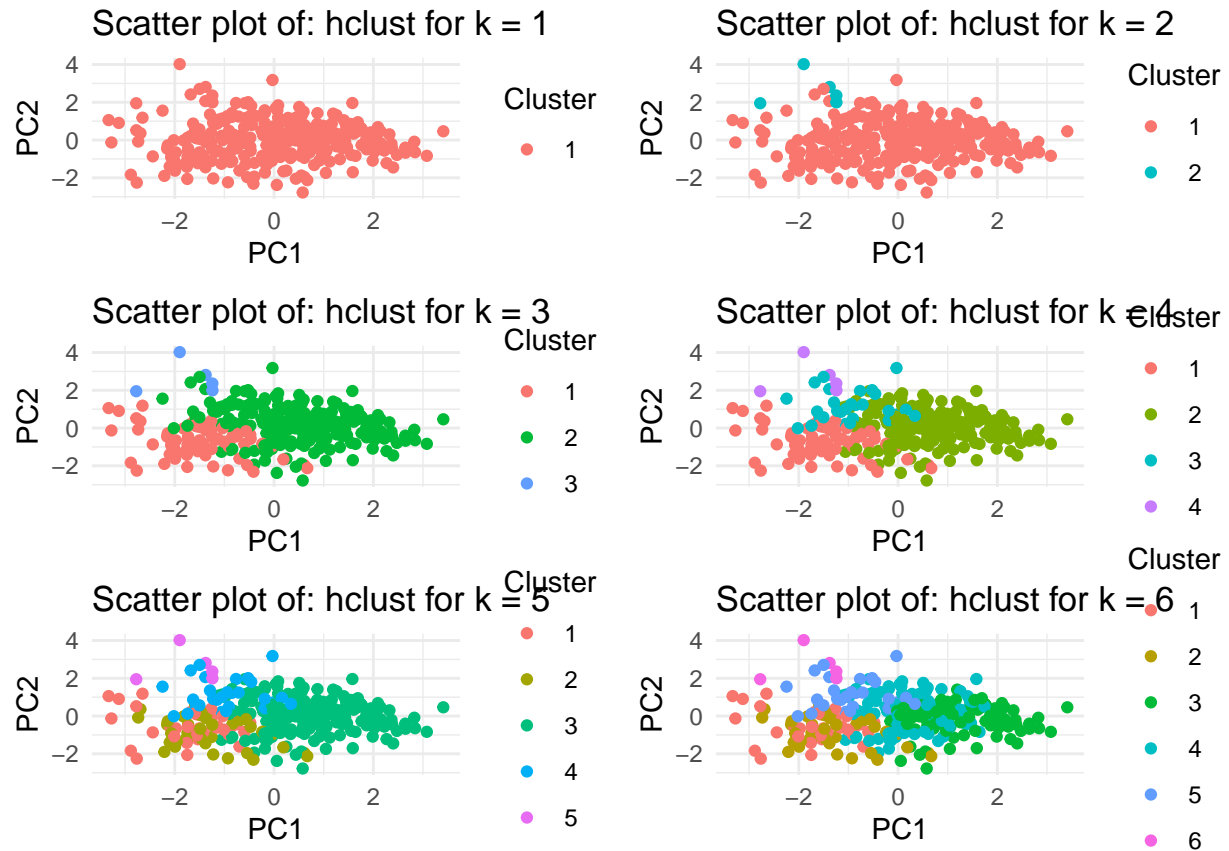
```



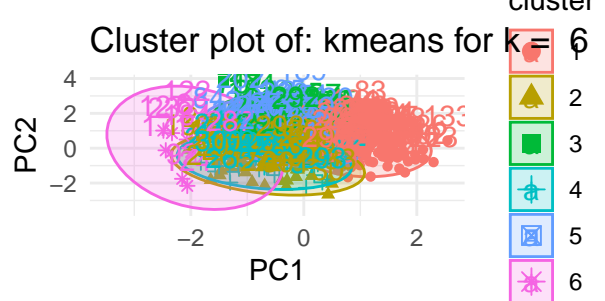
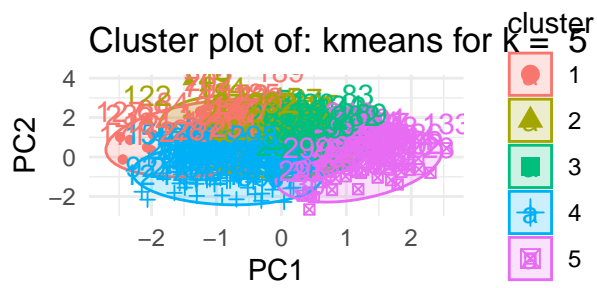
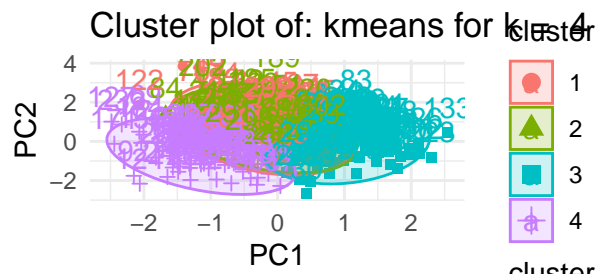
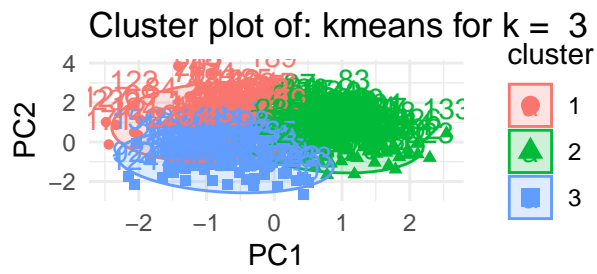
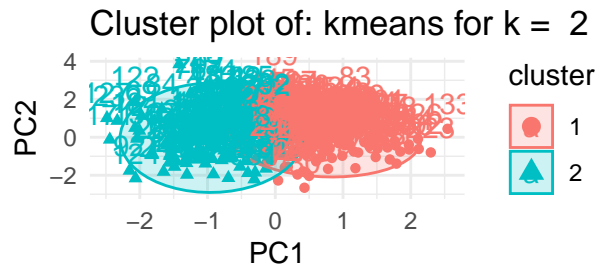






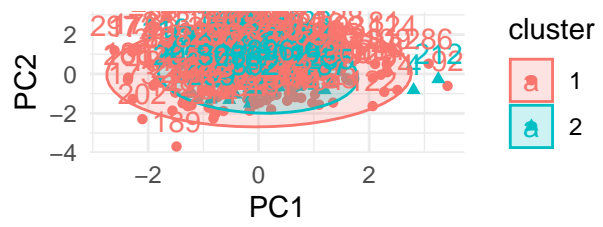


```
for( i in 1:(m-1)) {
  create_grid_plot_results(data,lists_of_cluster_results,i,n)
}
```

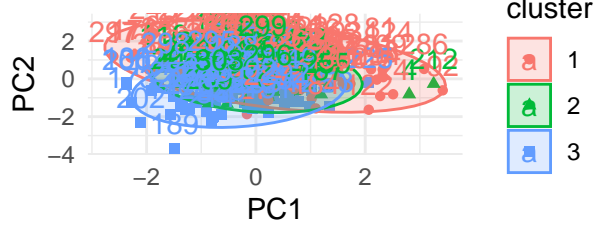




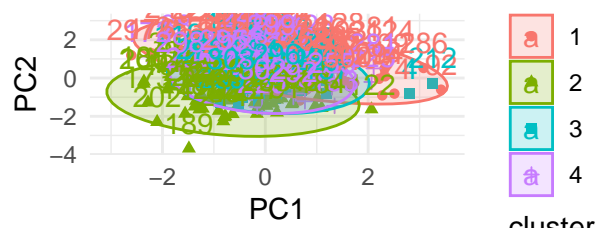
Cluster plot of: kmedoids for k = 2



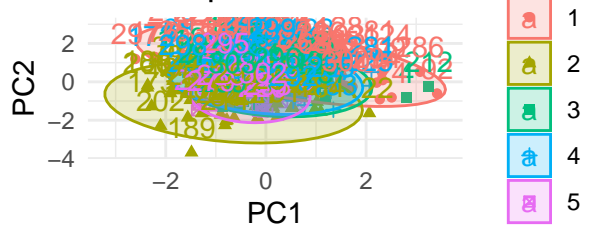
Cluster plot of: kmedoids for k = 3



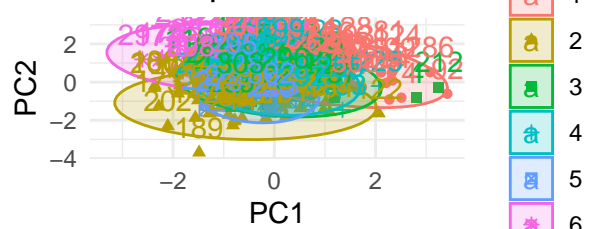
Cluster plot of: kmedoids for k = 4



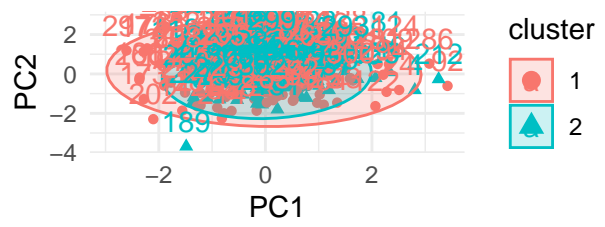
Cluster plot of: kmedoids for k = 5



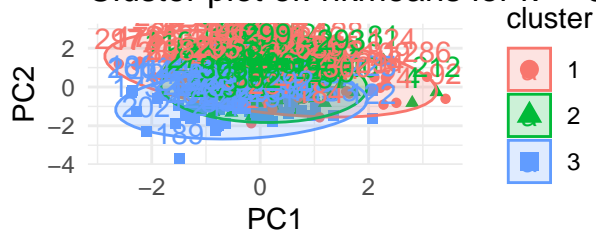
Cluster plot of: kmedoids for k = 6



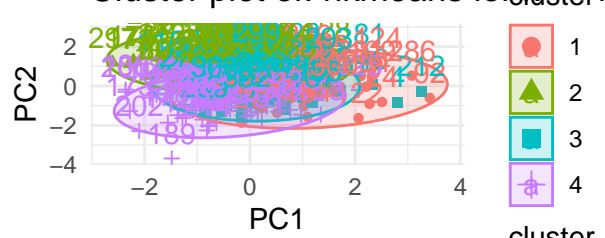
Cluster plot of: hkmeans for k = 2



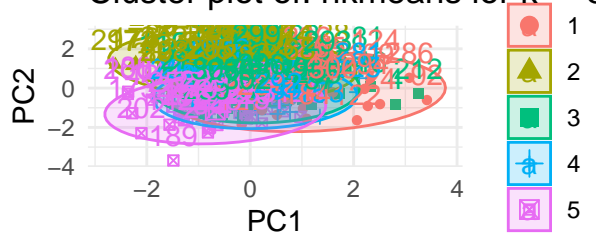
Cluster plot of: hkmeans for k = 3



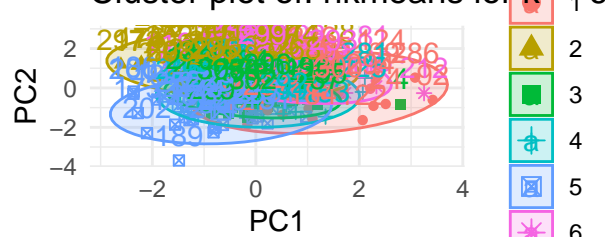
Cluster plot of: hkmeans for k = 4

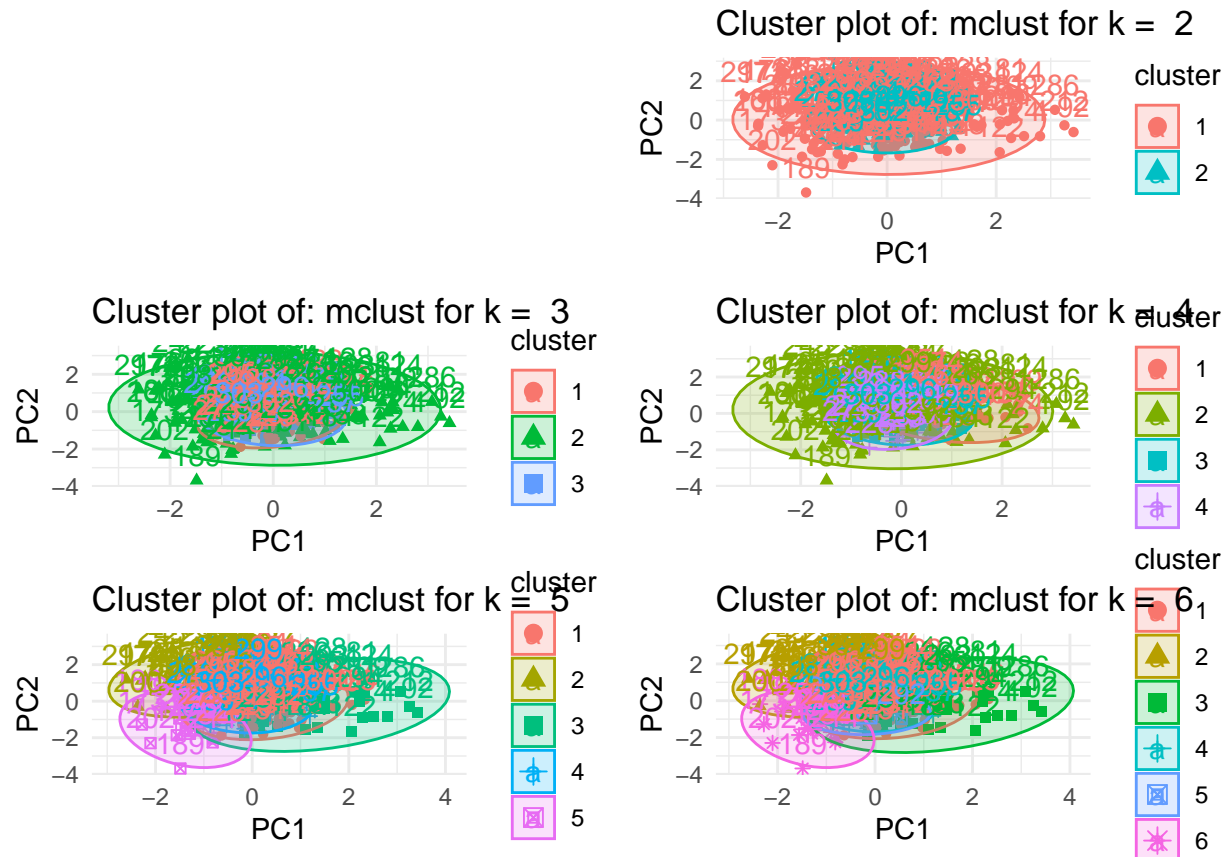


Cluster plot of: hkmeans for k = 5



Cluster plot of: hkmeans for k = 6





```
# DBSCAN Clustering
```

```
# Perform DBSCAN clustering
```

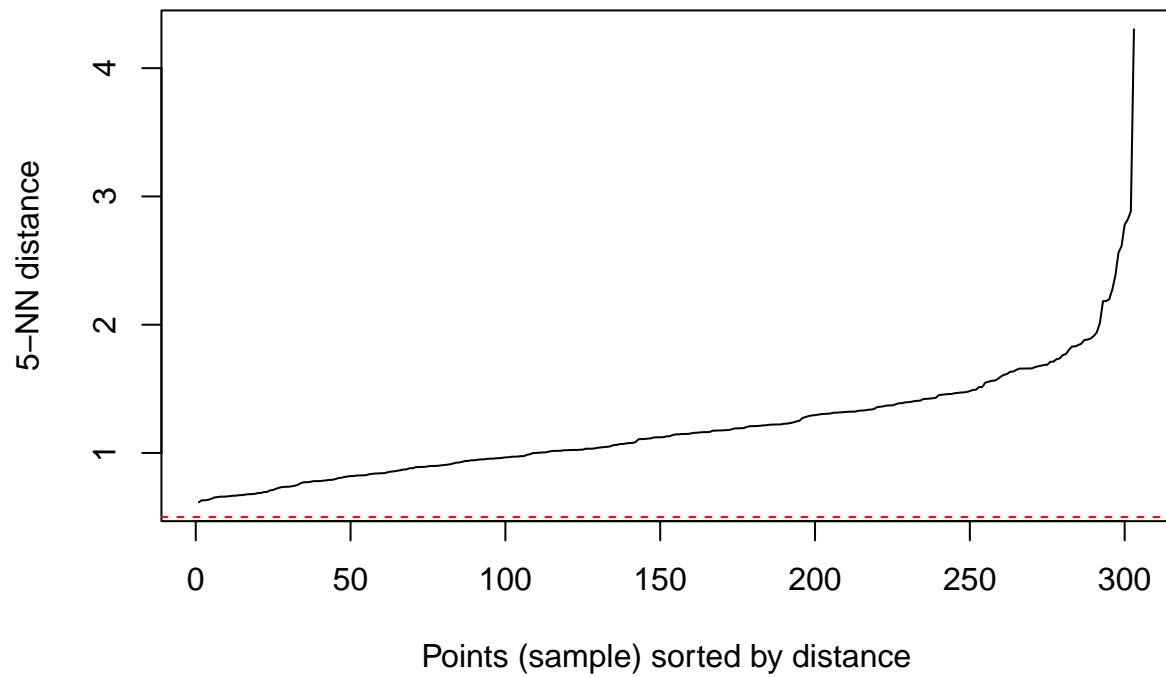
```
dbscan_result <- dbscan(data, eps = 1, minPts = 4)
dbscan_result$cluster
```

```
## [1] 1 0 0 0 2 1 0 1 1 0 1 1 1 1 0 0 1 1 1 1 1 0 0 1 1 1 0 1 0 1 1 1 1 1 1 0
## [38] 4 0 1 0 2 0 1 1 1 0 0 0 1 1 1 1 1 3 1 1 1 1 1 0 1 1 0 0 0 1 0 0 1 1 3 1
## [75] 1 0 1 1 1 0 0 1 0 0 1 1 1 0 1 1 0 0 1 0 1 1 1 1 1 1 1 1 1 0 1 1 1 1 0 1
## [112] 1 1 0 3 1 1 2 5 1 1 0 0 0 1 1 0 0 1 1 1 1 0 1 1 1 0 3 0 1 1 0 1 5 1 1 4 1
## [149] 1 1 0 1 0 0 3 0 0 1 1 1 1 0 1 0 1 1 1 1 1 1 0 0 0 0 1 0 1 0 0 1 1 0 1 0 0
## [186] 1 1 4 0 1 1 0 1 3 0 0 0 1 1 4 1 0 0 0 1 0 1 0 1 1 1 0 2 0 1 1 1 1 5 1 1 1
## [223] 1 3 0 0 1 0 0 1 1 0 0 0 1 1 3 1 1 1 1 1 1 1 0 0 1 1 1 1 0 0 3 1 0 1 1 0 0
## [260] 1 1 1 1 1 0 0 1 1 1 1 1 0 0 0 1 0 1 1 1 1 0 1 1 1 1 0 0 1 1 1 1 1 0 0 1 1
## [297] 0 0 1 0 0 1 0
```

```
# Parameteroptimierung für DBSCAN
```

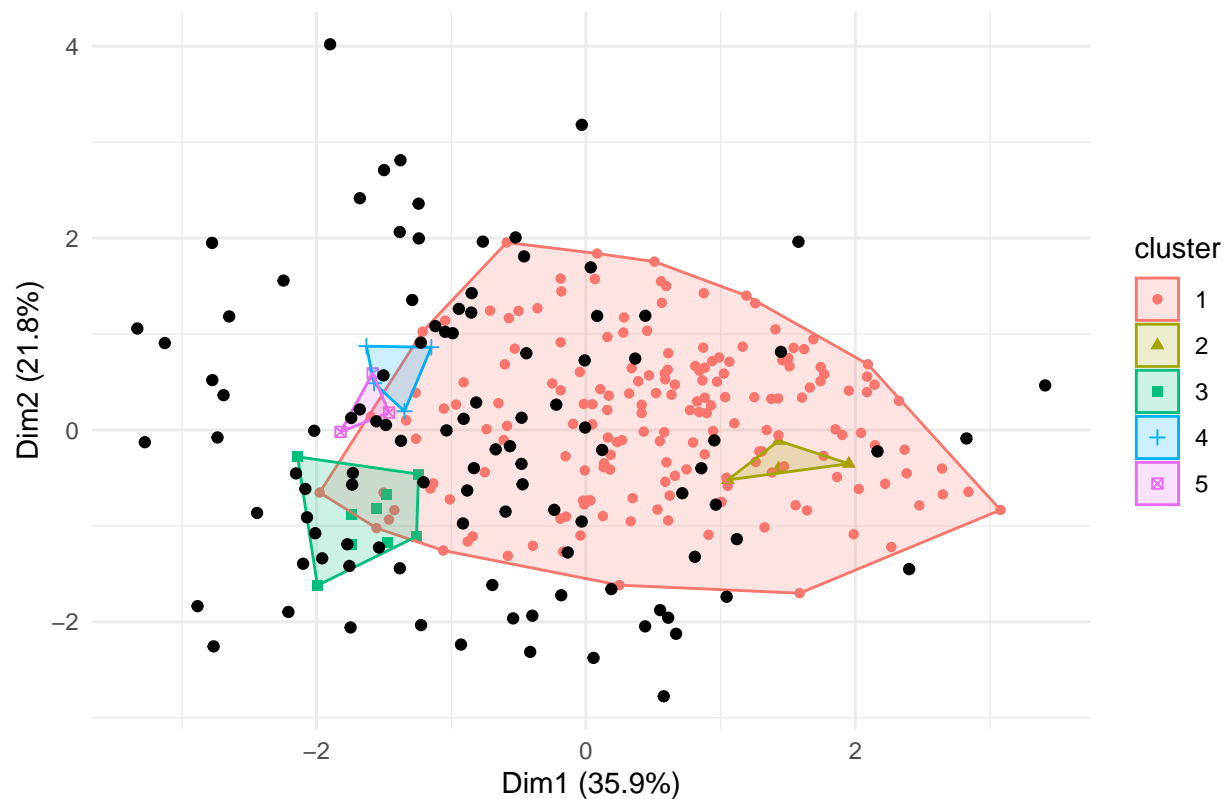
```
kNNdistplot(data, k = 5)
```

```
abline(h = 0.5, col = "red", lty = 2)
```



```
# DBSCAN Cluster Visualisierung
fviz_cluster(dbscan_result, data = data, geom = "point", stand = FALSE, show.clust.cent = FALSE) +
  theme_minimal() +
  ggtitle("DBSCAN Clustering of Heart Disease Patients")
```

## DBSCAN Clustering of Heart Disease Patients



```
# Anzahl der Cluster und Rauschpunkte
```

```
cat("Anzahl der ermittelten cluster: ", length(unique(dbscan_result$cluster[dbscan_result$cluster > 0])))
```

```
## Anzahl der ermittelten cluster: 5
```

```
cat("Anzahl der als Rauschen erkannten Punkte: ", sum(dbscan_result$cluster == 0), "\n")
```

```
## Anzahl der als Rauschen erkannten Punkte: 104
```

```
# Clusterzuweisungen in Dataframe umwandeln
```

```
dbscan_clusters <- data.frame(Cluster = factor(dbscan_result$cluster))
```

```
# Originaldaten hinzufügen
```

```
dbscan_clusters <- cbind(dbscan_clusters, heart_disease_patients)
```

```
# Clusterzusammenfassung erstellen
```

```
cluster_summary <- aggregate(. ~ Cluster, data = dbscan_clusters, mean)
```

```
print(cluster_summary)
```

```
##   Cluster    age trestbps    chol  thalach  oldpeak
## 1      0 56.74038 138.9904 252.0288 140.2596 1.685577
## 2      1 52.89385 127.1285 241.7095 157.9050 0.624581
## 3      2 39.25000 134.5000 200.0000 175.0000 1.550000
## 4      3 59.88889 125.8889 271.6667 101.4444 1.611111
```

```
## 5      4 59.75000 158.7500 271.0000 120.2500 0.400000
## 6      5 63.66667 128.3333 314.0000 128.3333 1.866667
```

In diesem Abschnitt wird das DBSCAN-Clustering durchgeführt. Die Parameter wurden so gewählt, dass die Anzahl der Cluster als sinnvoll erschienen und möglichst wenig Rauschen in den Daten erkannt wurde. Die Visualisierung zeigt, dass die Cluster mittels DBSCAN Verfahren nicht wirklich gut voneinander getrennt sind. Die Clusterzusammenfassung zeigt die Mittelwerte der Cluster für jede numerische Variable.

Auch die visuelle Analyse zeigt, dass keine eindeutige Clusterstruktur erkennbar ist. Die Cluster sind nicht klar voneinander getrennt und die Clusterzusammenfassung zeigt, dass die Mittelwerte der Cluster für jede numerische Variable sehr ähnlich sind. Vielmehr scheinen die Daten eine unstrukturierte Punktwolke zu bilden. Dies könnte darauf hindeuten, dass die Daten nicht in klar abgegrenzte Cluster unterteilt werden können.