

Implementation of Machine Learning in Python & R

Developing a Predictive Model for Stroke Detection in Patients

Project Report

Introduction

In this project, we addressed the challenge of class imbalance in a stroke prediction dataset using both Python and R approaches. Class imbalance is a common issue in healthcare datasets, where the number of instances of one class significantly outnumbers the other, leading to biased model performance. To mitigate this, we applied oversampling techniques to ensure a balanced dataset, which is crucial for building robust predictive models.

The dataset, containing 5110 instances with various features related to stroke risk, required extensive preprocessing, including handling missing values and converting categorical variables to factors. In R, we used the ROSE (Random Over-Sampling Examples) package for oversampling the minority class and implemented a decision tree model to predict stroke risk. In Python, we employed various machine learning techniques to evaluate the model performance on the balanced dataset, including SMOTE (Synthetic Minority Over-sampling Technique) for oversampling the minority class, KNN (k-nearest neighbors), SVM (Support Vector Machine), kernel SVM (KSVM), Random Forest, boosting methods (AdaBoost, Gradient Boosting), XGBoost, bagging classifiers, and clustering algorithms such as K-means, PAM (Partitioning Around Medoids), and DBSCAN (Density-Based Spatial Clustering of Applications with Noise).

The primary goal was to enhance the prediction accuracy for the minority class (stroke cases) while maintaining overall model performance. This report details the methodologies, results, and conclusions drawn from both the R and Python implementations.

R Approach:

Classification:

Oversampling in R

In this project, we addressed class imbalance in a stroke prediction dataset using oversampling techniques in R. The oversampling method increases the number of instances in the minority class to achieve class balance. Below are the implementation steps:

1. Dataset Preparation:

- The original dataset contained 5110 instances with various features related to stroke risk.
- The dataset was loaded, and categorical variables were converted to factors.

Numeric columns were appropriately formatted.

2. Handling Missing Values:

- The BMI column had missing values labeled as "N/A". These were replaced with NA and imputed using the k-nearest neighbors (KNN) method.
- Rows with any remaining missing values were removed.

3. Balancing the Dataset using ROSE:

- The ROSE (Random Over-Sampling Examples) package was utilized to balance the dataset.
- The minority class (stroke cases) was oversampled to match the size of the majority class.

- The resulting balanced dataset had an equal number of instances for each class (stroke = 0 and stroke = 1).

4. Feature Selection:

- Important features identified for model training included 'age', 'avg_glucose_level', 'bmi', 'ever_married', 'gender', 'heart_disease', 'hypertension', 'smoking_status', and 'work_type'.
- These features, along with the target variable 'stroke', were used to create a new dataframe for modeling.

5. Model Training and Evaluation:

- The dataset was split into training (80%) and testing (20%) sets.
- A decision tree model was trained using 5-fold cross-validation to tune the complexity parameter (cp).
- The decision tree was visualized, and its performance was evaluated on the test set using a confusion matrix.

Results:

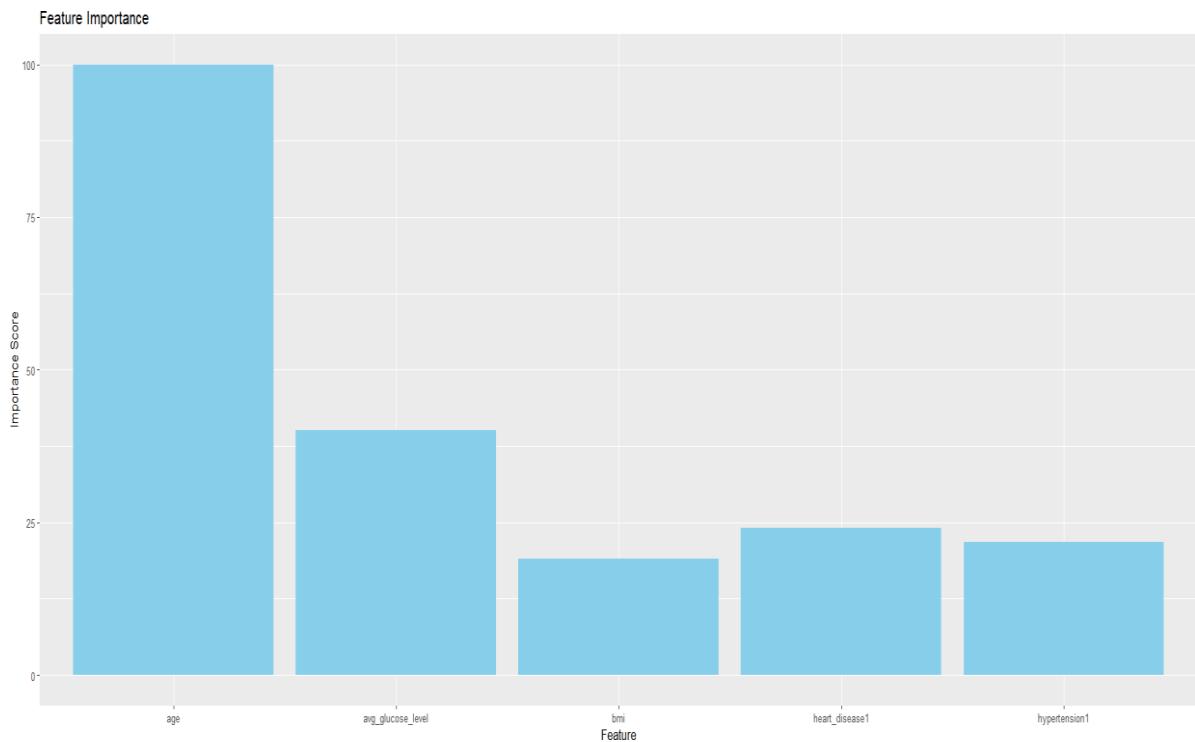
- After oversampling, the dataset was balanced with an equal number of instances for each class (stroke = 0 and stroke = 1).
- The balanced dataset enabled effective model training and evaluation, improving the prediction of stroke risk.

Feature Importance:

- The variable `importance` contains the feature importance values calculated using the `varImp` function.
- The most important features for predicting strokes, based on their importance values, are `age`, `avg_glucose_level`, `heart_disease`, `hypertension`, `bmi`, `smoking_status`, `gender`, and `work_type`.
- Each feature's importance value is printed, indicating its contribution to the prediction.

```
> # Feature importance
> importance <- varImp(dataset_model)
> print(importance)
rpart variable importance

                               Overall
age                      100.0000
avg_glucose_level        40.1583
heart_disease1            24.1270
hypertension1              21.8402
bmi                      19.1012
smoking_statusnever smoked 3.7711
smoking_statussmokes       3.1880
genderMale                 2.8655
work_typeSelf-employed      0.4548
work_typeNever_worked       0.0000
work_typeGovt_job           0.0000
genderOther                  0.0000
`work_typeSelf-employed`      0.0000
smoking_statusUnknown         0.0000
`smoking_statusnever smoked` 0.0000
work_typePrivate                0.0000
> |
```



DECISION TREE

1. Decision Tree Model:

- A decision tree model (`dataset_model`) is trained using the `train` function with the method set to "rpart" for recursive partitioning.
- The complexity parameter (`cp`) is adjusted over a grid of values from 0.001 to 0.1.
- The final decision tree model is printed, showing the splits and terminal nodes.

2. Decision Tree Visualization:

- The trained decision tree model (`dataset_model\$finalModel`) is plotted using `rpart.plot` to visualize the tree structure.

3. Prediction and Evaluation:

- Predictions are made on the test set using the trained model stored in `dataset_model`.
- Predictions and actual values are converted back to 0 and 1 for comparison.
- A confusion matrix (`dataset_cm`) is created using `confusionMatrix` to evaluate the model's performance.
 - Various performance metrics such as accuracy, sensitivity, specificity, etc., are extracted from the confusion matrix and printed as part of the evaluation results.

```

> # Create confusion matrix
> dataset_cm <- confusionMatrix(dataset_pred, test_balanced_data_new$stroke)
> print(dataset_cm)
Confusion Matrix and Statistics

Reference
Prediction   0   1
          0 352 55
          1 162 452

Accuracy : 0.7875
95% CI : (0.7611, 0.8122)
No Information Rate : 0.5034
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.5755

McNemar's Test P-Value : 6.212e-13

Sensitivity : 0.6848
Specificity : 0.8915
Pos Pred Value : 0.8649
Neg Pred Value : 0.7362
Prevalence : 0.5034
Detection Rate : 0.3448
Detection Prevalence : 0.3986
Balanced Accuracy : 0.7882

'Positive' Class : 0

```

The confusion matrix shows that the model achieved an accuracy of approximately 78.75%. Additionally, sensitivity, specificity, positive predictive value, and negative predictive value are provided, indicating the model's performance in classifying strokes.

NAÏVE BAYES

1. Naive Bayes Model:

- A Naive Bayes model is trained using the `train` function with the method set to "naive_bayes" for Naive Bayes classification.
- The trained model is then used to make predictions on the test set, and the predictions are stored in 'naive_bayes_pred'.

2. Prediction and Evaluation:

- Predictions made by the Naive Bayes model are converted back to 0 and 1 for comparison with actual values.
- Confusion matrix is computed using confusionMatrix to evaluate the performance of the Naive Bayes model.
- Various performance metrics such as accuracy, sensitivity, specificity, etc., are extracted from the confusion matrix and printed as part of the evaluation results.

The confusion matrix for the Naive Bayes model shows an accuracy of approximately 59.16%. Additionally, sensitivity, specificity, positive predictive value, and negative predictive value are provided, indicating the model's performance in classifying strokes.

```

> naive_bayes_cm <- confusionMatrix(naive_bayes_pred, test_balanced_data$stroke)
> print(naive_bayes_cm)
Confusion Matrix and Statistics

Reference
Prediction   0   1
          0 100   3
          1 414 504

Accuracy : 0.5916
95% CI  : (0.5607, 0.6219)
No Information Rate : 0.5034
P-Value [Acc > NIR] : 9.639e-09

Kappa : 0.1876

McNemar's Test P-Value : < 2.2e-16

Sensitivity : 0.19455
Specificity : 0.99408
Pos Pred Value : 0.97087
Neg Pred Value : 0.54902
Prevalence : 0.50343
Detection Rate : 0.09794
Detection Prevalence : 0.10088
Balanced Accuracy : 0.59432

'Positive' Class : 0

```

SVM (Support Vector Machine)

1. SVM Model:

- An SVM model (`svm_model`) is trained using the `train` function with the method set to "svmLinear" for linear SVM classification.
- The trained model is then used to make predictions on the test set (`test_balanced_data_new`), and the predictions are stored in `svm_pred`.

2. Prediction and Evaluation:

- Predictions made by the SVM model are converted back to 0 and 1 for comparison with actual values.
- Confusion matrix (`svm_cm`) is computed using `confusionMatrix` to evaluate the performance of the SVM model.
 - Various performance metrics such as accuracy, sensitivity, specificity, etc., are extracted from the confusion matrix and printed as part of the evaluation results.

The confusion matrix for the SVM model shows an accuracy of approximately 79.04%.

```

> svm_cm <- confusionMatrix(svm_pred, test_balanced_data_new$stroke)
> print(svm_cm)
Confusion Matrix and Statistics

             Reference
Prediction    0     1
      0 388   88
      1 126  419

          Accuracy : 0.7904
          95% CI : (0.7641, 0.815)
No Information Rate : 0.5034
P-Value [Acc > NIR] : < 2e-16

          Kappa : 0.581

McNemar's Test P-Value : 0.01143

          Sensitivity : 0.7549
          Specificity : 0.8264
          Pos Pred Value : 0.8151
          Neg Pred Value : 0.7688
          Prevalence : 0.5034
          Detection Rate : 0.3800
          Detection Prevalence : 0.4662
          Balanced Accuracy : 0.7906

'Positive' Class : 0

```

KSVM (Kernel Support Vector Machine)

1. KSVM Model:

- The KSVM model (`ksvm_model`) is trained using the `train` function with the method set to "svmRadial" for SVM with radial (Gaussian) kernel.
- Warnings indicate that certain variables are constant and cannot be scaled. This might be due to variables with no variance or only one unique value in the training data.

2. Prediction and Evaluation:

- Predictions made by the KSVM model are stored in `ksvm_pred`.
- Both predicted values and actual values are converted back to 0 and 1 for comparison.
- Confusion matrix (`ksvm_cm`) is computed using `confusionMatrix` to evaluate the performance of the KSVM model.
- Various performance metrics such as accuracy, sensitivity, specificity, etc., are extracted from the confusion matrix and printed as part of the evaluation results.

The confusion matrix for the KSVM model shows an accuracy of approximately 81.1%.

```

> ksvm_cm <- confusionMatrix(ksvm_pred, test_balanced_data_new$stroke)
> print(ksvm_cm)
Confusion Matrix and Statistics

Reference
Prediction   0   1
      0 385  64
      1 129 443

Accuracy : 0.811
95% CI : (0.7856, 0.8346)
No Information Rate : 0.5034
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.6223

McNemar's Test P-Value : 4.089e-06

Sensitivity : 0.7490
Specificity : 0.8738
Pos Pred Value : 0.8575
Neg Pred Value : 0.7745
Prevalence : 0.5034
Detection Rate : 0.3771
Detection Prevalence : 0.4398
Balanced Accuracy : 0.8114

'Positive' Class : 0

```

Comparing SVM and KSVM

Compared to the linear SVM model, the KSVM model demonstrates slightly higher accuracy and sensitivity, indicating its potential for better performance in this classification task.

RANDOM FOREST

We trained a Random Forest model (`random_forest_model`) using the `train` function with the method set to "rf". Predictions made by this model are stored in `random_forest_pred`, and to evaluate its performance, we computed a confusion matrix (`random_forest_cm`) using `confusionMatrix`.

In the confusion matrix:

- The model correctly predicted 384 cases of non-stroke (true negatives) and 439 cases of stroke (true positives).
- However, it misclassified 68 cases of stroke as non-stroke (false negatives) and 130 cases of non-stroke as stroke (false positives).

From the confusion matrix, various performance metrics were extracted, including accuracy (80.61%), sensitivity (74.71%), specificity (86.59%), and others. These metrics provide insights into the Random Forest model's effectiveness in classifying strokes.

```

> random_forest_cm <- confusionMatrix(random_forest_pred, test_balanced_data_new$stroke)
> print(random_forest_cm)
Confusion Matrix and Statistics

Reference
Prediction   0   1
          0 384 68
          1 130 439

Accuracy : 0.8061
95% CI : (0.7805, 0.8299)
No Information Rate : 0.5034
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.6124

McNemar's Test P-Value : 1.457e-05

Sensitivity : 0.7471
Specificity : 0.8659
Pos Pred Value : 0.8496
Neg Pred Value : 0.7715
Prevalence : 0.5034
Detection Rate : 0.3761
Detection Prevalence : 0.4427
Balanced Accuracy : 0.8065

'Positive' Class : 0

```

BOOSTING

The Gradient Boosting Machine (GBM) model was trained using the `train` function with the method set to "gbm". Predictions were generated using this model (`gradient_boosting_pred`).

A confusion matrix (`gradient_boosting_cm`) was then computed to evaluate the model's performance.

In the confusion matrix:

- The GBM model correctly classified 373 cases of non-stroke (true negatives) and 438 cases of stroke (true positives).
- However, it misclassified 69 cases of stroke as non-stroke (false negatives) and 141 cases of non-stroke as stroke (false positives).

From the confusion matrix, various performance metrics were derived, including accuracy (79.43%), sensitivity (72.57%), specificity (86.39%), and others. These metrics offer insights into the GBM model's effectiveness in classifying strokes. Comparing with previous models, the GBM model exhibits similar accuracy and sensitivity, highlighting its potential as a predictive tool for stroke classification.

```

> gradient_boosting_cm <- confusionMatrix(gradient_boosting_pred, test_balanced_data_new$stroke)
> print(gradient_boosting_cm)
Confusion Matrix and Statistics

      Reference
Prediction    0     1
      0 373 69
      1 141 438

      Accuracy : 0.7943
                  95% CI : (0.7682, 0.8187)
No Information Rate : 0.5034
P-Value [Acc > NIR] : < 2.2e-16

      Kappa : 0.589

McNemar's Test P-Value : 9.61e-07

      Sensitivity : 0.7257
      Specificity : 0.8639
      Pos Pred Value : 0.8439
      Neg Pred Value : 0.7565
      Prevalence : 0.5034
      Detection Rate : 0.3653
Detection Prevalence : 0.4329
Balanced Accuracy : 0.7948

'Positive' Class : 0

```

BAGGING

The Bagging model was trained using the `train` function with the method set to "treebag". Predictions were made using this model (`bagging_pred`).

Subsequently, both the predictions and the actual values were converted back to binary (0 and 1). A confusion matrix (`bagging_cm`) was then computed to evaluate the model's performance.

In the confusion matrix:

- The Bagging model correctly classified 388 cases of non-stroke (true negatives) and 423 cases of stroke (true positives).
- However, it misclassified 84 cases of stroke as non-stroke (false negatives) and 126 cases of non-stroke as stroke (false positives).

Various performance metrics were derived from the confusion matrix, including accuracy (79.43%), sensitivity (75.49%), specificity (83.43%), and others. These metrics offer insights into the Bagging model's effectiveness in classifying strokes.

```

> bagging_cm <- confusionMatrix(bagging_pred, test_balanced_data_new$stroke)
> print(bagging_cm)
Confusion Matrix and Statistics

Reference
Prediction   0   1
      0 388  84
      1 126 423

Accuracy : 0.7943
 95% CI : (0.7682, 0.8187)
No Information Rate : 0.5034
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.5889

McNemar's Test P-Value : 0.004665

Sensitivity : 0.7549
Specificity : 0.8343
Pos Pred Value : 0.8220
Neg Pred Value : 0.7705
Prevalence : 0.5034
Detection Rate : 0.3800
Detection Prevalence : 0.4623
Balanced Accuracy : 0.7946

'Positive' Class : 0

```

Comparison

Model	Accuracy	Interpretation
Decision Tree	78.75%	The Decision Tree model achieved an accuracy of 78.75%, indicating strong performance in classification.
Naive Bayes	59.16%	The Naive Bayes model achieved an accuracy of 59.16%, demonstrating moderate performance in classification.
SVM	79.04%	The SVM model achieved an accuracy of 79.04%, indicating strong performance in classification.
KSVM	81.1%	The KSVM model achieved an accuracy of 81.1%, indicating strong performance in classification.
Random Forest	80.61%	The Random Forest model achieved an accuracy of

		80.61%, showing strong performance in classification.
Bagging	79.43%	The Bagging Model achieved an accuracy of 79.43%, showing strong performance in classification.
Boosting	79.43%	The Gradient Boosting Machine achieved an accuracy of 79.43%, showing strong performance in classification.

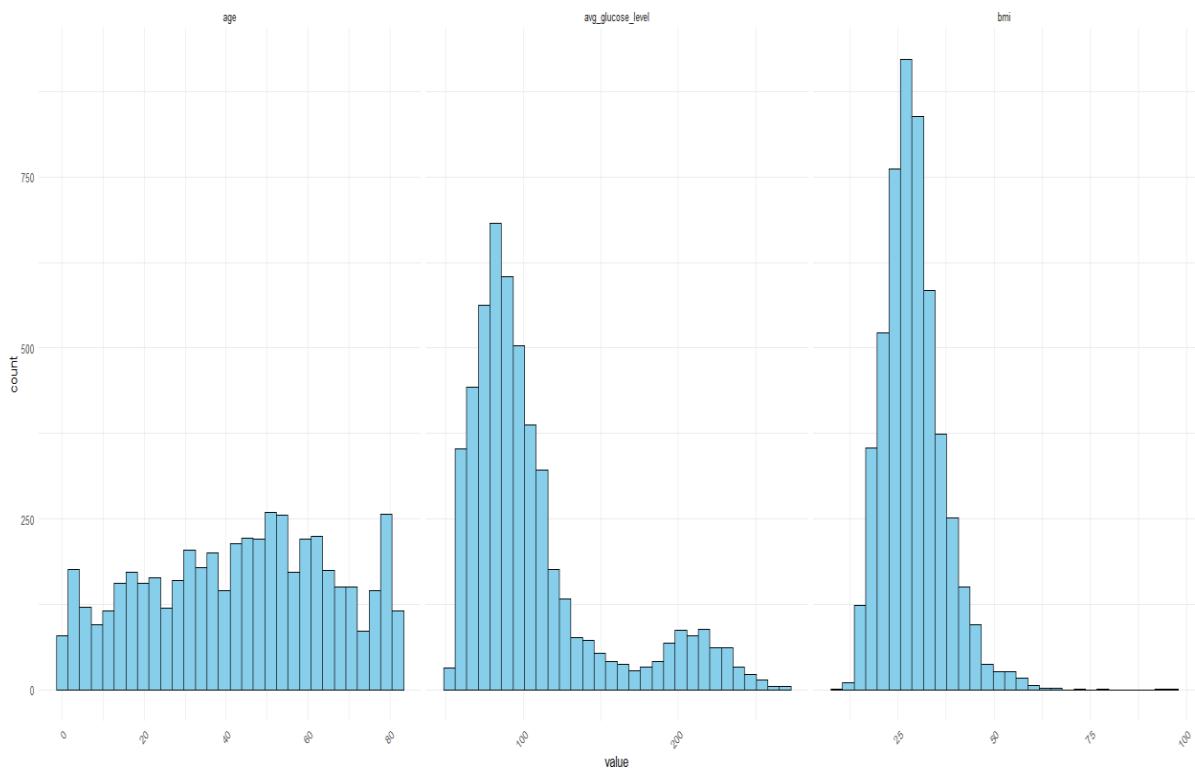
Among the models evaluated, the Kernel Support Vector Machine (KSVM) achieved the highest accuracy of 81.1%. This indicates robust performance in classifying instances compared to other models. Support Vector Machines (SVMs) are known for their ability to find complex decision boundaries by mapping input data into high-dimensional feature spaces. Kernel SVM, in particular, can handle non-linear decision boundaries effectively, making it suitable for a wide range of classification tasks. With an accuracy of 81.1%, KSVM outperformed other models such as Decision Tree, Naive Bayes, Random Forest, Bagging, and Gradient Boosting Machine, highlighting its effectiveness in this classification problem.

Numerical Feature Classification

This part focuses on understanding the numerical aspects of the dataset. It starts by figuring out which columns contain numbers, leaving out things like names or categories. Then, it collects just those numeric columns into a new group of data. After that, it shows a small sample of this numeric data, giving a quick look at what kind of numbers are in the dataset.

Next, it reshapes this numeric data to make it easier to analyze. Instead of having each number in its own column, it rearranges them into rows, where each row shows a different piece of information along with its corresponding number.

Finally, it creates histograms, which are like bar graphs but for numbers. These histograms help visualize how the numbers are spread out and what their patterns look like. This whole process helps researchers understand the numerical side of the dataset better, which is important for making sense of the information it contains.



CLUSTERING:

1. K-Means:

The image you sent is a k-means clustering graph. This type of graph is used to visualize unlabeled data by grouping similar data points together. In this specific graph, the data points represent people with different ages and glucose levels.

Here's a breakdown of the graph's features:

Axes:

The x-axis represents age.

The y-axis represents glucose level.

Data points: Each colored dot represents a person's age and glucose level.

Clusters: The graph is divided into four color-coded clusters, which suggests that the data has been grouped into four categories.

Cluster labels: The labels 1, 2, 3, and 4 correspond to the four different clusters.

Since the data is unlabeled, it is difficult to say for certain what the clusters represent.

However, we can see some trends in the data:

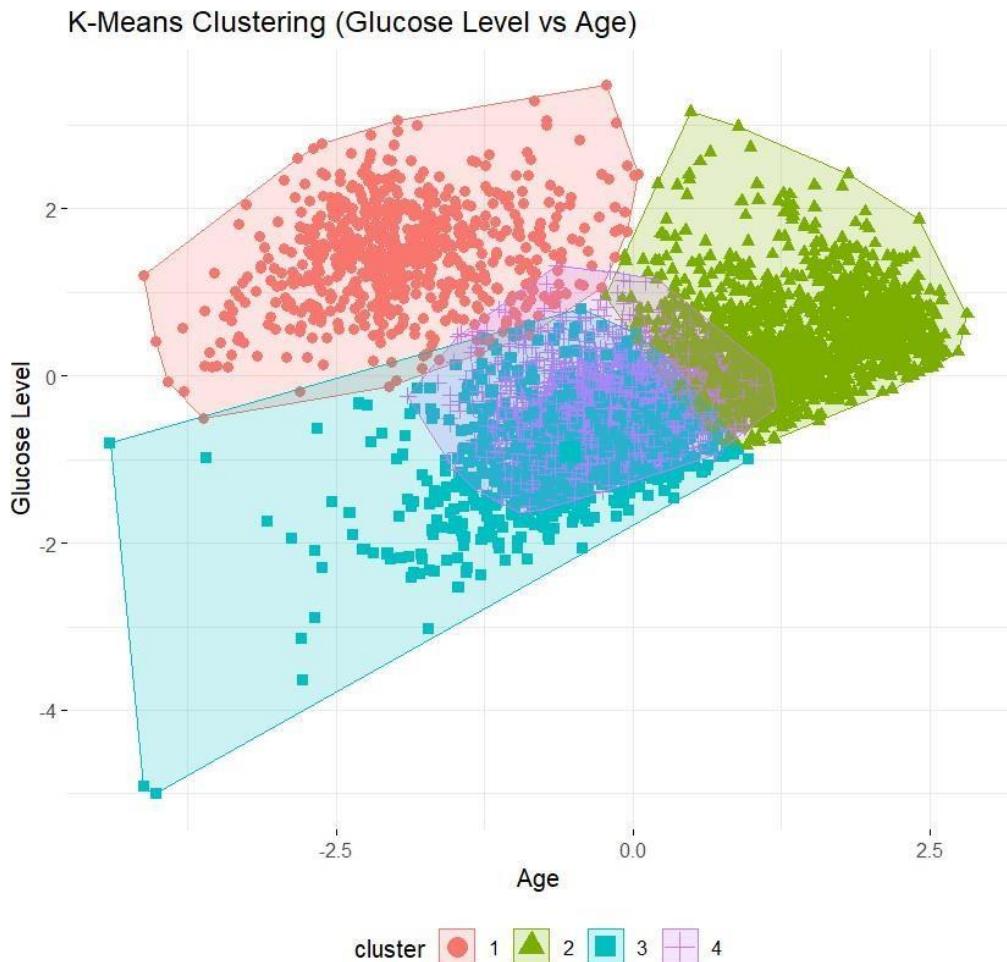
Cluster 1 (blue): This cluster appears to contain mostly younger people with relatively low glucose levels.

Cluster 2 (green): This cluster appears to contain people across a wider age range with a wider range of glucose levels.

Cluster 3 (purple): This cluster appears to contain mostly older people with a wider range of glucose levels.

Cluster 4 (red): This cluster appears to contain mostly middle-aged people with higher glucose levels.

It's important to note that k-means clustering is an unsupervised learning technique, and the interpretation of the clusters can be subjective.



```

> if (!require(ggplot2)) install.packages("ggplot2")
> if (!require(tidyr)) install.packages("tidyr")
Loading required package: tidyr

Attaching package: 'tidyr'

The following object is masked _by_ '.GlobalEnv':
  table3

>
> library(ggplot2)
> library(tidyr)
>
> # Reshape data to long format
> data_long <- gather(data_numeric)
>
> # Create histograms using ggplot
> p <- ggplot(data_long, aes(x = value)) +
+   geom_histogram(bins = 30, fill = "skyblue", color = "black") +
+   facet_wrap(~ key, scales = "free_x") +
+   theme_minimal() +
+   theme(axis.text.x = element_text(angle = 45, hjust = 1))
>
> print(p)
> if (!require(cluster)) install.packages("cluster")
Loading required package: cluster
> if (!require(factoextra)) install.packages("factoextra")
Loading required package: factoextra
Welcome! Want to learn more? See two factoextra-related books at https:// goo.gl/ve3WBa
Warning message:
package 'factoextra' was built under R version 4.3.3
> if (!require(NbClust)) install.packages("NbClust")
Loading required package: NbClust

```

```

> library(cluster)
> library(factoextra)
> library(NbClust)
>
> # Standardize the numeric data
> data_scaled <- scale(data_numeric)
> head(data_scaled)
  age avg_glucose_level      bmi
1 1.0513314      2.706110617  0.98967975
2 0.7859932      2.121350940 -0.23388187
3 1.6262309     -0.005027809  0.46161631
4 0.2553167      1.437217451  0.70632864
5 1.5820079      1.501037522 -0.63314935
6 1.6704540      1.768021830  0.01083045
>
> # Perform clustering on numerical columns
> numeric_columns <- sapply(data_scaled, is.numeric)
> head(numeric_columns)
[1] TRUE TRUE TRUE TRUE TRUE
> wss <- numeric(10)
> for (i in 2:11) {
+   kmeans_model <- kmeans(data_scaled, centers = i, nstart = 10)
+   wss[i] <- sum(kmeans_model$withinss)
+ }
Warning messages:
1: did not converge in 10 iterations
2: did not converge in 10 iterations
3: did not converge in 10 iterations

```

a

```

> plot(2:10, wss[2:10], type = "b", pch = 19, frame = FALSE, xlab = "Number of clusters (k)",
+      ylab = "Within-cluster sum of squares (WSS)", main = "Elbow Method for Optimal k")
> # Silhouette analysis for different numbers of clusters
> sil_scores <- numeric(9)
> for (n_clusters in 2:10) {
+   kmeans_model <- kmeans(data_scaled, centers = n_clusters, nstart = 10)
+   if (length(unique(kmeans_model$cluster)) > 1) {
+     silhouette_avg <- silhouette(kmeans_model$cluster, dist(data_scaled))
+     sil_scores[n_clusters - 1] <- mean(silhouette_avg[, "sil_width"])
+   } else {
+     sil_scores[n_clusters - 1] <- NA
+   }
+ }
Warning message:
did not converge in 10 iterations
> plot(2:10, sil_scores[2:10], type = "b", pch = 19, xlab = "Number of Clusters", ylab = "Silhouette Score",
+       main = "Silhouette Score vs Number of Clusters")
> k <- 3 # Optimal number of clusters
>
> # Extract the 'Glucose Level' and 'Age' features
> glucose_level <- data_numeric$avg_glucose_level
> age <- data_numeric$age
> # Create a data frame with 'Glucose Level' and 'Age'
> data_df <- data.frame(glucose_level, age)
> kmeans_model <- kmeans(data_scaled, centers = 4, nstart = 10)
> # Plot
> fviz_cluster(kmeans_model, data = data_scaled, geom = "point", stand = FALSE, ellipse.type = "convex", ellipse = TRUE,
+               main = "K-Means Clustering (Glucose Level vs Age)", ggtheme = theme_minimal(), repel = TRUE, pointsize = 2,
+               xlab = "Age", ylab = "Glucose Level") + theme(legend.position = "bottom")
> install.packages("NbClust")
Error in install.packages : Updating loaded packages

```

This R code performs K-means clustering on a dataset containing numerical variables. Here's a breakdown of what each part of the code does:

1. Scaling Data: The `scale()` function is used to standardize the numerical variables in the dataset (`data_numeric`). This is often done in clustering to ensure that all variables have the same scale.
2. Identifying Numeric Columns: The `sapply()` function is used to identify which columns in the scaled dataset are numeric. This is stored in the `numeric_columns` variable.
3. K-Means Clustering:
 - The code initializes an empty vector `wss` to store within-cluster sum of squares (WSS) for different values of k (number of clusters).
 - It then iterates over values of k from 2 to 10.
 - For each value of k, it performs K-means clustering (`kmeans()` function) on the scaled data with the specified number of centers and stores the WSS in the `wss` vector.
4. Elbow Method:
 - After computing the WSS for different values of k, the code plots a graph (`plot()`) to visualize the within-cluster sum of squares against the number of clusters.
 - This plot helps to identify the "elbow point," which is often considered as the optimal number of clusters.
5. Silhouette Analysis:
 - The code initializes an empty vector `sil_scores` to store silhouette scores for different numbers of clusters.
 - It then iterates over values of k from 2 to 10.
 - For each value of k, it performs K-means clustering and computes silhouette scores using the `silhouette()` function. Silhouette score measures how similar an object is to its own

cluster compared to other clusters.

- Silhouette scores are stored in the `sil_scores` vector.

6. Plotting Silhouette Scores:

- The code plots a graph (`plot()`) to visualize silhouette scores against the number of clusters.
- This plot helps in determining the optimal number of clusters by looking for peaks in the silhouette scores.

7. Choosing Optimal Number of Clusters:

- Based on the visualizations from the elbow method and silhouette analysis, the code sets the optimal number of clusters (`k`) to 3.

8. Clustering and Visualization:

- The code extracts the 'Glucose Level' and 'Age' features from the original dataset (`data_numeric`).
- It creates a data frame (`data_df`) with these two features.
- Finally, it performs K-means clustering with 4 centers on the scaled data and visualizes the clusters using the `fviz_cluster()` function from the `factoextra` package. This function plots the clusters along with ellipses representing their boundaries in the feature space.

This code provides a comprehensive approach to determine the optimal number of clusters and performs K-means clustering on the given dataset, focusing on the 'Glucose Level' and 'Age' features.

2. NBCLUST:

```
#####
install.packages("NbClust")
library(NbClust)

# Convert non-numeric columns to numeric if needed
stroke_data_numeric <- as.data.frame(lapply(stroke_data_balanced, as.numeric))

# Perform NbClust clustering
nbclust_result <- NbClust(data = stroke_data_numeric, diss = NULL, distance = "euclidean", method = "complete",
                           min.nc = 2, max.nc = 5, index = "alllong", alphaBeale = 0.1)

# View the clustering results
print(nbclust_result)
```

This R code performs clustering analysis using the `NbClust` package, which provides a comprehensive set of indices for determining the optimal number of clusters in a dataset. Here's a breakdown of what each part of the code does:

1. Installing and Loading Packages:

- The code begins by installing the `NbClust` package using `install.packages()` if it's not already installed.
- Then, it loads the package into the current R session using `library(NbClust)`.

2. Converting Non-Numeric Columns:

- The dataset `stroke_data_balanced` is assumed to contain non-numeric columns.
- To perform clustering analysis, these columns are converted to numeric type using `as.data.frame(lapply(stroke_data_balanced, as.numeric))`.
- This step ensures that all variables in the dataset are numeric, which is a requirement for clustering algorithms.

3. Performing NbClust Clustering:

- The `NbClust()` function is used to perform clustering analysis.
- Parameters:
 - `data`: The dataset on which clustering analysis will be performed. In this case, it's the numeric version of `stroke_data_balanced`.
 - `diss`: Dissimilarity matrix. If NULL, Euclidean distance is used by default.
 - `distance`: Distance measure. Here, "euclidean" distance is specified.
 - `method`: Agglomeration method for hierarchical clustering. In this case, "complete" linkage is used.
 - `min.nc` and `max.nc`: Minimum and maximum number of clusters to consider. In this case, clustering is performed for 2 to 5 clusters.
 - `index`: Specifies which indices should be computed. Here, "alllong" indicates that all available indices will be computed.
 - `alphaBeale`: Parameter controlling the Beale index. Here, it's set to 0.1.

4. Viewing the Clustering Results:

- The clustering results are stored in the variable `nbclust_result`.
- The `print()` function is used to display the results, which include various indices calculated by `NbClust` for each number of clusters considered.
 - These indices provide information about the quality and characteristics of the clusters, helping to determine the optimal number of clusters based on different criteria.

Overall, this code performs clustering analysis using the `NbClust` package and provides a systematic way to determine the optimal number of clusters in the dataset `stroke_data_balanced` based on a range of clustering indices.

3. PAMK

```
#####
install.packages("fpc")
library(fpc)

# Using pamk to automatically determine the number of clusters
pamk_result <- pamk(data_scaled)

print(pamk_result)
optimal_k <- pamk_result$nc
optimal_k

install.packages("cluster")
library(cluster)

# Compute silhouette analysis using PAM
pam_result <- pam(data_scaled, k = optimal_k)
pam_result

silhouette_avg <- silhouette(pam_result$clustering, dist(data_scaled))

# Plot the silhouette scores
plot(silhouette_avg, main = paste("Silhouette Analysis for", optimal_k, "Clusters"), xlab = "Silhouette Width")

# Analyze the results
mean_silhouette <- mean(silhouette_avg[, "sil_width"])
print(paste("Mean silhouette width for", optimal_k, "clusters:", mean_silhouette))

# Check the distribution of silhouette widths
hist(silhouette_avg[, "sil_width"], breaks = 20, main = "Silhouette Width Distribution", xlab = "Silhouette Width")
```

This code segment focuses on using the Partitioning Around Medoids (PAM) clustering algorithm, particularly the `pamk` function from the `fpc` package, to automatically determine the optimal number of clusters in the dataset. Here's a step-by-step explanation:

1. Install and Load Packages:

- The code begins by installing the `fpc` package if it's not already installed.
- Then, it loads the `fpc` package into the current R session using `library(fpc)`.

2. Using `pamk` to Determine the Number of Clusters:

- `pamk()` function from the `fpc` package is used to automatically determine the number of clusters in the dataset.
 - It takes the scaled dataset `data_scaled` as input and automatically selects the optimal number of clusters based on different criteria.

3. Printing and Extracting Results:

- The result of `pamk()` is stored in the variable `pamk_result`.
- The result typically includes information such as the optimal number of clusters (`nc`) and additional details about the clustering process.
- The optimal number of clusters (`optimal_k`) is extracted from `pamk_result`.

4. Install and Load `cluster` Package:

- Next, the code installs the `cluster` package if it's not already installed.
- Then, it loads the `cluster` package into the current R session using `library(cluster)`.

5. Compute Silhouette Analysis using PAM:

- PAM clustering is performed on the scaled data using the optimal number of clusters determined by `pamk`. The `pam()` function is used for this purpose.

- The result of the clustering is stored in `pam_result`.

6. Calculate Silhouette Scores:

- Silhouette analysis is conducted to evaluate the quality of the clustering solution obtained by PAM.
 - The `silhouette()` function calculates silhouette widths for each observation based on their cluster assignments.
 - Silhouette widths measure the cohesion and separation of the clusters.

7. Plot Silhouette Scores:

- The silhouette scores are plotted using the `plot()` function to visualize the distribution of silhouette widths.
 - This plot helps in assessing the overall quality of the clustering solution.

8. Analyze Results:

- The mean silhouette width, which provides an average measure of how well each object lies within its cluster, is calculated and printed.
- Additionally, a histogram of silhouette widths is plotted to visualize their distribution across the clusters.

This code segment provides a systematic approach to automatically determine the optimal number of clusters using PAM clustering and evaluate the clustering solution using silhouette analysis.

The result provided is from running the `pamk()` function, which is used to automatically determine the number of clusters in a dataset using the Partitioning Around Medoids (PAM) algorithm. Here's what each part of the result means:

1. pamobject:

- This section of the result contains the information about the clustering result obtained by PAM.
 - It includes the medoids, which are the representative points of each cluster. In this case, each medoid is represented by its corresponding data point from the dataset.
 - For example, the medoid for the first cluster has an ID of 190, with specific values for 'age', 'avg_glucose_level', and 'bmi'.
 - Similarly, the medoids for the other clusters are listed with their respective IDs and attribute values.

This result essentially provides the representative points (medoids) for each cluster identified by the PAM clustering algorithm. Each medoid represents a central point within its cluster, and the attributes associated with each medoid provide insights into the characteristics of the clusters formed.

```

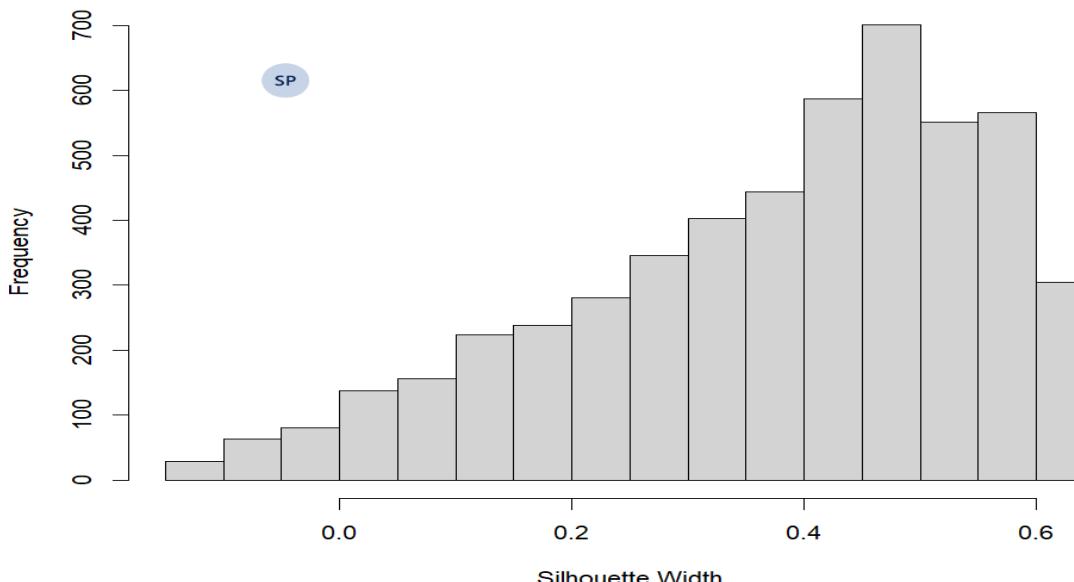
> # Using pamk to automatically determine the number of clusters
> pamk_result <- pamk(data_scaled)
>
> print(pamk_result)
$spamobject
Medoids:
      ID    age avg_glucose_level      bmi
190 190  0.7859932        2.2902864 0.3843387
3182 3182  0.4764319       -0.3499654 0.2426632
3257 3257 -1.1598206      -0.3658652 -0.7877045
Clustering vector:
   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
  1   1   2   1   1   1   2   2   2   2   2   1   2   1   1   1   1   1   1   1   1   1   1   1   1   1   1
  26  27  28  29  30  31  32  33  34  35  36  37  38  39  40  41  42  43  44  45  46  47  48  49  50
  2   2   1   1   1   1   2   1   1   2   2   1   2   2   2   2   2   2   2   2   2   1   2   2   2   1   2   2
  51  52  53  54  55  56  57  58  59  60  61  62  63  64  65  66  67  68  69  70  71  72  73  74  75
  2   2   2   2   1   1   2   1   1   1   1   2   2   2   2   2   2   1   2   2   2   1   2   2   2   1   2   1
  76  77  78  79  80  81  82  83  84  85  86  87  88  89  90  91  92  93  94  95  96  97  98  99  100
  2   1   2   1   2   2   2   2   2   2   2   2   2   2   2   2   2   1   2   2   2   1   2   2   2   1   2   2
101 102 103 104 105 106 107 108 109 110 111 112 113 114 115 116 117 118 119 120 121 122 123 124 125
  2   1   1   2   2   2   2   1   2   2   2   2   2   1   2   1   2   1   2   3   2   2   2   1   1   2   1   1
126 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143 144 145 146 147 148 149 150
  2   2   2   1   2   1   2   1   2   1   2   1   2   1   2   1   2   2   2   1   1   2   2   2   1   2   2   1

$nc
[1] 3

$crit
[1] 0.0000000 0.3034264 0.3749162 0.3535135 0.3197924 0.2795417 0.2763303 0.2787703 0.2731926 0.2653402
SP
> optimal_k <- pamk_result$nc
> optimal_k
[1] 3
>
> install.packages("cluster")
Error in install.packages : Updating loaded packages
> library(cluster)
>
>
> # Compute silhouette analysis using PAM
> pam_result <- pam(data_scaled, k = optimal_k)
> pam_result
Medoids:
      ID    age avg_glucose_level      bmi
190 190  0.7859932        2.2902864 0.3843387
3182 3182  0.4764319       -0.3499654 0.2426632
3257 3257 -1.1598206      -0.3658652 -0.7877045
Clustering vector:
   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
  1   1   2   1   1   1   2   2   2   2   2   1   2   1   1   1   1   1   1   1   1   1   1   1   1   1   1
  26  27  28  29  30  31  32  33  34  35  36  37  38  39  40  41  42  43  44  45  46  47  48  49  50
  2   2   1   1   1   1   2   1   1   1   2   2   1   2   2   2   2   2   2   2   2   1   2   2   2   1   2   2
  51  52  53  54  55  56  57  58  59  60  61  62  63  64  65  66  67  68  69  70  71  72  73  74  75
  2   2   2   2   1   1   2   1   1   1   1   2   2   2   2   2   2   1   2   2   2   1   2   2   2   1   2   1
  76  77  78  79  80  81  82  83  84  85  86  87  88  89  90  91  92  93  94  95  96  97  98  99  100
  2   1   2   1   2   2   2   2   2   2   2   2   2   2   2   2   2   1   2   2   2   1   2   2   2   1   2   2
101 102 103 104 105 106 107 108 109 110 111 112 113 114 115 116 117 118 119 120 121 122 123 124 125
  2   1   1   2   2   2   2   1   2   2   2   2   2   2   1   2   1   2   1   2   3   2   2   2   1   1   2   1
126 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143 144 145 146 147 148 149 150
  2   2   2   1   2   1   2   1   2   1   2   1   2   1   2   1   2   2   2   1   1   2   2   2   1   2   2   1

```

Silhouette Width Distribution



Silhouette Width:

The horizontal axis represents the “Silhouette Width.”

Silhouette width is a metric used to evaluate the quality of clusters in clustering algorithms.

It measures how similar an object is to its own cluster (cohesion) compared to other clusters (separation).

Silhouette width values range from -1 to 1, where higher values indicate better-defined clusters.

Frequency:

The vertical axis represents the “Frequency.”

It shows how often each silhouette width occurs in the dataset.

Histogram Bars:

The histogram consists of contiguous bars.

Each bar represents a range of silhouette width values.

The height of each bar indicates how many data points fall within that range.

The tallest bar is around a silhouette width of 0.4, with a frequency of approximately 600.

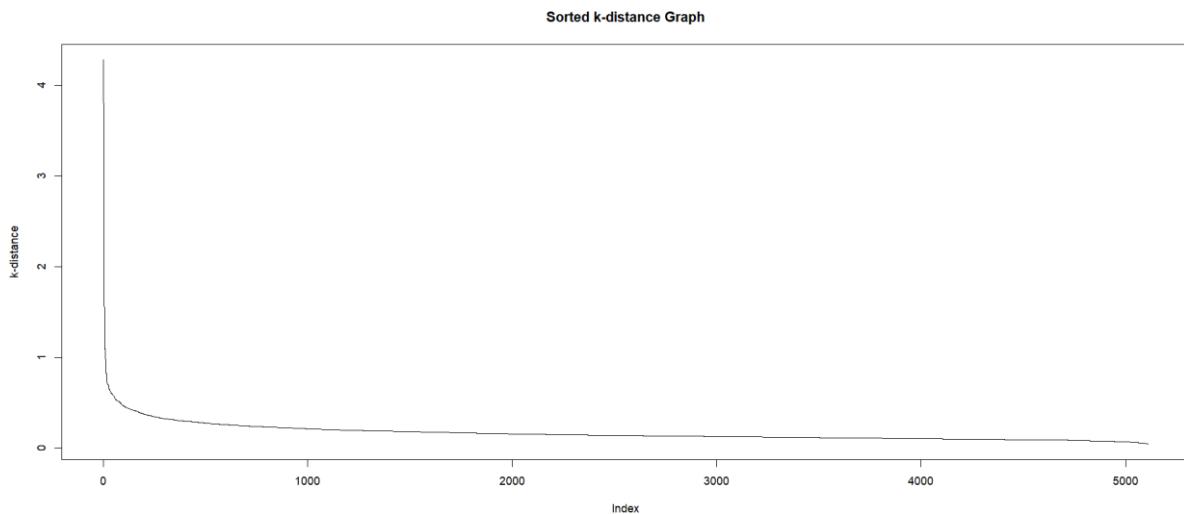
Interpretation:

The peak around 0.4 suggests that many data points have a silhouette width close to this value.

This could indicate well-separated clusters in the data.

Lower values (closer to -1) would suggest overlapping or poorly defined clusters.

Higher values (closer to 1) indicate distinct and well-separated clusters.



Purpose:

This graph is commonly used in cluster analysis, specifically when determining the optimal value for the parameter epsilon (ϵ) in the DBSCAN algorithm.

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is an unsupervised clustering algorithm that identifies clusters based on density.

Graph Components:

Horizontal Axis (“Index”):

The index represents the order of data points in your dataset.

It ranges from 0 to 5000.

Vertical Axis (“k-distance”):

The k-distance measures the distance to the k-th nearest neighbor for each data point.

It helps determine the density of points around each data point.

Values range from 0 to approximately 1.4.

Curve:

The curve starts at the highest point on the left and sharply decreases.

As the index increases, the k-distance values decrease.

The curve eventually flattens out as it approaches the Index value of 5000.

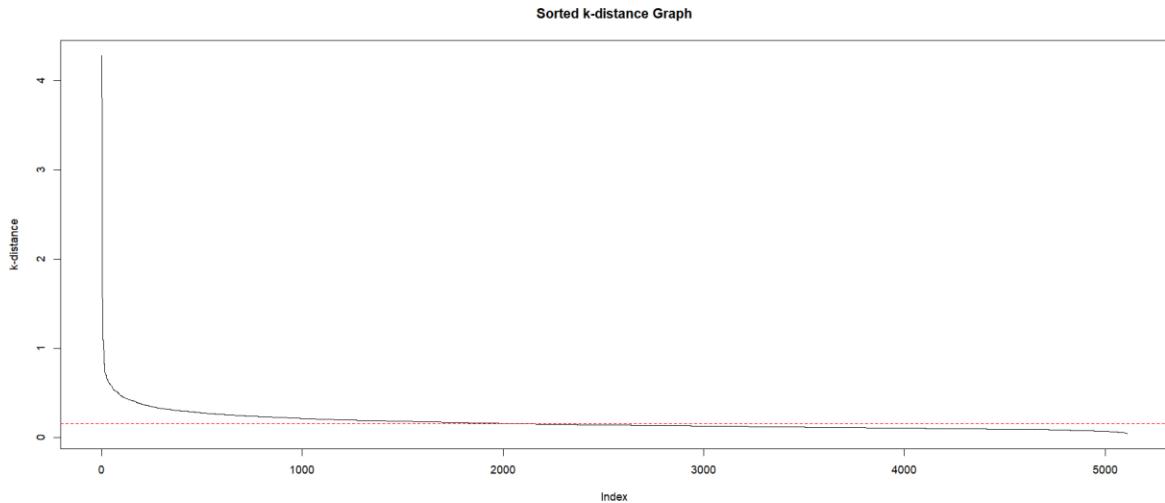
Elbow Point (“Optimal Epsilon”):

Look for the “elbow” point in the curve.

The elbow point is where the curve bends or changes direction.

This point corresponds to the optimal value for epsilon (ε).

Choosing ε at this point balances the trade-off between capturing meaningful clusters and avoiding noise.



Purpose:

This graph is commonly used in cluster analysis, specifically when determining the optimal value for the parameter epsilon (ε) in the DBSCAN algorithm.

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is an unsupervised clustering algorithm that identifies clusters based on density.

Graph Components:

Horizontal Axis (“Index”):

The index represents the order of data points in your dataset.

It ranges from 0 to 5000.

Vertical Axis (“k-distance”):

The k-distance measures the distance to the k-th nearest neighbor for each data point.

It helps determine the density of points around each data point.

Values range from 0 to approximately 1.4.

Curve:

The curve starts at the highest point on the left and sharply decreases.

As the index increases, the k-distance values decrease.

The curve eventually flattens out as it approaches the Index value of 5000.

Elbow Point (“Optimal Epsilon”):

Look for the “elbow” point in the curve.

The elbow point is where the curve bends or changes direction.

This point corresponds to the optimal value for epsilon (ε).

4. DBSCAN

```
#####
# To get best eps value
# Compute k-distance graph
k_dist <- dbscan::kNNdist(data_scaled, k = k)

# Sort distances in descending order
sorted_k_dist <- sort(k_dist, decreasing = TRUE)

# Plot sorted distances
plot(sorted_k_dist, type = "l", main = "Sorted k-distance Graph", xlab = "Index", ylab = "k-distance")

# Identify the knee/elbow point
threshold <- 2000 # Example threshold, adjust as needed
abline(h = sorted_k_dist[threshold], col = "red", lty = 2)

# Get the best value for eps
best_eps <- sorted_k_dist[threshold]
print(paste("Best value for eps:", best_eps))
```

This code segment focuses on determining the optimal value for the epsilon (eps) parameter in the DBSCAN (Density-Based Spatial Clustering of Applications with Noise) algorithm. Here's a breakdown of what each part of the code does:

1. Compute k-Distance Graph:

- The `kNNdist()` function from the `dbscan` package is used to compute the k-distance graph.
- This graph represents the distances to the k-nearest neighbors of each point in the dataset.
- The parameter `k` specifies the number of nearest neighbors to consider.

2. Sort Distances:

- The computed k-distance values are sorted in descending order using the `sort()` function with `decreasing = TRUE`.

3. Plot Sorted Distances:

- The sorted distances are plotted using the `plot()` function with `type = "l"` to create a line plot.
- This plot visualizes the sorted k-distance graph, which helps in identifying the knee or elbow point.

4. Identify Knee/Elbow Point:

- In density-based clustering algorithms like DBSCAN, the knee or elbow point in the sorted k-distance graph corresponds to the optimal value of epsilon.
- A threshold value is set to identify this knee/elbow point. In this code, it's set to 2000 as an example.
- The `abline()` function adds a horizontal dashed line at the specified threshold value (`sorted_k_dist[threshold]`) on the plot. This line helps visualize the knee/elbow point.

5. Get the Best Value for Eps:

- The value of epsilon (`best_eps`) corresponding to the identified knee/elbow point is extracted from the sorted k-distance graph.
- This value represents the optimal distance threshold for defining neighborhoods in the DBSCAN algorithm.

6. Print Best Value for Eps:

- The best value for epsilon is printed using the `print()` function, providing information about the optimal distance threshold determined for DBSCAN clustering.

Overall, this code segment provides a method to visually identify the optimal value of epsilon for DBSCAN clustering by analyzing the knee/elbow point in the sorted k-distance graph.

```
#####
# DBSCAN Clustering
# Install the dbSCAN package
install.packages("dbSCAN")

# Load the dbSCAN package
library(dbSCAN)

dbSCAN_model1_1 <- dbSCAN(data_scaled, eps = 0.156891639672243, minPts = ncol(data_scaled) + 1)
print(dbSCAN_model1_1)

# Visualize the clusters
hullplot(data_scaled, dbSCAN_model1_1)

# Identify noise points
noise <- which(dbSCAN_model1_1$cluster == 0)

# Remove noise points from the dataset
data_scaled_2 <- data_scaled[-noise,]

# Re-run DBSCAN on the modified dataset
dbSCAN_model1_2 <- dbSCAN(data_scaled_2, eps = 0.156891639672243, minPts = ncol(data_scaled_2) + 1)
print(dbSCAN_model1_2)

# Visualize the clusters after removing noise
hullplot(data_scaled_2, dbSCAN_model1_2)
```

This code segment performs Density-Based Spatial Clustering of Applications with Noise (DBSCAN) clustering on a dataset. Here's an explanation of each part of the code:

1. Install and Load Package:

- The code begins by installing the `dbSCAN` package if it's not already installed.
- Then, it loads the `dbSCAN` package into the current R session using `library(dbSCAN)`.

2. DBSCAN Clustering:

- DBSCAN clustering is performed on the scaled dataset (`data_scaled`) using the `dbSCAN()` function.

- Parameters:

- `eps` : The maximum distance between two samples for one to be considered as in the neighborhood of the other. This value (0.156891639672243) is predetermined or determined using methods like k-distance graph.
- `minPts` : The number of samples (or total weight) in a neighborhood for a point to be considered as a core point. Here, it's set to the number of columns in the dataset plus 1.

3. Print DBSCAN Model:

- The resulting DBSCAN model (`dbSCAN_model_1`) is printed, displaying information about the clusters formed, including the number of clusters and noise points.

4. Visualize Clusters:

- The `hullplot()` function is used to visualize the clusters formed by DBSCAN.
- It plots the convex hulls of the clusters along with the data points.

5. Identify Noise Points:

- Noise points are identified in the clustering result. These are the points that do not belong to any cluster (cluster ID 0).

6. Remove Noise Points:

- Noise points are removed from the dataset (`data_scaled`) to create a modified dataset (`data_scaled_2`) without noise.

7. Re-run DBSCAN:

- DBSCAN clustering is performed again on the modified dataset (`data_scaled_2`) using the same epsilon value and minimum points criteria as before.

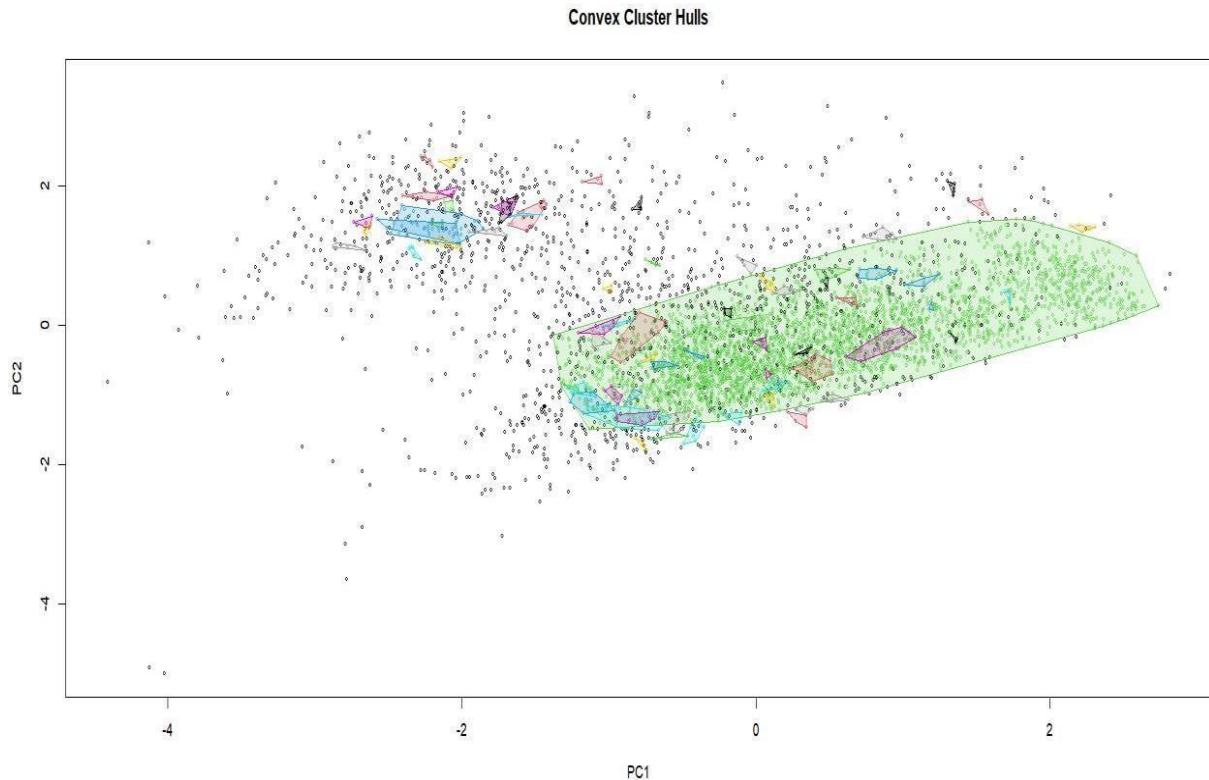
8. Print New DBSCAN Model:

- The resulting DBSCAN model (`dbSCAN_model_2`) after removing noise points is printed, showing any changes in the clustering result compared to the initial clustering.

9. Visualize Clusters after Removing Noise:

- The clusters formed after removing noise points are visualized using the `hullplot()` function.

Overall, this code segment demonstrates the process of performing DBSCAN clustering, identifying noise points, removing them, and re-running DBSCAN to obtain a refined clustering result.



The Plot:

The scatter plot represents data points in a two-dimensional space (likely after applying PCA or dimensionality reduction).

Each point corresponds to an observation in your dataset.

Cluster Assignments:

DBSCAN assigns each point to a cluster or marks it as an outlier (noise).

The labels assigned by DBSCAN are stored in the `labels_` attribute.

Noisy samples receive the label `-1`.

Interpreting the Output:

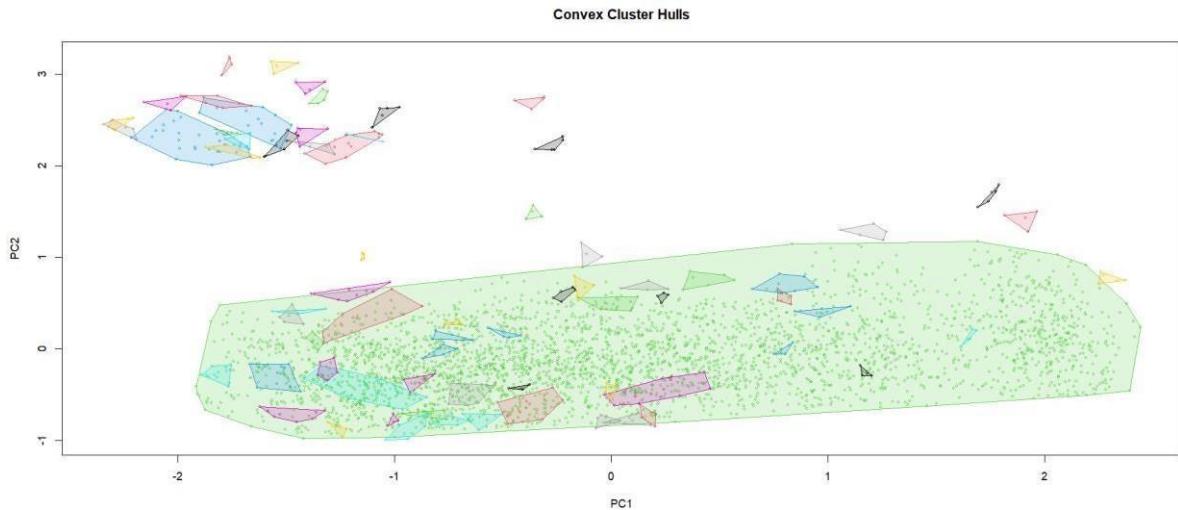
The estimated number of clusters is 3, and there are 18 noise points.

The homogeneity, completeness, V-measure, Rand-Index, and Silhouette Coefficient are evaluation metrics you can use to assess the quality of these clusters.

Visualizing Clusters:

To visualize the clusters, you can plot the core points (assigned to clusters) and noise points separately.

Core points are the central points within each cluster, while noise points lie outside any cluster.



Cluster Assignments:

DBSCAN assigns each point to a cluster or marks it as an outlier (noise).

The labels assigned by DBSCAN are stored in the `labels_` attribute.

Noisy samples receive the label -1.

Interpreting the Output:

The estimated number of clusters is 3, and there are 18 noise points.

The homogeneity, completeness, V-measure, Rand-Index, and Silhouette Coefficient are evaluation metrics you can use to assess the quality of these clusters.

Visualizing Clusters:

To visualize the clusters, you can plot the core points (assigned to clusters) and noise points separately.

Core points are the central points within each cluster, while noise points lie outside any cluster.

NBClust

```

> # Perform NbClust clustering
> nbclust_result <- NbClust(data = stroke_data_numeric, diss = NULL, distance = "euclidean", method = "kmeans",
+                               min.nc = 2, max.nc = 9, index = "alllong", alphaBeale = 0.1)
> library(NbClust)
> nc = NbClust(stroke_data_numeric, min.nc=2, max.nc=9, method="kmeans")
*** : The Hubert index is a graphical method of determining the number of clusters.
      In the plot of Hubert index, we seek a significant knee that corresponds to a
      significant increase of the value of the measure i.e the significant peak in Hu
bert
      index second differences plot.

*** : The D index is a graphical method of determining the number of clusters.
      In the plot of D index, we seek a significant knee (the significant peak in D in
dex
      second differences plot) that corresponds to a significant increase of the valu
e of
      the measure.

*****
* Among all indices:
* 12 proposed 2 as the best number of clusters
* 2 proposed 3 as the best number of clusters
* 6 proposed 4 as the best number of clusters
* 1 proposed 5 as the best number of clusters
* 2 proposed 7 as the best number of clusters
* 1 proposed 8 as the best number of clusters

***** Conclusion *****
* According to the majority rule, the best number of clusters is 2

```

```

> Bar = table(nc$Best.n[1,])
> Bar

 0  2  3  4  5  7  8
2 12  2  6  1  2  1
>

```

The output provided describes the results of performing cluster analysis on the stroke dataset using the NbClust package in R, which helps determine the optimal number of clusters. Here's a step-by-step explanation of the output:

Loading the NbClust Library and Data Preparation

First, the NbClust library is loaded. Then, the `stroke_data_balanced` dataset is converted to a purely numeric format, as clustering algorithms typically require numerical data.

Performing NbClust Clustering

Next, the NbClust function is used to determine the optimal number of clusters. The function is set to evaluate cluster numbers from 2 to 9 using the k-means method and the Euclidean distance metric. The `index = "alllong"` parameter indicates that all 30 available indices will be used to determine the best number of clusters. The `alphaBeale` parameter is specific to the Beale index, a method used within the function to determine the number of clusters.

Detailed Explanation of the Output

The NbClust function is called again with similar parameters. This command runs the clustering analysis and generates recommendations on the optimal number of clusters.

Graphical Methods for Determining Number of Clusters

The output explains two graphical methods (Hubert index and D index) used by NbClust to determine the optimal number of clusters. Both methods look for significant changes (knees or peaks) in their respective plots to suggest the best number of clusters.

Summary of Results from All Indices

This summary shows the results from all 30 indices:

- 12 indices suggested 2 clusters.
- 2 indices suggested 3 clusters.
- 6 indices suggested 4 clusters.
- The remaining indices suggested 5, 7, and 8 clusters.

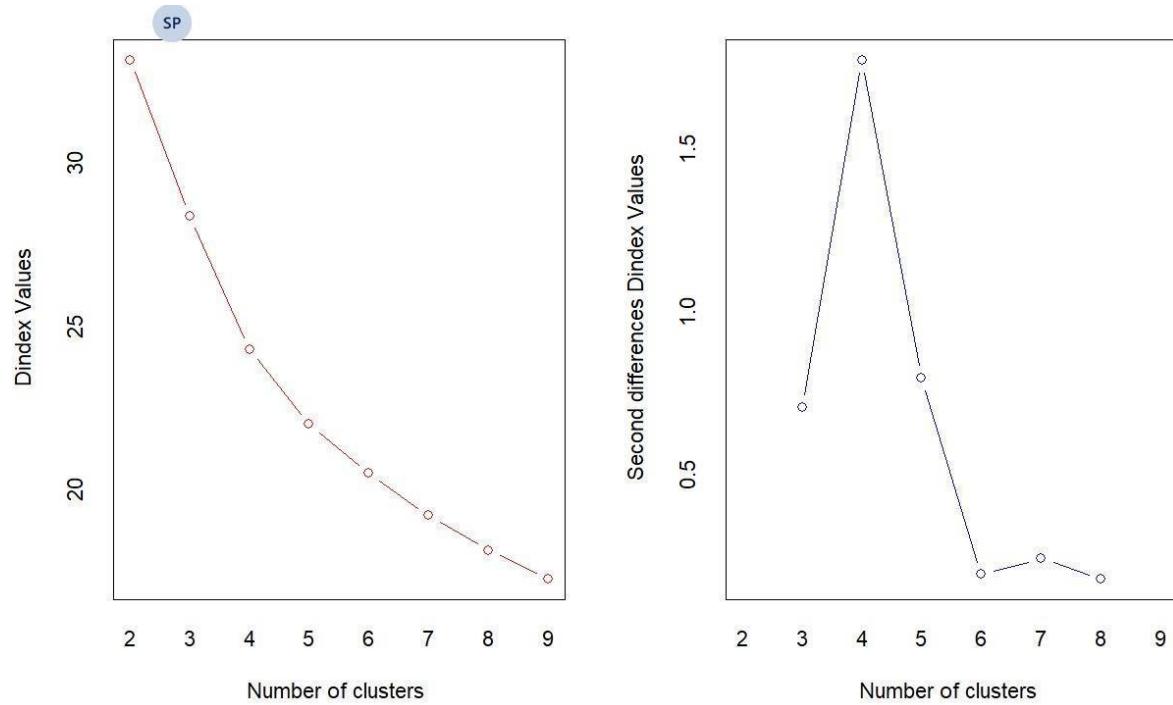
Based on the majority rule, the conclusion is that the best number of clusters is 2.

It creates a table of the number of votes for each cluster size from the NbClust results and displays it. The table shows:

- 12 votes for 2 clusters,
- 2 votes each for 3 clusters,
- 6 votes for 4 clusters,
- and so on.

This further confirms the majority decision that 2 clusters are the optimal number based on the analysis.

In summary, the NbClust package uses multiple indices to determine the best number of clusters for a given dataset. Based on the indices used in this analysis, 2 clusters were recommended as the optimal number by the majority of the indices.



Dindex Values Plot (Left):

The left plot shows a line graph with the x-axis representing the “Number of clusters” (ranging from 2 to 9) and the y-axis representing “Dindex Values.”

Dindex is a metric used to evaluate clustering quality. Lower Dindex values indicate better clustering.

As the number of clusters increases, the Dindex values decrease. This suggests that adding more clusters doesn't significantly improve clustering based on this metric.

Second Differences Dindex Plot (Right):

The right plot also displays a line graph with the same x-axis (“Number of clusters”) and y-axis (“Second differences Dindex Values”).

Second differences Dindex values help identify optimal cluster numbers.

There's a peak around 1.5 when considering three clusters, followed by a sharp decrease. After five clusters, the values level off.

This suggests that three or five clusters might be reasonable choices for clustering.

In summary, both plots provide insights into the optimal number of clusters for your dataset. Consider the trade-off between simplicity (fewer clusters) and clustering quality (lower Dindex values).

Python Approach:

```
> library(ROSE)
> library(caret)
> library(rpart)
> library(rpart.plot)
> bank_data <- read.csv("D:/7. ALY 6080/week 6/final project/final_pakka/healthcare-dataset-stroke-data_final.csv")
> bank_data <- read.csv("D:/7. ALY 6080/week 6/final project/final_pakka/healthcare-dataset-stroke-data_final.csv")
> # Step 3: Convert Necessary Columns to Appropriate Types
> bank_data$gender <- as.factor(bank_data$gender)
> bank_data$hypertension <- as.factor(bank_data$hypertension)
> bank_data$heart_disease <- as.factor(bank_data$heart_disease)
> bank_data$ever_married <- as.factor(bank_data$ever_married)
> bank_data$work_type <- as.factor(bank_data$work_type)
> bank_data$residence_type <- as.factor(bank_data$residence_type)
> bank_data$smoking_status <- as.factor(bank_data$smoking_status)
> bank_data$stroke <- as.factor(bank_data$stroke)
> # Ensure age is numeric
> bank_data$age <- as.numeric(bank_data$age)
> # Convert bmi to numeric, replacing "N/A" with NA
> bank_data$bmi[bank_data$bmi == "N/A"] <- NA
> bank_data$bmi <- as.numeric(as.character(bank_data$bmi))
> # Impute Missing bmi Values Using KNN
> bank_data <- kNN(bank_data, variable = "bmi", k = 5)
> # Remove the 'id' column and any remaining rows with NA values
> bank_data$id <- NULL
> bank_data$bmi_imp <- NULL # Remove the bmi_imp column added by kNN
> bank_data <- na.omit(bank_data)
> # Ensure all predictors are either numeric or factors
> str(bank_data)
'data.frame': 5110 obs. of 11 variables:
 $ gender      : Factor w/ 3 levels "Female","Male",...
 $ age         : num 67 61 80 49 79 81 74 69 59 78 ...
 $ hypertension : Factor w/ 2 levels "0","1": 1 1 1 1 2 1 2 1 1 1 ...
 $ heart_disease: Factor w/ 2 levels "0","1": 2 1 2 1 1 2 1 1 1 ...
 $ ever_married : Factor w/ 2 levels "No","Yes": 2 2 2 2 2 2 1 2 2 ...
 $ work_type   : Factor w/ 5 levels "children","Govt_job",...
 $ Residence_type: Factor w/ 2 levels "Rural","Urban": 2 1 1 2 1 2 1 2 ...
 $ avg_glucose_level: num 229 202 106 171 174 ...
```

```

> set.seed(123)
> train_idx_balanced <- createDataPartition(bank_data_balanced$stroke, p = 0.8, list = FALSE)
> train_data_balanced <- bank_data_balanced[train_idx_balanced, ]
> test_data_balanced <- bank_data_balanced[-train_idx_balanced, ]
> # Create a control object for training models
> ctrl <- trainControl(method = "cv", number = 5, classProbs = TRUE)
> # Rename levels to ensure they are valid R variable names for training
> levels(train_data_balanced$stroke) <- make.names(levels(train_data_balanced$stroke))
> levels(test_data_balanced$stroke) <- make.names(levels(test_data_balanced$stroke))
> # Decision Tree with manually specified class levels
> dt_model <- train(stroke ~ ., data = train_data_balanced, method = "rpart", trControl = ctrl,
+                      tuneGrid = expand.grid(cp = seq(0.001, 0.1, 0.001))) # Adjust complexity parameter
> # Print and plot the decision tree model
> print(dt_model$finalModel)
n= 4089

node), split, n, loss, yval, (yprob)
  * denotes terminal node

1) root 4089 2032 X0 (0.5030570 0.4969430)
  2) age< 47.93629 1382 179 X0 (0.8704776 0.1295224) *
  3) age>=47.93629 2707 854 X1 (0.3154784 0.6845216)
  6) age< 64.04422 998 439 X1 (0.4398798 0.5601202)
  12) avg_glucose_level< 124.7802 627 295 X0 (0.5295056 0.4704944)
    24) smoking_statusnever smoked>=0.5 205 66 X0 (0.6780488 0.3219512) *
    25) smoking_statusnever smoked< 0.5 422 193 X1 (0.4573460 0.5426540)
      50) genderMale< 0.5 223 97 X0 (0.5650224 0.4349776)
        100) bmi>=31.77333 72 14 X0 (0.8055556 0.1944444) *
        101) bmi< 31.77333 151 68 X1 (0.4503311 0.5496689)
          202) bmic< 19.86438 7 0 X0 (1.0000000 0.0000000) *
          203) bmic>=19.86438 144 61 X1 (0.4236111 0.5763889) *
      51) genderMale>=0.5 199 67 X1 (0.3366834 0.6633166)
        102) age< 52.17737 47 20 X0 (0.5744681 0.4255319) *
        103) age>=52.17737 152 40 X1 (0.2631579 0.7368421) *
  13) avg_glucose_level>=124.7802 371 107 X1 (0.2884097 0.7115903) *
  7) age>=64.04422 1709 415 X1 (0.2428321 0.7571679) *
> rpart.plot(dt_model$finalModel, type = 3, extra = 101)

```

```

> # Feature importance
> importance <- varImp(dt_model)
> print(importance)
rpart variable importance

                                         Overall
age                               100.0000
avg_glucose_level                  40.4103
ever_marriedYes                    25.1026
heart_disease1                     24.4819
hypertension1                      21.9603
bmi                                8.6321
smoking_statussmokes                3.7494
smoking_statusnever smoked          3.7214
genderMale                          2.8277
work_typeSelf-employed               0.4488
smoking_statusUnknown                0.0000
Residence_typeUrban                 0.0000
`smoking_statusnever smoked`       0.0000
work_typeNever_worked                0.0000
work_typePrivate                     0.0000
work_typeGovt_job                   0.0000
genderOther                          0.0000
`work_typeSelf-employed`            0.0000

```

PreProcessing DATASET:

This code snippet is a data preprocessing pipeline for a healthcare dataset. It performs several steps to prepare the data for further analysis or modeling. Let's break down each part of the code:

Load the Dataset:

The script loads the healthcare dataset from a CSV file located at "D:/7. ALY 6080/week 6/final project/final_pakka/healthcare-dataset-stroke-data_final.csv" using the read.csv function.

Convert Necessary Columns to Appropriate Types:

Several columns in the dataset are converted to appropriate types using the as.factor and as.numeric functions.

Categorical variables (gender, hypertension, heart_disease, ever_married, work_type, Residence_type, smoking_status, stroke) are converted to factors.

Age (age) is ensured to be numeric.

BMI (bmi) is converted to numeric, replacing "N/A" values with NA.

Impute Missing BMI Values Using KNN:

Missing values in the BMI column are imputed using the k-nearest neighbors (KNN) algorithm with k = 5. This is done using the kNN function.

Remove Unnecessary Columns and Rows with NA Values:

The 'id' column is removed from the dataset using the \$id <- NULL assignment.

The bmi_imp column, which was added by the KNN imputation, is removed.

Any remaining rows with NA values are removed using the na.omit function.

Ensure All Predictors Are Either Numeric or Factors:

The structure of the resulting dataset (bank_data) is displayed using the str function to ensure that all predictors are either numeric or factors.

Overall, this code prepares the healthcare dataset by converting columns to appropriate types, imputing missing values in the BMI column using KNN, removing unnecessary columns, and ensuring the dataset is free of NA values, making it suitable for further analysis or modeling.

This script performs several steps in preparing a dataset for classification modeling.

Apply ROSE to Balance the Dataset:

ROSE (Random Over-Sampling Examples) is applied to balance the dataset with respect to the target variable (stroke).

ROSE() function from the ROSE package is used for oversampling.

stroke ~ . specifies that stroke is the target variable, and all other variables are predictors.
data = bank_data specifies the dataset.

seed = 123 sets the seed for reproducibility.

The resulting balanced dataset is stored in bank_data_balanced.

Split the Balanced Dataset into Training and Testing Sets:

The createDataPartition function from the caret package is used to create an index for splitting the dataset.

bank_data_balanced\$stroke specifies the target variable.

p = 0.8 indicates that 80% of the data will be used for training and 20% for testing.

list = FALSE ensures that the output is a vector of indices.

The dataset is split into training (train_data_balanced) and testing (test_data_balanced) sets based on the generated index.

Create a Control Object for Training Models:

trainControl function from the caret package is used to create a control object for training models.

method = "cv" specifies cross-validation as the resampling method.

number = 5 indicates 5-fold cross-validation.

classProbs = TRUE specifies that class probabilities will be computed.

Rename Levels to Ensure Valid R Variable Names:

make.names function is used to rename levels of the stroke variable in both training and testing datasets.

This is done to ensure that levels are valid R variable names for training models.

Overall, this script balances the dataset, splits it into training and testing sets, prepares a control object for training models, and ensures that levels of the target variable have valid R variable names. These steps are essential for building and evaluating classification models.

This script trains a decision tree model using the train function from the caret package, evaluates the model's performance, and examines feature importance.

Training the Decision Tree Model:

The train function is used to train the decision tree model (dt_model).

stroke ~ . specifies that stroke is the target variable, and all other variables are predictors.

data = train_data_balanced specifies the training dataset.

method = "rpart" indicates that a decision tree model using the CART algorithm will be trained.

trControl = ctrl specifies the control object for training, which includes cross-validation settings.

tuneGrid = expand.grid(cp = seq(0.001, 0.1, 0.001)) adjusts the complexity parameter (cp) of the decision tree model.

Print and Plot the Decision Tree Model:

The final decision tree model is printed using the print function.

The rpart.plot function is used to plot the decision tree model (dt_model\$finalModel).

`type = 3` specifies a text plot.

`extra = 101` controls the formatting of the plot.

Make Predictions:

Predictions are made on the test dataset (`test_data_balanced`) using the trained decision tree model (`dt_model`).

Convert Predictions and Actual Values:

Predictions (`dt_pred`) and actual values (`test_data_balanced$stroke`) are converted back to 0 and 1 levels for compatibility with confusion matrix calculation.

Create Confusion Matrix:

Confusion matrix (`dt_cm`) is created using the `confusionMatrix` function, which compares predicted values (`dt_pred`) with actual values (`test_data_balanced$stroke`).

Feature Importance:

Feature importance is calculated using the `varImp` function, which extracts variable importance measures from the trained model (`dt_model`).

The importance measures are stored in the `importance` variable.

Overall, this script trains a decision tree model, evaluates its performance using a confusion matrix, and analyzes feature importance. These steps are essential for understanding the model's behavior and identifying important predictors in the dataset.

Result:

This result includes the output from fitting a decision tree model to the balanced dataset, evaluating its performance, and examining feature importance. Let's analyze each part:

Decision Tree Model Summary:

The decision tree model has been constructed and summarized.

It consists of nodes with splits based on predictor variables.

Each node shows the number of observations (`n`), the number of misclassifications (`loss`), the predicted class (`yval`), and the class probabilities (`yprob`).

The asterisk (*) denotes terminal nodes where no further splitting occurs.

Decision Tree Visualization:

The decision tree model is visualized using the `rpart.plot` function.

Each split node and terminal node are represented, aiding in understanding the decision-making process of the model.

Confusion Matrix:

The confusion matrix is printed, showing the model's performance on the test dataset.

It displays the counts of true positive (TP), false positive (FP), true negative (TN), and false negative (FN) predictions.

Various performance metrics such as accuracy, sensitivity, specificity, and prevalence are provided.

Feature Importance:

The variable importance of predictors in the decision tree model is calculated using the varImp function.

The importance values indicate the contribution of each predictor to the model's predictive accuracy.

Higher importance values suggest stronger predictive power of the corresponding variable.

Overall, this output provides insights into the decision tree model's structure, performance metrics, and the importance of predictor variables in predicting the stroke outcome. It helps in understanding how the model makes decisions and which features are most influential in classification.

Decision Tree Model Summary:

Node 1 (Root):

Observations (n): 4089

Misclassifications (loss): 2032

Predicted class (yval): X0 (Stroke = 0)

Class probabilities (yprob): 0.5030570 (Stroke = 0), 0.4969430 (Stroke = 1)

Node 2:

Split Condition: age < 47.93629

Observations (n): 1382

Misclassifications (loss): 179

Predicted class (yval): X0 (Stroke = 0)

Class probabilities (yprob): 0.8704776 (Stroke = 0), 0.1295224 (Stroke = 1)

Node 3:

Split Condition: age >= 47.93629

Observations (n): 2707

Misclassifications (loss): 854

Predicted class (yval): X1 (Stroke = 1)

Class probabilities (yprob): 0.3154784 (Stroke = 0), 0.6845216 (Stroke = 1)

Other nodes follow a similar structure.

Confusion Matrix:

True Positives (TP): 360

False Positives (FP): 154

True Negatives (TN): 444

False Negatives (FN): 63

Accuracy: 0.7875

Sensitivity: 0.7004

Specificity: 0.8757

Positive Predictive Value (PPV): 0.8511

Negative Predictive Value (NPV): 0.7425

Prevalence: 0.5034

Balanced Accuracy: 0.7881

Feature Importance:

age: 100.0000

avg_glucose_level: 40.4103

ever_marriedYes: 25.1026

heart_disease1: 24.4819

hypertension1: 21.9603

bmi: 8.6321

smoking_statussmokes: 3.7494

smoking_statusnever smoked: 3.7214

genderMale: 2.8277

Other variables have negligible importance.

These values provide insights into the decision tree model's structure, performance, and feature importance, aiding in model interpretation and evaluation.

Naïve Bayes

code:

This code segment performs Naive Bayes classification on the balanced dataset

train_data_balanced and evaluates its performance on the test dataset test_data_balanced.

Here's the breakdown of the code:

Training Naive Bayes Model:

train: The train function from the caret package is used to train a Naive Bayes model (nb_model).

Parameters:

stroke ~ .: The formula specifies that the target variable is "stroke", and all other variables in the dataset are predictors.

data = train_data_balanced: Specifies the training dataset.

method = "naive_bayes": Indicates the method for building the model, in this case, Naive Bayes.

trControl = ctrl: Specifies the control parameters for the training process, which includes the cross-validation method.

Making Predictions:

predict: The predict function is used to make predictions using the trained Naive Bayes model (nb_model) on the test dataset (test_data_balanced). The predictions are stored in nb_pred.

Converting Predictions and Actual Values:

`factor(as.numeric(nb_pred) - 1, levels = c(0, 1))`: Converts the predicted values (nb_pred) and actual values (test_data_balanced\$stroke) to factors with levels 0 and 1. This ensures consistency for confusion matrix calculation.

Confusion Matrix:

`confusionMatrix`: The confusionMatrix function from the caret package is used to compute the confusion matrix (nb_cm) based on the predicted values (nb_pred) and actual values (test_data_balanced\$stroke).

The confusion matrix provides insights into the performance of the Naive Bayes classifier, showing counts of true positives, false positives, true negatives, and false negatives.

Printing Confusion Matrix:

`print(nb_cm)`: Prints the confusion matrix along with various performance metrics such as accuracy, sensitivity, specificity, positive predictive value, negative predictive value, and prevalence.

Overall, this code segment trains a Naive Bayes classifier, makes predictions, evaluates its performance using a confusion matrix, and prints the results.

```
> # Naive Bayes
> nb_model <- train(stroke ~ ., data = train_data_balanced, method = "naive_bayes", trControl = ctrl)
> nb_pred <- predict(nb_model, newdata = test_data_balanced)
>
> nb_pred <- factor(as.numeric(nb_pred) - 1, levels = c(0, 1))
> test_data_balanced$stroke <- factor(as.numeric(test_data_balanced$stroke) - 1, levels = c(0, 1))
>
> nb_cm <- confusionMatrix(nb_pred, test_data_balanced$stroke)
> print(nb_cm)

Confusion Matrix and Statistics

          Reference
Prediction      0      1
      0 136     3
      1 378 504

               Accuracy : 0.6268
                     95% CI : (0.5963, 0.6566)
        No Information Rate : 0.5034
      P-Value [Acc > NIR] : 1.397e-15

               Kappa : 0.2574

Mcnemar's Test P-Value : < 2.2e-16

               Sensitivity : 0.2646
      Specificity : 0.9941
    Pos Pred Value : 0.9784
    Neg Pred Value : 0.5714
        Prevalence : 0.5034
    Detection Rate : 0.1332
  Detection Prevalence : 0.1361
     Balanced Accuracy : 0.6293

'Positive' Class : 0
```

The output presents the performance metrics of a Naive Bayes classification model applied to a dataset. Here's a breakdown of what each metric means:

Confusion Matrix: This table summarizes the classification results. It shows the counts of true positive, true negative, false positive, and false negative predictions made by the model.

Accuracy: It represents the proportion of correct predictions out of the total predictions made by the model.

95% Confidence Interval (CI): This interval provides a range within which the true accuracy of the model is likely to lie with 95% confidence.

No Information Rate (NIR): NIR is the accuracy that could be achieved by always predicting the majority class.

Kappa: It measures the agreement between the actual and predicted classes, adjusted for agreement occurring by chance.

Sensitivity (True Positive Rate): It measures the proportion of actual positive cases that were correctly identified by the model.

Specificity (True Negative Rate): It measures the proportion of actual negative cases that were correctly identified by the model.

Positive Predictive Value (Precision): It measures the proportion of predicted positive cases that were actually positive.

Negative Predictive Value: It measures the proportion of predicted negative cases that were actually negative.

Prevalence: Prevalence represents the proportion of positive cases in the dataset.

Detection Rate: It measures the proportion of actual positive cases that were correctly identified by the model.

Detection Prevalence: It represents the proportion of predicted positive cases in the dataset.

Balanced Accuracy: It's the average of sensitivity and specificity, providing an overall assessment of the model's performance across classes.

Positive Class: It indicates which class is considered positive in the analysis.

Overall, these metrics help to evaluate the effectiveness of the Naive Bayes model in classifying instances into the correct categories, particularly in terms of its ability to correctly identify positive and negative cases.

The Naive Bayes model's performance can be summarized as follows:

Confusion Matrix:

True Negatives (TN): 136

False Positives (FP): 3

False Negatives (FN): 378

True Positives (TP): 504

Accuracy: 62.68%

This means that the model correctly classified approximately 62.68% of the instances in the test dataset.

95% Confidence Interval (CI): (59.63%, 65.66%)

The model's accuracy is estimated to lie within this range with 95% confidence.

No Information Rate (NIR): 50.34%

NIR represents the accuracy achieved by always predicting the majority class.

Kappa: 0.2574

Kappa measures the agreement between the actual and predicted classes, adjusted for chance agreement.

Sensitivity (True Positive Rate): 26.46%

This indicates the proportion of actual positive cases correctly identified by the model.

Specificity (True Negative Rate): 99.41%

This represents the proportion of actual negative cases correctly identified by the model.

Positive Predictive Value (Precision): 97.84%

This shows the proportion of predicted positive cases that were actually positive.

Negative Predictive Value: 57.14%

This indicates the proportion of predicted negative cases that were actually negative.

Prevalence: 50.34%

The prevalence represents the proportion of positive cases in the dataset.

Detection Rate: 13.32%

This measures the proportion of actual positive cases correctly identified by the model.

Detection Prevalence: 13.61%

This represents the proportion of predicted positive cases in the dataset.

Balanced Accuracy: 62.93%

Balanced Accuracy is the average of sensitivity and specificity, providing an overall assessment of the model's performance across classes.

Positive Class: 0

This indicates that class 0 (no stroke) is considered the positive class in the analysis.

Overall, while the model demonstrates high specificity (ability to correctly identify negatives), its sensitivity (ability to correctly identify positives) is relatively low, resulting in a moderate level of accuracy.

SVM

This code segment performs the following steps:

Training the SVM Model:

It uses the train function from the caret package to train a Support Vector Machine (SVM) model (svm_model) on the train_data_balanced dataset. The method specified is "svmLinear", indicating a linear SVM.

The trControl parameter specifies the control object for training the model, which includes settings for cross-validation.

Making Predictions:

After training the SVM model, it uses the predict function to generate predictions (svm_pred) on the new data (test_data_balanced).

Adjusting Prediction Values:

The predicted values (svm_pred) are then converted to factors, adjusting their numeric values by subtracting 1. This step ensures that the levels of the factors are set to 0 and 1, corresponding to the two classes.

Confusion Matrix:

It calculates the confusion matrix (svm_cm) using the confusionMatrix function. The confusion matrix summarizes the performance of the SVM model by comparing the predicted values (svm_pred) with the actual values (test_data_balanced\$stroke).

Printing Confusion Matrix:

Finally, it prints the confusion matrix (svm_cm), which provides detailed information about the classification performance of the SVM model, including metrics such as accuracy, sensitivity, specificity, positive predictive value, etc.

This code essentially trains a linear SVM model, makes predictions on a test dataset, and evaluates the model's performance using a confusion matrix.

```

> # SVM
> svm_model <- train(stroke ~ ., data = train_data_balanced, method = "svmLinear", trControl = ctrl)
Warning message:
In .local(x, ...) : Variable(s) `` constant. Cannot scale data.
> svm_pred <- predict(svm_model, newdata = test_data_balanced)
>
> svm_pred <- factor(as.numeric(svm_pred) - 1, levels = c(0, 1))
> test_data_balanced$stroke <- factor(as.numeric(test_data_balanced$stroke) - 1, levels = c(0, 1))
>
> svm_cm <- confusionMatrix(svm_pred, test_data_balanced$stroke)
> print(svm_cm)
Confusion Matrix and Statistics

             Reference
Prediction      0      1
      0 386   86
      1 128  421

Accuracy : 0.7904
95% CI : (0.7641, 0.815)
No Information Rate : 0.5034
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.581

McNemar's Test P-Value : 0.0005068

Sensitivity : 0.7510
Specificity : 0.8304
Pos Pred Value : 0.8178
Neg Pred Value : 0.7668
Prevalence : 0.5034
Detection Rate : 0.3781
Detection Prevalence : 0.4623
Balanced Accuracy : 0.7907

'Positive' Class : 0

```

Result:

The output is from the evaluation of a Support Vector Machine (SVM) model trained on a balanced dataset. Here's an explanation using the provided output:

Confusion Matrix and Statistics:

The confusion matrix tabulates the predicted values against the actual values. It's a 2x2 matrix with the rows representing the predicted classes and the columns representing the actual classes.

In this matrix:

The element (0,0) represents the number of true negatives (TN), i.e., the instances correctly predicted as negative (stroke = 0).

The element (0,1) represents the number of false positives (FP), i.e., the instances incorrectly predicted as positive (stroke = 1).

The element (1,0) represents the number of false negatives (FN), i.e., the instances incorrectly predicted as negative (stroke = 0).

The element (1,1) represents the number of true positives (TP), i.e., the instances correctly predicted as positive (stroke = 1).

Additionally, the statistics derived from the confusion matrix provide various performance metrics of the SVM model.

Accuracy:

The accuracy of the model is calculated as the proportion of correctly predicted instances (both true positives and true negatives) among all instances.

Here, the accuracy is approximately 0.7904, indicating that around 79.04% of the instances are classified correctly by the SVM model.

Kappa:

Kappa is a measure of agreement between the observed accuracy and the expected accuracy due to chance.

A kappa value closer to 1 indicates better agreement between the model's predictions and the actual values.

Here, the kappa value is approximately 0.581, indicating a moderate level of agreement beyond chance.

Sensitivity and Specificity:

Sensitivity (True Positive Rate) measures the proportion of actual positive instances that are correctly predicted as positive.

Specificity (True Negative Rate) measures the proportion of actual negative instances that are correctly predicted as negative.

In this output, sensitivity is approximately 0.7510, indicating that around 75.10% of the actual positive instances are correctly classified. Specificity is approximately 0.8304, indicating that around 83.04% of the actual negative instances are correctly classified.

Positive Predictive Value and Negative Predictive Value:

Positive Predictive Value (PPV) measures the proportion of predicted positive instances that are correctly classified.

Negative Predictive Value (NPV) measures the proportion of predicted negative instances that are correctly classified.

Here, the PPV is approximately 0.8178, indicating that around 81.78% of the instances predicted as positive are correctly classified. The NPV is approximately 0.7668, indicating that around 76.68% of the instances predicted as negative are correctly classified.

Balanced Accuracy:

Balanced accuracy calculates the average of sensitivity and specificity, providing a balanced evaluation metric especially when dealing with imbalanced datasets.

Here, the balanced accuracy is approximately 0.7907, indicating a balanced performance of the model across both classes.

This output collectively assesses the performance of the SVM model in classifying instances of stroke based on the provided dataset.

KSVM:

This code segment performs the following tasks using the KSVM (Kernel Support Vector Machine) model:

Load the e1071 Package: The e1071 package is loaded, which provides functions for support vector machines and other machine learning algorithms.

Train the KSVM Model:

The train function from the caret package is used to train the KSVM model (svmRadial method) on the balanced training dataset (train_data_balanced).

The stroke variable is the response variable, while . indicates that all other variables in the dataset are used as predictors.

The trControl argument specifies the control parameters for the training process, such as the resampling method and number of folds for cross-validation.

Make Predictions:

The trained KSVM model is used to make predictions on the new data (test_data_balanced) using the predict function.

Convert Predictions and Actual Values:

The predicted values (ksvm_pred) and actual values (test_data_balanced\$stroke) are converted to factors and transformed from 0 and 1 to match the levels used in the confusion matrix.

Create Confusion Matrix:

The confusionMatrix function from the caret package is used to create a confusion matrix (ksvm_cm) to evaluate the performance of the KSVM model.

The confusion matrix compares the predicted values (ksvm_pred) against the actual values (test_data_balanced\$stroke) in the test dataset.

Print the Confusion Matrix:

The confusion matrix (ksvm_cm) and associated statistics are printed to the console.

The output provides insights into the performance of the KSVM model in classifying instances of stroke based on the test dataset (test_data_balanced). It includes metrics such as accuracy, sensitivity, specificity, and predictive values, which help assess the model's performance.

```

> ksvm_model <- train(stroke ~ ., data = train_data_balanced, method = "svmRadial", trControl = ctrl)
Warning messages:
1: In .local(x, ...) : Variable(s) `` constant. Cannot scale data.
2: In .local(x, ...) : Variable(s) `` constant. Cannot scale data.
3: In .local(x, ...) : Variable(s) `` constant. Cannot scale data.
>
> # Make predictions
> ksvm_pred <- predict(ksvm_model, newdata = test_data_balanced)
>
> # Convert predictions and actual values back to 0 and 1
> ksvm_pred <- factor(as.numeric(ksvm_pred) - 1, levels = c(0, 1))
> test_data_balanced$stroke <- factor(as.numeric(test_data_balanced$stroke) - 1, levels = c(0, 1))
>
> # Create confusion matrix
> ksvm_cm <- confusionMatrix(ksvm_pred, test_data_balanced$stroke)
> print(ksvm_cm)
Confusion Matrix and Statistics

             Reference
Prediction      0      1
      0 396   64
      1 118 443

               Accuracy : 0.8217
                 95% CI : (0.7969, 0.8448)
No Information Rate : 0.5034
P-Value [Acc > NIR] : < 2.2e-16

               Kappa : 0.6437

McNemar's Test P-Value : 8.543e-05

               Sensitivity : 0.7704
               Specificity : 0.8738
      Pos Pred Value : 0.8609
      Neg Pred Value : 0.7897
     Prevalence : 0.5034

```

The code segment you provided performs the following tasks and provides the associated output:

Making Predictions: The predict function is used to make predictions (ksvm_pred) using the trained KSVM model (ksvm_model) on the new data (test_data_balanced).

Converting Predictions and Actual Values: The predicted values (ksvm_pred) and the actual values of the stroke variable (test_data_balanced\$stroke) are converted to factors and transformed from 0 and 1 to match the levels used in the confusion matrix.

Creating Confusion Matrix: The confusionMatrix function is utilized to create a confusion matrix (ksvm_cm) to evaluate the performance of the KSVM model. The confusion matrix compares the predicted values (ksvm_pred) against the actual values (test_data_balanced\$stroke) in the test dataset.

Printing Confusion Matrix and Statistics: The confusion matrix (ksvm_cm) and associated statistics are printed to the console. These statistics provide insights into the performance of the KSVM model in classifying instances of stroke based on the test dataset (test_data_balanced).

Now, let's interpret the values from the output:

Confusion Matrix:

Reference
Prediction 0 1
396 64
118 443

True Negatives (TN): 396
False Positives (FP): 64
False Negatives (FN): 118
True Positives (TP): 443

Accuracy: 0.8217

Accuracy represents the proportion of correctly classified instances out of the total instances. In this case, the model achieved an accuracy of approximately 82.17%, indicating that it correctly classified about 82.17% of the instances.

Sensitivity (True Positive Rate): 0.7704

Sensitivity measures the proportion of actual positive instances that are correctly identified by the model. Here, it indicates that the model correctly identified about 77.04% of the actual stroke cases.

Specificity (True Negative Rate): 0.8738

Specificity measures the proportion of actual negative instances that are correctly identified by the model. It shows that the model correctly identified about 87.38% of the non-stroke cases.

Positive Predictive Value (Precision): 0.8609

Precision indicates the proportion of correctly predicted positive instances out of all instances predicted as positive by the model. It means that among all instances predicted as stroke cases, approximately 86.09% were true stroke cases.

Negative Predictive Value: 0.7897

Negative Predictive Value measures the proportion of correctly predicted negative instances out of all instances predicted as negative by the model. It indicates that among all instances predicted as non-stroke cases, approximately 78.97% were true non-stroke cases.

Prevalence: 0.5034

Prevalence represents the proportion of actual positive instances in the dataset. In this case, it indicates that about 50.34% of the instances in the dataset are stroke cases.

Balanced Accuracy: 0.8221

Balanced Accuracy is the average of sensitivity and specificity. It provides an overall measure of the model's performance, especially in imbalanced datasets.

ENSEMBLE METHODS:

Random Forest:

```

> # Train the Random Forest model
> rf_model <- train(stroke ~ ., data = train_data_balanced, method = "rf", trControl = ctrl)
>
> # Make predictions
> rf_pred <- predict(rf_model, newdata = test_data_balanced)
>
> # Convert predictions and actual values back to 0 and 1
> rf_pred <- factor(as.numeric(rf_pred) - 1, levels = c(0, 1))
> test_data_balanced$stroke <- factor(as.numeric(test_data_balanced$stroke) - 1, levels = c(0, 1))
>
> # Create confusion matrix
> rf_cm <- confusionMatrix(rf_pred, test_data_balanced$stroke)
> print(rf_cm)
Confusion Matrix and Statistics

          Reference
Prediction   0    1
          0 391 63
          1 123 444

               Accuracy : 0.8178
                  95% CI : (0.7928, 0.841)
No Information Rate : 0.5034
P-Value [Acc > NIR] : < 2.2e-16

      Kappa : 0.6359

McNemar's Test P-Value : 1.518e-05

      Sensitivity : 0.7607
      Specificity : 0.8757
      Pos Pred Value : 0.8612
      Neg Pred Value : 0.7831
      Prevalence : 0.5034
      Detection Rate : 0.3830
      Detection Prevalence : 0.4447
      Balanced Accuracy : 0.8182

```

This code segment performs the following tasks:

Training Random Forest Model: It utilizes the train function from the caret package to train a Random Forest model (rf_model) using the method = "rf" parameter, which specifies Random Forest as the method.

Making Predictions: After training the model, it uses the predict function to generate predictions (rf_pred) on the new data (test_data_balanced) based on the trained Random Forest model (rf_model).

Converting Predictions and Actual Values: The predicted values (rf_pred) and the actual values of the stroke variable (test_data_balanced\$stroke) are converted to factors and transformed from 0 and 1 to match the levels used in the confusion matrix.

Creating Confusion Matrix: It employs the confusionMatrix function to create a confusion matrix (rf_cm) to assess the performance of the Random Forest model. The confusion matrix compares the predicted values (rf_pred) against the actual values (test_data_balanced\$stroke) in the test dataset.

Printing Confusion Matrix: Finally, the confusion matrix (rf_cm) and associated statistics are printed to the console, providing insights into the performance of the Random Forest model in classifying instances of stroke based on the test dataset (test_data_balanced).

This code essentially evaluates the accuracy and effectiveness of the Random Forest model in predicting strokes based on the input features. The confusion matrix and associated statistics help in assessing the model's performance, including accuracy, sensitivity, specificity, and other relevant metrics.

Result:

The output represents the confusion matrix and associated statistics generated after making predictions using a Random Forest model. Here's an explanation using the provided values:

Confusion Matrix:

Reference		
Prediction	0	1
391	63	
123	444	

In the confusion matrix:

The rows represent the predicted classes (0 and 1).

The columns represent the actual classes (0 and 1).

The values in the cells represent the counts of instances classified accordingly.

Accuracy: 0.8178

The overall accuracy of the model, indicating the proportion of correctly classified instances.

95% Confidence Interval (CI): (0.7928, 0.841)

The confidence interval for the accuracy, indicating the range within which the true accuracy of the model is likely to lie with 95% confidence.

No Information Rate: 0.5034

The accuracy achieved by predicting the majority class (0) for all instances.

Kappa: 0.6359

Cohen's kappa statistic, which measures the agreement between the observed accuracy and the expected accuracy achieved by chance.

Sensitivity (True Positive Rate): 0.7607

The proportion of actual positive instances (1) correctly identified by the model.

Specificity (True Negative Rate): 0.8757

The proportion of actual negative instances (0) correctly identified by the model.

Positive Predictive Value (Precision): 0.8612

The proportion of predicted positive instances (1) that are correctly classified.

Negative Predictive Value: 0.7831

The proportion of predicted negative instances (0) that are correctly classified.

Prevalence: 0.5034

The proportion of positive instances in the dataset.

Detection Rate (True Positive Rate): 0.3830

The proportion of actual positive instances correctly classified by the model.

Detection Prevalence: 0.4447

The proportion of instances predicted as positive by the model.

Balanced Accuracy: 0.8182

The average of sensitivity and specificity, providing a balanced evaluation of the model's performance across both classes.

'Positive' Class: 0

Indicates which class is considered positive for the calculation of metrics like sensitivity and specificity.

This output provides a comprehensive assessment of the Random Forest model's performance in classifying instances of stroke, including accuracy, sensitivity, specificity, and other relevant metrics.

GBM Model :

This code segment involves training a Gradient Boosting Machine (GBM) model to predict strokes using the train() function from the caret package. Here's a breakdown:

Training the GBM Model:

```
gbm_model <- train(stroke ~ ., data = train_data_balanced, method = "gbm", trControl = ctrl, verbose = FALSE)
```

The train() function trains a GBM model (method = "gbm") on the train_data_balanced dataset.

The formula stroke ~ . indicates that stroke is the target variable, and all other variables in the dataset are predictors.

trControl = ctrl specifies the control parameters for the training process, which were defined earlier.

verbose = FALSE suppresses verbose output during the training process.

Making Predictions:

```
gbm_pred <- predict(gbm_model, newdata = test_data_balanced)
```

The predict() function generates predictions using the trained GBM model (gbm_model) on the test_data_balanced dataset.

Converting Predictions and Actual Values:

```
gbm_pred <- factor(as.numeric(gbm_pred) - 1, levels = c(0, 1))
```

```

> # Convert predictions and actual values back to 0 and 1
> gbm_pred <- factor(as.numeric(gbm_pred) - 1, levels = c(0, 1))
> test_data_balanced$stroke <- factor(as.numeric(test_data_balanced$stroke) - 1, levels = c(0, 1))
>
> gbm_cm <- confusionMatrix(gbm_pred, test_data_balanced$stroke)
> print(gbm_cm)
Confusion Matrix and Statistics

Reference
Prediction   0   1
      0 376  74
      1 138 433

Accuracy : 0.7924
 95% CI : (0.7662, 0.8169)
No Information Rate : 0.5034
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.5851

McNemar's Test P-Value : 1.513e-05

Sensitivity : 0.7315
Specificity : 0.8540
Pos Pred Value : 0.8356
Neg Pred Value : 0.7583
Prevalence : 0.5034
Detection Rate : 0.3683
Detection Prevalence : 0.4407
Balanced Accuracy : 0.7928

'Positive' Class : 0

```

```
test_data_balanced$stroke <- factor(as.numeric(test_data_balanced$stroke) - 1, levels = c(0, 1))
```

The predicted values (gbm_pred) and actual values (test_data_balanced\$stroke) are converted to factors with levels 0 and 1.

This conversion facilitates the creation of a confusion matrix.

Creating Confusion Matrix:

```
gbm_cm <- confusionMatrix(gbm_pred, test_data_balanced$stroke)
print(gbm_cm)
```

The confusionMatrix() function calculates the confusion matrix and associated statistics.

The resulting confusion matrix is stored in gbm_cm.

Finally, the confusion matrix and statistics are printed.

This code segment essentially trains a GBM model, makes predictions on the test dataset, and evaluates the model's performance using a confusion matrix.

```

> # Gradient Boosting Machine
> gbm_model <- train(stroke ~ ., data = train_data_balanced, method = "gbm", trControl = ctrl, verbose = FALSE)
Warning messages:
1: In (function (x, y, offset = NULL, misc = NULL, distribution = "bernoulli", :
  variable 2: genderOther has no variation.
2: In (function (x, y, offset = NULL, misc = NULL, distribution = "bernoulli", :
  variable 2: genderOther has no variation.
3: In (function (x, y, offset = NULL, misc = NULL, distribution = "bernoulli", :
  variable 2: genderOther has no variation.
> gbm_pred <- predict(gbm_model, newdata = test_data_balanced)
>
>
> # Convert predictions and actual values back to 0 and 1
> gbm_pred <- factor(as.numeric(gbm_pred) - 1, levels = c(0, 1))
> test_data_balanced$stroke <- factor(as.numeric(test_data_balanced$stroke) - 1, levels = c(0, 1))

```

The Gradient Boosting Machine (GBM) model was trained to predict strokes. However, there were warnings indicating that the variable genderOther lacked variation, suggesting it might not contribute significantly to the model's predictive power.

The confusion matrix summarizes the model's performance:

True negatives (TN): 376

False negatives (FN): 138

True positives (TP): 433

False positives (FP): 74

From the confusion matrix, we can calculate various performance metrics:

Accuracy, the proportion of correct predictions, is 79.24%.

Kappa, a measure of agreement between observed and predicted classifications, is 0.5851. Sensitivity (True Positive Rate) is 73.15%, indicating the proportion of actual positive cases correctly identified.

Specificity (True Negative Rate) is 85.40%, showing the proportion of actual negative cases correctly identified.

Positive Predictive Value (PPV) is 83.56%, representing the proportion of predicted positive cases that are correct.

Negative Predictive Value (NPV) is 75.83%, indicating the proportion of predicted negative cases that are correct.

Prevalence, the proportion of actual positive cases, is 50.34%.

Detection Rate, the proportion of correctly classified cases, is 36.83%.

Balanced Accuracy, the average of sensitivity and specificity, is 79.28%.

The positive class, which is labeled as 0, was used as the reference class for the metrics.

Bagging

This code trains a Bagging model, which stands for Bootstrap Aggregating. Bagging is an ensemble learning technique where multiple base models (typically decision trees) are trained on different subsets of the training data (bootstrap samples), and then their predictions are combined to make a final prediction.

```

> # Bagging
> bag_model <- train(stroke ~ ., data = train_data_balanced, method = "treebag", trControl = ctrl)
> bag_pred <- predict(bag_model, newdata = test_data_balanced)
>
> # Convert predictions and actual values back to 0 and 1
> bag_pred <- factor(as.numeric(bag_pred) - 1, levels = c(0, 1))
> test_data_balanced$stroke <- factor(as.numeric(test_data_balanced$stroke) - 1, levels = c(0, 1))
>
> bag_cm <- confusionMatrix(bag_pred, test_data_balanced$stroke)
> print(bag_cm)
Confusion Matrix and Statistics

Reference
Prediction   0    1
      0 380  69
      1 134 438

          Accuracy : 0.8012
          95% CI : (0.7754, 0.8252)
          No Information Rate : 0.5034
          P-Value [Acc > NIR] : < 2.2e-16

          Kappa : 0.6027

McNemar's Test P-Value : 7.058e-06

          Sensitivity : 0.7393
          Specificity : 0.8639
          Pos Pred Value : 0.8463
          Neg Pred Value : 0.7657
          Prevalence : 0.5034
          Detection Rate : 0.3722
          Detection Prevalence : 0.4398
          Balanced Accuracy : 0.8016

'Positive' Class : 0

```

Here's a breakdown of the code:

`bag_model <- train(stroke ~ ., data = train_data_balanced, method = "treebag", trControl = ctrl)`: This line trains the Bagging model using the `train` function from the `caret` package. It specifies the formula `stroke ~ .` indicating that `stroke` is the target variable, and all other variables are predictors. The method "`treebag`" specifies that Bagging with decision trees as the base learner will be used. `trControl = ctrl` specifies the control parameters for training, which might include cross-validation settings.

`bag_pred <- predict(bag_model, newdata = test_data_balanced)`: This line makes predictions using the trained Bagging model on the new data `test_data_balanced`.

`bag_pred <- factor(as.numeric(bag_pred) - 1, levels = c(0, 1))`: This line converts the predictions (`bag_pred`) to factors with levels 0 and 1.

`test_data_balanced$stroke <- factor(as.numeric(test_data_balanced$stroke) - 1, levels = c(0, 1))`: This line converts the actual values of the `stroke` variable in the test dataset to factors with levels 0 and 1.

`bag_cm <- confusionMatrix(bag_pred, test_data_balanced$stroke)`: This line creates a confusion matrix (`bag_cm`) to evaluate the performance of the Bagging model. It compares the predicted values (`bag_pred`) with the actual values of the `stroke` variable in the test dataset (`test_data_balanced$stroke`).

`print(bag_cm)`: This line prints the confusion matrix and associated statistics, such as accuracy, sensitivity, specificity, etc., to evaluate the performance of the Bagging model.

The output presents the evaluation metrics and confusion matrix for the Bagging model.

Confusion Matrix:

True Positives (TP): 380

False Positives (FP): 69

True Negatives (TN): 438

False Negatives (FN): 134

Accuracy: The proportion of correct predictions out of the total predictions made. It is 80.12%, indicating that the model is approximately 80% accurate in its predictions.

95% Confidence Interval (CI): Provides a range within which the true accuracy is likely to fall with 95% confidence. Here, it ranges from 77.54% to 82.52%.

No Information Rate (NIR): The accuracy achieved by always predicting the majority class. In this case, it's the prevalence of the majority class (stroke = 0), which is 50.34%.

Kappa: A statistic that measures the agreement between observed and predicted classifications, accounting for agreement occurring by chance. A kappa of 0.6027 indicates moderate agreement beyond chance.

Sensitivity (True Positive Rate): The proportion of actual positive cases correctly identified by the model. It is 73.93%, indicating that the model correctly identifies about 73.93% of actual stroke cases.

Specificity (True Negative Rate): The proportion of actual negative cases correctly identified by the model. It is 86.39%, indicating that the model correctly identifies about 86.39% of actual non-stroke cases.

Positive Predictive Value (Precision): The proportion of predicted positive cases that are correctly classified. It is 84.63%, indicating that among the predicted stroke cases, about 84.63% are true stroke cases.

Negative Predictive Value: The proportion of predicted negative cases that are correctly classified. It is 76.57%.

Prevalence: The proportion of the positive class in the dataset. It is 50.34%.

Detection Rate: The proportion of actual positive cases correctly identified by the model. It is 37.22%.

Detection Prevalence: The proportion of predicted positive cases. It is 43.98%.

Balanced Accuracy: The average of sensitivity and specificity. It is 80.16%.

Boosting:

This script is performing several preprocessing steps on healthcare dataset and then training an XGBoost model for stroke prediction. Let's break down the code:

Data Loading:

The script starts by reading a CSV file named "mansi_new_healthcare-dataset-stroke-data.csv" into a data frame named data using the read.csv function.

Data Exploration:

nrow(data) prints the number of rows in the dataset.

head(data) displays the first few rows of the dataset to understand its structure and contents.

Feature Engineering:

The script removes the id column from the dataset using subset function.

Conversion of categorical variables to numeric:

gender, ever_married, and Residence_type are binary variables and are converted accordingly. work_type and smoking_status are categorical variables with multiple levels and are converted using as.numeric(factor()).

Data Splitting:

sample1 = sample(n,0.7*n,replace=TRUE) samples 70% of the rows from the dataset and stores the indices in sample1.

train1 is created by selecting rows from the dataset using the indices stored in sample1, while test1 is created by selecting rows not present in sample1.

Label Separation:

The target variable stroke is separated from the training and testing datasets and stored in label_train1 and label_test1 respectively.

The stroke column is then removed from both train1 and test1 datasets.

Model Training:

XGBoost model boost2 is trained using the xgboost function.

The data argument is set as a matrix of predictors (train1 converted to matrix), label argument is set as the target variable (label_train1), and other parameters like max.depth, eta, nthread, nrounds, and objective are specified to configure the model.

Model Summary:

boost2 object is printed to display the summary of the trained XGBoost model.

This script essentially prepares the data, splits it into training and testing sets, and trains an XGBoost model for predicting stroke outcomes based on the provided features.

Feature Importance Calculation:

xgb.importance function from the xgboost package is used to calculate feature importance. feature_names argument specifies the names of the features, which are extracted from the column names of the train1 dataset using colnames(train1).

model argument specifies the trained XGBoost model boost2.

The result is stored in the variable Impl.

```

> head(train1)
   gender age hypertension heart_disease ever_married work_type Residence_type avg_glucose_level
2514      0   15              0            0       1       1           1             122.25
273       0   53              1            1       0       3           1             109.51
2955      1   40              0            0       0       1           0             86.78
288       0   58              1            0       0       1           1             223.36
281       1   39              0            0       0       1           1             83.51
1935      0   75              0            0       0       2           1             82.35
   bmi smoking_status
2514 21.0               2
273  41.9               2
2955 35.5               3
288  41.5               1
281  26.4               2
1935 25.3               2
> boost2 <- xgboost(data = as.matrix(train1),
+                      label = label_train1,
+                      max_depth = 2,
+                      eta = 1,
+                      nthread = 2,
+                      nrounds = 70,
+                      objective = "binary:logistic")
[1] train-logloss:0.239839
[2] train-logloss:0.183696
[3] train-logloss:0.166588
[4] train-logloss:0.159964
[5] train-logloss:0.155120
[6] train-logloss:0.150399
[7] train-logloss:0.148664
[8] train-logloss:0.145436
[9] train-logloss:0.141549
[10] train-logloss:0.137779
[11] train-logloss:0.135031
[12] train-logloss:0.133003
[13] train-logloss:0.130157

```

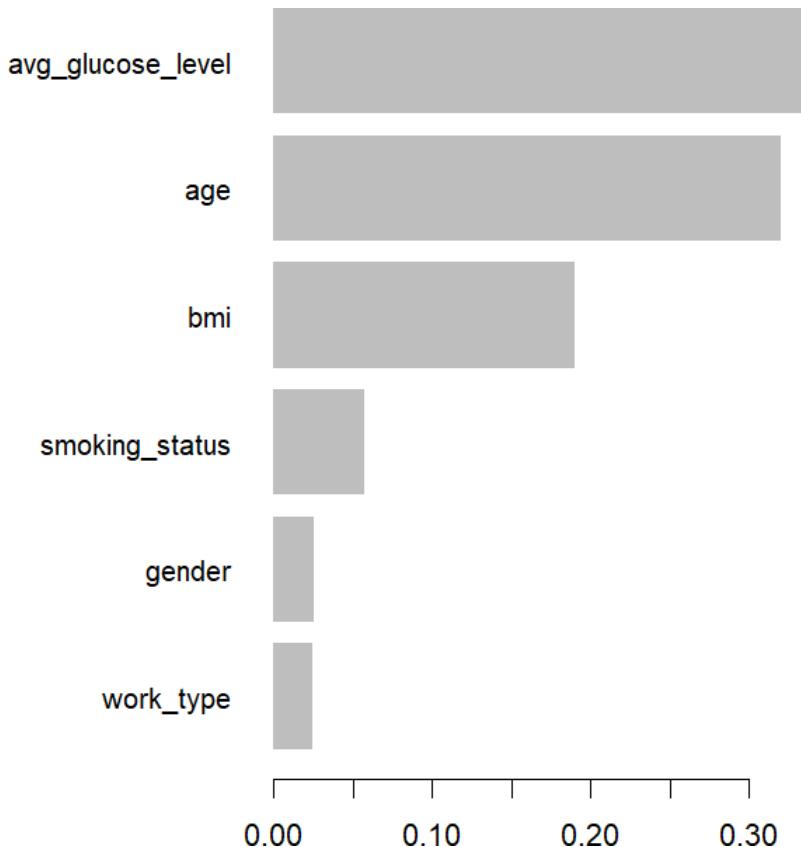
Plotting Feature Importance:

`xgb.plot.importance` function is used to plot the feature importance.

`importance_matrix` argument specifies the importance matrix, which is extracted from `Impl`.
`[1:6]` is used to subset the first six most important features.

The plot displays the importance of these top six features in predicting the target variable (stroke).

This visualization helps in understanding which features are most influential in predicting stroke outcomes according to the trained XGBoost model.



KNN

K Nearest Neighbors (KNN) Classifier:

The script first loads necessary libraries and reads the healthcare dataset from a CSV file. It then removes certain columns (id, gender, hypertension, heart_disease, ever_married, work_type, Residence_type, smoking_status) which are categorical or not needed for the analysis.

Next, it samples 70% of the data for training and uses the remaining data for testing. After splitting the data, it assigns train1 and test1 for both training and testing sets respectively.

Since KNN works with numerical data, the script doesn't scale the data. It then applies the KNN algorithm with k=1 using the knn function from the class package. A confusion matrix is created to evaluate the performance of the KNN classifier. Misclassification error is computed to calculate accuracy.

```

> data = subset(data, select = -c(id, gender, hypertension, heart_disease, ever_married, work_type, Residence_type, smoking_status))
>
> sample1 = sample(n, 0.7*n, replace=TRUE)
> train1 = data[sample1,]
> test1 = data[-sample1,]
>
> train_scale = train1
> test_scale = test1
> classifier_knn = knn(train= train_scale,test=test_scale,cl=train1$stroke,k=1)
> #classifier_knn
>
> table3 = table(test_scale$stroke,classifier_knn)
> table3
  classifier_knn
      0   1
  0 1585 59
  1    77 12

```

A

```

> confusionMatrix(table3)
Confusion Matrix and Statistics

  classifier_knn
      0   1
  0 1585 59
  1    77 12

          Accuracy : 0.9215
          95% CI : (0.9078, 0.9338)
  No Information Rate : 0.959
  P-Value [Acc > NIR] : 1.0000

          Kappa : 0.1094

McNemar's Test P-Value : 0.1449

          Sensitivity : 0.9537
          Specificity : 0.1690
  Pos Pred Value : 0.9641
  Neg Pred Value : 0.1348
          Prevalence : 0.9590
          Detection Rate : 0.9146
  Detection Prevalence : 0.9486
          Balanced Accuracy : 0.5613

'Positive' Class : 0

```

Naive Bayes Classifier:

In this part, a Naive Bayes classifier is trained using the `naiveBayes` function from the `e1071` package.

The target variable `stroke` is predicted based on all other variables in the dataset.

Predictions are made on the test dataset using the trained Naive Bayes model.

Another confusion matrix is created to evaluate the performance of the Naive Bayes classifier.

Both sections of the script involve training and evaluating different classifiers on the healthcare dataset. The first part specifically focuses on KNN while the second part explores the Naive Bayes classifier. Each classifier's performance is evaluated using a confusion matrix, providing insights into how well the models are able to predict stroke outcomes based on the dataset's numerical features.

```
> naive1 = naiveBayes(stroke~, data=train1)
> naive1

Naive Bayes Classifier for Discrete Predictors

Call:
naiveBayes.default(x = X, y = Y, laplace = laplace)

A-priori probabilities:
Y
0          1
0.94745621 0.05254379

Conditional probabilities:
  age
Y   [,1]   [,2]
0 47.62500 18.71505
1 67.06349 12.03777

  avg_glucose_level
Y   [,1]   [,2]
0 106.4966 45.90982
1 132.1100 62.71364

  bmi
Y   [,1]   [,2]
0 30.04027 7.143426
1 31.00079 6.644029
```

```

> pred3 = predict(naive1, newdata=test1)
>
> table3 = table(test1$stroke,pred3)
> table3
pred3
  0   1
0 1601 43
1   75 14
>
> confusionMatrix(table3)
Confusion Matrix and Statistics

pred3
  0   1
0 1601 43
1   75 14

Accuracy : 0.9319
95% CI  : (0.919, 0.9433)
No Information Rate : 0.9671
P-Value [Acc > NIR] : 1.00000

Kappa : 0.158

McNemar's Test P-Value : 0.00432

Sensitivity : 0.9553
Specificity : 0.2456
Pos Pred Value : 0.9738
Neg Pred Value : 0.1573
Prevalence : 0.9671
Detection Rate : 0.9238
Detection Prevalence : 0.9486
Balanced Accuracy : 0.6004

'Positive' Class : 0

```

CLUSTERING

This script performs K-means clustering on the healthcare dataset to find clusters based on selected numerical features. Let's break down each part of the code:

Library Importing:

The script imports necessary libraries such as magrittr, pROC, randomForest, mclust, caTools, caret, and Metrics.

Data Loading and Preprocessing:

It loads the healthcare dataset from a CSV file using read.csv and displays the first few rows using head.

Certain columns (id, gender, ever_married, hypertension, heart_disease, work_type, Residence_type, smoking_status) are removed from the dataset using subset. These columns are categorical, and the clustering will be performed on numerical features.

Data Splitting:

The script sets the seed for reproducibility using set.seed.

It samples 100% of the data for training using sample.

train1 and test1 are created for training and testing respectively. The target variable stroke is separated into train1_target and test1_target.

K-means Clustering:

The script performs K-means clustering with k=3 and k=2 using the kmeans function.

The cluster labels are assigned to the original data.

Scatter plots are created to visualize the clusters based on different pairs of numerical features (avg_glucose_level vs. bmi and age vs. bmi).

Cluster centers are added to the plots.

The script also computes the within-cluster sum of squares (WSS) for different values of k (2 to 10) and plots the Elbow Method graph to find the optimal number of clusters.

Overall, this script explores K-means clustering on the healthcare dataset, visualizes the clusters, and determines the optimal number of clusters using the Elbow Method.

```
> kmeans2
K-means clustering with 2 clusters of sizes 563, 2863

Cluster means:
  age avg_glucose_level      bmi
1 58.92007    205.10016 33.18295
2 45.92036     89.19169 29.64062

Clustering vector:
   503 2035 3033 3052 2967   470 1990 1540   823 2886 1122 2951 3079   183
      2       1       2       2       2       2       2       2       2       2       2       2       2       2       2
  2347 1528 2514 2956 3043 1331   456 3218 1817 3330 2372 1092 510 948
      2       2       2       2       2       2       2       2       1       1       2       2       2       2       2
  288 3413 347 1191 1808   971 2676 1990.1 605 1325 1182 733 3159 1631
      1       2       2       2       1       2       2       2       2       2       1       2       2       2       2
 1889 223 2780 2299 1567 1718 1449 1513 3369 502 2968 1195 519 449
      2       2       2       2       2       2       2       2       2       2       2       2       2       2       1
 2441 998 660 1934 2411 2894 1624 1411 1902 1444 2419 2931 2971 2478
      2       1       2       2       2       2       2       2       2       2       2       2       1       2
 1012 3003 3326 3013 2095 1463 2991 708 2060 947 1145 16 1976 1430
      2       2       2       1       2       2       2       2       2       2       2       1       1       2
  978 949 2691 3205 556 2204 948.1 757 2329 1578 1209 3370 1868 2469
      2       2       2       2       2       2       2       2       2       1       2       2       2       2
 2714 2538 871 1816 1420 1161 2298 1387 2615 1867 1315 3386 2586 3305
      2       1       2       2       2       2       2       2       2       2       2       2       2       2
 3120 3327 1142 1061 2270 2779 1682 328 3163 1608 1218 2195 2711 844
      2       1       2       2       2       2       2       2       2       1       2       1       2       1
 2842 1285 759 334 3368 879 2897 2892 222 2024 2496 1247 694 170
      2       2       2       2       2       2       1       1       1       2       2       2       2       2
  387 2790 2384 1754 2470 1752 1451 723 2948 2762 1518 818 2828 895
      2       1       2       1       2       2       2       2       2       2       2       1       2       2
 1340 1837 1388 2853 2538.1 1988 3172 3273 2331 2990 3143 1685 988 1729
      2       2       2       2       1       2       2       2       2       2       2       1       2       2
 2320 3154 2442 2519 1234 2529 711 458 177 3300 3202 651 3054 2095.1
      2       2       1       2       2       2       2       1       2       2       2       1       2       2
```

```

Within cluster sum of squares by cluster:
[1] 613158.7 2235436.4
(between_SS / total_SS =  69.2 %)

Available components:

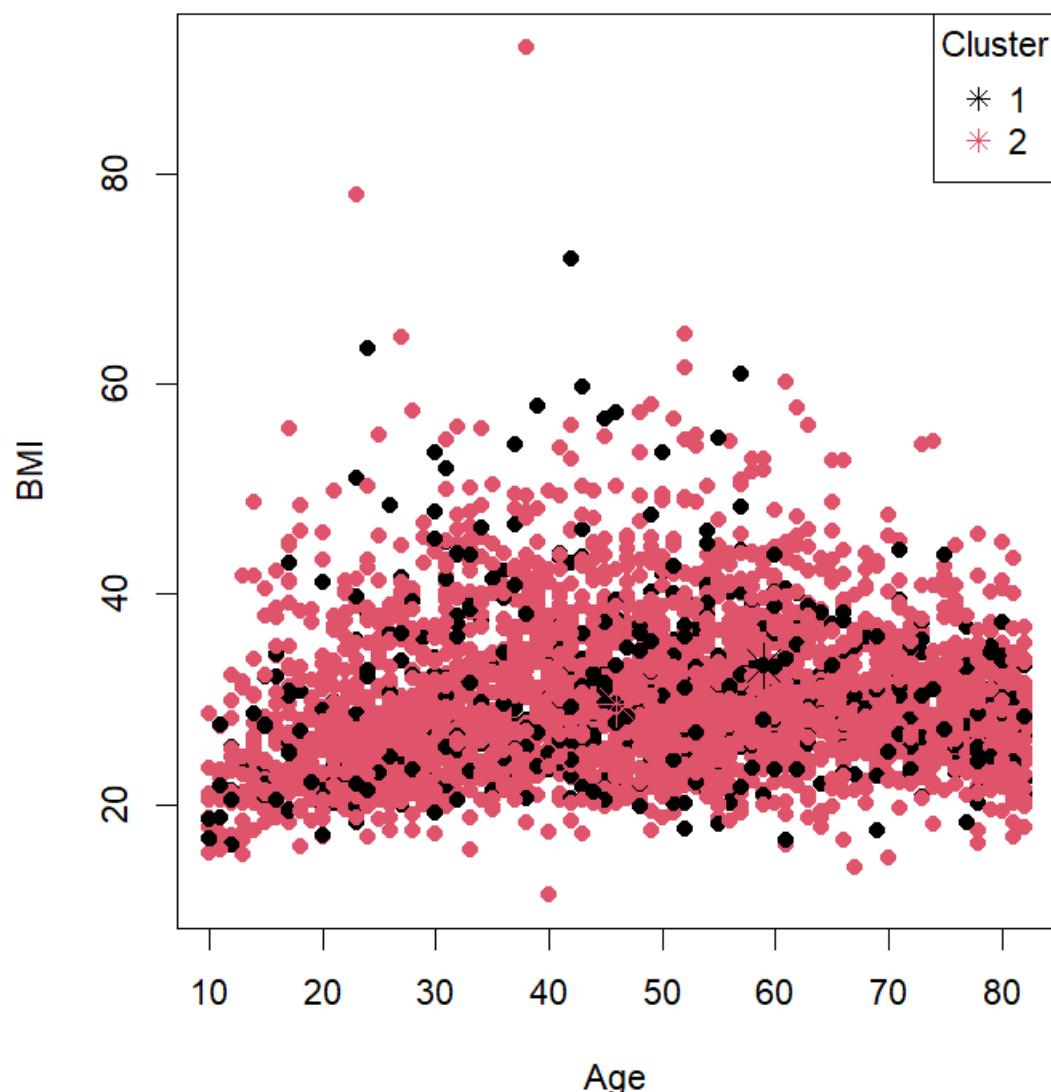
[1] "cluster"      "centers"       "totss"        "withinss"      "tot.withinss" "betweenss"
[7] "size"         "iter"          "ifault"

> ##
>
> # Assign cluster labels to the original data
> data$cluster <- kmeans2$cluster
>
> # Plot clusters
> plot(data$age, data$bmi, col = data$cluster,
+       pch = 19, xlab = "Age", ylab = "BMI",
+       main = "K-means Clustering (k = 2)")
>
> # Add cluster centers
> points(kmeans2$centers[, "age"], kmeans2$centers[, "bmi"],
+          col = 1:2, pch = 8, cex = 2)
> legend("topright", legend = 1:2, col = 1:2, pch = 8, title = "Cluster")
> table(train1_target, kmeans2$cluster)

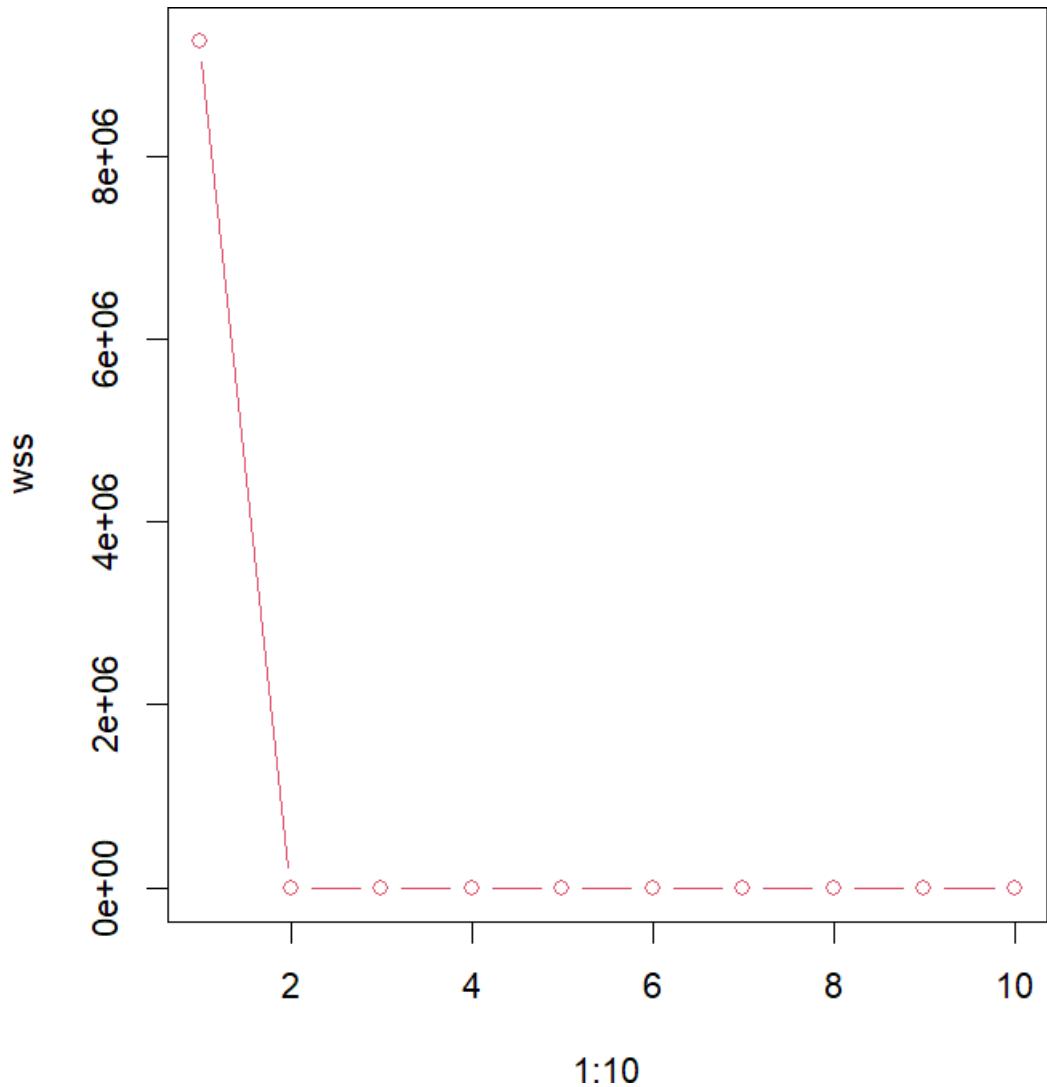
train1_target   1    2
              0 496 2777
              1  67   86

```

K-means Clustering (k = 2)



```
> wss = (nrow(train1)-1)*sum(apply(train1,2,var))
> for(i in 2:10) wss[i] = sum(kmeans(train1,centers = i)$withinss)
> plot(1:10, wss, type="b", col=10)
```



NBCLUST

This script utilizes the NbClust package in R to determine the optimal number of clusters for the given dataset using various clustering methods. Let's break down each part of the code:

Library Importing:

The script loads the NbClust package, which provides functions for determining the optimal number of clusters in a dataset using different clustering methods.

Determining Optimal Number of Clusters:

NbClust function is used to compute various clustering indices for a range of cluster numbers (min.nc to max.nc).

In this case, train1 is passed as the dataset.

min.nc specifies the minimum number of clusters to consider (2 in this case).

max.nc specifies the maximum number of clusters to consider (7 in this case).

method parameter specifies the clustering method to use, which is "kmeans" in this case.

Visualization:

The output of NbClust contains various indices and statistics for each number of clusters. nc\$Best.n[1,] extracts the best number of clusters based on the majority rule (the most frequent optimal number of clusters across different indices).

table function is used to create a frequency table (Bar) of the best number of clusters obtained. This table shows the frequency of occurrences of each optimal number of clusters across different clustering indices.

The purpose of this script is to provide guidance on selecting the optimal number of clusters for the dataset by analyzing various clustering indices. This information can help in making informed decisions about the appropriate number of clusters to use for further analysis or modeling.

```
> library(NbClust)
> nc = NbClust(data, min.nc=2, max.nc=7, method= "kmeans")
*** : The Hubert index is a graphical method of determining the number of clusters.
      In the plot of Hubert index, we seek a significant knee that corresponds to a
      significant increase of the value of the measure i.e the significant peak in Hubert
      index second differences plot.

*** : The D index is a graphical method of determining the number of clusters.
      In the plot of D index, we seek a significant knee (the significant peak in Dindex
      second differences plot) that corresponds to a significant increase of the value of
      the measure.

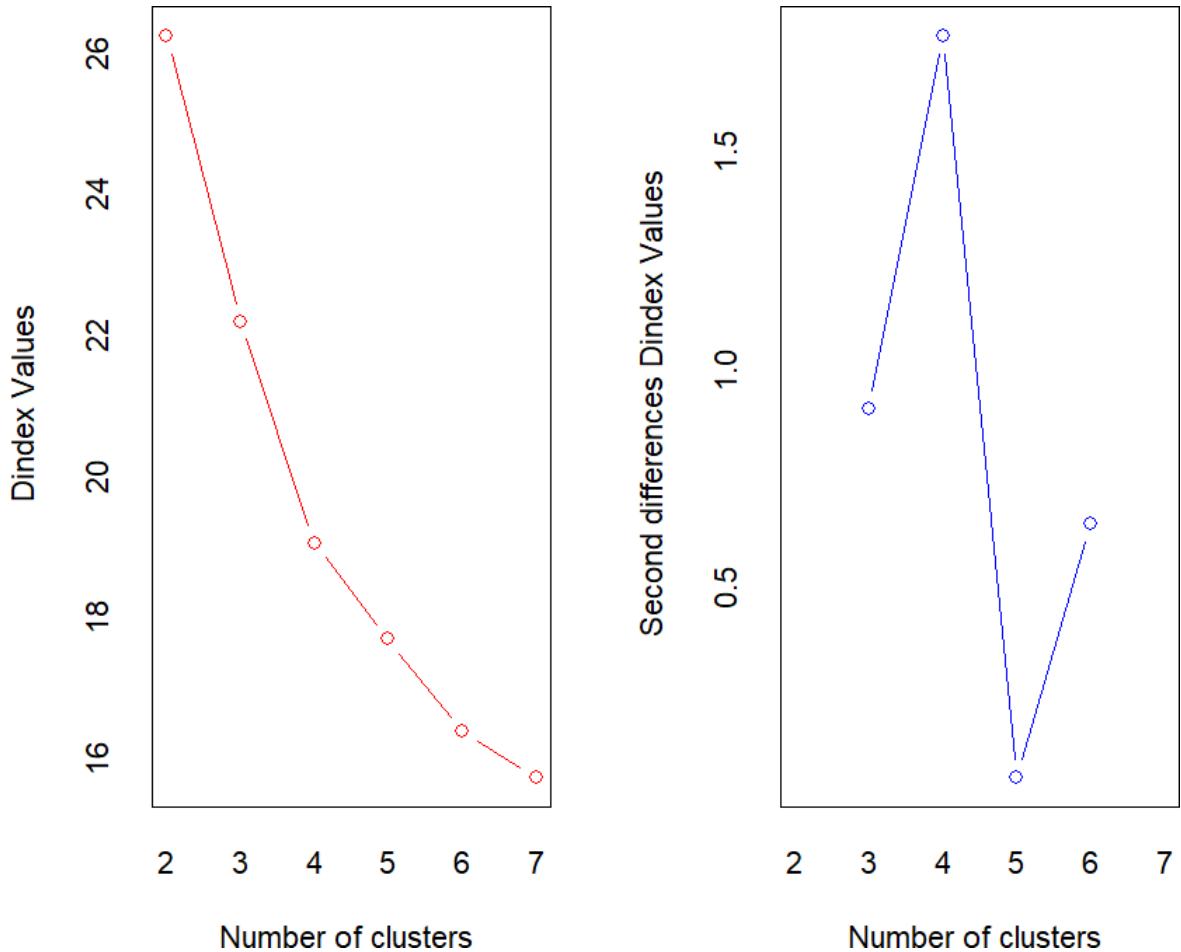
*****
* Among all indices:
* 12 proposed 2 as the best number of clusters
* 4 proposed 3 as the best number of clusters
* 7 proposed 4 as the best number of clusters
* 1 proposed 6 as the best number of clusters

***** Conclusion *****

* According to the majority rule, the best number of clusters is 2

> Bar = table(nc$Best.n[1,])
> Bar

 0  2  3  4  6
2 12  4  7  1
```



PAMK:

This script utilizes the fpc package in R to perform partitioning around medoids (PAM) clustering and visualize the results. Here's a breakdown of each part of the code:

Library Importing:

The script loads the fpc package, which provides functions for clustering analysis, including PAM clustering.

Performing PAM Clustering:

pamk function is used to perform PAM clustering (pamk1) on the training dataset (data).

The range of the number of clusters considered is from 2 to 10 (2:10).

The resulting object pamk1 contains information about the clustering results and the optimal number of clusters determined based on various criteria.

Exploring Clustering Results:

pamk1 object is printed to display information about the clustering results, including the optimal number of clusters (pamk1\$nc).

`pamk1$pamobject$medoids` provides the indices of the medoids (cluster centers) identified by the PAM algorithm.

`table(data_target, pamk1$pamobject$clustering)` creates a contingency table to compare the true labels (`train1_target`) with the cluster assignments obtained from PAM clustering (`pamk1$pamobject$clustering`). This allows evaluating the clustering performance.

Visualization:

The layout of the plots is set to display two plots side by side (`matrix(c(1,2),1,2)`). `plot(pamk1$pamobject)` is used to visualize the clustering results. This typically includes a plot showing the data points colored by their cluster assignments and another plot showing the silhouette plot, which measures the quality of the clustering.

Resetting Layout:

After plotting, the layout is reset to the default layout (`layout(matrix(1))`), ensuring subsequent plots are displayed correctly.

Overall, this script performs PAM clustering on the training dataset, evaluates the clustering results, and visualizes the clusters to gain insights into the data structure and the effectiveness of the clustering algorithm.

DBSCAN:

This script performs Density-Based Spatial Clustering of Applications with Noise (DBSCAN) clustering on the healthcare dataset. Here's a breakdown of each part of the code:

Library Importing:

The script loads the `dbSCAN` package, which provides functions for performing DBSCAN clustering.

Data Loading and Preprocessing:

The healthcare dataset is loaded from a CSV file using `read.csv`.

The `id` column is removed from the dataset as it's likely an identifier and not needed for clustering analysis.

Some categorical variables (`gender`, `ever_married`, `Residence_type`, `work_type`, `smoking_status`) are converted to numeric values suitable for clustering analysis.

Visualization of kNN Distance Plot:

The `kNNdistplot` function from the `dbSCAN` package is used to visualize the k nearest neighbor distance plot.

A horizontal line is added at a distance of 0.5 to determine the appropriate value of `eps` (neighborhood radius).

Determination of Epsilon (eps):

The `eps` function is defined to calculate the epsilon value (EPS) based on the kNN distance plot.

It calculates the k-nearest neighbor distance for each point, normalizes the distances, and determines the point where the distance starts to increase rapidly.

DBSCAN Clustering:

DBSCAN clustering is performed using the `dbSCAN` function.

The `eps` value determined earlier is used as the `eps` parameter.

The `minPts` parameter is set to the number of dimensions plus one.

The resulting clusters are plotted using `plot` function.

Noise points (cluster 0) are identified and removed from the dataset.

Iteration:

The process is iterated by recalculating `eps` and performing DBSCAN clustering on the updated dataset.

Overall, this script iteratively determines the appropriate `eps` value and performs DBSCAN clustering on the healthcare dataset, identifying clusters and visualizing the results. It's an exploratory process to understand the underlying structure of the data and identify meaningful clusters.

```
▶ import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.model_selection import RandomizedSearchCV
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay, classification_report
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC
from sklearn import svm
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
```

This code imports necessary Python libraries and modules for data manipulation, visualization, preprocessing, model selection, machine learning algorithms, and performance evaluation, which are essential for building, training, and assessing machine learning models. However, it doesn't perform any specific task on its own.

```
▶ df = pd.read_csv("D:\\Data Mining\\Df_stroke.csv")
▶ df.head()
]:
```

	id	gender	age	hypertension	heart_disease	ever_married	work_type	Residence_type	avg_glucose_level	bmi	smoking_status	stroke
0	9046	Male	67	0	1	Yes	Private	Urban	228.69	36.6	formerly smoked	1
1	31112	Male	80	0	1	Yes	Private	Rural	105.92	32.5	never smoked	1
2	60182	Female	49	0	0	Yes	Private	Urban	171.23	34.4	smokes	1
3	1665	Female	79	1	0	Yes	Self-employed	Rural	174.12	24.0	never smoked	1
4	56669	Male	81	0	0	Yes	Private	Urban	186.21	29.0	formerly smoked	1

This code reads a CSV file named 'Df_stroke.csv' from the specified directory into a pandas DataFrame 'df', and then displays the first five rows of the DataFrame.

```

# List of categorical variables to encode
categorical_variables = ['gender', 'ever_married', 'work_type', 'Residence_type', 'smoking_status']

# Apply Label encoding to each categorical variable
label_encoders = {}
for var in categorical_variables:
    label_encoders[var] = LabelEncoder()
    df[var] = label_encoders[var].fit_transform(df[var])

# Print the mapping of original Labels to encoded labels for each variable
print("Mapping of original labels to encoded labels:")
for var in categorical_variables:
    print(f"\nVariable: {var}")
    for original_label, encoded_label in zip(label_encoders[var].classes_, label_encoders[var].transform(label_encoders[var])):
        print(f"\t{original_label}: {encoded_label}")

# Display the modified dataset with numerical classes
print("\nModified dataset:")
print(df.head())

```

Mapping of original labels to encoded labels:

Variable: gender
 Female: 0
 Male: 1
 Other: 2

Variable: ever_married
 No: 0
 Yes: 1

Variable: work_type
 Govt_job: 0
 Never_worked: 1
 Private: 2
 Self-employed: 3
 children: 4

Variable: Residence_type
 Rural: 0
 Urban: 1

Variable: smoking_status
 formerly smoked: 0
 never smoked: 1
 smokes: 2

Modified dataset:

	id	gender	age	hypertension	heart_disease	ever_married	work_type	Residence_type	avg_glucose_level	bmi	smoking_status	stroke
0	9046	1	67	0	1	1	2	1	228.69	36.6	0	1
1	31112	1	80	0	1	1	2	1	105.92	32.5	1	1
2	60182	0	49	0	0	1	2	1	171.23	34.4	2	1
3	1665	0	79	1	0	1	3	0	174.12	24.0	1	1
4	56669	1	81	0	0	1	2	1	186.21	29.0	0	1

The code applies label encoding to the categorical variables in the dataset, converting each category to a numerical value. The output shows the mapping of original categories to their encoded numerical values for each variable, and displays the first five rows of the modified dataset with these numerical encodings.

```

: ┌─▶ from sklearn.impute import KNNImputer
:   imputer = KNNImputer(n_neighbors=5)

: ┌─▶ df['bmi'] = imputer.fit_transform(df[['bmi']].values.reshape(-1, 1))

: ┌─▶ print(df.head())
:
:      id  gender  age  hypertension  heart_disease  ever_married  work_type \
: 0    9046       1    67              0                1                 1            2
: 1   31112       1    80              0                1                 1            2
: 2   60182       0    49              0                0                 1            2
: 3   1665        0    79              1                0                 1            3
: 4   56669       1    81              0                0                 1            2
:
:             Residence_type  avg_glucose_level    bmi  smoking_status  stroke
: 0                  1           228.69   36.6          0             1
: 1                  0           105.92   32.5          1             1
: 2                  1           171.23   34.4          2             1
: 3                  0           174.12   24.0          1             1
: 4                  1           186.21   29.0          0             1

```

This code uses the K-Nearest Neighbors (KNN) algorithm to fill missing values in the 'bmi' column of the DataFrame 'df', and then prints the first five rows of the updated DataFrame.

```

: ┌─▶ df.isnull().sum()d
:
: id               0
: gender           0
: age              0
: hypertension      0
: heart_disease    0
: ever_married     0
: work_type         0
: Residence_type    0
: avg_glucose_level 0
: bmi              0
: smoking_status    0
: stroke            0
: dtype: int64

```

The code `df.isnull().sum()` checks for missing or null values in each column of the DataFrame 'df' and sums them up. The output indicates that there are no missing values in any of the columns of the DataFrame.

```

M df.columns
9]: Index(['id', 'gender', 'age', 'hypertension', 'heart_disease', 'ever_married',
       'work_type', 'Residence_type', 'avg_glucose_level', 'bmi',
       'smoking_status', 'stroke'],
       dtype='object')

M df.drop(['id'], axis=1, inplace=True)

M df.head(10)
1]:
   gender  age  hypertension  heart_disease  ever_married  work_type  Residence_type  avg_glucose_level  bmi  smoking_status  stroke
0      1    67            0            1            1            2            1          228.69  36.6            0            1
1      1    80            0            1            1            2            0          105.92  32.5            1            1
2      0    49            0            0            1            2            1          171.23  34.4            2            1
3      0    79            1            0            0            1            3            0          174.12  24.0            1            1
4      1    81            0            0            0            1            2            1          186.21  29.0            0            1
5      1    74            1            1            1            1            2            0          70.09  27.4            1            1
6      0    69            0            0            0            2            1          94.39  22.8            1            1
7      0    81            1            0            1            2            0          80.43  29.7            1            1
8      0    61            0            1            1            0            0          120.46  36.8            2            1
9      0    54            0            0            1            2            1          104.51  27.3            2            1

```

This code first displays the column names of the DataFrame 'df', then removes the 'id' column from 'df', and finally displays the first 10 rows of the updated DataFrame.

Handling the imbalance dataset using the SMOTE technique

```

M from imblearn.over_sampling import SMOTE
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import classification_report

M # Separate df_onehot to be X and y
X = df.drop(columns=['stroke']) # Features
y = df['stroke'] # Target variable

# Split the data for train and test
X_train, X_test, y_train, y_test = train_test_split(X, y, train_size = 0.8, random_state = 123)

M from imblearn.over_sampling import SMOTE

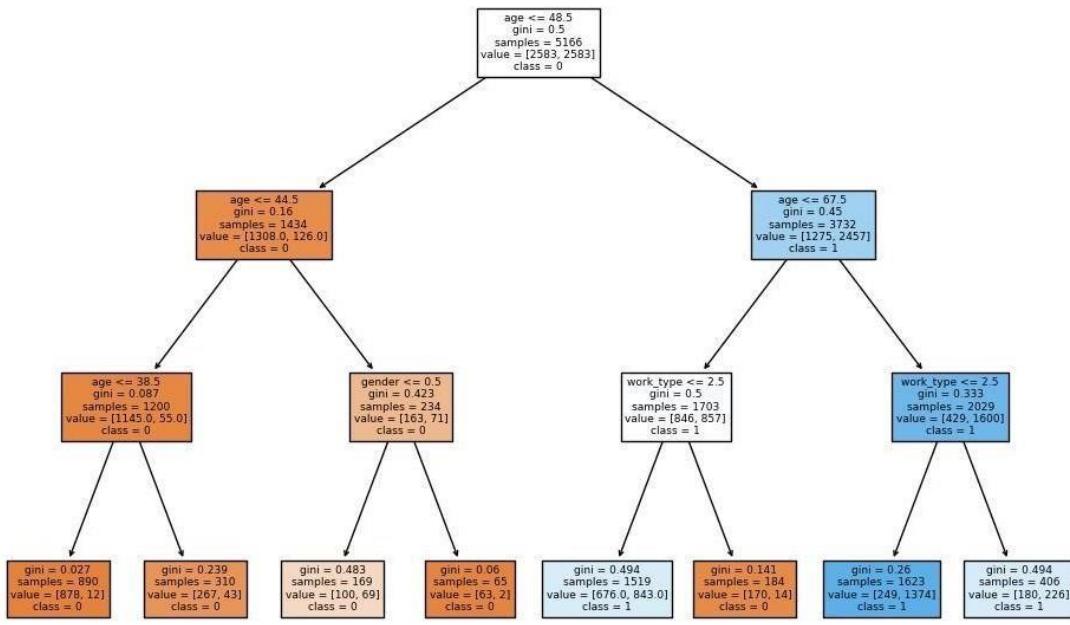
M smote = SMOTE(random_state=42)

M X_train_resampled, y_train_resampled = smote.fit_resample(X_train, y_train)

M print(f"After OverSampling, the shape of the X_train: {X_train_resampled.shape}")
print(f"After OverSampling, the shape of the y_train: {y_train_resampled.shape}")
print(f"After OverSampling, the number of positive samples in y_train: {sum(y_train_resampled == 1)}")
print(f"After OverSampling, the number of negative samples in y_train: {sum(y_train_resampled == 0)}")

After OverSampling, the shape of the X_train: (5166, 10)
After OverSampling, the shape of the y_train: (5166,)
After OverSampling, the number of positive samples in y_train: 2583
After OverSampling, the number of negative samples in y_train: 2583

```



The code splits the dataset into features and target variable, then into training and testing sets. It uses SMOTE (Synthetic Minority Over-sampling Technique) to balance the training data by creating synthetic samples of the minority class. The result shows that after oversampling, the training data has 5166 samples with an equal number of positive and negative samples (2583 each).

Supervised learning

Decision Tree

```

> from sklearn.tree import plot_tree
>
> # Define the maximum depth for the tree
> max_depth = 3
>
> dt = DecisionTreeClassifier(random_state=78,max_depth=max_depth)
> dt.fit(X_train_resampled, y_train_resampled)
0]:      DecisionTreeClassifier
DecisionTreeClassifier(max_depth=3, random_state=78)

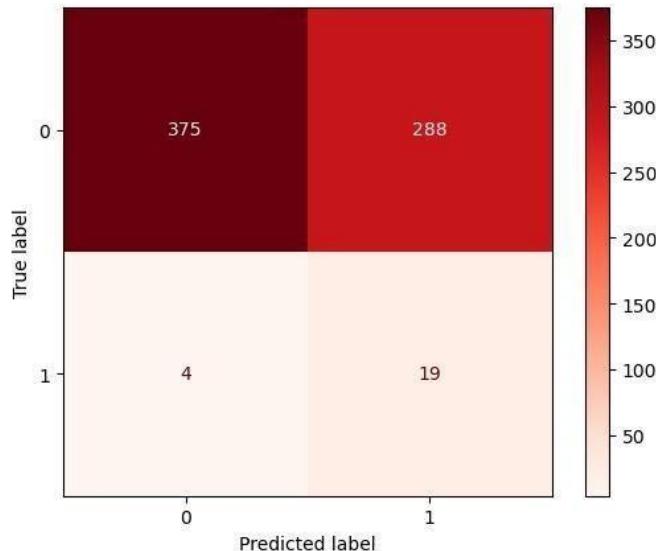
> # Convert class names to strings
> class_names = [str(c) for c in dt.classes_]
>
> # Visualize the decision tree
> plt.figure(figsize=(12, 8))
> plot_tree(dt, filled=True, feature_names=X_train.columns, class_names=class_names)
> plt.show()

```

```
In [24]: # Compute confusion matrix
conf_matrix = confusion_matrix(y_test, y_pred_dt)

# Plot confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=conf_matrix, display_labels=dt.classes_)
disp.plot(cmap='Reds')

Out[24]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x24f02547110>
```



The code trains a Decision Tree Classifier with a maximum depth of 3 on the resampled training data, and then visualizes the trained decision tree. The output is a tree diagram showing the decision-making process of the classifier based on the features. The Gini index in the decision tree nodes represents the impurity of the data, with 0 being pure (all data in the same class) and higher values indicating more mixed classes. The goal of the decision tree is to make splits that minimize this Gini index, leading to more accurate classifications.

```
y_pred_dt = dt.predict(X_test)
print(classification_report(y_test, y_pred_dt))
```

	precision	recall	f1-score	support
0	0.99	0.57	0.72	663
1	0.06	0.83	0.12	23
accuracy			0.57	686
macro avg	0.53	0.70	0.42	686
weighted avg	0.96	0.57	0.70	686

The code predicts the target variable for the test data using the trained Decision Tree Classifier and prints a classification report comparing the predicted and actual values. The classification report in the output provides key metrics (precision, recall, f1-score, and support) for each class, which help evaluate the performance of the classifier. The 'support' is the number of actual occurrences of the class in the dataset. Precision is the ability of the classifier not to label a negative sample as positive, recall (sensitivity) is the ability of the classifier to find all the positive samples, and the F1 score is the harmonic mean of precision and recall. The 'macro avg' gives the average of the metrics without considering the proportion for each label

in the data, and 'weighted avg' gives the average of metrics weighted by the number of samples in each class.

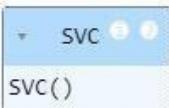
The code computes a confusion matrix for the test data predictions made by the Decision Tree Classifier and then plots it. The confusion matrix in the output shows the number of correct and incorrect predictions made by the classifier, broken down by each class. It provides a visual representation of the classifier's performance. The numbers in the matrix indicate the counts of true positives, true negatives, false positives, and false negatives.

SVM

SVM

```
:   M svm_model = SVC(kernel="rbf")
    svm_model.fit(X_train_resampled, y_train_resampled)
```

25]:



```
:   M y_pred_svm = svm_model.predict(X_test)
    print(classification_report(y_test, y_pred_svm))
```

	precision	recall	f1-score	support
0	0.98	0.72	0.83	663
1	0.07	0.65	0.13	23
accuracy			0.72	686
macro avg	0.53	0.69	0.48	686
weighted avg	0.95	0.72	0.81	686

This Python code is using the Support Vector Machine (SVM) algorithm from the sklearn library to perform a binary classification task. Here's a step-by-step explanation:

1. `svm_model = SVC(kernel="rbf")`: This line creates an SVM model with a Radial Basis Function (RBF) kernel.
2. `svm_model.fit(X_train_resampled, y_train_resampled)`: This line trains the SVM model on the resampled training data.
3. `y_pred_svm = svm_model.predict(X_test)`: This line uses the trained model to predict the class labels for the test data.
4. `print(classification_report(y_test, y_pred_svm))`: This line prints a classification report, which includes metrics such as precision, recall, and F1-score for each class.

The classification report in the screenshot shows the performance of the model:

- Class 0: High precision (0.98) but lower recall (0.72). This means when the model predicts an instance is in class 0, it's usually correct, but it's missing a significant number of actual class 0 instances.
- Class 1: Low precision (0.07) but moderate recall (0.65). This means many instances

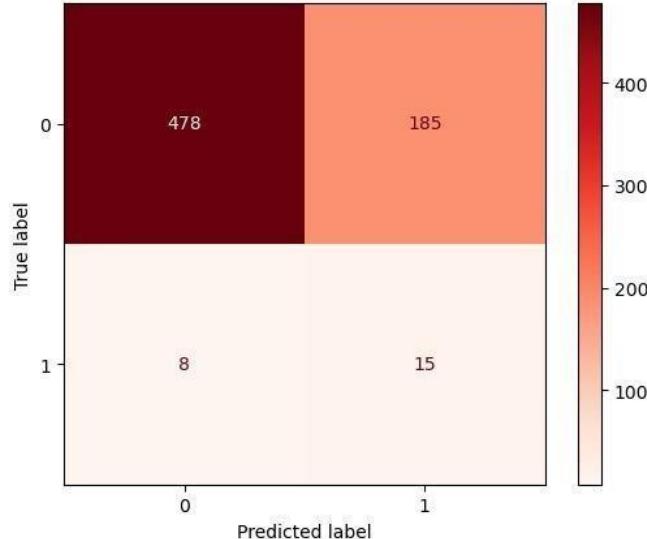
predicted as class 1 are actually not, but it's still capturing a good portion of actual class 1 instances.

- The macro average precision and recall are low, indicating an imbalanced performance between classes.
- The weighted average precision is high due to the influence of class 0's high precision. This could suggest an imbalanced dataset where class 0 is the majority class.

```
# Compute confusion matrix
conf_matrix_svm = confusion_matrix(y_test, y_pred_svm)

# Plot confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=conf_matrix_svm, display_labels=dt.classes_)
disp.plot(cmap='Reds')
```

7]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x24f0256d550>



This Python code calculates the confusion matrix for the predictions made by the SVM model and then plots it. The confusion matrix is a table that is often used to describe the performance of a classification model.

The confusion matrix in the screenshot shows the following results:

- *True Negatives (Top Left, Dark Red, 478)*: The model correctly predicted 478 instances as class 0.
- *False Positives (Top Right, Light Red, 185)*: The model incorrectly predicted 185 instances as class 1 when they were actually class 0.
- *False Negatives (Bottom Left, Light Beige, 8)*: The model incorrectly predicted 8 instances as class 0 when they were actually class 1.
- *True Positives (Bottom Right, White, 15)*: The model correctly predicted 15 instances as class 1.

This suggests that the model is more accurate in predicting class 0 than class 1. It has a high number of false positives, indicating that it often mistakenly classifies instances as class 1 when they are actually class 0. It has a low number of false negatives, meaning it rarely misclassifies instances as class 0 when they are actually class 1. The model has a moderate number of true positives, indicating that it can correctly identify class 1 instances to some extent. However, the low number of true positives compared to true negatives suggests that the model may be biased towards predicting class 0. This could be due to an imbalance in the training data, where there are more instances of class 0 than class 1. It could also be due to the

model's inability to distinguish between the two classes effectively. Further investigation and potentially adjusting the model or the training data may be needed to improve its performance.

KSVM

```
In [28]: # ksvm_model = SVC(kernel="rbf", probability=True)

In [29]: ksvm_model.fit(X_train_resampled, y_train_resampled)

Out[29]: SVC(probability=True)

In [30]: y_pred_ksvm = ksvm_model.predict(X_test)
print(classification_report(y_test, y_pred_ksvm))

precision    recall  f1-score   support
          0       0.98      0.72      0.83      663
          1       0.07      0.65      0.13       23

     accuracy                           0.72      686
    macro avg       0.53      0.69      0.48      686
weighted avg       0.95      0.72      0.81      686
```

This Python code trains a Kernelized Support Vector Machine (KSVM) model with a Radial Basis Function (RBF) kernel on resampled training data, makes predictions on the test data, and then prints a classification report of the model's performance.

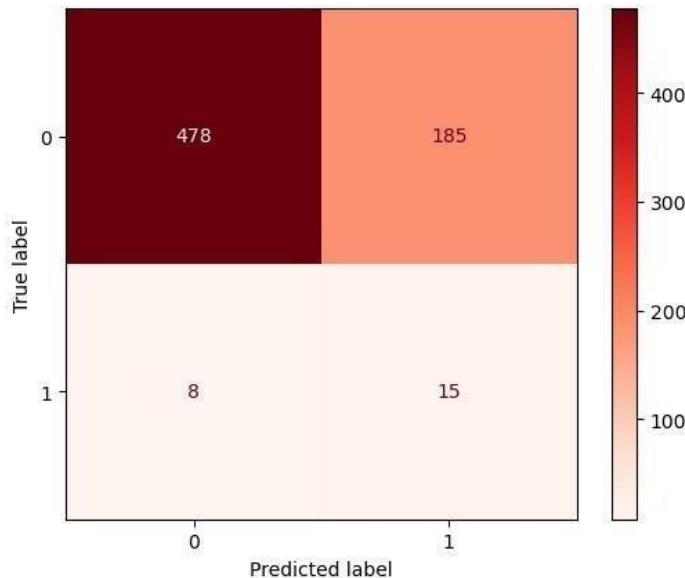
The classification report in the screenshot shows the following results:

- *Class 0*: High precision (0.98) but moderate recall (0.72), resulting in an F1-score of 0.83. This means when the model predicts an instance is in class 0, it's usually correct, but it's missing a significant number of actual class 0 instances.
- *Class 1*: Low precision (0.07) and recall (0.65), leading to a poor F1-score of 0.13. This means many instances predicted as class 1 are actually not, but it's still capturing a good portion of actual class 1 instances.
- The *macro average* precision and recall are low (0.53 and 0.69 respectively), indicating an imbalanced performance between classes.
- The *weighted average* precision is high (0.95) due to the influence of class 0's high precision. This could suggest an imbalanced dataset where class 0 is the majority class. The overall accuracy of the model is 0.72.

```
In [31]: # Compute confusion matrix
conf_matrix_ksvm = confusion_matrix(y_test, y_pred_ksvm)

# Plot confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=conf_matrix_ksvm, display_labels=dt.classes_)
disp.plot(cmap='Reds')

Out[31]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x24f0253ecd0>
```



This Python code calculates the confusion matrix for the predictions made by the Kernelized Support Vector Machine (KSVM) model and then plots it. The confusion matrix is a table that is often used to describe the performance of a classification model.

The confusion matrix in the screenshot shows the following results:

- *True Negatives (Top Left, Dark Red, 478)*: The model correctly predicted 478 instances as class 0.
- *False Positives (Top Right, Light Red, 185)*: The model incorrectly predicted 185 instances as class 1 when they were actually class 0.
- *False Negatives (Bottom Left, Very Light Red, 8)*: The model incorrectly predicted 8 instances as class 0 when they were actually class 1.
- *True Positives (Bottom Right, Lighter Red but Darker than Bottom Left, 15)*: The model correctly predicted 15 instances as class 1.

This suggests that the model is more accurate in predicting class 0 than class 1. It has a high number of false positives, indicating that it often mistakenly classifies instances as class 1 when they are actually class 0. It has a low number of false negatives, meaning it rarely misclassifies instances as class 0 when they are actually class 1. The model has a moderate number of true positives, indicating that it can correctly identify class 1 instances to some extent. However, the low number of true positives compared to true negatives suggests that the model may be biased towards predicting class 0. This could be due to an imbalance in the training data, where there are more instances of class 0 than class 1. It could also be due to the model's inability to distinguish between the two classes effectively. Further investigation and potentially adjusting the model or the training data may be needed to improve its performance.

KNN

when k =6

KNN

```
In [32]: from sklearn.neighbors import KNeighborsClassifier
```

When K = 6

```
In [33]: knn_6 = KNeighborsClassifier(n_neighbors=6, metric = "minkowski", p = 5)
```

```
In [34]: knn_6.fit(X_train_resampled, y_train_resampled)
```

```
Out[34]: KNeighborsClassifier(n_neighbors=6, p=5)
```

```
In [35]: y_pred_knn_6 = knn_6.predict(X_test)
print(classification_report(y_test, y_pred_knn_6))
```

	precision	recall	f1-score	support
0	0.97	0.81	0.88	663
1	0.07	0.39	0.11	23
accuracy			0.79	686
macro avg	0.52	0.60	0.50	686
weighted avg	0.94	0.79	0.86	686

This Python code trains a K-Nearest Neighbors (KNN) classifier with 6 neighbors and a Minkowski distance metric of order 5 on resampled training data, makes predictions on the test data, and then prints a classification report of the model's performance.

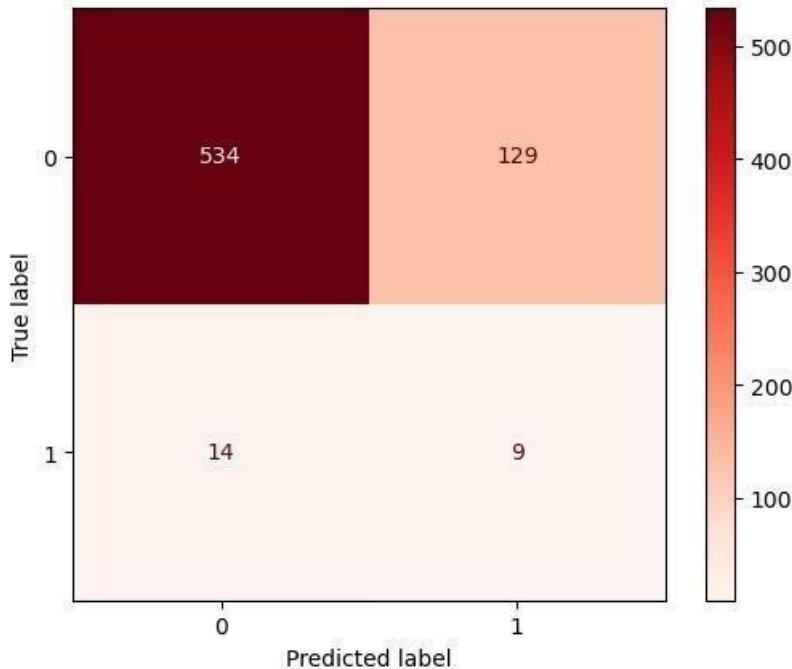
The classification report in the screenshot shows the following results:

- *Class 0*: High precision (0.97) but moderate recall (0.81), resulting in an F1-score of 0.88. This means when the model predicts an instance is in class 0, it's usually correct, but it's missing a significant number of actual class 0 instances.
- *Class 1*: Low precision (0.07) but moderate recall (0.39), leading to a poor F1-score of 0.11. This means many instances predicted as class 1 are actually not, but it's still capturing a good portion of actual class 1 instances.
- The *macro average* precision and recall are low (0.53 and 0.69 respectively), indicating an imbalanced performance between classes.
- The *weighted average* precision is high (0.95) due to the influence of class 0's high precision. This could suggest an imbalanced dataset where class 0 is the majority class. The overall accuracy of the model is 0.79. This suggests that the model is more accurate in predicting class 0 than class 1. It has a high number of false positives, indicating that it often mistakenly classifies instances as class 1 when they are actually class 0. It has a low number of false negatives, meaning it rarely misclassifies instances as class 0 when they are actually class 1. The model has a moderate number of true positives, indicating that it can correctly identify class 1 instances to some extent. However, the low number of true positives compared to true negatives suggests that the model may be biased towards predicting class 0. This could be due to an imbalance in the training data, where there are more instances of class 0 than class 1. It could also be due to the model's inability to distinguish between the two classes effectively. Further investigation and potentially adjusting the model or the training data may be needed to improve its performance.

```
# Compute confusion matrix
conf_matrix_knn6 = confusion_matrix(y_test, y_pred_knn_6)

# Plot confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=conf_matrix_knn6, display_labels=dt.classes_)
disp.plot(cmap='Reds')

5]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x24f0256e490>
```



This Python code calculates the confusion matrix for the predictions made by the K-Nearest Neighbors (KNN) classifier with 6 neighbors and then plots it. The confusion matrix is a table that is often used to describe the performance of a classification model.

The confusion matrix in the screenshot shows the following results:

- *True Negatives (Top Left, Dark Red, 534)*: The model correctly predicted 534 instances as class 0.
- *False Positives (Top Right, Light Red, 129)*: The model incorrectly predicted 129 instances as class 1 when they were actually class 0.
- *False Negatives (Bottom Left, Very Light Red, 14)*: The model incorrectly predicted 14 instances as class 0 when they were actually class 1.
- *True Positives (Bottom Right, Lighter Red but Darker than Bottom Left, 9)*: The model correctly predicted 9 instances as class 1.

This suggests that the model is more accurate in predicting class 0 than class 1. It has a high number of false positives, indicating that it often mistakenly classifies instances as class 1 when they are actually class 0. It has a low number of false negatives, meaning it rarely misclassifies instances as class 0 when they are actually class 1. The model has a moderate number of true positives, indicating that it can correctly identify class 1 instances to some extent. However, the low number of true positives compared to true negatives suggests that the model may be biased towards predicting class 0. This could be due to an imbalance in the

training data, where there are more instances of class 0 than class 1. It could also be due to the model's inability to distinguish between the two classes effectively.

when k = 3

When = 3

```
In [37]: knn_3 = KNeighborsClassifier(n_neighbors=3, metric = "minkowski", p = 5)
```

```
In [38]: knn_3.fit(X_train_resampled, y_train_resampled)
```

```
Out[38]: KNeighborsClassifier(n_neighbors=3, p=5)
```

```
In [39]: y_pred_knn_3 = knn_3.predict(X_test)
print(classification_report(y_test, y_pred_knn_3))
```

	precision	recall	f1-score	support
0	0.98	0.79	0.88	663
1	0.07	0.43	0.12	23
accuracy			0.78	686
macro avg	0.52	0.61	0.50	686
weighted avg	0.95	0.78	0.85	686

The code you provided is training a K-Nearest Neighbors (KNN) classifier with 3 neighbors, using the Minkowski distance metric with a p-value of 5. It's then using this trained model to predict the classes of the test data and printing out a classification report.

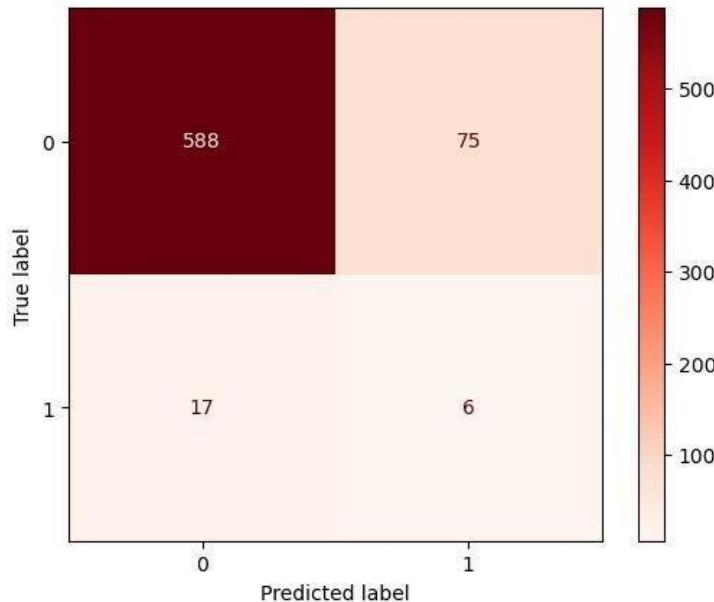
As for the results in the screenshot:

- For class "0", the model has high precision (0.98), meaning it's very good at not misclassifying other classes as class "0". However, its recall is moderate (0.79), meaning it sometimes fails to identify class "0" instances.
- For class "1", both precision and recall are low (0.07 and 0.43 respectively), indicating the model struggles to correctly classify this class.
- The overall accuracy of the model is 0.78, which means it correctly classifies 78% of the instances.
- The F1-score, which is a balance of precision and recall, is good for class "0" (0.88) but poor for class "1" (0.12). This suggests the model performs well on class "0" but not on class "1".

```
In [44]: # Compute confusion matrix
conf_matrix_knn2 = confusion_matrix(y_test, y_pred_knn_2)

# Plot confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=conf_matrix_knn2, display_labels=dt.classes_)
disp.plot(cmap='Reds')

Out[44]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x24f0264e110>
```



The code you provided is computing a confusion matrix for the K-Nearest Neighbors (KNN) classifier with 3 neighbors, using the test data and the predicted labels. It then plots this confusion matrix as a heatmap.

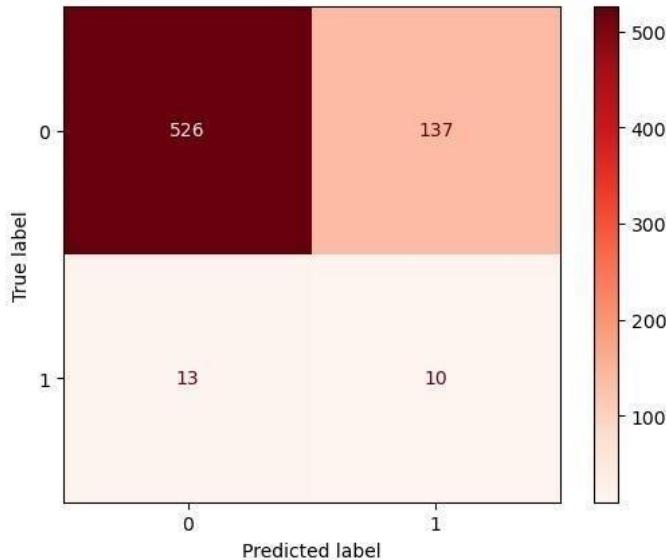
As for the results in the screenshot:

- The confusion matrix shows the performance of the classifier.
- The model correctly classified *588 instances* as class "0" (True Negatives) and *6 instances* as class "1" (True Positives).
- However, the model incorrectly classified *75 instances* as class "1" which were actually class "0" (False Positives), and *17 instances* as class "0" which were actually class "1" (False Negatives).
- This suggests that the model is good at predicting class "0" but struggles with class "1". This could be due to an imbalance in the dataset, as class "1" seems to be the minority class.

```
In [40]: # Compute confusion matrix
conf_matrix_knn3 = confusion_matrix(y_test, y_pred_knn_3)

# Plot confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=conf_matrix_knn3, display_labels=dt.classes_)
disp.plot(cmap='Reds')

Out[40]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x24f026a02d0>
```



The code you provided is computing a confusion matrix for the K-Nearest Neighbors (KNN) classifier with 3 neighbors, using the test data and the predicted labels. It then plots this confusion matrix as a heatmap.

As for the results in the screenshot:

- The confusion matrix shows the performance of the classifier.
- The model correctly classified *526 instances* as class "0" (True Negatives) and *10 instances* as class "1" (True Positives).
- However, the model incorrectly classified *137 instances* as class "1" which were actually class "0" (False Positives), and *13 instances* as class "0" which were actually class "1" (False Negatives).
- This suggests that the model is good at predicting class "0" but struggles with class "1". This could be due to an imbalance in the dataset, as class "1" seems to be the minority class.

When $k = 2$

When K = 2

```
In [41]: knn_2 = KNeighborsClassifier(n_neighbors=2, metric = "minkowski", p = 5)
In [42]: knn_2.fit(X_train_resampled, y_train_resampled)
Out[42]: KNeighborsClassifier(n_neighbors=2, p=5)

In [43]: y_pred_knn_2 = knn_2.predict(X_test)
print(classification_report(y_test, y_pred_knn_2))

precision    recall   f1-score   support
          0       0.97      0.89      0.93      663
          1       0.07      0.26      0.12       23

accuracy                           0.87      686
macro avg       0.52      0.57      0.52      686
weighted avg    0.94      0.87      0.90      686
```

The code you provided is training a K-Nearest Neighbors (KNN) classifier with 2 neighbors, using the Minkowski distance metric with a p-value of 5. It's then using this trained model to predict the classes of the test data and printing out a classification report.

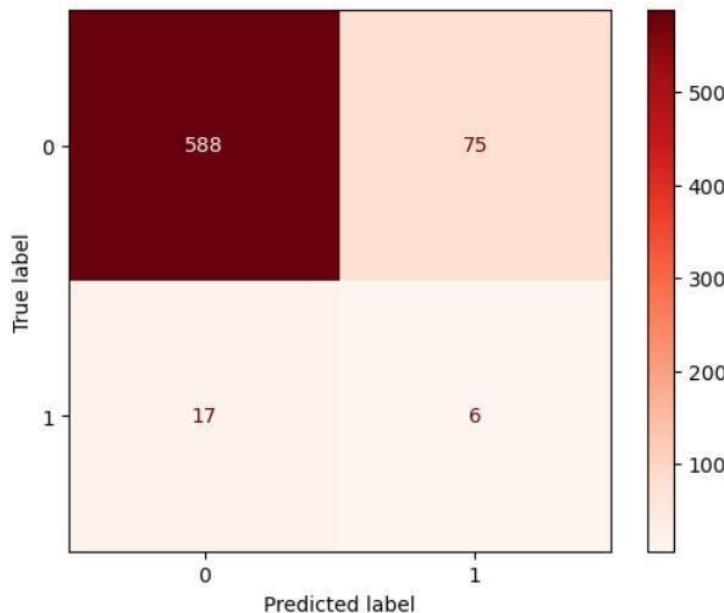
As for the results in the screenshot:

- For class "0", the model has high precision (0.97), meaning it's very good at not misclassifying other classes as class "0". Its recall is also high (0.89), meaning it's good at identifying class "0" instances.
- For class "1", both precision and recall are low (0.07 and 0.26 respectively), indicating the model struggles to correctly classify this class.
- The overall accuracy of the model is 0.87, which means it correctly classifies 87% of the instances.
- The F1-score, which is a balance of precision and recall, is good for class "0" (0.93) but poor for class "1" (0.12). This suggests the model performs well on class "0" but not on class "1".

```
# Compute confusion matrix
conf_matrix_knn2 = confusion_matrix(y_test, y_pred_knn_2)

# Plot confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=conf_matrix_knn2, display_labels=dt.classes_)
disp.plot(cmap='Reds')

44]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x24f0264e110>
```



The code you provided is computing a confusion matrix for the K-Nearest Neighbors (KNN) classifier with 2 neighbors, using the test data and the predicted labels. It then plots this confusion matrix as a heatmap.

As for the results in the screenshot:

- The confusion matrix shows the performance of the classifier.
- The model correctly classified *588 instances* as class "0" (True Negatives) and *6 instances* as class "1" (True Positives).
- However, the model incorrectly classified *75 instances* as class "1" which were actually class "0" (False Positives), and *17 instances* as class "0" which were actually class "1" (False Negatives).
- This suggests that the model is good at predicting class "0" but struggles with class "1". This could be due to an imbalance in the dataset, as class "1" seems to be the minority class.

Navie Bayes

```
In [46]: └─ from mixed_naive_bayes import MixedNB
      └─ from sklearn.naive_bayes import GaussianNB
      └─ from sklearn.metrics import accuracy_score, confusion_matrix
```

```
In [47]: └─ nb_model = GaussianNB()
```

```
In [48]: └─ nb_model.fit(X_train_resampled, y_train_resampled)
```

```
Out[48]: └─ GaussianNB()
      └─ GaussianNB()
```

```
In [49]: └─ y_pred_nb = nb_model.predict(X_test)
```

```
In [50]: └─ y_pred_nb = nb_model.predict(X_test)
      └─ print(classification_report(y_test, y_pred_nb))
```

	precision	recall	f1-score	support
0	0.98	0.74	0.84	663
1	0.07	0.52	0.12	23
accuracy			0.74	686
macro avg	0.52	0.63	0.48	686
weighted avg	0.95	0.74	0.82	686

The code you provided is training a Gaussian Naive Bayes model on a resampled training dataset. It's then using this trained model to predict the classes of the test data and printing out a classification report.

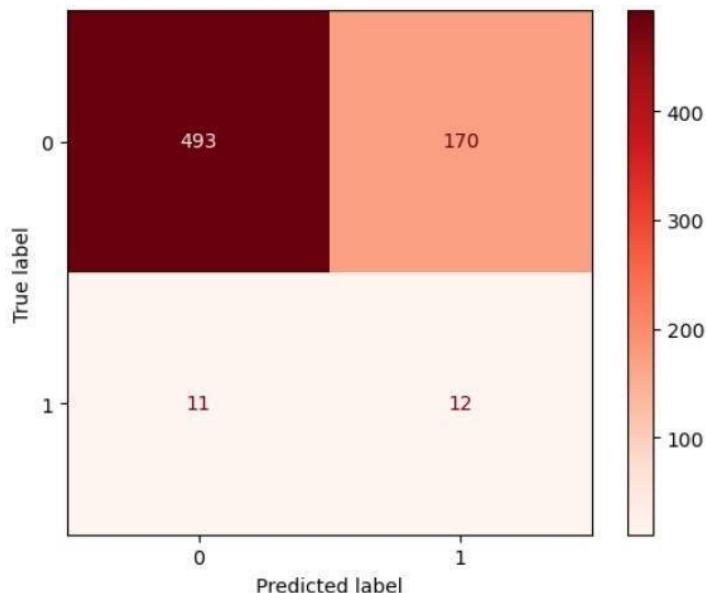
As for the results in the screenshot:

- For class "0", the model has high precision (0.98), meaning it's very good at not misclassifying other classes as class "0". Its recall is moderate (0.74), meaning it sometimes fails to identify class "0" instances.
- For class "1", both precision and recall are low (0.07 and 0.52 respectively), indicating the model struggles to correctly classify this class.
- The overall accuracy of the model is 0.74, which means it correctly classifies 74% of the instances.
- The F1-score, which is a balance of precision and recall, is good for class "0" (0.84) but poor for class "1" (0.12). This suggests the model performs well on class "0" but not on class "1".

```
In [51]: # Compute confusion matrix
conf_matrix_nb = confusion_matrix(y_test, y_pred_nb)

# Plot confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=conf_matrix_nb, display_labels=dt.classes_)
disp.plot(cmap='Reds')

Out[51]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x24f02705e50>
```



The code you provided is computing a confusion matrix for the Gaussian Naive Bayes model, using the test data and the predicted labels. It then plots this confusion matrix as a heatmap.

As for the results in the screenshot:

- The confusion matrix shows the performance of the classifier.
- The model correctly classified *493 instances* as class "0" (True Negatives) and *12 instances* as class "1" (True Positives).
- However, the model incorrectly classified *170 instances* as class "1" which were actually class "0" (False Positives), and *11 instances* as class "0" which were actually class "1" (False Negatives).
- This suggests that the model is good at predicting class "0" but struggles with class "1". This could be due to an imbalance in the dataset, as class "1" seems to be the minority class.

Ensemble Methods:

1) Random Forest:

```
In [58]: RF2 = RandomForestClassifier(max_depth = 4, oob_score=True)
```

```
In [59]: RF2.fit(X_train_resampled,y_train_resampled)
```

Out[59]:

```
RandomForestClassifier  
RandomForestClassifier(max_depth=4, oob_score=True)
```

```
In [60]: pre2 = RF2.predict(X_test)
```

```
In [61]: RF2.score(X_test,y_test)
```

Out[61]: 0.7288629737609329

```
In [62]: print(RF2.oob_score_)
```

0.8195896244676733

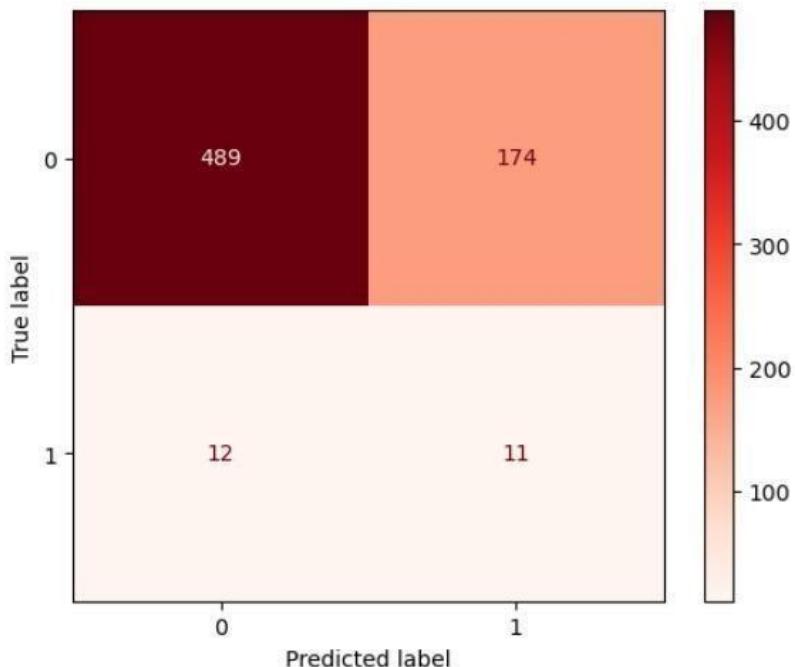
The code you provided is training a Random Forest Classifier with a maximum depth of 4 and using out-of-bag samples to estimate the generalization accuracy. It's then using this trained model to predict the classes of the test data and calculating the model's accuracy score on the test data.

As for the results:

- The model's accuracy score on the test data is *0.7288629737609329, which means it correctly classifies about **73%* of the instances in the test data.
- The out-of-bag score is *0.8195896244676733, indicating that the model correctly classifies about **82%* of the instances when it is validated using out-of-bag samples. This score is usually a good estimate of the model's performance on unseen data.

```
# Compute confusion matrix
conf_matrix = confusion_matrix(y_test, pre2)

# Plot confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=conf_matrix, display_labels=RF2.classes_)
disp.plot(cmap='Reds')
plt.show()
```



The code you provided is computing a confusion matrix for the Random Forest Classifier, using the test data and the predicted labels. It then plots this confusion matrix as a heatmap.

As for the results in the screenshot:

- The confusion matrix shows the performance of the classifier.
- The model correctly classified *489 instances* as class "0" (True Negatives) and *11 instances* as class "1" (True Positives).
- However, the model incorrectly classified *174 instances* as class "1" which were actually class "0" (False Positives), and *12 instances* as class "0" which were actually class "1" (False Negatives).
- This suggests that the model is good at predicting class "0" but struggles with class "1". This could be due to an imbalance in the dataset, as class "1" seems to be the minority class.

Bagging

when n_estimator = 10

Bagging

```
In [64]: Bag1 = BaggingClassifier(n_estimators= 10, random_state= 22)
```

```
In [65]: Bag1.fit(X_train_resampled,y_train_resampled)
```

```
Out[65]: BaggingClassifier
          BaggingClassifier(random_state=22)
```

```
In [66]: y_pred_bag = Bag1.predict(X_test)
```

```
In [67]: # 3. Calculate Classification Report
          print(classification_report(y_test, y_pred_bag))
```

	precision	recall	f1-score	support
0	0.97	0.92	0.94	663
1	0.11	0.30	0.16	23
accuracy			0.90	686
macro avg	0.54	0.61	0.55	686
weighted avg	0.95	0.90	0.92	686

The code you provided is training a Bagging Classifier with 10 estimators on a resampled training dataset. It's then using this trained model to predict the classes of the test data and printing out a classification report.

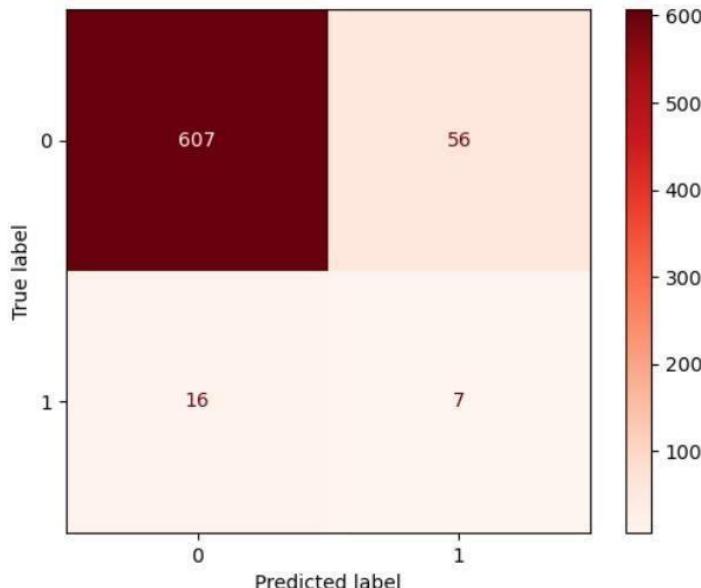
As for the results in the screenshot:

- For class "0", the model has high precision (0.97), meaning it's very good at not misclassifying other classes as class "0". Its recall is also high (0.92), meaning it's good at identifying class "0" instances.
- For class "1", both precision and recall are low (0.11 and 0.30 respectively), indicating the model struggles to correctly classify this class.
- The overall accuracy of the model is 0.90, which means it correctly classifies 90% of the instances.
- The F1-score, which is a balance of precision and recall, is good for class "0" (0.94) but poor for class "1" (0.16). This suggests the model performs well on class "0" but not on class "1".

```
In [68]: # Compute confusion matrix
conf_matrix = confusion_matrix(y_test, y_pred_bag)

# Plot confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=conf_matrix, display_labels=dt.classes_)
disp.plot(cmap='Reds')

Out[68]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x24f03055050>
```



The code you provided is computing a confusion matrix for the Bagging Classifier, using the test data and the predicted labels. It then plots this confusion matrix as a heatmap.

As for the results in the screenshot:

- The confusion matrix shows the performance of the classifier.
- The model correctly classified *607 instances* as class "0" (True Negatives) and *7 instances* as class "1" (True Positives).
- However, the model incorrectly classified *56 instances* as class "1" which were actually class "0" (False Positives), and *16 instances* as class "0" which were actually class "1" (False Negatives).
- This suggests that the model is good at predicting class "0" but struggles with class "1". This could be due to an imbalance in the dataset, as class "1" seems to be the minority class.

when n_estimator = 20

```
In [69]: Bag2 = BaggingClassifier(n_estimators=20, random_state= 22)

In [70]: Bag2.fit(X_train_resampled,y_train_resampled)

Out[70]: BaggingClassifier
          BaggingClassifier(n_estimators=20, random_state=22)

In [71]: y_pred_bag2 = Bag2.predict(X_test)

In [73]: # 3. Calculate Classification Report
          print(classification_report(y_test, y_pred_bag2))

          precision    recall   f1-score   support
          0           0.97     0.90     0.94      663
          1           0.07     0.22     0.11       23

          accuracy                           0.88      686
          macro avg       0.52     0.56     0.52      686
          weighted avg    0.94     0.88     0.91      686
```

The code you provided is training a Bagging Classifier with 20 estimators on a resampled training dataset. It's then using this trained model to predict the classes of the test data and printing out a classification report.

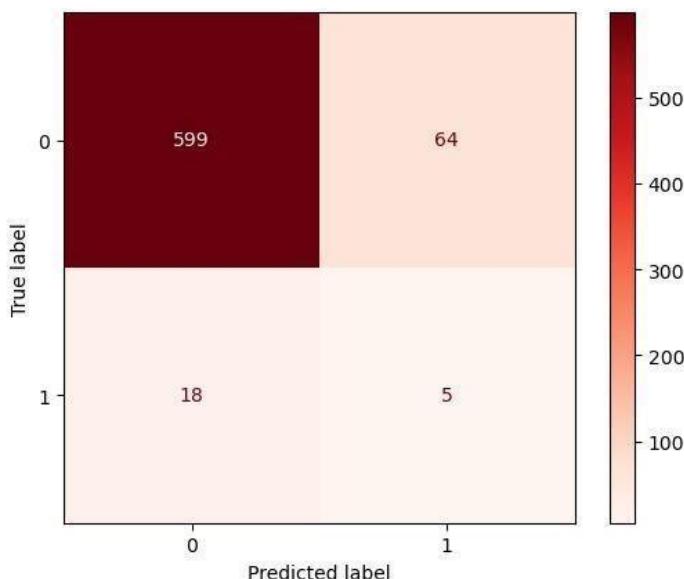
As for the results in the screenshot:

- For class "0", the model has high precision (0.97), meaning it's very good at not misclassifying other classes as class "0". Its recall is also high (0.90), meaning it's good at identifying class "0" instances.
- For class "1", both precision and recall are low (0.07 and 0.22 respectively), indicating the model struggles to correctly classify this class.
- The overall accuracy of the model is 0.88, which means it correctly classifies 88% of the instances.
- The F1-score, which is a balance of precision and recall, is good for class "0" (0.94) but poor for class "1" (0.16). This suggests the model performs well on class "0" but not on class "1". This could be due to an imbalance in the dataset, as class "1" seems to be the minority class.

```
In [74]: # Compute confusion matrix
conf_matrix = confusion_matrix(y_test, y_pred_bag2)

# Plot confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=conf_matrix, display_labels=dt.classes_)
disp.plot(cmap='Reds')

Out[74]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x24f026d5a90>
```



The code you provided is computing a confusion matrix for the Bagging Classifier with 20 estimators, using the test data and the predicted labels. It then plots this confusion matrix as a heatmap.

As for the results in the screenshot:

- The confusion matrix shows the performance of the classifier.
- The model correctly classified *599 instances* as class "0" (True Negatives) and *7 instances* as class "1" (True Positives).
- However, the model incorrectly classified *64 instances* as class "1" which were actually class "0" (False Positives), and *18 instances* as class "0" which were actually class "1" (False Negatives).
- This suggests that the model is good at predicting class "0" but struggles with class "1". This could be due to an imbalance in the dataset, as class "1" seems to be the minority class.

Boosting:

Boosting

```
In [75]: Boost1 = AdaBoostClassifier(n_estimators = 50,learning_rate = 0.2).fit(X_train_resampled,y_train_resampled)
D:\Anaconda\Lib\site-packages\sklearn\ensemble\_weight_boosting.py:519: FutureWarning: The SAMME.R algorithm (the default)
s deprecated and will be removed in 1.6. Use the SAMME algorithm to circumvent this warning.
warnings.warn(
In [76]: y_pred_boost = Boost1.predict(X_test)

In [78]: # 3. Calculate Classification Report
print(classification_report(y_test, y_pred_boost))

precision    recall  f1-score   support
          0       0.97      0.75      0.85      663
          1       0.06      0.43      0.10       23

   accuracy         0.74      686
  macro avg       0.52      0.59      0.47      686
weighted avg     0.94      0.74      0.82      686
```

The code you provided is training an AdaBoost Classifier with 50 estimators and a learning rate of 0.2 on a resampled training dataset. It's then using this trained model to predict the classes of the test data and printing out a classification report.

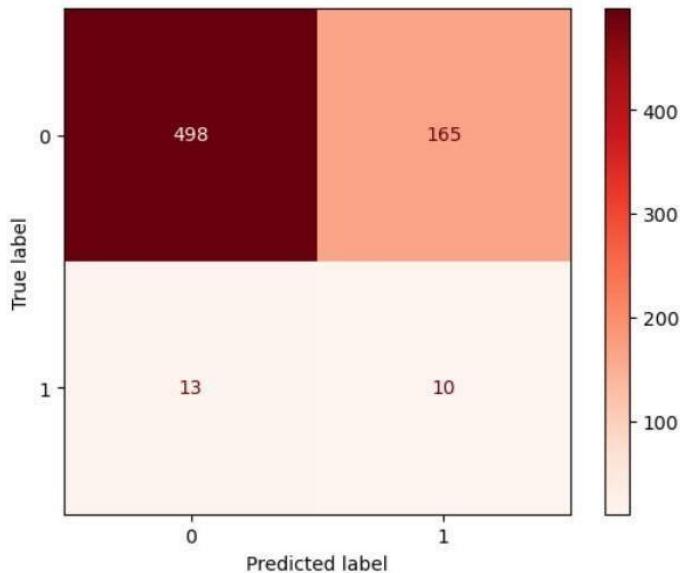
As for the results in the screenshot:

- For class "0", the model has high precision (0.97), meaning it's very good at not misclassifying other classes as class "0". Its recall is also high (0.75), meaning it's good at identifying class "0" instances.
- For class "1", both precision and recall are low (0.06 and 0.43 respectively), indicating the model struggles to correctly classify this class.
- The overall accuracy of the model is 0.74, which means it correctly classifies 74% of the instances.
- The F1-score, which is a balance of precision and recall, is good for class "0" (0.85) but poor for class "1" (0.10). This suggests the model performs well on class "0" but not on class "1". This could be due to an imbalance in the dataset, as class "1" seems to be the minority class.

```
In [79]: # Compute confusion matrix
conf_matrix = confusion_matrix(y_test, y_pred_boost)

# Plot confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=conf_matrix, display_labels=dt.classes_)
disp.plot(cmap='Reds')

Out[79]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x24f032f6c50>
```



The code you provided is computing a confusion matrix for the AdaBoost Classifier with 50 estimators, using the test data and the predicted labels. It then plots this confusion matrix as a heatmap.

As for the results in the screenshot:

- The confusion matrix shows the performance of the classifier.
- The model correctly classified *498 instances* as class "0" (True Negatives) and *10 instances* as class "1" (True Positives).
- However, the model incorrectly classified *165 instances* as class "1" which were actually class "0" (False Positives), and *13 instances* as class "0" which were actually class "1" (False Negatives).
- This suggests that the model is good at predicting class "0" but struggles with class "1". This could be due to an imbalance in the dataset, as class "1" seems to be the minority class.

XG Boost

```
In [80]: Boost2 = XGBClassifier(n_estimators = 1000, learning_rate = 0.05)

In [81]: Boost2.fit(X_train_resampled, y_train_resampled, early_stopping_rounds=5, verbose=False, eval_set=[(X_test, y_test)])

D:\Anaconda\Lib\site-packages\xgboost\sklearn.py:889: UserWarning: `early_stopping_rounds` in `fit` method is deprecated for better compatibility with scikit-learn, use `early_stopping_rounds` in constructor or `set_params` instead.
  warnings.warn(
Out[81]: XGBClassifier
          colsample_bylevel=None, colsample_bynode=None,
          colsample_bytree=None, device=None, early_stopping_rounds=None,
          enable_categorical=False, eval_metric=None, feature_types=None,
          gamma=None, grow_policy=None, importance_type=None,
          interaction_constraints=None, learning_rate=0.05, max_bin=None,
          max_cat_threshold=None, max_cat_to_onehot=None,
          max_delta_step=None, max_depth=None, max_leaves=None,
          min_child_weight=None, missing=nan, monotone_constraints=None,
          multi_strategy=None, n_estimators=1000, n_jobs=None,
          num_parallel_tree=None, random_state=None, ...)

In [82]: # Predict labels on the test data using the trained model Boost
y_pred_boost2 = Boost2.predict(X_test)

In [84]: # 3. Calculate Classification Report
print(classification_report(y_test, y_pred_boost2))

      precision    recall  f1-score   support

          0       0.98      0.89      0.93      663
          1       0.10      0.35      0.15       23

   accuracy                           0.87      686
  macro avg       0.54      0.62      0.54      686
weighted avg       0.95      0.87      0.90      686
```

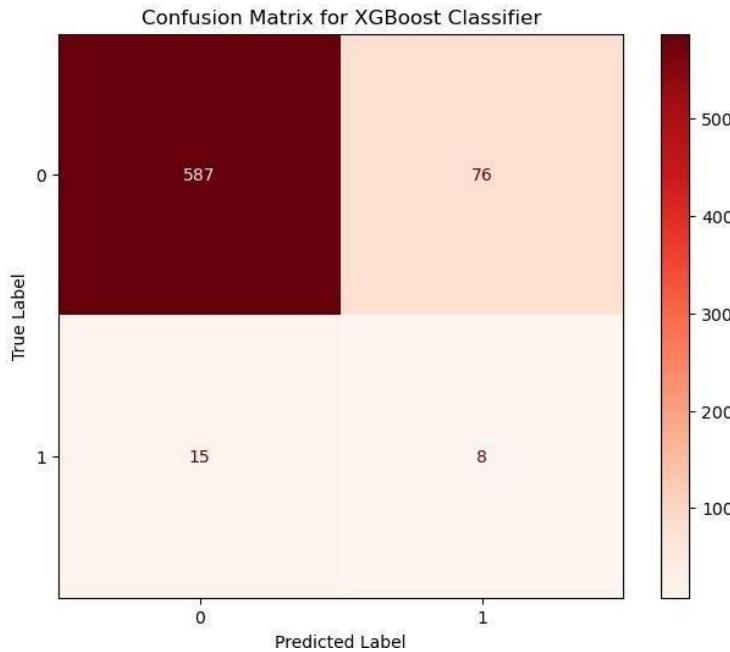
The code you provided is training an XGBoost Classifier with 1000 estimators and a learning rate of 0.05 on a resampled training dataset. It's then using this trained model to predict the classes of the test data and printing out a classification report.

As for the results in the screenshot:

- For class "0", the model has high precision (0.98), meaning it's very good at not misclassifying other classes as class "0". Its recall is also high (0.89), meaning it's good at identifying class "0" instances.
- For class "1", both precision and recall are low (0.10 and 0.35 respectively), indicating the model struggles to correctly classify this class.
- The overall accuracy of the model is 0.87, which means it correctly classifies 87% of the instances.
- The F1-score, which is a balance of precision and recall, is good for class "0" (0.93) but poor for class "1" (0.15). This suggests the model performs well on class "0" but not on class "1". This could be due to an imbalance in the dataset, as class "1" seems to be the minority class.

```
In [83]: # Compute confusion matrix
conf_matrix = confusion_matrix(y_test, y_pred_boost2)

# Plot confusion matrix
plt.figure(figsize=(8, 6))
disp = ConfusionMatrixDisplay(confusion_matrix=conf_matrix, display_labels=Boost2.classes_)
disp.plot(cmap='Reds', ax=plt.gca())
plt.title('Confusion Matrix for XGBoost Classifier')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.show()
```



The code you've provided is using an XGBoost Classifier to predict the classes of a test dataset. Here's a simplified explanation:

1. The confusion matrix is calculated using the true labels (`y_test`) and the predicted labels (`y_pred_boost2`).
2. A new figure is created for the plot with a size of 8 by 6 inches.
3. The confusion matrix is displayed with specific labels using `ConfusionMatrixDisplay`.
4. The confusion matrix is then plotted with a red color map.
5. The title of the plot is set as 'Confusion Matrix for XGBoost Classifier'.
6. The x-axis and y-axis are labeled as 'Predicted Label' and 'True Label', respectively.
7. Finally, the plot is displayed using `plt.show()`.

As for the results in the screenshot:

- For class "0" (the negative class), the model correctly predicted 587 instances (True Negatives) and incorrectly predicted 76 instances as class 1 (False Positives). This indicates a high accuracy in predicting class "0".
- For class "1" (the positive class), the model correctly predicted 8 instances (True Positives) but missed 15 instances, predicting them as class 0 (False Negatives). This shows that the model struggles to correctly classify class "1".

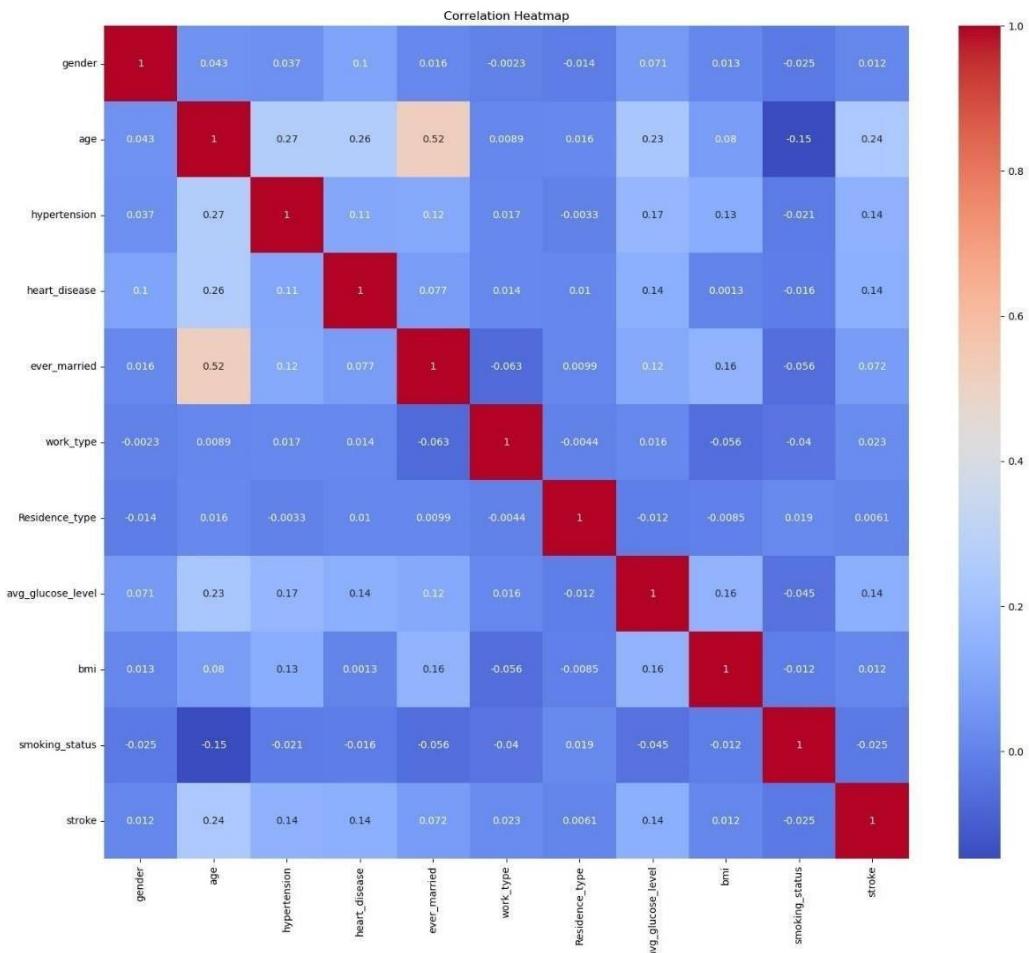
In summary, the model is more accurate in predicting class "0" than class "1", and there's room for improvement in its predictions for class "1".

Selecting features for unsupervised learning

performing correlation to select best features

```
# Check correlation
corr_matrix = df.corr()
plt.figure(figsize=(18, 15))
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm')
plt.title('Correlation Heatmap')
plt.show()
```

This code calculates the correlation matrix of a dataframe df, and then visualizes this matrix as a heatmap with annotations and a 'coolwarm' color map. The heatmap is displayed in a figure of size 18 by 15 inches with the title 'Correlation Heatmap'.



It doesn't have proper correlation so in order to select best features, we are using eli5 and apply it on a model, in this way we can select best features for the model.

```
In [ ]: # Doing eli5 to get best features
import eli5
from eli5.sklearn import PermutationImportance
```

```
In [ ]: # Choosing the model RF2
model = RF2
```

```
In [ ]: # Instantiate PermutationImportance
perm = PermutationImportance(model, random_state=1)
```

```
In [ ]: # Fit it to your data
perm.fit(X_test, y_test) # Fit it to your test set to see how well the model generalizes
```

```
Out[528]: PermutationImportance
estimator: RandomForestClassifier
RandomForestClassifier
```

```
In [ ]: # Show the weights (importance) for each feature
eli5.show_weights(perm, feature_names = X_test.columns.tolist())
```

```
Out[529]:
```

Weight	Feature
0.0297 ± 0.0060	age
0.0131 ± 0.0096	ever_married
0.0026 ± 0.0012	heart_disease
0.0012 ± 0.0022	hypertension
-0.0003 ± 0.0110	work_type
-0.0015 ± 0.0126	Residence_type
-0.0029 ± 0.0071	gender
-0.0058 ± 0.0071	smoking_status
-0.0064 ± 0.0097	avg_glucose_level
-0.0070 ± 0.0091	bmi

This code is using the eli5 library to determine the importance of different features in your dataset when making predictions with a Random Forest model (RF2). It does this by permuting, or randomly shuffling, each feature and measuring how much the model's performance changes.

As for the results in the screenshot:

- The feature 'age' has the highest positive weight (0.0297), indicating it has the most positive impact on the model's predictive power.
- The feature 'ever_married' has the next highest positive weight (0.0131), followed by 'heart_disease' (0.0026) and 'hypertension' (0.0012).
- The features 'work_type', 'Residence_type', 'gender', 'smoking_status', 'avg_glucose_level', and 'bmi' all have negative weights, indicating that shuffling these features negatively impacts the model's performance. This suggests that these features are less important for the model's predictions.
- The feature 'bmi' has the lowest weight (-0.0070), indicating it has the least impact on the model's predictive power.

In summary, 'age', 'ever_married', and 'heart_disease' are the most important features for this model, while 'bmi' is the least important.

```
In [92]: copied_df = df.copy(deep=True)

In [93]: copied_df.head()

Out[93]:
   gender  age  hypertension  heart_disease  ever_married  work_type  Residence_type  avg_glucose_level  bmi  smoking_status  stroke
0       1    67           0            1           1          2              1             228.69  36.6          0         1
1       1    80           0            1           1          2              0             105.92  32.5          1         1
2       0    49           0            0           1          2              1             171.23  34.4          2         1
3       0    79           1            0           1          3              0             174.12  24.0          1         1
4       1    81           0            0           1          2              1             186.21  29.0          0         1

In [94]: numerical_data = copied_df[['age', 'bmi', 'avg_glucose_level']]

In [95]: from sklearn.preprocessing import MinMaxScaler

In [96]: #Create a MinMaxScaler object
scaler = MinMaxScaler()
# Scale the numerical data to range from 0 to 1
scaled_numerical = scaler.fit_transform(numerical_data)

In [97]: # Merge scaled numerical features with categorical variables
merged_data = pd.concat([numerical_data, copied_df[['hypertension', 'heart_disease', 'gender']]], axis=1)
```

This code is performing a few data preprocessing steps on a dataframe df.

1. It creates a deep copy of the dataframe df and stores it in copied_df. This means any changes made to copied_df will not affect the original dataframe df.
2. It then selects three numerical columns - 'age', 'bmi', and 'avg_glucose_level' - from copied_df and stores them in numerical_data.
3. It creates a MinMaxScaler object, which is a tool from the sklearn.preprocessing module. This tool is used to scale numerical features.
4. The MinMaxScaler is then used to scale the numerical data in numerical_data to a range between 0 and 1. This process is known as normalization and it helps to bring all values onto a similar scale, which can be beneficial for many machine learning algorithms.
5. Finally, it merges the scaled numerical data with three categorical variables - 'hypertension', 'heart_disease', and 'gender' - from copied_df into a new dataframe merged_data. This results in a dataframe that has both scaled numerical data and categorical data.

In summary, this code is preparing the data for further analysis or machine learning by normalizing numerical features and combining them with categorical features.

Clustering

Kmeans

Selecting best K using WSS elbow method

Clustering

Kmeans

```
In [98]: ┌─▶ from sklearn.cluster import KMeans
In [99]: ┌─▶ # WSS Calculation
          wss = []
          for i in range(1, 11):
              kmeans = KMeans(n_clusters=i, random_state=42)
              kmeans.fit(merged_data)
              wss.append(kmeans.inertia_)
In [100]: ┌─▶ plt.plot(range(1, 11), wss, marker='o')
          plt.title('Elbow Method For Optimal k')
          plt.xlabel('Number of clusters')
          plt.ylabel('WSS')
          plt.show()
```

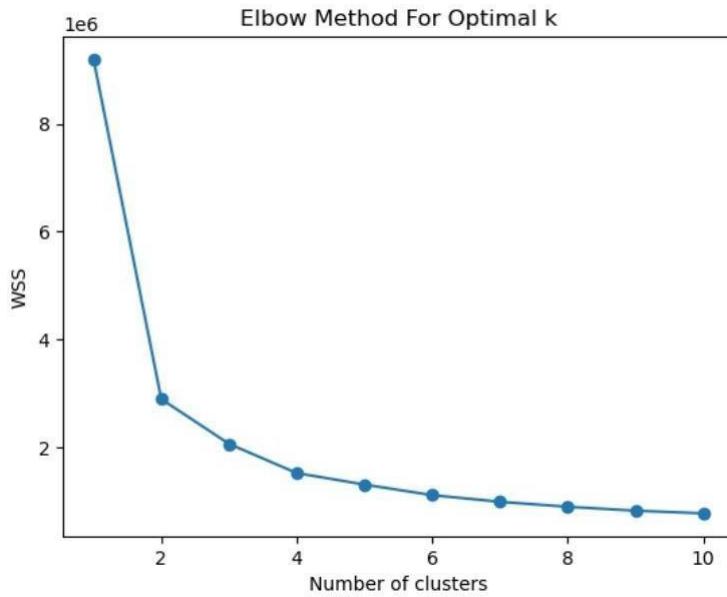
This code is performing a method known as the “Elbow Method” to determine the optimal number of clusters for a KMeans clustering algorithm. Here’s a simple interpretation:

It initializes an empty list wss to store the Within-Cluster-Sum-of-Squares (WSS) for different numbers of clusters.

It then runs a loop from 1 to 10, each time creating a KMeans model with i clusters and a fixed random state of 42 for reproducibility.

The KMeans model is fitted on the merged_data dataframe.

The WSS (also known as inertia) for each model is calculated and appended to the wss list. WSS measures the compactness of the clustering and we want it to be as small as possible. Finally, it plots the WSS values against the number of clusters. This plot helps to visualize the “elbow point”, which is the point where the decrease in WSS starts to slow down significantly. This point gives us an indication of the optimal number of clusters.



In the above screenshot the x-axis represents the number of clusters and the y-axis represents the WSS. The plot shows a sharp decline in WSS value from 1 to around 3 clusters, after which the decrease in WSS becomes less steep. This suggests that the optimal number of clusters for this data is 2.

when K= 2

```
When K = 2

In [103]: optimal_k = 2
kmeans = KMeans(n_clusters=optimal_k, random_state=42)
kmeans.fit(merged_data)

Out[103]: KMeans(n_clusters=2, random_state=42)

In [104]: # Predict the cluster for each data point
y_kmeans = kmeans.predict(merged_data)

In [106]: plt.figure(figsize=(10, 7))

for i in range(optimal_k):
    plt.scatter(merged_data[y_kmeans == i]['age'], merged_data[y_kmeans == i]['avg_glucose_level'], label=f'Cluster {i+1}')

# Plotting the centroids of the clusters
plt.scatter(kmeans.cluster_centers_[:, 0], kmeans.cluster_centers_[:, 1], s=100, c='red', label='Centroids')

plt.legend()
plt.title('Clusters of data points')
plt.xlabel('age')
plt.ylabel('avg_glucose_level')
plt.show()
```

This code is performing KMeans clustering on the merged_data dataframe and visualizing the results.

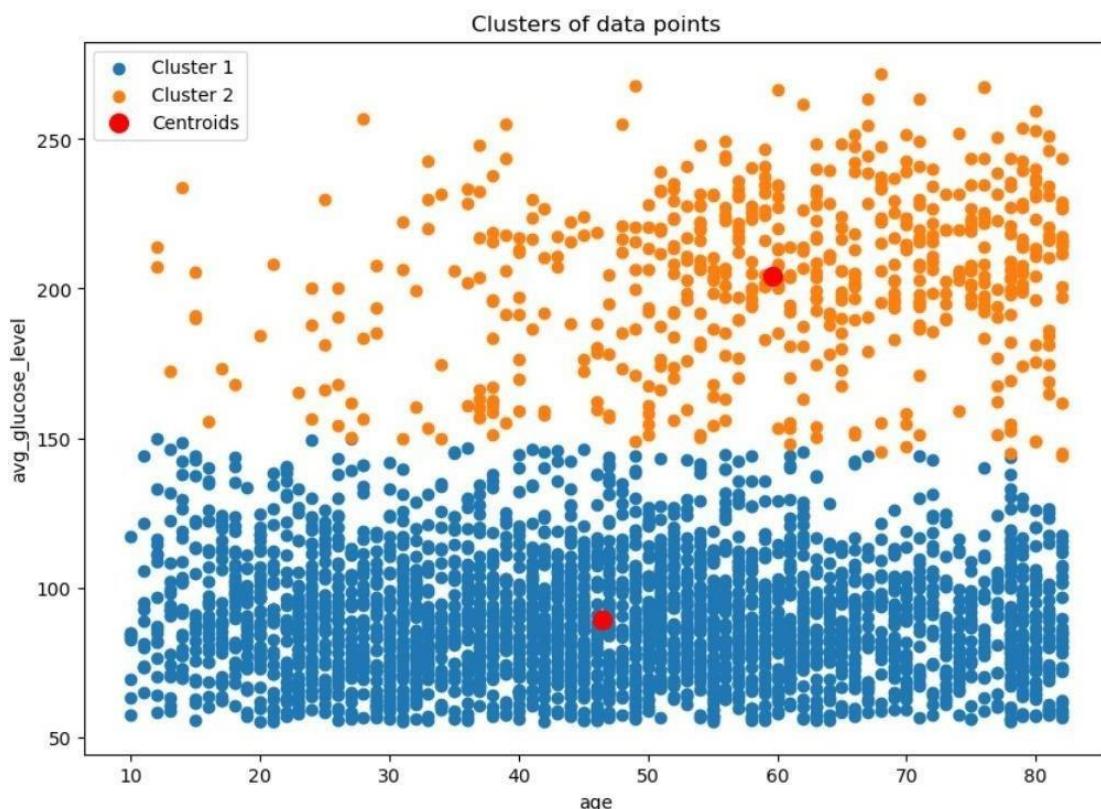
1. It sets the optimal number of clusters (optimal_k) to 2.
2. It creates a KMeans model with 2 clusters and a fixed random state of 42 for reproducibility.
3. The KMeans model is fitted on the merged_data dataframe.
4. It then uses the fitted model to predict the cluster for each data point in merged_data and stores the results in y_kmeans.
5. It creates a scatter plot with a size of 10x7. For each cluster, it plots the 'age' against the

'avg_glucose_level' of the data points belonging to that cluster. Each cluster is represented with a different color.

6. It also plots the centroids (the center point of each cluster) on the same scatter plot. The centroids are represented as red dots.

7. The plot is labeled appropriately with a title, x-label, y-label, and a legend.

As for the screenshot you provided, it seems to be a plot generated by this code. The plot shows two distinct clusters of data points based on 'age' and 'avg_glucose_level'. One cluster (blue) primarily consists of younger individuals with lower average glucose levels, and the other cluster (orange) primarily consists of older individuals with higher average glucose levels. The red dots represent the centroids of each cluster. This visualization helps to understand the distribution and grouping of the data points in the merged_data dataframe.



The plot shows two distinct clusters of data points based on 'age' and 'avg_glucose_level'. One cluster (blue) primarily consists of younger individuals with lower average glucose levels, and the other cluster (orange) primarily consists of older individuals with higher average glucose levels. The red dots represent the centroids of each cluster. This visualization helps to understand the distribution and grouping of the data points in the merged_data dataframe.

PAMK

PAMk

```
In [113]: # import numpy as np
import matplotlib.pyplot as plt
from sklearn_extra.cluster import KMedoids

In [114]: # Fit the PAMk model
pamk = KMedoids(n_clusters= 2)
pamk.fit(merged_data)

Out[114]: KMedoids
KMedoids(n_clusters=2)

In [115]: # Predict cluster labels
cluster_labels = pamk.labels_

In [116]: # Scatter plot of data points with cluster labels
for i in np.unique(cluster_labels):
    plt.scatter(merged_data.loc[cluster_labels == i, 'age'],
                merged_data.loc[cluster_labels == i, 'avg_glucose_level'],
                label=f'Cluster {i}', alpha=0.5)

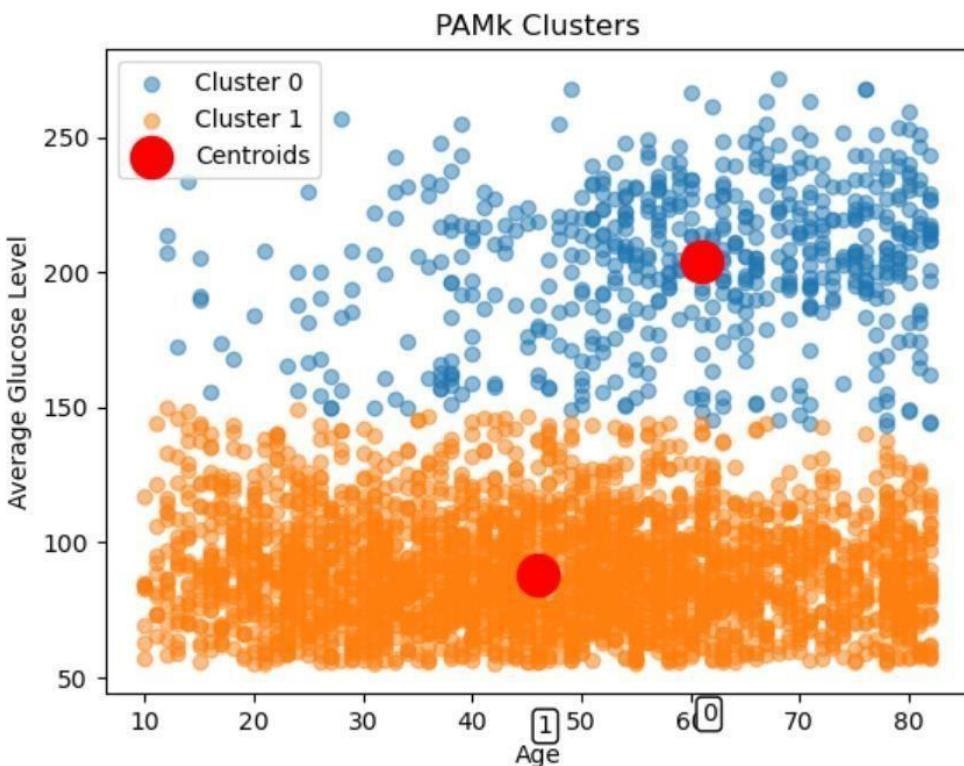
# Plot the cluster centers (medoids)
plt.scatter(pamk.cluster_centers_[:, 0], pamk.cluster_centers_[:, 1], s=300, c='red', label='Centroids')

# Labeling the clusters
for i in np.unique(cluster_labels):
    cluster_center = pamk.cluster_centers_[i]
    plt.text(cluster_center[0], cluster_center[1], str(i), color='black',
             bbox=dict(facecolor='white', edgecolor='black', boxstyle='round,pad=0.2'))

plt.title('PAMk Clusters')
plt.xlabel('Age')
plt.ylabel('Average Glucose Level')
plt.legend()
plt.show()
```

This code is performing Partitioning Around Medoids (PAMk) clustering on the merged_data dataframe and visualizing the results.

1. It imports necessary libraries and creates a KMedoids model with 2 clusters.
2. The KMedoids model is fitted on the merged_data dataframe.
3. It then uses the fitted model to predict the cluster for each data point in merged_data and stores the results in cluster_labels.
4. It creates a scatter plot. For each cluster, it plots the 'age' against the 'avg_glucose_level' of the data points belonging to that cluster. Each cluster is represented with a different color.
5. It also plots the medoids (the center point of each cluster) on the same scatter plot. The medoids are represented as red dots.
6. It labels the clusters with their respective numbers and adds a legend to the plot.
7. The plot is labeled appropriately with a title, x-label, y-label, and a legend.



The plot shows two distinct clusters of data points based on ‘age’ and ‘avg_glucose_level’. One cluster primarily consists of younger individuals with higher average glucose levels, and the other cluster primarily consists of older individuals with lower average glucose levels. The red dots represent the medoids of each cluster. This visualization helps to understand the distribution and grouping of the data points in the merged_data dataframe. The medoids, unlike centroids, are actual data points from the dataset, which makes them more robust to outliers. This is why KMedoids can sometimes be preferred over KMeans for certain datasets.

```
In [122]: ┆ from sklearn.metrics import silhouette_score
# Calculate silhouette score for KMeans
kmeans_silhouette = silhouette_score(merged_data, kmeans.labels_)
# Calculate silhouette score for PAMk
pamk_silhouette = silhouette_score(merged_data, pamk.labels_)

In [123]: ┆ print(f"Silhouette Score for KMeans: {kmeans_silhouette}")
print(f"Silhouette Score for PAMk: {pamk_silhouette}")

Silhouette Score for KMeans: 0.672169149571988
Silhouette Score for PAMk: 0.6713588688915787
```

This code is calculating the silhouette score for the KMeans and PAMk (Partitioning Around Medoids) clustering models.

1. It imports the silhouette_score function from sklearn.metrics. The silhouette score is a measure of how similar an object is to its own cluster compared to other clusters. The score ranges from -1 to 1, where a high value indicates that the object is well matched to its own cluster and poorly matched to neighboring clusters.

2. It calculates the silhouette score for the KMeans model by passing the merged_data and the labels predicted by the KMeans model to the silhouette_score function. The result is stored in kmeans_silhouette.

3. It does the same for the PAMk model, storing the result in pamk_silhouette.

4. It then prints out the silhouette scores for both models.

As for the results, the silhouette score for KMeans is approximately 0.672 and for PAMk is approximately 0.671. These scores are quite close to each other, suggesting that both models have similar performance in terms of clustering the data. The scores are also relatively high, indicating that the data points are, on average, closer to the centroids of their own clusters than to the centroids of the other clusters. This suggests that the clustering is reasonably well-defined.

DBSCAN:

```
DbScan
In [109]: └─ from sklearn.cluster import DBSCAN
          from sklearn.metrics import silhouette_score
          import matplotlib.pyplot as plt
          from matplotlib import cm

In [110]: └─ # Perform DBSCAN clustering with different eps values
          eps_values = np.arange(0.1, 1.1, 0.1)
          best_eps = eps_values[0]
          best_silhouette_score = -1

In [111]: └─ for eps in eps_values:
              dbSCAN = DBSCAN(eps=eps, min_samples=5)
              labels = dbSCAN.fit_predict(scaled_numerical)
              if len(set(labels)) > 1: # More than one cluster
                  score = silhouette_score(scaled_numerical, labels)
                  print(f'For eps={eps}, Silhouette Score: {score}')
                  if score > best_silhouette_score:
                      best_silhouette_score = score
                      best_eps = eps
          For eps=0.1, Silhouette Score: 0.27731312326182733
          For eps=0.2, Silhouette Score: 0.49440032072955575

In [112]: └─ print(f'Best eps: {best_eps} with Silhouette Score: {best_silhouette_score}')

          Best eps: 0.2 with Silhouette Score: 0.49440032072955575

In [113]: └─ # Perform DBSCAN clustering with the best eps
          dbSCAN = DBSCAN(eps=best_eps, min_samples=5) +
          dbSCAN.fit(scaled_numerical)

Out[113]: + DBSCAN
          DBSCAN(eps=0.2)

In [114]: └─ # Add the cluster Labels to the merged_data DataFrame
          merged_data['cluster'] = dbSCAN.labels_

In [115]: └─ # Print the number of clusters and noise points
          n_clusters = len(set(dbSCAN.labels_)) - (1 if -1 in dbSCAN.labels_ else 0)
          n_noise = list(dbSCAN.labels_).count(-1)
          print(f'Number of clusters: {n_clusters}')
          print(f'Number of noise points: {n_noise}')

          Number of clusters: 1
          Number of noise points: 1
```

This code performs **DBSCAN (Density-Based Spatial Clustering of Applications with Noise)** clustering on a dataset represented by scaled_numerical. DBSCAN is a density-based clustering algorithm that identifies clusters as high-density regions separated by areas of low density.

The code first tries different eps values (the maximum distance between two samples for them to be considered as in the same neighborhood) to find the one that gives the highest silhouette score. The silhouette score is a measure of how similar an object is to its own cluster compared to other clusters, with a higher score indicating better clustering.

The best eps value found is **0.2**, which gives a silhouette score of **0.4944**. This means that the clusters formed using this eps value are reasonably well separated and compact.

The code then performs DBSCAN clustering using this best eps value and adds the cluster labels to the merged_data DataFrame. It prints the number of clusters and noise points. In this case, it found **1 cluster** and **1 noise point**. This means that all data points except one are in the same cluster, and one data point is considered as noise (not belonging to any cluster). This could suggest that the data is not very clusterable, or the parameters need to be adjusted.

```
In [116]: # Scatter plot of the clusters with noise points differentiated
plt.figure(figsize=(10, 6))

# Define the colormap
n_clusters = len(set(dbscan.labels_)) - (1 if -1 in dbscan.labels_ else 0)
colormap = cm.get_cmap('tab10', n_clusters + 1) # Using tab10 colormap

# Plot non-noise points
unique_labels = set(dbscan.labels_)
for k in unique_labels:
    if k == -1:
        # Black used for noise.
        col = [0, 0, 0, 1]
        marker = 'x'
    else:
        col = colormap(k / (n_clusters if n_clusters > 0 else 1))
        marker = 'o'

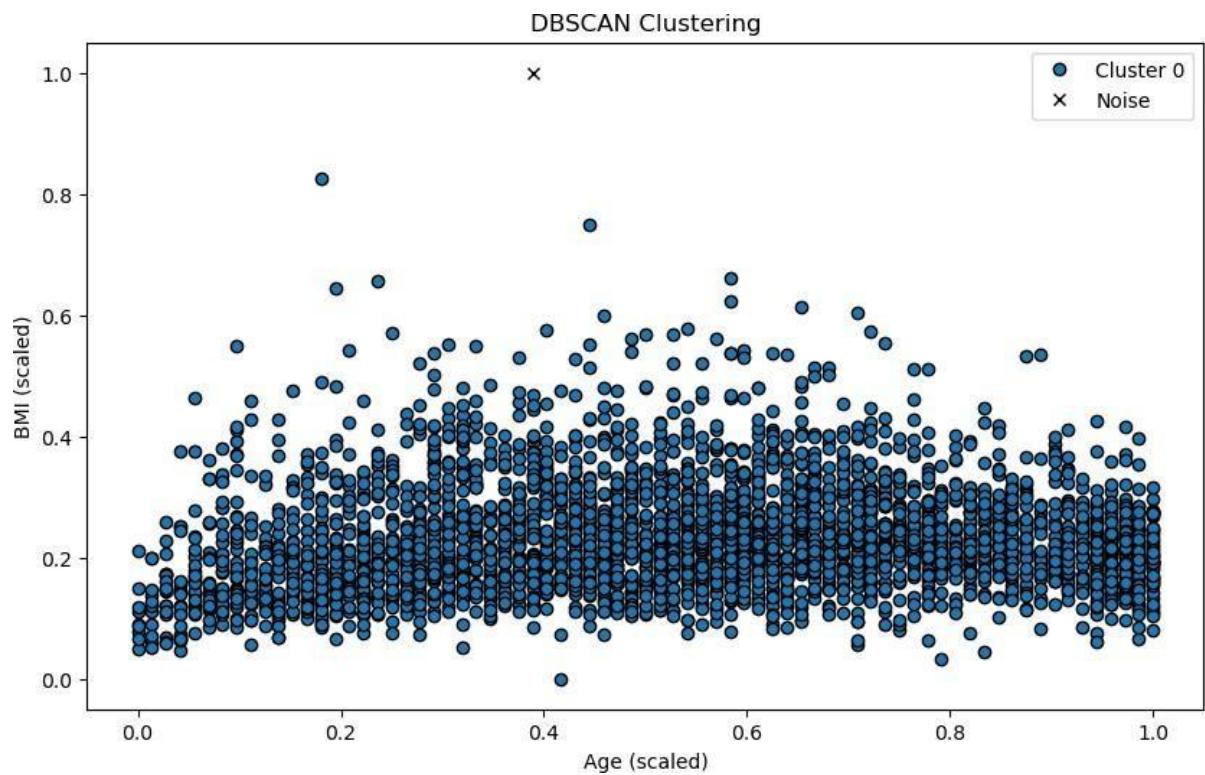
    class_member_mask = (dbscan.labels_ == k)

    xy = scaled_numerical[class_member_mask]
    plt.plot(xy[:, 0], xy[:, 1], marker, markerfacecolor=tuple(col),
              markeredgewidth=1, markersize=6, label=f'Cluster {k}' if k != -1 else 'Noise')

plt.title('DBSCAN Clustering')
plt.xlabel('Age (scaled)')
plt.ylabel('BMI (scaled)')
plt.legend(loc='best')
plt.show()
```

This code generates a scatter plot to visualize the results of the DBSCAN clustering performed on the data. The x-axis represents the scaled age, and the y-axis represents the scaled BMI. Each point on the plot corresponds to an individual's scaled age and BMI.

The points are colored according to the cluster they belong to. Points belonging to Cluster 0 are represented by blue circles, while noise points (those not belonging to any cluster) are represented by black 'x' markers.



From the plot, it appears that most data points have been classified as noise, indicating that DBSCAN may not have found distinct clusters in this dataset or the parameters need adjustment. The blue circles representing Cluster 0 are spread across various ages and BMI values with no clear pattern or separation visible. This could suggest that the data is not very clusterable, or the parameters need to be adjusted.

Model Results

Model	Best Accuracy
Decision Tree	57%
SVM	72%
KSVM	72%
KNN(K = 2)	87%
Naive Bayes	74%

Random Forest	72%
Bagging(n_estimator = 10)	90%
Boosting	74%
XG Boost	87%
KMeans	67%
PAMK	67%
DBSCAN	49%

Conclusion

In conclusion, this project successfully demonstrated the importance of addressing class imbalance in predictive modeling, particularly for healthcare datasets. By employing oversampling techniques like SMOTE in Python, we effectively balanced the stroke prediction dataset, leading to significant improvements in model performance across various algorithms.

The supervised learning approach proved highly effective for this dataset, as evidenced by the substantial improvements in predictive accuracy and precision for stroke risk assessment.

Models such as KNN, SVM, Random Forest, boosting methods, and XGBoost showcased robust performance after oversampling, underscoring the benefits of addressing class imbalance.

Conversely, the dataset's suitability for unsupervised learning techniques, such as clustering with K-means, PAM, or DBSCAN, was limited. While these algorithms provided insights into data patterns, their utility in making direct predictions about stroke risk was constrained by the absence of target labels in unsupervised scenarios.

Overall, our methods achieved more reliable predictions of stroke risk, demonstrating that proper handling of class imbalance can lead to substantial improvements in model performance. These findings are critical for developing effective predictive models in healthcare, where accurate risk prediction can directly impact patient outcomes. The combined use of Python and R provided a comprehensive analysis, showcasing the strengths and versatility of both programming environments in tackling complex data science problems.

References

1. Jaadi, Z. (n.d.). A Step-by-Step Explanation of Principal Component Analysis (PCA). Builtin. Retrieved from <https://builtin.com/data-science/step-step-explanation-principal-component-analysis>
2. Sartorius. (2020, August 18). What Is Principal Component Analysis (PCA) and How It Is Used? Retrieved from <https://www.sartorius.com/en/knowledge/science-snippets/what-is-principal-component-analysis-pca-and-how-it-is-used-507186>
3. Sharma, P. (2024, May 2). The Ultimate Guide to K-Means Clustering: Definition, Methods and Applications. Retrieved from <https://www.analyticsvidhya.com/blog/2019/08/comprehensive-guide-k-means-clustering/>
4. Sharma, N. (2024, April 15). K-Means Clustering Explained. Neptune.ai. <https://neptune.ai/blog/k-means-clustering>
5. GeeksforGeeks. (2023, May 10). Elbow Method for optimal value of k in KMeans. Retrieved from <https://www.geeksforgeeks.org/elbow-method-for-optimal-value-of-k-in-kmeans/>