

AMERICAN INTERNATIONAL UNIVERSITY-BANGLADESH (AIUB) FACULTY OF SCIENCE AND TECHNOLOGY

FINAL TERM PROJECT

INTRODUCTION TO DATA SCIENCE

Summer 2022-2023

Section: B

Submitted By

Name: Md. Nasifur Rahman

ID: 20-43651-2

Submitted To

Touhedul Islam

Assistant Professor

Department of Computer Science

Date of Submission: August 16,2023

Source: https://www.kaggle.com/datasets/uciml/mushroom-classification

Introduction:

This dataset includes descriptions of hypothetical samples corresponding to 23 species of gilled mushrooms in the Agaricus and Lepiota Family Mushroom drawn from The Audubon Society Field Guide to North American Mushrooms Each species is identified as definitely edible, poisonous, or of unknown edibility and not recommended. This latter class was combined with the poisonous one. The Guide clearly states that there is no simple rule for determining the edibility of a mushroom; no rule like "leaflets three, let it be" for Poisonous Oak and Ivy.

Attribute Information:

Class: Tt define is it edible or poisonous.

Cap Shape: The shape of the mushroom's cap (e.g., convex, flat, bell, etc.).

Cap Surface: The texture of the mushroom's cap (e.g., fibrous, smooth, scaly, etc.).

Cap Color: The color of the mushroom's cap (e.g., brown, yellow, red, etc.).

Bruises: Presence or absence of bruises on the mushroom (bruises, no bruises).

Odor: The odor of the mushroom (e.g., almond, foul, anise, etc.).

Gill Attachment: How the gills are attached to the cap (e.g., free, attached).

Gill Spacing: The spacing between gills (e.g., close, crowded).

Gill Size: The size of the gills (e.g., broad, narrow).

Gill Color: The color of the gills (e.g., black, pink, white, etc.).

Stalk Shape: The shape of the mushroom's stalk (e.g., enlarging, tapering).

Stalk Root: The root type of the stalk (e.g., bulbous, club, equal, rooted).

Stalk Surface Above Ring: The texture of the stalk above the ring (e.g., fibrous, smooth).

Stalk Surface Below Ring: The texture of the stalk below the ring (e.g., fibrous, smooth).

Stalk Color Above Ring: The color of the stalk above the ring (e.g., brown, white).

Stalk Color Below Ring: The color of the stalk below the ring (e.g., brown, white).

Veil Type: The type of veil on the mushroom (partial, universal).

Veil Color: The color of the veil (e.g., brown, white).

Ring Number: The number of rings on the stalk (e.g., none, one, two).

Ring Type: The type of ring on the stalk (e.g., cobwebby, evanescent).

Spore Print Color: The color of the spore print (e.g., black, brown, white, etc.).

Population: The abundance of mushrooms in a specific population (e.g., abundant, clustered).

Habitat: The habitat where the mushrooms are typically found (e.g., grasses, woods, meadows, etc.).

Import Dataset:

```
1 library(caret)
2 library(ggplot2)
3 library(lattice)
4 library(readr)
5 library(class)
6
7 data<- read.csv("E:/mushrooms.csv",header = TRUE,sep=',')|
8 str(data)</pre>
```

Output:

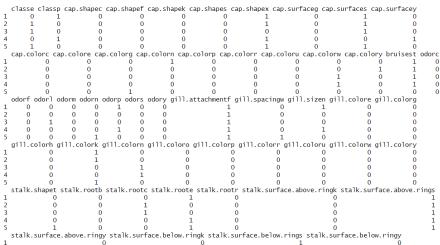
Explanation:

- **caret:** A comprehensive package for training and evaluating machine learning models.
- **ggplot2:** A popular package for creating data visualizations.
- **lattice:** A package for creating lattice plots, which are similar to traditional scatterplots but offer more flexibility.
- readr: A package for reading and parsing data efficiently.
- **class:** A package that provides various classification algorithms, including the knearest neighbors (k-NN) algorithm.
- **str():** The str() function is used to display the structure of the data object. It provides information about the data frame's structure, including the names of its columns, data types of each column, and a preview of the data.

Data Pre-processing:

• There is no missing value.

• Converted all values in numeric.



Calculated Z-Score normalization

```
zscore_normalized_data <- scale(encoded_data)</pre>
           head(zscore_normalized_data)

      Classe
      classp
      cap.shapec
      cap.shapef
      cap.shapek
      cap.shapes
      cap.shapex
      cap.shapex

0.9646211 -0.9646211 -0.02219347 -0.7961608 -0.3368573 -0.06288113 -1.0945233 -0.02219347 

4 -1.0365489 1.0365489 -0.02219347 -0.7961608 -0.3368573 -0.06288113 1.1054186 -0.02219347 

5 0.9646211 -0.9646211 -0.02219347 -0.7961608 -0.3368573 -0.06288113 1.1054186 -0.02219347
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          1.4758508
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                -0.6774919
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   1.4758508
                   0.9646211 -0.9646211 -0.02219347 -0.7961608 -0.3368573 -0.06288113 1.1054186
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               -0.02219347
             cap.surfacey cap.colorc cap.colore cap.colorg cap.colorn cap.colorn cap.colorr cap.colorr cap.coloru cap.colorw cap.color
                           -0.8152744 -0.07378939 -0.4758376 -0.5410833 -0.6253381 -0.1343238 -0.04441978 -0.04441978 
1.2264300 -0.07378939 -0.4758376 -0.5410833 -0.6253381 -0.1343238 -0.04441978 -0.04441978
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               2.609732 -0.3898652
2.609732 -0.3898652
                          1.2264300 -0.07378939 -0.4758376 -0.5410833 -0.6253381 -0.1343238 -0.04441978 -0.04441978 -0.381314 -0.38898652 
1.2264300 -0.07378939 -0.4758376 -0.5410833 -0.6253381 -0.1343238 -0.04441978 -0.04441978 -0.383134 -0.38388652
                1.1858436 -0.1555724 -0.6017711 -0.2275528 -0.066712 -0.8760876 -0.1803687 -0.2762286 -0.2762286 -0.2762286 -0.1858436 -0.1555724 -0.6017711 -0.2275528 -0.066712 -0.8760876 -0.1803687 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.2762286 -0.
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       odory gill.attachmentf
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              0.1628864
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              0.1628864
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            0.1628864
                   1.1858436 \ -0.1555724 \ -0.6017711 \ -0.2275528 \ -0.066712 \ -0.8760876 \ \ 5.5435180 \ -0.2762286 \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.276286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.2762286 \ \ -0.276286 \ \ -0.276286 \ \ -0.276286 \ \ -0.276286 \ \ -0.276286 \
             0.1628864
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              0.1628864
             gill.spacingw gill.sizen gill.colore gill.
                                    -0.4388366 -0.6689971 -0.1093466 -0.3193666 -0.3146646 -0.2299361
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               2.5982835
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         -0.0891037
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           -0.4742807
                                                                                                                                                                                       -0.1093466 -0.3193666 -0.3146646
                                    -0.4388366 1.4945907
                                                                                                                                                                                                                                                                                                                                                                                                                                      -0.2299361
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     2.5982835
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          -0.0891037
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             -0.4742807
                                    2.2784719 -0.6689971 -0.1093466 -0.3193666 -0.3146646 4.3484982 -0.4388366 -0.6689971 -0.1093466 -0.3193666 -0.3146646 -0.2299361
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            -0.3848221
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             -0.0891037
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               -0.4742807
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              2.5982835 -0.0891037
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         -0.4742807
gill.colorr gill.coloru gill.coloru gill.colory stalk.shapet stalk.rootb stalk.rootc stalk.roott 1-0.05442976 -0.2538848 -0.4166868 -0.1034305 -1.1447353 -0.9318467 -0.2710317 2.5005603 -0.1555724
```

Correlation:

> round(cor(encoded_data),digits = 2)							
	classe classp	cap.shapec	cap.shapef	cap.shapek	cap.shapes o	ap.shapex (cap.surfaceg	
classe	1.00 -1.00	-0.02	-0.02	-0.16	0.06	0.03	-0.02	
classp	-1.00 1.00	0.02	0.02	0.16	-0.06	-0.03	0.02	
cap.shapec	-0.02 0.02	1.00	-0.02	-0.01	0.00	-0.02	0.25	
cap.shapef	-0.02 0.02	-0.02	1.00	-0.27	-0.05	-0.72	-0.01	
cap.shapek	-0.16 0.16	-0.01	-0.27	1.00	-0.02	-0.30	0.01	
cap.shapes	0.06 -0.06	0.00	-0.05	-0.02	1.00	-0.06	0.00	
cap.shapex	0.03 -0.03	-0.02	-0.72	-0.30	-0.06	1.00	-0.02	
cap.surfaceg	-0.02 0.02	0.25	-0.01	0.01	0.00	-0.02	1.00	
cap.surfaces	-0.10 0.10		-0.09	0.14	-0.04	-0.04	-0.02	
cap.surfacey	-0.09 0.09		0.03	0.02	-0.05	-0.02	-0.02	
		cap.surfacey	cap.coloro	cap.colore	cap.colorg		cap.colorp cap	o.colorr
classe	-0.10	-0.09				0.04	-0.03	0.04
classp	0.10	0.09				-0.04	0.03	-0.04
cap.shapec	-0.02	0.02				-0.01	0.00	0.00
cap.shapef	-0.09	0.03				0.02	-0.05	0.01
cap.shapek	0.14	0.02				0.13	0.00	-0.01
cap.shapes	-0.04	-0.05				0.03	-0.01	0.00
cap.shapex	-0.04	-0.02				-0.06	0.04	0.00
cap.surfaceg	-0.02	-0.02				-0.01	0.00	0.00
cap.surfaces	1.00	-0.55				0.08	0.05	-0.03
cap.surfacey	-0.55	1.00				0.02	-0.03	0.05
							rn odorp odors	
classe	0.04	0.13	-0.11	0.50 -0.16			79 -0.19 -0.29	
classp	-0.04	-0.13	0.11	-0.50 0.16				
cap.shapec	0.00	0.02	0.02		-0.01 -0.01			
cap.shapef	0.01	-0.09	-0.01		0.08 -0.10			
cap.shapek	-0.01	-0.07	-0.12	-0.23 -0.05				0.21
cap.shapes	0.00	-0.02	-0.02	-0.05 -0.01	-0.04 -0.0	0.00 0.0	07 -0.01 -0.02	-0.02

Explanation:

- The cor() function is used to compute the correlation matrix of the data frame.
- round() function is applied to the correlation matrix calculated in the previous step. The digits parameter is set to 2, indicating that the correlation coefficients should be rounded to two decimal places.

KNN:

Spited data into training and test sets and defined target column index.

```
23  set.seed(123)
24  train_indices <- sample(1:nrow(encoded_data), 0.7 * nrow(encoded_data))
25  training_data <- encoded_data[train_indices, ]
26  test_data <- encoded_data[-train_indices, ]
27  target_column_index <- 1</pre>
```

Explanation:

- The set.seed() function is used to set the random number generator's seed. This ensures that the random sampling and other random operations performed later in the code will produce the same results each time the code is run.
- The sample() function is used to randomly sample indices from 1 to the total number of rows in the encoded_data data frame. The second argument 0.7 * nrow(encoded_data) specifies that approximately 70% of the rows will be sampled as training data. The sampled indices are stored in the train_indices variable.
- training_data will be used to train a machine learning model and test_data ill be used to evaluate the performance of the trained model.

Pearson's correlation:

```
correlation_matrix <- cor(training_data[,-target_column_index], training_data[, target_column_index])</pre>
  important_attributes <- colnames(training_data)[abs(correlation_matrix) > 0.1]
> correlation_matrix <- cor(training_data[,-target_column_index], training_data[, target_column_index])
> print(correlation_matrix)
classp
                           -1.000000000
cap.shapec
                           -0.019286750
cap.shapef
                           -0.016268908
cap.shapek
                           -0.170760033
                            0.061982566
cap, shapes
cap.shapex
                            0.031250830
cap.surfacea
                            -0.019286750
                           -0.096512452
-0.084785007
cap.surfaces
cap.surfacey
cap.colorc
                            0.030118100
                            -0.097945624
cap.colore
cap.colorg
                            0.044894826
cap.colorn
                            0.044972971
                            -0.035714883
cap, colorp
cap.colorr
                            0.040823263
cap.coloru
                            0.040823263
cap.colorw
                            0.137216698
                            -0.112370785
cap.colory
bruisest
                            0.502135130
                            -0.163956964
odorf
                           -0.619386698
```

Explanation:

- Calculate coefficient and select important attributes
- abs(correlation_matrix) > 0.1: This creates a logical matrix where each element is TRUE if the absolute value of the corresponding correlation coefficient is greater than 0.1, and FALSE otherwise.

KNN model:

Accuracy test:

```
> accuracy_test <- sum(knn_model == test_dataset[, target_column_index]) / nrow(test_dataset)
> cat("Accuracy (Test Set):", accuracy_test, "\n")
Accuracy (Test Set): 0.9983593
```

10-fold cross validation:

```
> num_folds <- 10</pre>
> cv_folds <- createFolds(training_dataset[, target_column_index], k = num_folds)</pre>
  cv_accuracy <- numeric(num_folds)</pre>
  for (fold in 1:num_folds) {
     train_indices <- cv_folds[[fold]]
train_data <- training_dataset[train_indices, ]</pre>
     test_data <- training_dataset[-train_indices, ]</pre>
     knn model <- knn(
        train = train_data[, -target_column_index],
        test = test_data[, -target_column_index],
        cl = train_data[, target_column_index],
        k = best_k
     accuracy <- sum(knn_model == test_data[, target_column_index]) / nrow(test_data)</pre>
     cv_accuracy[fold] <- accuracy</pre>
> cat("Accuracy (", num_folds, "-fold CV):\n")
Accuracy ( 10 -fold CV):
> print(cv_accuracy)
  \begin{smallmatrix} 1 \end{smallmatrix} ] \hspace{0.1cm} 0.9953107 \hspace{0.1cm} 0.9892515 \hspace{0.1cm} 0.9984369 \hspace{0.1cm} 0.9906195 \hspace{0.1cm} 0.9931601 \hspace{0.1cm} 0.9972646 \hspace{0.1cm} 0.9972640 \hspace{0.1cm} 0.9992183 \hspace{0.1cm} 0.9986320 \hspace{0.1cm} 0.9878859 
> average_accuracy <- mean(cv_accuracy)
> cat("Average Accuracy:", average_accuracy, "\n")
Average Accuracy: 0.9947043
```

Confusion matrix:

```
xtab <- table(knn_model, test_dataset[, target_column_index])</pre>
 cm <- caret::confusionMatrix(xtab)</pre>
 xtab <- table(knn_model, test_dataset[, target_column_index])</pre>
cm <- caret::confusionMatrix(xtab)</pre>
print(cm)
onfusion Matrix and Statistics
nn_model
      0 1148
           4 1286
              Accuracy: 0.9984
  95% CI : (0.9958, 0.9996)
No Information Rate : 0.5275
  P-Value [Acc > NIR] : <2e-16
                 Kappa: 0.9967
Mcnemar's Test P-Value : 0.1336
           Sensitivity: 0.9965
           Specificity: 1.0000
        Pos Pred Value: 1.0000
        Neg Pred Value : 0.9969
            Prevalence: 0.4725
        Detection Rate : 0.4709
 Detection Prevalence : 0.4709
     Balanced Accuracy: 0.9983
      'Positive' Class : 0
```

Explanation:

- Confusion matrix calculated using the table() function. xtab is a table containing the counts of true positive, true negative, false positive, and false negative predictions.
- confusionMatrix() function from the "caret" package used to create a confusion matrix object that contains various metrics for evaluating the model's performance.

Recall and Precision:

```
> library(Metrics)
> Metrics::recall(knn_model, test_dataset[, target_column_index])
[1] 0.9968992
> |
> precision_value <- Metrics::precision(predicted_labels, true_labels)
> print(precision_value)
[1] 1
```

Explanation:

- The Metrics package in R provides functions for computing various metrics commonly used in evaluating the performance of machine learning models. The Metrics::recall() function specifically calculates the recall metric for evaluating the performance of a binary classification model.
- The Metrics::precision() function calculates the precision as the ratio of true positive
 predictions to the total number of instances that were predicted as positive by the model.