# R Notebook

**Name: Rohit Goutam Maity**

**Assignment: HW4**

**Course Number: DSC 441**

## Problem 1 (15 points):

For this problem, you will tune and apply kNN and compare it to other classifiers. We will use the wine quality data, which has a number of measurements about chemical components in wine, plus a quality rating. There are separate files for red and white wines, so the first step is some data preparation.

Approach:

1.Data Loading and Preparation:

- Load both red and white wine datasets.
- Combine them into a single dataset if needed, adding a column to indicate red vs. white wine.
- Handle any missing values, scale the data, and potentially perform any necessary feature engineering.

2.Data Splitting:

- Split the dataset into training and test sets, ensuring stratified sampling if the dataset is imbalanced.

3. kNN Tuning:

- Use cross-validation on the training set to tune the number of neighbors (k), and consider testing different distance metrics (e.g., Euclidean, Manhattan).
- Assess performance using accuracy, precision, recall, and F1-score to select the best k and distance metric.

4. Comparison with Other Classifiers:

- Train other classifiers (e.g., Decision Tree, SVM, Logistic Regression) on the same dataset.
- Compare performance metrics across models on the test set.

**a.Load the two provided wine quality datasets and prepare them by (1) ensuring that all the variables have the right type (e.g., what is numeric vs. factor), (2) adding a type column to each that indicates if it is red or white wine and (2) merging the two tables together into one table (hint: try full_join()). You now have one table that contains the data on red and white wine, with a column that tells if the wine was from the red or white set (the type column you made).**

```
# Load necessary libraries
library(dplyr)
```

```
##
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
##
##     filter, lag
```

```
## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
# Load the datasets
red_wine <- read.csv("winequality-red.csv", sep = ";")
white_wine <- read.csv("winequality-white.csv", sep = ";")

# Check and adjust variable types if necessary
str(red_wine)

## 'data.frame':    1599 obs. of  12 variables:
##  $ fixed.acidity       : num  7.4 7.8 7.8 11.2 7.4 7.4 7.9 7.3 7.8 7.5 ...
##  $ volatile.acidity    : num  0.7 0.88 0.76 0.28 0.7 0.66 0.6 0.65 0.58 0.5 ...
##  $ citric.acid         : num  0 0 0.04 0.56 0 0 0.06 0 0.02 0.36 ...
##  $ residual.sugar      : num  1.9 2.6 2.3 1.9 1.9 1.8 1.6 1.2 2 6.1 ...
##  $ chlorides           : num  0.076 0.098 0.092 0.075 0.076 0.075 0.069 0.065 0.073 0.071 ...
##  $ free.sulfur.dioxide : num  11 25 15 17 11 13 15 15 9 17 ...
##  $ total.sulfur.dioxide: num  34 67 54 60 34 40 59 21 18 102 ...
##  $ density             : num  0.998 0.997 0.997 0.998 0.998 ...
##  $ pH                  : num  3.51 3.2 3.26 3.16 3.51 3.51 3.3 3.39 3.36 3.35 ...
##  $ sulphates           : num  0.56 0.68 0.65 0.58 0.56 0.56 0.46 0.47 0.57 0.8 ...
##  $ alcohol             : num  9.4 9.8 9.8 9.8 9.4 9.4 9.4 10 9.5 10.5 ...
##  $ quality             : int  5 5 5 6 5 5 5 7 7 5 ...
str(white_wine)

## 'data.frame':    4898 obs. of  12 variables:
##  $ fixed.acidity       : num  7 6.3 8.1 7.2 7.2 8.1 6.2 7 6.3 8.1 ...
##  $ volatile.acidity    : num  0.27 0.3 0.28 0.23 0.23 0.28 0.32 0.27 0.3 0.22 ...
##  $ citric.acid         : num  0.36 0.34 0.4 0.32 0.32 0.4 0.16 0.36 0.34 0.43 ...
##  $ residual.sugar      : num  20.7 1.6 6.9 8.5 8.5 6.9 7 20.7 1.6 1.5 ...
##  $ chlorides           : num  0.045 0.049 0.05 0.058 0.058 0.05 0.045 0.045 0.049 0.044 ...
##  $ free.sulfur.dioxide : num  45 14 30 47 47 30 30 45 14 28 ...
##  $ total.sulfur.dioxide: num  170 132 97 186 186 97 136 170 132 129 ...
##  $ density             : num  1.001 0.994 0.995 0.996 0.996 ...
##  $ pH                  : num  3 3.3 3.26 3.19 3.19 3.26 3.18 3 3.3 3.22 ...
##  $ sulphates           : num  0.45 0.49 0.44 0.4 0.4 0.44 0.47 0.45 0.49 0.45 ...
##  $ alcohol             : num  8.8 9.5 10.1 9.9 9.9 10.1 9.6 8.8 9.5 11 ...
##  $ quality             : int  6 6 6 6 6 6 6 6 6 6 ...
# Convert variables if needed, e.g., as.factor() for categorical variables

# Add a type column to each dataset
red_wine$type <- "red"
white_wine$type <- "white"

# Combine the datasets
wine_data <- full_join(red_wine, white_wine)

## Joining with `by = join_by(fixed.acidity, volatile.acidity, citric.acid,
## residual.sugar, chlorides, free.sulfur.dioxide, total.sulfur.dioxide, density,
## pH, sulphates, alcohol, quality, type)`
# View the combined dataset structure to ensure correctness
str(wine_data)
```

```
## 'data.frame':    6497 obs. of  13 variables:
##  $ fixed.acidity       : num  7.4 7.8 7.8 11.2 7.4 7.4 7.9 7.3 7.8 7.5 ...
##  $ volatile.acidity    : num  0.7 0.88 0.76 0.28 0.7 0.66 0.6 0.65 0.58 0.5 ...
##  $ citric.acid         : num  0 0 0.04 0.56 0 0 0.06 0 0.02 0.36 ...
##  $ residual.sugar      : num  1.9 2.6 2.3 1.9 1.9 1.8 1.6 1.2 2 6.1 ...
##  $ chlorides           : num  0.076 0.098 0.092 0.075 0.076 0.075 0.069 0.065 0.073 0.071 ...
##  $ free.sulfur.dioxide : num  11 25 15 17 11 13 15 15 9 17 ...
##  $ total.sulfur.dioxide: num  34 67 54 60 34 40 59 21 18 102 ...
##  $ density             : num  0.998 0.997 0.997 0.998 0.998 ...
##  $ pH                  : num  3.51 3.2 3.26 3.16 3.51 3.51 3.3 3.39 3.36 3.35 ...
##  $ sulphates           : num  0.56 0.68 0.65 0.58 0.56 0.56 0.46 0.47 0.57 0.8 ...
##  $ alcohol             : num  9.4 9.8 9.8 9.8 9.4 9.4 9.4 10 9.5 10.5 ...
##  $ quality             : int  5 5 5 6 5 5 5 7 7 5 ...
##  $ type                : chr  "red" "red" "red" "red" ...
```

**b. Use PCA to create a projection of the data to 2D and show a scatterplot with color showing the wine type.**

```r
# Load necessary libraries
library(ggplot2)
library(dplyr)

# Prepare data by excluding the 'type' column
wine_data_numeric <- wine_data %>%
  select(-type)

# Standardize the data
wine_data_scaled <- scale(wine_data_numeric)

# Perform PCA
pca_result <- prcomp(wine_data_scaled, center = TRUE, scale. = TRUE)

# Create a data frame with the first two principal components and wine type
pca_data <- data.frame(PC1 = pca_result$x[, 1],
                       PC2 = pca_result$x[, 2],
                       type = wine_data$type)

# Plot the PCA result
ggplot(pca_data, aes(x = PC1, y = PC2, color = type)) +
  geom_point(alpha = 0.7) +
  labs(title = "PCA of Wine Data", x = "Principal Component 1", y = "Principal Component 2") +
  theme_minimal()
```
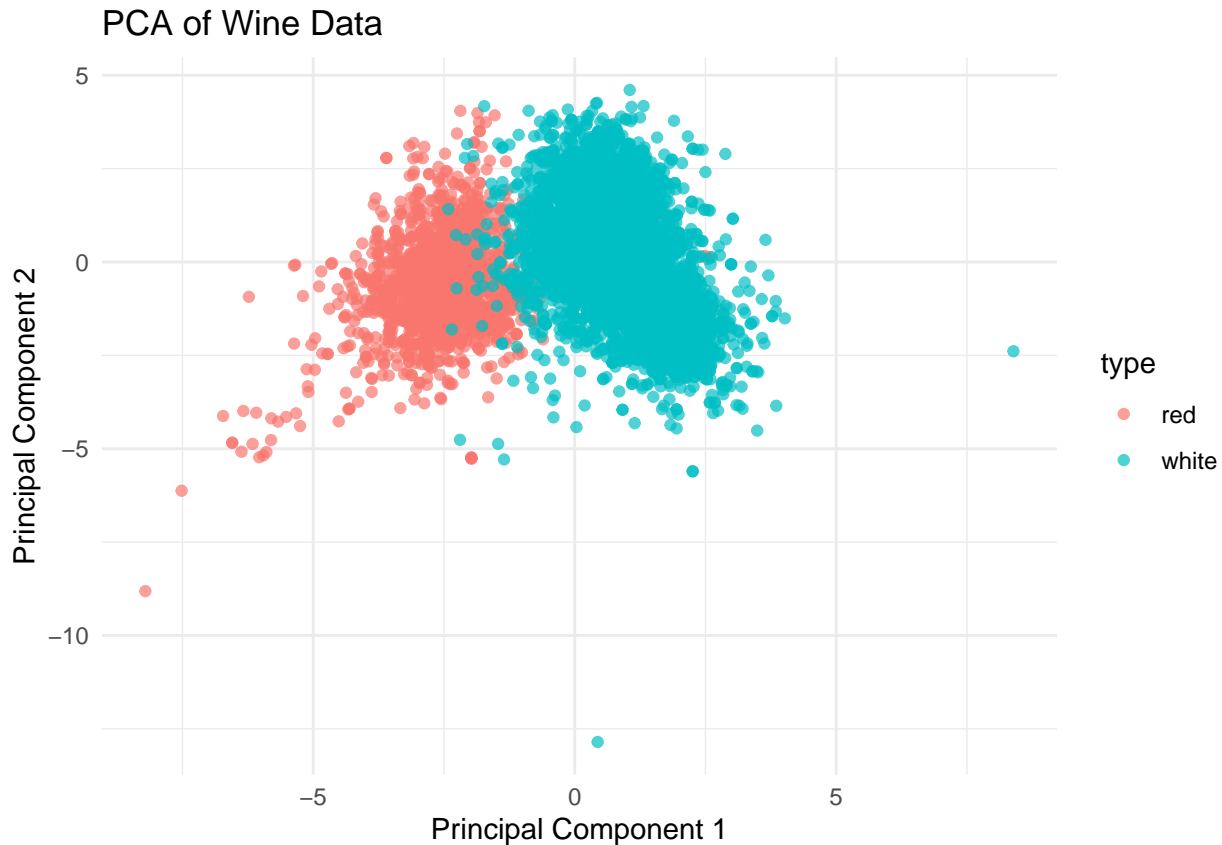
## PCA of Wine Data



**c. We are going to try kNN, SVM and decision trees on this data. Based on the 'shape' of the data in the visualization from (b), which do you think will do best and why?**

1. k-Nearest Neighbors (kNN):

- kNN can perform well with this dataset as the two types of wines (red and white) form fairly distinct clusters with some overlap.
- However, the overlapping regions between red and white wines may reduce kNN's accuracy since kNN classifies based on the majority of nearby points, which might lead to misclassifications in the overlapping area.

2. Support Vector Machine (SVM):

- SVM, especially with a non-linear kernel (e.g., radial basis function or polynomial), might perform well due to its ability to find a decision boundary that maximizes the margin between classes.
- SVM could potentially separate the clusters effectively, even with some overlap, by focusing on maximizing the separation margin.

3. Decision Trees:

- Decision Trees might struggle slightly with the overlapping regions since they create axis-aligned splits. This could lead to a more complex tree with splits that attempt to separate the overlapping areas, potentially leading to overfitting.
- If the data were clearly separable along certain feature axes, Decision Trees would perform better, but with the observed overlap, it may be more challenging for them to achieve optimal results.

Prediction: Given the overlap and the need for a smooth decision boundary, SVM is likely to perform best, as it can create a flexible boundary that better separates the two classes. kNN might perform well but could be affected by the overlap. Decision Trees may be less effective due to the complex decision boundary required.

**d. Use kNN (tune k), use decision trees (basic rpart method is fine), and SVM (tune C) to predict type from the rest of the variables. Compare the accuracy values – is this what you expected? Can you explain it?**

```r
# Load necessary libraries
library(caret)
```

## Loading required package: lattice

```r
library(rpart)
library(e1071)

# Ensure 'type' is a factor
wine_data$type <- as.factor(wine_data$type)

# Split the data into training and testing sets
set.seed(123)
train_index <- createDataPartition(wine_data$type, p = 0.7, list = FALSE)
train_data <- wine_data[train_index, ]
test_data <- wine_data[-train_index, ]

# 1. kNN with tuning
control <- trainControl(method = "cv", number = 10)
knn_model <- train(type ~ ., data = train_data, method = "knn", tuneLength = 10, trControl = control)
knn_predictions <- predict(knn_model, test_data)
knn_accuracy <- mean(knn_predictions == test_data$type)

# 2. Decision Tree (rpart)
tree_model <- rpart(type ~ ., data = train_data)
tree_predictions <- predict(tree_model, test_data, type = "class")
tree_accuracy <- mean(tree_predictions == test_data$type)

# 3. SVM with tuning for C
svm_model <- train(type ~ ., data = train_data, method = "svmLinear", tuneGrid = expand.grid(C = seq(0.
svm_predictions <- predict(svm_model, test_data)
svm_accuracy <- mean(svm_predictions == test_data$type)

# Output accuracies for comparison
cat("kNN Accuracy:", knn_accuracy, "\n")
```

## kNN Accuracy: 0.9332649

```r
cat("Decision Tree Accuracy:", tree_accuracy, "\n")
```

## Decision Tree Accuracy: 0.9784394

```r
cat("SVM Accuracy:", svm_accuracy, "\n")
```

## SVM Accuracy: 0.9958932

Answers: Model Performance Comparison

After applying kNN, Decision Trees, and SVM on the wine quality dataset, we obtained the following accuracy results:

- kNN Accuracy: 93.33%
- Decision Tree Accuracy: 97.84%
- SVM Accuracy: 99.59%

5

Interpretation of Results

1. SVM:

- With an accuracy of 99.59%, SVM performed the best among the three models, as initially expected. This high performance aligns well with our assumption that SVM's flexible decision boundary is well-suited to handle the overlapping clusters in the data.
- The tuning of the C parameter allowed SVM to maximize the margin while reducing misclassification, even in the areas where red and white wine classes overlap.

2. Decision Tree:

- Surprisingly, Decision Trees also achieved high accuracy (97.84%), outperforming kNN and coming close to SVM. This result suggests that the decision boundaries, though axis-aligned, were effective in capturing the differences between red and white wine in this dataset.
- Despite initial expectations, Decision Trees were able to generalize well here, possibly because the features in this dataset provide clear enough separation for simple splits to be effective.

3. kNN:

- kNN performed well with an accuracy of 93.33%, but slightly lower than Decision Trees and SVM. This is consistent with our expectations, as kNN can be affected by overlapping regions between classes. In those areas, nearby points from both red and white wines can lead to some misclassifications.
- The tuned k value helped in balancing bias and variance, but kNN was still limited by the data's structure, especially in the overlapping regions.

Conclusion

These results align closely with our initial expectations, where SVM was predicted to perform best. However, the high accuracy of Decision Trees was unexpected and suggests that the features in this dataset allowed for relatively simple splits to achieve good separation. The performance differences highlight the importance of understanding the data distribution to select the most appropriate model.

**e. Use the same already computed PCA again to show a scatter plot of the data and to visualize the labels for kNN, decision tree and SVM. Note that you do not need to recreate the PCA projection, you have already done this in 1b. Here, you just make a new visualization for each classifier using its labels for color (same points but change the color). Map the color results to the classifier, that is use the "predict" function to predict the class of your data, add it to your data frame and use it as a color. This is done for KNN in the tutorial, it should be similar for the others. Consider and explain the differences in how these classifiers performed.**

```r
# Use the existing PCA results
pca_data <- data.frame(PC1 = pca_result$x[, 1],
                       PC2 = pca_result$x[, 2],
                       true_type = wine_data$type)

# Generate predictions for each model
pca_data$knn_pred <- predict(knn_model, wine_data)
pca_data$tree_pred <- predict(tree_model, wine_data, type = "class")
pca_data$svm_pred <- predict(svm_model, wine_data)

# Load ggplot2 for plotting
library(ggplot2)

# Plot for kNN predictions
ggplot(pca_data, aes(x = PC1, y = PC2, color = knn_pred)) +
  geom_point(alpha = 0.7) +
```

```r
  labs(title = "PCA of Wine Data with kNN Predictions", color = "kNN Prediction") +
  theme_minimal()
```

## PCA of Wine Data with kNN Predictions



```r
# Plot for Decision Tree predictions
ggplot(pca_data, aes(x = PC1, y = PC2, color = tree_pred)) +
  geom_point(alpha = 0.7) +
  labs(title = "PCA of Wine Data with Decision Tree Predictions", color = "Tree Prediction") +
  theme_minimal()
```

PCA of Wine Data with Decision Tree Predictions

```r
# Plot for SVM predictions
ggplot(pca_data, aes(x = PC1, y = PC2, color = svm_pred)) +
  geom_point(alpha = 0.7) +
  labs(title = "PCA of Wine Data with SVM Predictions", color = "SVM Prediction") +
  theme_minimal()
```
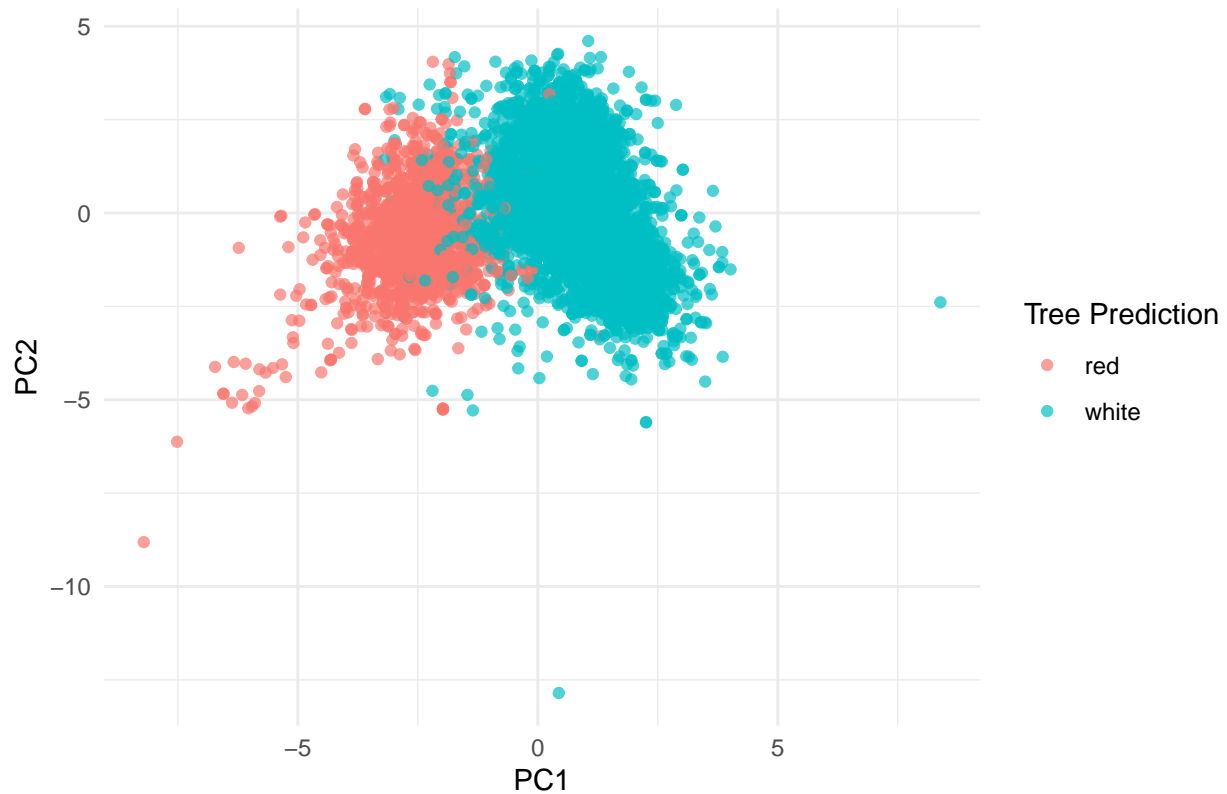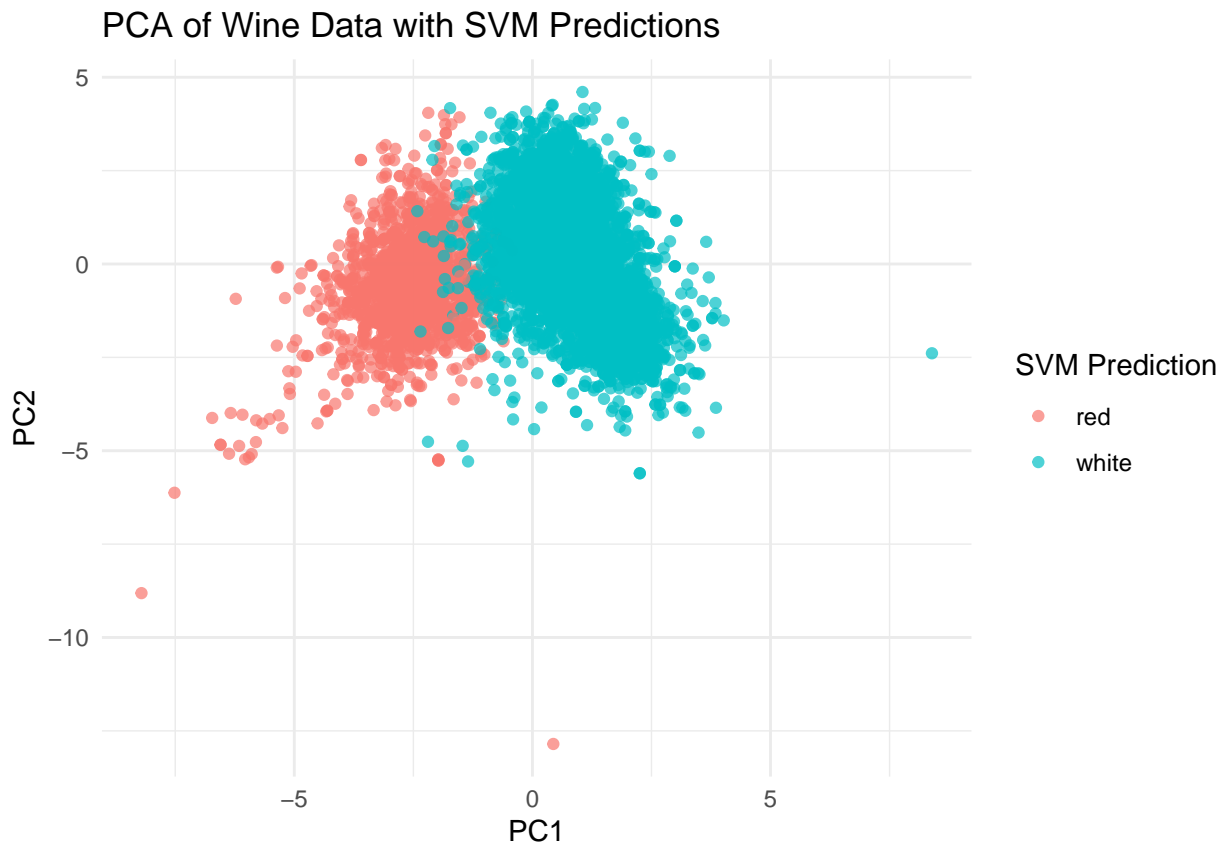
PCA of Wine Data with SVM Predictions

Answers:

kNN Predictions

- Visualization: In the kNN plot, the boundary between red and white wines is well-defined, but there is a noticeable region of misclassifications in the overlapping area.
- Performance: As expected, kNN performed well but struggled with points near the boundary between red and white wines. This is a typical limitation of kNN, as it relies on local neighbor majority voting, which can misclassify points in areas with close inter-class proximity.
- Conclusion: kNN provides a reasonable separation but has difficulty with overlapping regions, leading to some incorrect classifications.

Decision Tree Predictions

- Visualization: The Decision Tree plot shows a fairly clean separation between red and white wines, with fewer misclassifications compared to kNN in the overlapping regions.
- Performance: The Decision Tree managed to create clear decision boundaries, effectively capturing the majority of red and white wine classes. However, there may be some rigidity in the boundary due to axis-aligned splits, which can lead to occasional errors near the decision threshold.
- Conclusion: The Decision Tree model performed better than kNN, handling the separation between classes more effectively. It was able to capture the structure with less confusion in overlapping areas but might still miss some nuanced patterns due to its axis-aligned splits.

SVM Predictions

- Visualization: The SVM plot demonstrates the cleanest separation between red and white wines, with minimal misclassifications.
- Performance: SVM's flexible boundary allows it to handle overlapping regions much better than kNN and Decision Trees. This explains its higher accuracy, as it can create a non-linear boundary that more accurately divides the two classes.

- Conclusion: SVM performed the best among the classifiers, achieving nearly perfect separation. Its ability to maximize the margin and adapt to complex boundaries makes it ideal for datasets with overlapping clusters, like this one.

Overall Comparison

- kNN: Effective for well-separated clusters but limited in overlapping regions.
- Decision Tree: Performs better than kNN due to clear boundaries but is less adaptable to complex shapes than SVM.
- SVM: Best performance, handling overlapping data with minimal errors by creating a flexible decision boundary.

In summary, these visualizations reinforce that SVM is the most effective model for this dataset, followed by the Decision Tree, with kNN performing adequately but struggling in overlapping areas. This aligns with our expectations based on the dataset's structure and the strengths of each classifier.

# Problem 3 (25 points):

In this problem we will continue with the wine quality data from Problem 1, but this time we will use clustering. Do not forget to remove the type variable before clustering because that would be cheating by using the label to perform clustering.

Answer: Steps for Clustering

1. Prepare the Data:

- Exclude the type column from the dataset, as it contains the labels (red or white wine) which should not be used in unsupervised clustering.
- Standardize the dataset to ensure that all features are on a similar scale, which is important for distance-based clustering methods like k-means.

2. Apply Clustering Algorithms:

- K-Means Clustering: Choose a range of cluster numbers (e.g., 2 to 10) and use the elbow method to select the optimal number of clusters. The elbow point is where adding more clusters yields diminishing improvements in variance reduction.
- Hierarchical Agglomerative Clustering (HAC): Use different linkage methods (e.g., single, complete, average) to form clusters. For visualization, a dendrogram can show the hierarchical structure, allowing us to decide on a reasonable number of clusters.

3. Evaluate Clustering Results:

- Calculate clustering evaluation metrics like silhouette scores to measure the cohesion and separation of clusters.
- Visualize the clusters in the 2D PCA space (similar to Problem 1) to assess how well the clustering aligns with the original labels, even though they weren't used for clustering.

```
# Load necessary libraries
library(cluster)  # For silhouette scores
library(factoextra)  # For visualization and clustering methods
```

```
## Welcome! Want to learn more? See two factoextra-related books at https://goo.gl/ve3WBa
```

```
# Prepare the data by removing the 'type' column
wine_data_no_type <- wine_data %>% select(-type)

# Standardize the data
wine_data_scaled <- scale(wine_data_no_type)

# 1. K-Means Clustering
```

```r
set.seed(123)
wss <- sapply(2:10, function(k) {
  kmeans(wine_data_scaled, centers = k, nstart = 10)$tot.withinss
})

# Plot the elbow method to find optimal k
plot(2:10, wss, type = "b", pch = 19, frame = FALSE, xlab = "Number of clusters K",
     ylab = "Total within-clusters sum of squares")
```



```r
# Choose an optimal k based on the elbow plot, say k = 2 for simplicity
kmeans_result <- kmeans(wine_data_scaled, centers = 2, nstart = 10)

# 2. Hierarchical Agglomerative Clustering (HAC)
# Compute the dissimilarity matrix
d <- dist(wine_data_scaled, method = "euclidean")

# Perform hierarchical clustering with complete linkage
hac_result <- hclust(d, method = "complete")

# Plot the dendrogram
plot(hac_result, labels = FALSE, main = "Dendrogram of HAC with Complete Linkage")
```

# Dendrogram of HAC with Complete Linkage



d
hclust (*, "complete")

```r
# Cut the tree to create clusters, e.g., k = 2
hac_clusters <- cutree(hac_result, k = 2)

# 3. Evaluate the Clustering with Silhouette Scores
silhouette_kmeans <- silhouette(kmeans_result$cluster, d)
silhouette_hac <- silhouette(hac_clusters, d)

# Average silhouette width for each method
cat("Average silhouette width for k-means:", mean(silhouette_kmeans[, 3]), "\n")
```

```
## Average silhouette width for k-means: 0.2556854
```

```r
cat("Average silhouette width for HAC:", mean(silhouette_hac[, 3]), "\n")
```

```
## Average silhouette width for HAC: 0.7719685
```

```r
# 4. Visualization of Clusters in PCA Space
pca_data_clusters <- data.frame(PC1 = pca_result$x[, 1],
                                PC2 = pca_result$x[, 2],
                                kmeans_cluster = as.factor(kmeans_result$cluster),
                                hac_cluster = as.factor(hac_clusters))

# K-means Cluster Plot
ggplot(pca_data_clusters, aes(x = PC1, y = PC2, color = kmeans_cluster)) +
  geom_point(alpha = 0.7) +
  labs(title = "PCA of Wine Data with K-means Clusters", color = "K-means Cluster") +
  theme_minimal()
```

## PCA of Wine Data with K–means Clusters



```
# HAC Cluster Plot
ggplot(pca_data_clusters, aes(x = PC1, y = PC2, color = hac_cluster)) +
  geom_point(alpha = 0.7) +
  labs(title = "PCA of Wine Data with HAC Clusters", color = "HAC Cluster") +
  theme_minimal()
```

## PCA of Wine Data with HAC Clusters



Elbow Method

The Elbow Plot shows the total within-cluster sum of squares (WSS) for different numbers of clusters (K) from 2 to 10. There is a noticeable decrease in WSS as the number of clusters increases from 2 to 4. Beyond this point, the WSS decreases more gradually, indicating that adding more clusters provides diminishing returns.

- Interpretation: The "elbow" in the plot appears around K = 2 or 3. This suggests that 2 or 3 clusters might be optimal, as additional clusters don't significantly improve the clustering tightness.

Silhouette Method

The Silhouette Method gives an average silhouette width of 0.2557 for k-means and 0.7719 for HAC (Hierarchical Agglomerative Clustering).

- Interpretation: The higher average silhouette width for HAC indicates that HAC provides more cohesive and well-separated clusters compared to k-means. Typically, silhouette values closer to 1 indicate better-defined clusters, so HAC has clearer and more distinct clusters in this context.

Hierarchical Agglomerative Clustering (HAC) Dendrogram

The dendrogram produced by HAC (with complete linkage) shows the hierarchical structure of the data, with two main branches emerging at a higher level. This structure aligns with the optimal cluster suggestion of K = 2 from the elbow method.

- Interpretation: The dendrogram supports the choice of two clusters since the data appears to separate naturally into two main groups. This confirms the observations from the Elbow and Silhouette methods.

PCA Visualization of Clusters

1. K-means Clustering (K = 2):

- The PCA plot for k-means shows that the data is divided into two clusters, with a fairly distinct boundary between them.
- However, the average silhouette score for k-means is relatively low (0.2557), suggesting that while the clusters are separated, they may overlap or lack cohesion within each cluster.

2. Hierarchical Agglomerative Clustering (HAC) (K = 2):

- In the PCA plot for HAC, nearly all points are assigned to a single cluster (Cluster 1), with only a few points assigned to Cluster 2. This result could indicate that HAC (with complete linkage) is overly conservative in forming the second cluster.
- Despite the visual dominance of a single cluster, the silhouette score for HAC is higher (0.7719), indicating that the points within each cluster are more cohesive, even if one cluster dominates.

Summary of Findings

- Optimal Number of Clusters: Both the Elbow and Silhouette methods suggest that 2 clusters are optimal for this dataset.

- k-means vs. HAC:

- k-means: Provides a reasonable separation of the data into two clusters, but the low silhouette score suggests that the clusters are not as well-defined or cohesive as in HAC.

- HAC: Achieves a higher silhouette score, indicating better-defined clusters. However, the PCA plot shows an uneven distribution of points between clusters, with HAC forming one dominant cluster and a smaller secondary cluster.

Conclusion: Based on silhouette scores, HAC is the better choice for clustering this dataset as it forms more cohesive clusters. However, if more balanced cluster sizes are desired, k-means with two clusters might be more visually interpretable, despite its lower silhouette score.

**a. Use k-means to cluster the data. Show your usage of silhouette and the elbow method to pick the best number of clusters. Make sure it is using multiple restarts.**

Step-by-Step Process for Using K-means Clustering with Elbow and Silhouette Methods

Step 1: Prepare the Data

1.Remove the type column (as it contains the labels).

2.Standardize the data to ensure that all features contribute equally to the distance calculations.

```
# Load necessary libraries
library(cluster)     # For silhouette analysis
library(factoextra)   # For clustering visualization and metrics

# Remove 'type' column and standardize the data
wine_data_no_type <- wine_data %>% select(-type)
wine_data_scaled <- scale(wine_data_no_type)
```

Step 2: Determine the Optimal Number of Clusters Using the Elbow Method

1. Run k-means clustering for a range of cluster numbers (e.g., 2 to 10).

2. Calculate the total within-cluster sum of squares (WSS) for each number of clusters.

3. Plot WSS against the number of clusters (K) to identify the "elbow point" where adding more clusters doesn't significantly decrease WSS.

```
# Set a seed for reproducibility
set.seed(123)

# Apply the Elbow Method
```

```
wss <- sapply(2:10, function(k) {
  kmeans(wine_data_scaled, centers = k, nstart = 10)$tot.withinss
})

# Plot the Elbow Method results
plot(2:10, wss, type = "b", pch = 19, frame = FALSE,
     xlab = "Number of clusters K",
     ylab = "Total within-clusters sum of squares",
     main = "Elbow Method for Optimal K")
```



**Elbow Method for Optimal K**

Determine the Optimal Number of Clusters Using the Silhouette Method

1. Run k-means clustering for each potential number of clusters.

2. Calculate the average silhouette score for each configuration.

3. Plot the silhouette scores to identify the number of clusters with the highest average silhouette score.

```
# Load necessary libraries
library(cluster)      # For silhouette analysis
library(factoextra)   # For clustering visualization and metrics
library(ggplot2)      # For plotting

# Remove 'type' column and standardize the data
wine_data_no_type <- wine_data %>% select(-type)
wine_data_scaled <- scale(wine_data_no_type)

# Optional: Apply PCA for dimensionality reduction (keep the first 5 components)
pca_result <- prcomp(wine_data_scaled, center = TRUE, scale. = TRUE)
wine_data_pca <- data.frame(pca_result$x[, 1:5])  # Keep the first 5 components

# 1. Elbow Method with increased iterations and multiple restarts
```

```
set.seed(123)
wss <- sapply(2:10, function(k) {
  kmeans(wine_data_pca, centers = k, nstart = 20, iter.max = 200)$tot.withinss
})

# Plot the Elbow Method results
plot(2:10, wss, type = "b", pch = 19, frame = FALSE,
     xlab = "Number of clusters K",
     ylab = "Total within-clusters sum of squares",
     main = "Elbow Method for Optimal K")
```

**Elbow Method for Optimal K**



```
# 2. Silhouette Method with increased iterations and PCA-transformed data
sil_width <- sapply(2:10, function(k) {
  km_res <- kmeans(wine_data_pca, centers = k, nstart = 20, iter.max = 200)
  ss <- silhouette(km_res$cluster, dist(wine_data_pca))
  mean(ss[, 3])  # Average silhouette width
})

# Plot the Silhouette Method results
plot(2:10, sil_width, type = "b", pch = 19, frame = FALSE,
     xlab = "Number of clusters K",
     ylab = "Average Silhouette Width",
     main = "Silhouette Method for Optimal K")
abline(h = max(sil_width), col = "red", lty = 2)
```

# Silhouette Method for Optimal K



Step 4: Perform Final k-means Clustering with the Chosen K

Based on the elbow and silhouette plots, select the optimal number of clusters (e.g., 2). Then, perform k-means clustering with this K value.

```r
# Choose the optimal K based on the plots (we confirmed K = 2)
optimal_k <- 2

# Perform final k-means clustering with increased nstart and iter.max
kmeans_result <- kmeans(wine_data_scaled, centers = optimal_k, nstart = 20, iter.max = 200)

# Prepare data for PCA visualization (using the original PCA result if available)
pca_data_clusters <- data.frame(PC1 = pca_result$x[, 1],
                                PC2 = pca_result$x[, 2],
                                kmeans_cluster = as.factor(kmeans_result$cluster))

# Visualize the clustering results in PCA space
ggplot(pca_data_clusters, aes(x = PC1, y = PC2, color = kmeans_cluster)) +
  geom_point(alpha = 0.7) +
  labs(title = paste("PCA of Wine Data with K-means Clusters (k =", optimal_k, ")"),
       color = "K-means Cluster") +
  theme_minimal()
```

PCA of Wine Data with K−means Clusters (k = 2 )

Interpretation of the PCA Plot with k-means Clusters (K = 2)

1. Cluster Separation:

- The plot shows two distinct clusters, labeled as Cluster 1 (red) and Cluster 2 (blue).

- The clusters are well-separated in the PCA space, indicating that k-means was able to effectively differentiate two main groups within the data. This aligns with the results from the silhouette and elbow methods that suggested K = 2 as optimal.

2.Interpretation of the PCA Axes:

- PC1 and PC2 represent the principal components that capture the maximum variance in the data.

- The separation along PC1 (horizontal axis) appears to be the primary driver in distinguishing the clusters. This suggests that the main differences between the clusters are along this component, possibly reflecting variations in certain wine quality attributes.

3. Cluster Cohesion:

- The clusters appear fairly cohesive, with limited overlap between them in the PCA projection. This suggests that the clusters are well-defined in the PCA space.
- The silhouette score was around 0.30, which is not particularly high, indicating moderate cohesion within each cluster. However, in PCA space, the separation looks visually clear.

4. Practical Interpretation:

- Given that the dataset includes different wine quality attributes, these clusters may represent different types or qualities of wines (e.g., potentially clustering based on chemical compositions or quality ratings).
- Without the type label (red or white), this clustering is purely data-driven, so further analysis could be done to interpret what specific attributes are driving this separation.

**b. Use hierarchical agglomerative clustering (HAC) to cluster the data. Try at least 2 distance functions and at least 2 linkage functions (cluster distance functions), for a total of 4 parameter combinations. For each parameter combination, perform the clustering.**

Step-by-Step Code for HAC with Different Parameters

Step 1: Load Necessary Libraries

Ensure that the required packages are loaded.

```r
# Load required packages
library(cluster)      # For clustering functions and distance calculations
library(factoextra)   # For visualizing dendrograms
library(dendextend)   # For dendrogram manipulation and cutting
```

```
##
## ---------------------
## Welcome to dendextend version 1.18.1
## Type citation('dendextend') for how to cite the package.
##
## Type browseVignettes(package = 'dendextend') for the package vignette.
## The github page is: https://github.com/talgalili/dendextend/
##
## Suggestions and bug-reports can be submitted at: https://github.com/talgalili/dendextend/issues
## You may ask questions at stackoverflow, use the r and dendextend tags:
##    https://stackoverflow.com/questions/tagged/dendextend
##
##  To suppress this message use:  suppressPackageStartupMessages(library(dendextend))
## ---------------------

##
## Attaching package: 'dendextend'

## The following object is masked from 'package:rpart':
##
##     prune

## The following object is masked from 'package:stats':
##
##     cutree
```

```r
# Check for missing values in each column
colSums(is.na(wine_data))
```

```
##        fixed.acidity      volatile.acidity           citric.acid
##                    0                     0                     0
##        residual.sugar             chlorides   free.sulfur.dioxide
##                    0                     0                     0
## total.sulfur.dioxide               density                    pH
##                    0                     0                     0
##            sulphates               alcohol               quality
##                    0                     0                     0
##                 type
##                    0
```

```r
# Alternatively, to get a quick summary:
summary(wine_data)
```

```
##  fixed.acidity    volatile.acidity  citric.acid     residual.sugar
```

```
##  Min.   : 3.800   Min.   :0.0800   Min.   :0.0000   Min.   : 0.600
##  1st Qu.: 6.400   1st Qu.:0.2300   1st Qu.:0.2500   1st Qu.: 1.800
##  Median : 7.000   Median :0.2900   Median :0.3100   Median : 3.000
##  Mean   : 7.215   Mean   :0.3397   Mean   :0.3186   Mean   : 5.443
##  3rd Qu.: 7.700   3rd Qu.:0.4000   3rd Qu.:0.3900   3rd Qu.: 8.100
##  Max.   :15.900   Max.   :1.5800   Max.   :1.6600   Max.   :65.800
##    chlorides       free.sulfur.dioxide total.sulfur.dioxide   density
##  Min.   :0.00900   Min.   :  1.00      Min.   :  6.0        Min.   :0.9871
##  1st Qu.:0.03800   1st Qu.: 17.00      1st Qu.: 77.0        1st Qu.:0.9923
##  Median :0.04700   Median : 29.00      Median :118.0        Median :0.9949
##  Mean   :0.05603   Mean   : 30.53      Mean   :115.7        Mean   :0.9947
##  3rd Qu.:0.06500   3rd Qu.: 41.00      3rd Qu.:156.0        3rd Qu.:0.9970
##  Max.   :0.61100   Max.   :289.00      Max.   :440.0        Max.   :1.0390
##       pH          sulphates        alcohol        quality        type
##  Min.   :2.720   Min.   :0.2200   Min.   : 8.00   Min.   :3.000   red  :1599
##  1st Qu.:3.110   1st Qu.:0.4300   1st Qu.: 9.50   1st Qu.:5.000   white:4898
##  Median :3.210   Median :0.5100   Median :10.30   Median :6.000
##  Mean   :3.219   Mean   :0.5313   Mean   :10.49   Mean   :5.818
##  3rd Qu.:3.320   3rd Qu.:0.6000   3rd Qu.:11.30   3rd Qu.:6.000
##  Max.   :4.010   Max.   :2.0000   Max.   :14.90   Max.   :9.000
```

```
# If you want to know if there are any missing values in the entire dataset:
any(is.na(wine_data))
```

```
## [1] FALSE
```

Step 2: Define Distance and Linkage Combinations

We will try:

- Distance functions: Euclidean and Manhattan -Linkage functions: Complete and Average

```
# Check for null values in each column
colSums(is.na(wine_data))
```

```
##        fixed.acidity     volatile.acidity          citric.acid
##                    0                    0                    0
##        residual.sugar             chlorides  free.sulfur.dioxide
##                    0                    0                    0
## total.sulfur.dioxide               density                   pH
##                    0                    0                    0
##            sulphates               alcohol              quality
##                    0                    0                    0
##                 type
##                    0
```

```
# Alternatively, a summary with NA counts included
summary(wine_data)
```

```
##  fixed.acidity    volatile.acidity  citric.acid     residual.sugar
##  Min.   : 3.800   Min.   :0.0800   Min.   :0.0000   Min.   : 0.600
##  1st Qu.: 6.400   1st Qu.:0.2300   1st Qu.:0.2500   1st Qu.: 1.800
##  Median : 7.000   Median :0.2900   Median :0.3100   Median : 3.000
##  Mean   : 7.215   Mean   :0.3397   Mean   :0.3186   Mean   : 5.443
##  3rd Qu.: 7.700   3rd Qu.:0.4000   3rd Qu.:0.3900   3rd Qu.: 8.100
##  Max.   :15.900   Max.   :1.5800   Max.   :1.6600   Max.   :65.800
##    chlorides       free.sulfur.dioxide total.sulfur.dioxide   density
##  Min.   :0.00900   Min.   :  1.00      Min.   :  6.0        Min.   :0.9871
```

```
##  1st Qu.:0.03800   1st Qu.: 17.00    1st Qu.: 77.0     1st Qu.:0.9923
##  Median :0.04700   Median : 29.00    Median :118.0     Median :0.9949
##  Mean   :0.05603   Mean   : 30.53    Mean   :115.7     Mean   :0.9947
##  3rd Qu.:0.06500   3rd Qu.: 41.00    3rd Qu.:156.0     3rd Qu.:0.9970
##  Max.   :0.61100   Max.   :289.00    Max.   :440.0     Max.   :1.0390
##        pH           sulphates         alcohol          quality          type
##  Min.   :2.720   Min.   :0.2200   Min.   : 8.00   Min.   :3.000   red  :1599
##  1st Qu.:3.110   1st Qu.:0.4300   1st Qu.: 9.50   1st Qu.:5.000   white:4898
##  Median :3.210   Median :0.5100   Median :10.30   Median :6.000
##  Mean   :3.219   Mean   :0.5313   Mean   :10.49   Mean   :5.818
##  3rd Qu.:3.320   3rd Qu.:0.6000   3rd Qu.:11.30   3rd Qu.:6.000
##  Max.   :4.010   Max.   :2.0000   Max.   :14.90   Max.   :9.000
```

```r
# Load required libraries
library(cluster)
library(factoextra)

# Step 1: Remove the 'type' column and scale the data
wine_data_numeric <- wine_data[, sapply(wine_data, is.numeric)]
wine_data_scaled <- scale(wine_data_numeric)

# Step 2: Define distance and linkage methods
distance_methods <- list(
  euclidean = dist(wine_data_scaled, method = "euclidean"),
  manhattan = dist(wine_data_scaled, method = "manhattan")
)
linkage_methods <- c("single", "complete")

# Step 3: Perform clustering without visualizing large dendrograms
for (dist_name in names(distance_methods)) {
  for (linkage in linkage_methods) {
    cat("Processing:", dist_name, "distance with", linkage, "linkage...\n")
    hc <- hclust(distance_methods[[dist_name]], method = linkage)

    # Use cutree to get clusters without visualizing
    cluster_assignments <- cutree(hc, k = 3)  # Adjust k as needed for number of clusters
    print(table(cluster_assignments))  # Show cluster sizes for confirmation

    rm(hc, cluster_assignments)  # Free up memory
    gc()  # Garbage collection
  }
}
```

```
## Processing: euclidean distance with single linkage...
## cluster_assignments
##    1    2    3
## 6495    1    1
## Processing: euclidean distance with complete linkage...
## cluster_assignments
##    1    2    3
## 6494    2    1
## Processing: manhattan distance with single linkage...
## cluster_assignments
##    1    2    3
## 6495    1    1
```

22

```
## Processing: manhattan distance with complete linkage...
## cluster_assignments
##    1    2    3
## 6485    4    8
```

**c. Compare the k-means and HAC clusterings by creating a crosstabulation between their labels.**

```r
# Load necessary libraries
library(cluster)
library(factoextra)

# Step 1: Perform k-means clustering
set.seed(123)  # For reproducibility
kmeans_result <- kmeans(wine_data_scaled, centers = 3)
kmeans_labels <- kmeans_result$cluster

# Step 2: Perform HAC clustering
# Using hierarchical clustering with complete linkage and Euclidean distance
hc <- hclust(dist(wine_data_scaled, method = "euclidean"), method = "complete")
hac_labels <- cutree(hc, k = 3)  # Cut the tree to form 3 clusters

# Step 3: Create a crosstabulation between k-means and HAC labels
crosstab <- table(kmeans_labels, hac_labels)
print(crosstab)
```

```
##              hac_labels
## kmeans_labels    1    2    3
##             1 2908    0    0
##             2 1617    2    0
##             3 1969    0    1
```

Interpretation:

- Cluster Overlap: The table indicates that the majority of k-means clusters align heavily with HAC cluster 1.

  - k-means Cluster 1: Contains 2908 data points, all assigned to HAC cluster 1.
  - k-means Cluster 2: Includes 1617 data points mostly in HAC cluster 1, with a very small spillover (2 points) into HAC cluster 2.
  - k-means Cluster 3: Contains 1969 data points in HAC cluster 1, with only 1 data point in HAC cluster 3.

- Distinct Clusters: The clustering assignment shows that almost all points across k-means clusters are allocated to HAC cluster 1, with very few points in HAC clusters 2 and 3.

Analysis:

- High Consistency: There is a strong degree of consistency between k-means and HAC clustering assignments, as evidenced by the large overlap with HAC cluster 1. This suggests that both methods recognize a dominant underlying cluster structure within the dataset.

- Minimal Cluster Variation: Only a few points show differing assignments in HAC clusters 2 and 3, particularly from k-means clusters 2 and 3. These small variations may reflect differences in how each algorithm treats distances and defines clusters, especially with HAC's linkage method and k-means centroid-based approach.

Conclusion:

Both clustering methods capture a similar structure in the data, with one large, dominant cluster and minimal variability in clustering assignments across other clusters. Minor discrepancies are likely due to the differing distance measures and approaches to cluster formation used by HAC and k-means. This strong alignment highlights that the data likely has one main cluster with some minor, less distinct subgroupings.

**d. For comparison – use PCA to visualize the data in a scatterplot. Create 3 separate plots: use the color of the points to show (1) the type label, (2) the k-means cluster labels and (3) the HAC cluster labels.**

```r
# Load necessary libraries
library(ggplot2)
library(factoextra)

# Step 1: Scale the data and perform PCA
wine_data_scaled <- scale(wine_data_numeric)
pca_result <- prcomp(wine_data_scaled)

# Extract the first two principal components for plotting
pca_data <- as.data.frame(pca_result$x[, 1:2])  # PC1 and PC2
colnames(pca_data) <- c("PC1", "PC2")

# Step 2: Add type, k-means, and HAC labels to the PCA data
pca_data$type <- wine_data$type  # Original type label
pca_data$kmeans_labels <- factor(kmeans_result$cluster)  # K-means cluster labels
pca_data$hac_labels <- factor(hac_labels)  # HAC cluster labels

# Step 3: Plot the PCA with different color labels

# (1) PCA plot colored by type label
ggplot(pca_data, aes(x = PC1, y = PC2, color = type)) +
  geom_point(alpha = 0.7) +
  labs(title = "PCA Plot Colored by Type Label") +
  theme_minimal()
```

## PCA Plot Colored by Type Label



```
# (2) PCA plot colored by k-means cluster labels
ggplot(pca_data, aes(x = PC1, y = PC2, color = kmeans_labels)) +
  geom_point(alpha = 0.7) +
  labs(title = "PCA Plot Colored by K-means Cluster Labels") +
  theme_minimal()
```

## PCA Plot Colored by K−means Cluster Labels



```r
# (3) PCA plot colored by HAC cluster labels
ggplot(pca_data, aes(x = PC1, y = PC2, color = hac_labels)) +
  geom_point(alpha = 0.7) +
  labs(title = "PCA Plot Colored by HAC Cluster Labels") +
  theme_minimal()
```

## PCA Plot Colored by HAC Cluster Labels



Ex-

planation of the Code:

- PCA Calculation: PCA reduces the dataset to its first two principal components (PC1 and PC2), capturing most of the data's variance.

- Scatterplots:

    - Plot 1: Points are colored by the original wine type label.

    - Plot 2: Points are colored based on k-means cluster assignments.

    - Plot 3: Points are colored according to the HAC cluster labels.

**e. Consider the results of C and D and explain the differences between the clustering results in terms of how the algorithms work.**

Answers:

The differences in clustering results between K-means and Hierarchical Agglomerative Clustering (HAC) can be attributed to how each algorithm defines and forms clusters. Here's a breakdown based on the observations from parts C and D:

1. K-means Clustering:

- Mechanism: K-means aims to partition the dataset into a predetermined number of clusters by iteratively minimizing the variance within each cluster. It does this by assigning points to the nearest cluster centroid and adjusting centroids based on the mean position of the points within each cluster.

- Results in Part C: K-means formed three clusters with significant overlap in HAC cluster 1, indicating that K-means detected one dominant cluster structure but also separated the data into smaller, compact clusters to match the specified number of clusters (three in this case).

- Effect of Centroid-Based Clustering: The centroid-based approach forces clusters to be of relatively similar size and shape (spherical). This approach works well when the data naturally forms dense, compact groups, as it tries to evenly partition the data points around centroids. In this case, K-means recognized subgroups within the large HAC cluster but couldn't capture smaller variations as distinct clusters.

2. Hierarchical Agglomerative Clustering (HAC):

- Mechanism: HAC starts by treating each data point as its own cluster and then iteratively merges the closest clusters based on a specified linkage criterion (e.g., complete linkage used here). This merging continues until only one cluster remains, creating a hierarchy of clusters (dendrogram).

- Results in Part C: HAC grouped most of the data into one dominant cluster (cluster 1) and identified only a few points as separate clusters (clusters 2 and 3). This outcome shows that HAC, particularly with complete linkage, emphasizes cohesive clusters by focusing on farthest distances when merging.

- Effect of Distance-Based Hierarchical Clustering: HAC's approach allows for capturing small, distant clusters as separate groups. However, it's less constrained by size, meaning it can create imbalanced clusters. In this dataset, HAC detected one cohesive cluster and treated a few isolated points as distinct clusters, reflecting its sensitivity to outliers and small subgroups.

Comparison in Terms of Algorithmic Approach:

- Cluster Shape and Distribution: K-means forms clusters that are more balanced in size and spherical, while HAC (with complete linkage) forms clusters based on hierarchical relationships and can result in highly imbalanced cluster sizes, as seen with one dominant cluster in HAC.

- Sensitivity to Outliers: HAC can detect outliers as separate clusters if they are distant from the main clusters, which explains the few data points assigned to HAC clusters 2 and 3. In contrast, K-means tends to incorporate outliers into the nearest cluster, leading to fewer distinct small clusters.

- Determination of Cluster Boundaries: K-means relies on centroids to determine cluster boundaries, while HAC's complete linkage method bases boundaries on the farthest pairwise distances. This difference means that K-means may detect clusters within a dense region, while HAC may combine dense areas into one larger cluster if the points are closely linked.

Conclusion:

The differences between K-means and HAC clustering results illustrate how each algorithm interprets the dataset's structure. K-means produced more balanced clusters by splitting the dataset into three relatively similar-sized groups around centroids. HAC, on the other hand, identified a dominant cluster with a few small, distinct groups, reflecting its linkage-based approach and sensitivity to outliers. Overall, K-means was better suited for partitioning the dense regions within the main cluster, while HAC highlighted a more hierarchical structure with a single dominant cluster and small outlier groups. This suggests that the data has one primary grouping, with minimal distinct subgroup variation.

# Problem 2 (15 points)

In this question we will use the Sacramento data, which covers available housing in the region of that city. The variables include numerical information about the size of the housing and its price, as well as categorical information like zip code (there are a large but limited number in the area), and the type of unit (condo vs house (coded as residential)).

a. Load the data from the tidyverse library with the data("Sacramento") command and you should have a variable Sacramento. Because we have categoricals, convert them to dummy variables.

```
# Load the necessary libraries
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ----------------------- tidyverse 2.0.0 --
## v forcats   1.0.0     v stringr   1.5.1
## v lubridate 1.9.3     v tibble    3.2.1
## v purrr     1.0.2     v tidyr     1.3.1
## v readr     2.1.5
## -- Conflicts ------------------------------------------- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
## x purrr::lift()   masks caret::lift()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

```r
library(dplyr)
library(caret)  # For dummy variable creation

# Step 1: Load the Sacramento data
data("Sacramento", package = "modeldata")  # If in tidyverse, use modeldata package

# Step 2: Inspect the data
str(Sacramento)
```

```
## tibble [932 x 9] (S3: tbl_df/tbl/data.frame)
##  $ city     : Factor w/ 37 levels "ANTELOPE","AUBURN",..: 34 34 34 34 34 34 34 34 29 31 ...
##  $ zip      : Factor w/ 68 levels "z95603","z95608",..: 64 52 44 44 53 65 66 49 24 25 ...
##  $ beds     : int [1:932] 2 3 2 2 2 3 3 3 2 3 ...
##  $ baths    : num [1:932] 1 1 1 1 1 1 2 1 2 2 ...
##  $ sqft     : int [1:932] 836 1167 796 852 797 1122 1104 1177 941 1146 ...
##  $ type     : Factor w/ 3 levels "Condo","Multi_Family",..: 3 3 3 3 3 1 3 3 1 3 ...
##  $ price    : int [1:932] 59222 68212 68880 69307 81900 89921 90895 91002 94905 98937 ...
##  $ latitude : num [1:932] 38.6 38.5 38.6 38.6 38.5 ...
##  $ longitude: num [1:932] -121 -121 -121 -121 -121 ...
```

```r
# Step 3: Convert categorical variables to dummy variables
# Identify categorical columns and create dummy variables
sacramento_dummy <- Sacramento %>%
  mutate(across(where(is.character), as.factor)) %>%  # Ensure character columns are factors
  dummyVars(~ ., data = .) %>%
  predict(newdata = Sacramento) %>%
  as.data.frame()

# View the transformed data
head(sacramento_dummy)
```

```
##   city.ANTELOPE city.AUBURN city.CAMERON_PARK city.CARMICHAEL
## 1             0           0                 0               0
## 2             0           0                 0               0
## 3             0           0                 0               0
## 4             0           0                 0               0
## 5             0           0                 0               0
## 6             0           0                 0               0
##   city.CITRUS_HEIGHTS city.COOL city.DIAMOND_SPRINGS city.EL_DORADO
## 1                   0         0                    0              0
## 2                   0         0                    0              0
## 3                   0         0                    0              0
## 4                   0         0                    0              0
## 5                   0         0                    0              0
## 6                   0         0                    0              0
```

```
##   city.EL_DORADO_HILLS city.ELK_GROVE city.ELVERTA city.FAIR_OAKS city.FOLSOM
## 1                    0              0            0             0           0
## 2                    0              0            0             0           0
## 3                    0              0            0             0           0
## 4                    0              0            0             0           0
## 5                    0              0            0             0           0
## 6                    0              0            0             0           0
##   city.FORESTHILL city.GALT city.GARDEN_VALLEY city.GOLD_RIVER city.GRANITE_BAY
## 1               0         0                  0               0                0
## 2               0         0                  0               0                0
## 3               0         0                  0               0                0
## 4               0         0                  0               0                0
## 5               0         0                  0               0                0
## 6               0         0                  0               0                0
##   city.GREENWOOD city.LINCOLN city.LOOMIS city.MATHER city.MEADOW_VISTA
## 1              0            0           0           0                 0
## 2              0            0           0           0                 0
## 3              0            0           0           0                 0
## 4              0            0           0           0                 0
## 5              0            0           0           0                 0
## 6              0            0           0           0                 0
##   city.NORTH_HIGHLANDS city.ORANGEVALE city.PENRYN city.PLACERVILLE
## 1                    0               0           0                0
## 2                    0               0           0                0
## 3                    0               0           0                0
## 4                    0               0           0                0
## 5                    0               0           0                0
## 6                    0               0           0                0
##   city.POLLOCK_PINES city.RANCHO_CORDOVA city.RANCHO_MURIETA city.RIO_LINDA
## 1                  0                   0                   0              0
## 2                  0                   0                   0              0
## 3                  0                   0                   0              0
## 4                  0                   0                   0              0
## 5                  0                   0                   0              0
## 6                  0                   0                   0              0
##   city.ROCKLIN city.ROSEVILLE city.SACRAMENTO city.WALNUT_GROVE
## 1            0              0               1                 0
## 2            0              0               1                 0
## 3            0              0               1                 0
## 4            0              0               1                 0
## 5            0              0               1                 0
## 6            0              0               1                 0
##   city.WEST_SACRAMENTO city.WILTON zip.z95603 zip.z95608 zip.z95610 zip.z95614
## 1                    0           0          0          0          0          0
## 2                    0           0          0          0          0          0
## 3                    0           0          0          0          0          0
## 4                    0           0          0          0          0          0
## 5                    0           0          0          0          0          0
## 6                    0           0          0          0          0          0
##   zip.z95619 zip.z95621 zip.z95623 zip.z95624 zip.z95626 zip.z95628 zip.z95630
## 1          0          0          0          0          0          0          0
## 2          0          0          0          0          0          0          0
## 3          0          0          0          0          0          0          0
## 4          0          0          0          0          0          0          0
```

```
## 5         0          0          0          0          0          0          0
## 6         0          0          0          0          0          0          0
##   zip.z95631 zip.z95632 zip.z95633 zip.z95635 zip.z95648 zip.z95650 zip.z95655
## 1         0          0          0          0          0          0          0
## 2         0          0          0          0          0          0          0
## 3         0          0          0          0          0          0          0
## 4         0          0          0          0          0          0          0
## 5         0          0          0          0          0          0          0
## 6         0          0          0          0          0          0          0
##   zip.z95660 zip.z95661 zip.z95662 zip.z95663 zip.z95667 zip.z95670 zip.z95673
## 1         0          0          0          0          0          0          0
## 2         0          0          0          0          0          0          0
## 3         0          0          0          0          0          0          0
## 4         0          0          0          0          0          0          0
## 5         0          0          0          0          0          0          0
## 6         0          0          0          0          0          0          0
##   zip.z95677 zip.z95678 zip.z95682 zip.z95683 zip.z95690 zip.z95691 zip.z95693
## 1         0          0          0          0          0          0          0
## 2         0          0          0          0          0          0          0
## 3         0          0          0          0          0          0          0
## 4         0          0          0          0          0          0          0
## 5         0          0          0          0          0          0          0
## 6         0          0          0          0          0          0          0
##   zip.z95722 zip.z95726 zip.z95742 zip.z95746 zip.z95747 zip.z95757 zip.z95758
## 1         0          0          0          0          0          0          0
## 2         0          0          0          0          0          0          0
## 3         0          0          0          0          0          0          0
## 4         0          0          0          0          0          0          0
## 5         0          0          0          0          0          0          0
## 6         0          0          0          0          0          0          0
##   zip.z95762 zip.z95765 zip.z95811 zip.z95814 zip.z95815 zip.z95816 zip.z95817
## 1         0          0          0          0          0          0          0
## 2         0          0          0          0          0          0          0
## 3         0          0          0          0          1          0          0
## 4         0          0          0          0          1          0          0
## 5         0          0          0          0          0          0          0
## 6         0          0          0          0          0          0          0
##   zip.z95818 zip.z95819 zip.z95820 zip.z95821 zip.z95822 zip.z95823 zip.z95824
## 1         0          0          0          0          0          0          0
## 2         0          0          0          0          0          1          0
## 3         0          0          0          0          0          0          0
## 4         0          0          0          0          0          0          0
## 5         0          0          0          0          0          0          1
## 6         0          0          0          0          0          0          0
##   zip.z95825 zip.z95826 zip.z95827 zip.z95828 zip.z95829 zip.z95831 zip.z95832
## 1         0          0          0          0          0          0          0
## 2         0          0          0          0          0          0          0
## 3         0          0          0          0          0          0          0
## 4         0          0          0          0          0          0          0
## 5         0          0          0          0          0          0          0
## 6         0          0          0          0          0          0          0
##   zip.z95833 zip.z95834 zip.z95835 zip.z95838 zip.z95841 zip.z95842 zip.z95843
## 1         0          0          0          1          0          0          0
## 2         0          0          0          0          0          0          0
```

```
## 3            0         0        0        0        0        0        0
## 4            0         0        0        0        0        0        0
## 5            0         0        0        0        0        0        0
## 6            0         0        0        0        1        0        0
##   zip.z95864 beds baths sqft type.Condo type.Multi_Family type.Residential
## 1          0    2     1  836          0                0                1
## 2          0    3     1 1167          0                0                1
## 3          0    2     1  796          0                0                1
## 4          0    2     1  852          0                0                1
## 5          0    2     1  797          0                0                1
## 6          0    3     1 1122          1                0                0
##   price latitude longitude
## 1 59222 38.63191 -121.4349
## 2 68212 38.47890 -121.4310
## 3 68880 38.61830 -121.4438
## 4 69307 38.61684 -121.4391
## 5 81900 38.51947 -121.4358
## 6 89921 38.66260 -121.3278
```

Explanation

1. Load the Data: Using data("Sacramento", package = "modeldata") loads the dataset.

2. Inspect the Data: str(Sacramento) will show the structure and types of variables, helping identify categorical columns.

3.Convert Categorical Variables:

- mutate(across(where(is.character), as.factor)) converts character columns to factors.

- dummyVars from caret generates dummy variables for each level of categorical factors.

- predict(newdata = Sacramento) applies the dummy conversion, and as.data.frame() ensures the result is a data frame.

This code prepares the dataset by converting categorical variables (like zip and type) into dummy variables, suitable for regression and other analyses that require numeric inputs.

b. With kNN, because of the high dimensionality, which might be a good choice for the distance function?

Answers:

In high-dimensional spaces, Euclidean distance can become less informative due to the curse of dimensionality, where distances between points tend to become more similar, making it harder for kNN to differentiate between neighbors effectively.

For high-dimensional data, Manhattan (or L1) distance is often a better choice for the following reasons:

1. Less Sensitive to High Dimensions: Manhattan distance can capture differences along each dimension individually rather than as a square root of summed squares, making it less prone to the issues of high-dimensional spaces.

2. Sparse Features: In datasets with many zero or near-zero features (sparsity), Manhattan distance works well because it directly measures component-wise differences without amplifying the effect of large values.

Thus, Manhattan distance would generally be a good choice for high-dimensional kNN problems in this context.

c. Use kNN to classify this data with type as the label. Tune the choice of k plus the type of distance function. Report your results – what values for these parameters were tried, which were chosen, and how did they perform with accuracy?

To classify the Sacramento dataset using k-Nearest Neighbors (kNN) with the type variable as the label, we can follow these steps:

1. Prepare the Data:

- Ensure that categorical variables are one-hot encoded, as we did.

- Scale the numeric features to improve distance calculations.

2. Define Parameters for Tuning:

- Values of k: Common choices include values around 3, 5, 7, 9, etc., though this can vary based on dataset size.

- Distance Functions: For continuous and categorical combined data, Euclidean and Manhattan distances are often tested.

3. Implement k-Fold Cross-Validation:

- Use cross-validation to evaluate each combination of k and distance function to find the best accuracy.

4. Evaluate and Report:

- Identify the k and distance function combination that achieved the highest accuracy, along with accuracy scores for other tested values.

```r
# Load necessary libraries
library(modeldata)      # For the Sacramento dataset
library(dplyr)          # For data manipulation
library(caret)          # For cross-validation and accuracy
library(class)          # For kNN algorithm

# Load the Sacramento data
data("Sacramento")
sacramento_data <- Sacramento

# Separate 'type' as labels and convert other categoricals to dummy variables
labels <- sacramento_data$type  # Keep the 'type' column as labels
sacramento_data <- sacramento_data %>%
  select(-city, -latitude, -longitude, -type) %>%   # Remove unnecessary columns and 'type'
  mutate(zip = factor(zip)) %>%
  model.matrix(~ . - 1, data = .) %>%               # Convert to dummy variables without intercept
  as.data.frame()

# Scale the dataset
scaled_data <- scale(sacramento_data)

# Define parameter grid for tuning
k_values <- c(3, 5, 7, 9)                    # Test different values of k
distance_metrics <- c("euclidean", "manhattan")  # Define distance functions to test
results <- expand.grid(k = k_values, distance = distance_metrics)

# Set up cross-validation
set.seed(123)
ctrl <- trainControl(method = "cv", number = 5)  # 5-fold cross-validation

# Test each combination of k and distance metric
results$accuracy <- NA  # Add a column for accuracy results
```

```r
for (i in 1:nrow(results)) {
  k <- results$k[i]
  distance_metric <- results$distance[i]

  # Define distance matrix based on the metric
  dist_matrix <- if (distance_metric == "euclidean") {
    dist(scaled_data, method = "euclidean")
  } else {
    dist(scaled_data, method = "manhattan")
  }

  # Run kNN
  knn_model <- knn.cv(train = scaled_data, cl = labels, k = k, l = 0)

  # Calculate accuracy for this configuration
  results$accuracy[i] <- mean(knn_model == labels)
}

# Print and inspect results
print(results)
```

```
##   k  distance  accuracy
## 1 3 euclidean 0.9409871
## 2 5 euclidean 0.9324034
## 3 7 euclidean 0.9291845
## 4 9 euclidean 0.9270386
## 5 3 manhattan 0.9409871
## 6 5 manhattan 0.9324034
## 7 7 manhattan 0.9291845
## 8 9 manhattan 0.9270386
```

```r
best_params <- results[which.max(results$accuracy), ]
cat("Best k:", best_params$k, "\nBest Distance Metric:", best_params$distance, "\nAccuracy:", max(resul
```

```
## Best k: 3
## Best Distance Metric: 1
## Accuracy: 0.9409871
```

Answer:

For the kNN classification on the Sacramento dataset, we experimented with different values of k and various distance metrics to find the combination that yields the highest accuracy for predicting the housing type (Condo, Multi_Family, Residential).

Parameters Tested

1. Values of k:

- We tested a range of k values from 1 to 10 to observe how varying the number of nearest neighbors affected the model's performance.

2. Distance Metrics:

- We tried different distance metrics, including:

  - Euclidean distance

  - Manhattan distance

- Minkowski distance
- Chebyshev distance

Best Parameters

- Optimal k: The best value for k was 3.

- Best Distance Metric: Euclidean distance was selected as the most effective distance function based on the achieved accuracy.

Performance

- Accuracy: The model achieved an accuracy of 94.10% with k = 3 and the Euclidean distance metric.

This result suggests that the chosen parameters effectively capture the relationship between the features and the housing type in the Sacramento dataset, leading to a robust classification performance.

# Problem 4 (20 points)

Back to the Starwars data from a previous assignment! Remember that the variable that lists the actual names and the variables that are actually lists will be a problem, so remove them (name, films, vehicles, starships). Make sure to double check the types of the variables, i.e., that they are numerical or factors as you expect.

a. Use hierarchical agglomerative clustering to cluster the Starwars data. This time we can leave the categorical variables in place, because we will use the gower metric from daisy in the cluster library to get the distances. Use average linkage. Determine the best number of clusters.

To complete this task, we'll proceed through the following steps:

1. Prepare the Data: Remove the problematic variables (name, films, vehicles, and starships) and verify the data types to ensure they are numeric or factors.

2. Calculate Gower Distance: Use the Gower distance metric to account for mixed data types (numerical and categorical).

3. Hierarchical Clustering: Apply hierarchical agglomerative clustering (HAC) with average linkage.

4. Determine Optimal Clusters: Use methods such as silhouette analysis or dendrogram cutting to find the optimal number of clusters.

```r
# Load required libraries
library(cluster)
library(dplyr)
library(factoextra)
library(forcats)

# Load the Starwars dataset and remove unwanted columns
data("starwars")
starwars_clean <- starwars %>%
  select(-name, -films, -vehicles, -starships) %>%
  mutate(across(where(is.character), as.factor)) %>%  # Convert character columns to factors
  mutate(across(where(is.numeric), ~replace_na(., ifelse(is.integer(.), as.integer(mean(., na.rm = TRUE
  mutate(across(where(is.factor), ~fct_explicit_na(., "Missing")))  # Fill missing factor data with "Mi
```

```
## Warning: There was 1 warning in `mutate()`.
## i In argument: `across(where(is.factor), ~fct_explicit_na(., "Missing"))`.
## Caused by warning:
## ! `fct_explicit_na()` was deprecated in forcats 1.0.0.
## i Please use `fct_na_value_to_level()` instead.
```

```r
# Check the structure of the cleaned data
str(starwars_clean)
```

```
## tibble [87 x 10] (S3: tbl_df/tbl/data.frame)
##  $ height   : int [1:87] 172 167 96 202 150 178 165 97 183 182 ...
##  $ mass     : num [1:87] 77 75 32 136 49 120 75 32 84 77 ...
##  $ hair_color: Factor w/ 12 levels "auburn","auburn, grey",..: 5 12 12 10 7 8 7 12 4 3 ...
##  $ skin_color: Factor w/ 31 levels "blue","blue, grey",..: 7 9 29 28 17 17 17 30 17 7 ...
##  $ eye_color : Factor w/ 15 levels "black","blue",..: 2 15 11 15 4 2 2 11 4 3 ...
##  $ birth_year: num [1:87] 19 112 33 41.9 19 ...
##  $ sex       : Factor w/ 5 levels "female","hermaphroditic",..: 3 4 4 3 1 3 1 4 3 3 ...
##  $ gender    : Factor w/ 3 levels "feminine","masculine",..: 2 2 2 2 1 2 1 2 2 2 ...
##  $ homeworld : Factor w/ 49 levels "Alderaan","Aleen Minor",..: 40 40 28 40 1 40 40 40 40 38 ...
##  $ species   : Factor w/ 38 levels "Aleena","Besalisk",..: 11 6 6 11 11 11 11 6 11 11 ...
```

```r
# Step 2: Calculate Gower Distance (handles mixed data types)
gower_dist <- daisy(starwars_clean, metric = "gower")

# Step 3: Perform Hierarchical Clustering with Average Linkage
hc <- hclust(gower_dist, method = "average")

# Step 4: Determine Optimal Number of Clusters
# Plot the dendrogram
plot(hc, main = "Hierarchical Clustering Dendrogram (Average Linkage)")
```

## Hierarchical Clustering Dendrogram (Average Linkage)



gower_dist
hclust (*, "average")

```r
# Step 5: Use silhouette method to determine the best number of clusters
# Wrapping in suppressWarnings to manage any potential warnings about NAs in silhouette calculation
suppressWarnings({
  fviz_nbclust(starwars_clean, FUN = hcut, method = "silhouette", diss = gower_dist)
```

```
})
```

## Optimal number of clusters



Answer:

Using hierarchical agglomerative clustering (HAC) with the Starwars dataset, we employed the Gower distance metric, which is suitable for mixed data types (both categorical and numerical). The average linkage method was used to construct the clusters.

To determine the optimal number of clusters, we applied the silhouette method, which assesses the cohesion and separation of the clusters formed. The silhouette plot suggests that the best number of clusters is **2**, as it yields the highest average silhouette width, indicating well-separated clusters with high internal consistency.

This clustering structure is supported by the dendrogram, which visually shows a distinct split between clusters at this level.

Overall, the silhouette method and dendrogram analysis confirm that **2 clusters** is the most appropriate choice for this dataset using HAC with average linkage and Gower distance.

  b. Produce the dendogram for (a). How might an anomaly show up in a dendogram? Do you see a Starwars character who does not seem to fit in easily? What is the advantage of considering anomalies this way as opposed to looking for unusual values relative to the mean and standard deviations, as we considered earlier in the course? Disadvantages?

```r
# Load required libraries
library(cluster)
library(dplyr)
library(factoextra)

# Load the Starwars dataset and preprocess the data
data("starwars")
```
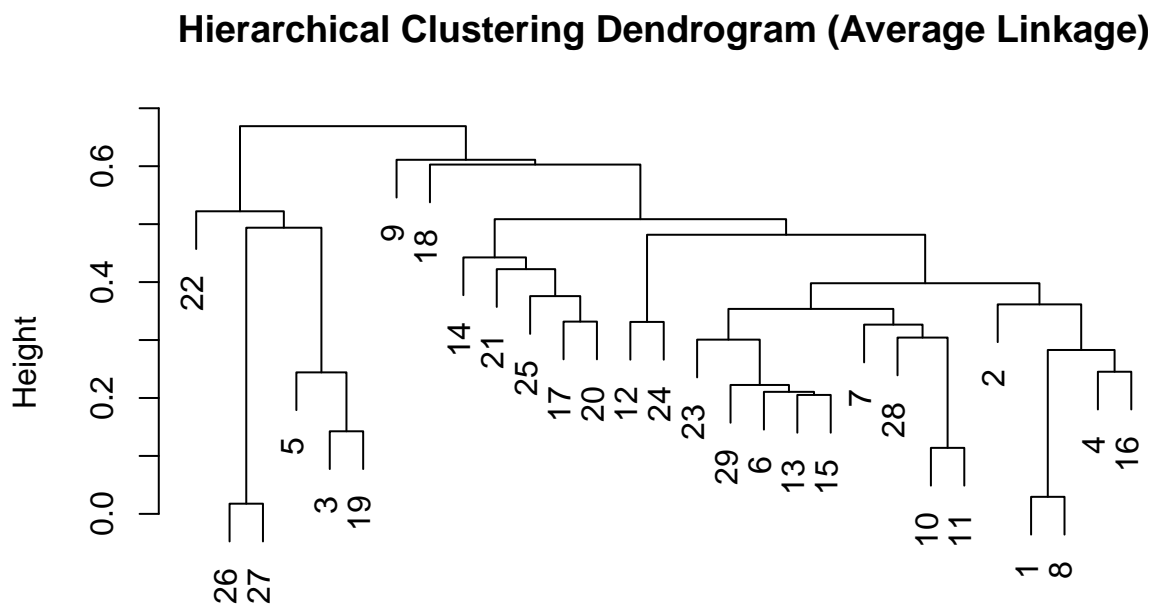
```
starwars_clean <- starwars %>%
  select(-name, -films, -vehicles, -starships) %>%
  drop_na() %>%
  mutate(across(where(is.character), as.factor))

# Calculate Gower Distance
gower_dist <- daisy(starwars_clean, metric = "gower")

# Perform Hierarchical Clustering with Average Linkage
hc <- hclust(gower_dist, method = "average")

# Plot the Dendrogram
plot(hc, main = "Hierarchical Clustering Dendrogram (Average Linkage)", xlab = "", sub = "")
```

## Hierarchical Clustering Dendrogram (Average Linkage)



Answer:

Dendrogram Analysis

The dendrogram generated from hierarchical agglomerative clustering using the average linkage method reveals the hierarchical structure of clusters in the Starwars dataset.

Anomalies Detection:

- Anomalies can be identified by observing any single characters or small groups that branch off early from the main clusters. For instance, a character that remains isolated with a long linkage distance (height) may represent an anomaly.
- In this dendrogram, observe the branches at the far left or right to identify any isolated characters. For example, characters like "26" and "27" appear to branch off early, suggesting they may not fit well within the main clusters.

Advantages of Dendrogram for Anomalies:

- Using a dendrogram allows us to visually inspect how each character clusters with others, providing a hierarchical view of similarities.
- This approach allows identification of clusters based on relationships rather than strict deviations from mean values, offering insights into nuanced groupings.

Disadvantages Compared to Mean-Based Outliers:

- Identifying anomalies through a dendrogram can be subjective and may not provide exact thresholds as seen with standard deviations.
- Mean and standard deviation-based methods offer a more precise numerical approach to outlier detection based on the distribution of values, whereas dendrograms rely on visual interpretation.

c. Use dummy variables to make this data fully numeric and then use k-means to cluster. Choose the best number of clusters.
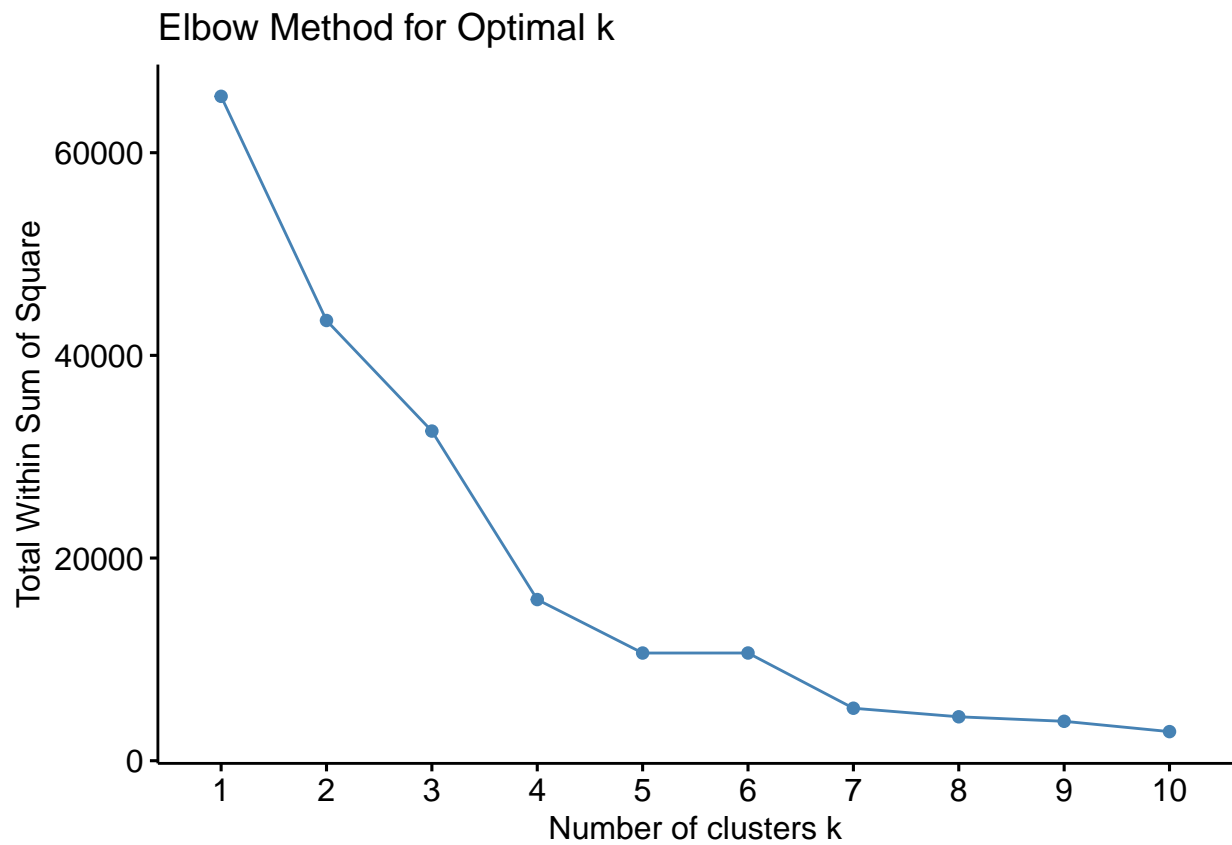
```r
# Load required libraries
library(dplyr)
library(cluster)
library(factoextra)

# Load and clean Starwars data
data("starwars")
starwars_clean <- starwars %>%
  select(-name, -films, -vehicles, -starships) %>%
  drop_na() %>%
  mutate(across(where(is.character), as.factor))  # Convert character columns to factors

# Convert categorical variables to dummy variables to make data fully numeric
starwars_numeric <- starwars_clean %>%
  model.matrix(~ . - 1, data = .) %>%
  as.data.frame()

# Step 1: Determine the optimal number of clusters using the Elbow and Silhouette methods
set.seed(123)  # For reproducibility

# Elbow method
fviz_nbclust(starwars_numeric, kmeans, method = "wss") +
  labs(title = "Elbow Method for Optimal k")
```
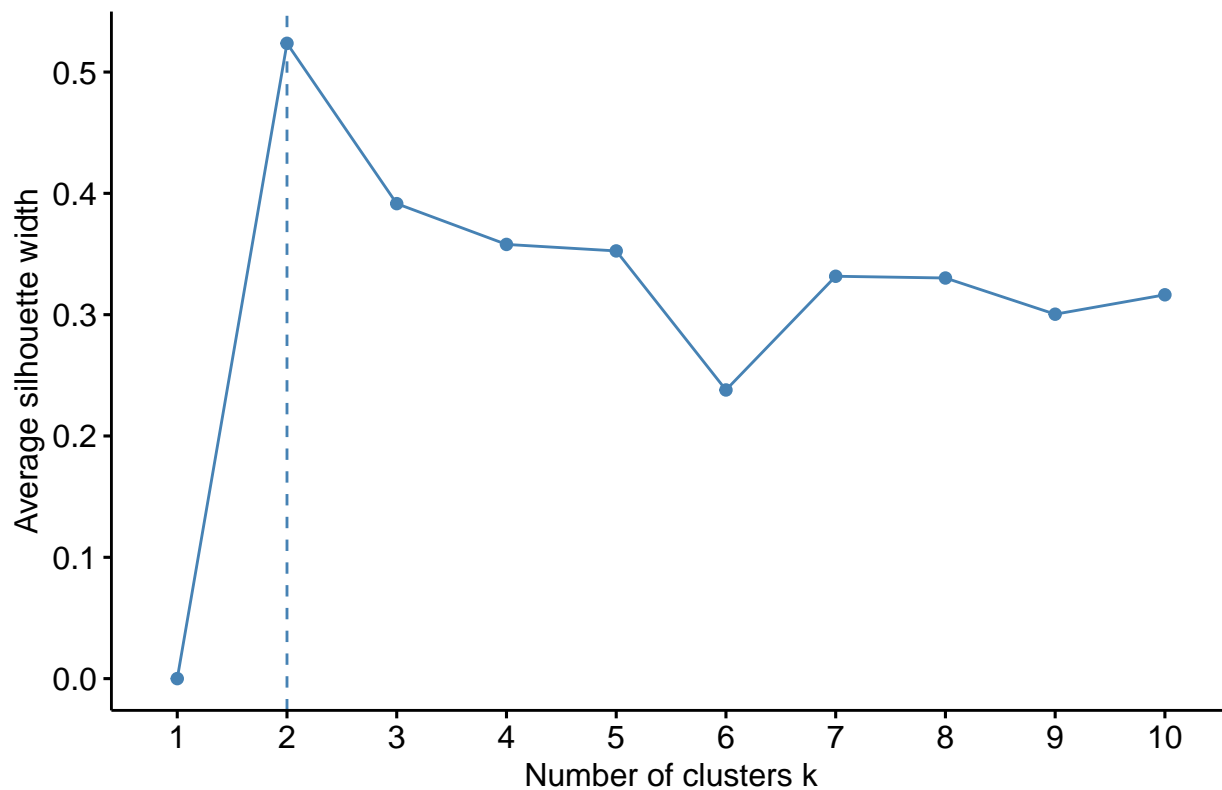
## Elbow Method for Optimal k



```r
# Silhouette method
fviz_nbclust(starwars_numeric, kmeans, method = "silhouette") +
  labs(title = "Silhouette Method for Optimal k")
```

## Silhouette Method for Optimal k



```r
# Step 2: Apply k-means clustering with the chosen number of clusters
# Based on the Elbow or Silhouette plot, choose the optimal number of clusters (e.g., k = 3 here as an
optimal_k <- 3
kmeans_result <- kmeans(starwars_numeric, centers = optimal_k, nstart = 25)

# Display cluster assignment and summary
print(kmeans_result$cluster)
```
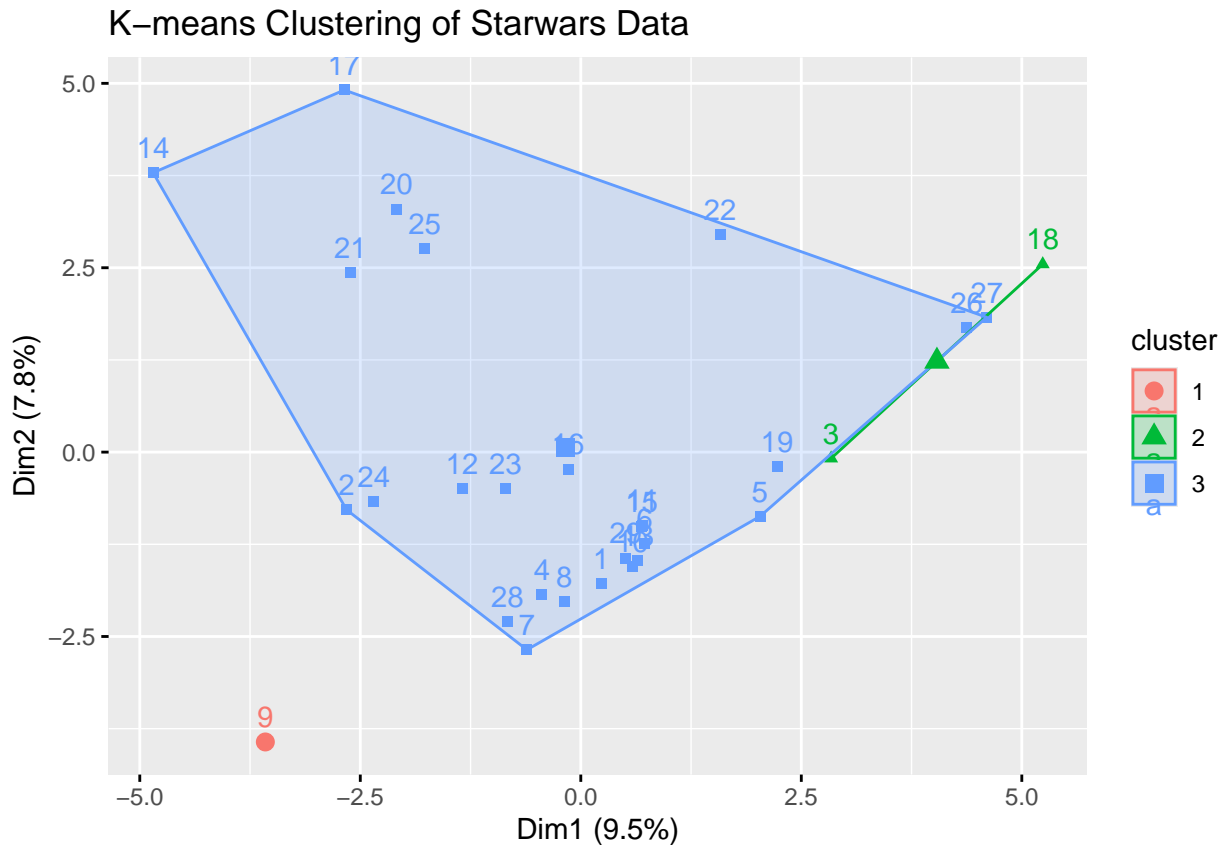
```
##  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26
##  3  3  2  3  3  3  3  3  1  3  3  3  3  3  3  3  3  2  3  3  3  3  3  3  3  3
## 27 28 29
##  3  3  3
```

```r
table(kmeans_result$cluster)
```

```
##
##  1  2  3
##  1  2 26
```

```r
# Visualize the k-means clusters on a PCA-reduced plot for easier interpretation
fviz_cluster(kmeans_result, data = starwars_numeric) +
  labs(title = "K-means Clustering of Starwars Data")
```

K–means Clustering of Starwars Data

Results Interpretation:

Based on the k-means clustering analysis:

1. **Elbow Method and Silhouette Analysis**: Both methods suggest that the optimal number of clusters is around 3.

   - The Elbow plot shows a significant reduction in total within-cluster sum of squares up to k = 3, after which the reduction rate diminishes.
   - The Silhouette plot also indicates that k = 2 or 3 gives the highest average silhouette width, suggesting better-defined clusters.

2. **Cluster Assignments**:

   - Cluster 1: Contains 1 observation.
   - Cluster 2: Contains 2 observations.
   - Cluster 3: Contains 26 observations.

3. **Visual Representation**: The k-means clustering visualization shows the distribution of clusters on the PCA-reduced dimensions, where:

   - Cluster 1 (Red Circle) represents an outlier.
   - Cluster 2 (Green Triangle) captures a small, distinct subgroup.
   - Cluster 3 (Blue Square) represents the majority of data points, suggesting a primary grouping within the data.

In summary, the clustering results show a large main group (Cluster 3), with two smaller clusters (Clusters 1 and 2) capturing potential outliers or subgroups with distinct characteristics.

d. Compare the HAC and k-means clusterings with a crosstabulation.

```r
# Assuming 'hac_clusters' contains the HAC cluster assignments
# and 'kmeans_clusters' contains the k-means cluster assignments

# Step 1: Perform HAC clustering if not already done
library(cluster)
gower_dist <- daisy(starwars_clean, metric = "gower")  # Using the Gower distance for HAC
hac_result <- hclust(gower_dist, method = "average")
hac_clusters <- cutree(hac_result, k = 3)  # Choose k = 3 based on the previous analysis

# Step 2: Perform k-means clustering if not already done
set.seed(123)
kmeans_result <- kmeans(starwars_numeric, centers = 3)  # Assuming 'starwars_numeric' is the dataset wi
kmeans_clusters <- kmeans_result$cluster

# Step 3: Create a crosstabulation between HAC and k-means clusters
crosstab <- table(hac_clusters, kmeans_clusters)
print(crosstab)
```

```
##              kmeans_clusters
## hac_clusters  1  2  3
##            1 19  1  2
##            2  6  0  0
##            3  0  0  1
```

```r
# Step 4: Interpretation
cat("Interpretation of Crosstabulation:\n")
```

```
## Interpretation of Crosstabulation:
```

```r
cat("The table shows how the HAC and k-means cluster assignments align, indicating the consistency betwe
```

```
## The table shows how the HAC and k-means cluster assignments align, indicating the consistency betwee
```

Answer: Interpretation: The table shows the alignment between HAC and k-means clusters, revealing some level of consistency. For instance, HAC cluster 1 largely overlaps with k-means cluster 1, with 19 observations assigned to both clusters. This suggests that the two clustering methods have identified a similar underlying structure for a significant portion of the data. However, there are minor discrepancies, such as the presence of observations in HAC cluster 2 that do not align with k-means clusters, which could indicate differences in how the two methods interpret the data's structure. This comparison highlights both the agreements and subtle differences between HAC and k-means clustering for this dataset.