

ASSIGNMENT 9

Rajasekhar Reddy Karna

2020-11-01

In this problem, you will use the nearest neighbors algorithm to fit a model on two simplified datasets. The first dataset (found in `binary-classifier-data.csv`) contains three variables; label, x, and y. The label variable is either 0 or 1 and is the output we want to predict using the x and y variables. The second dataset (found in `trinary-classifier-data.csv`) is similar to the first dataset except that the label variable can be 0, 1, or 2.

Note that in real-world datasets, your labels are usually not numbers, but text-based descriptions of the categories (e.g. spam or ham). In practice, you will encode categorical variables into numeric values.

a. Plot the data from each dataset using a scatter plot.

```
options(warn=-1)
library(ggplot2)
library(readr)
library(foreign)
library(caTools)
library(class)
library(caret)
```

```
## Loading required package: lattice
```

```
setwd("C:/Users/vahin/Documents/GitHub/dsc520/")
bi_classifier_df <- read.csv("data/binary-classifier-data.csv")
head(bi_classifier_df)
```

```
##   label      x      y
## 1     0 70.88469 83.17702
## 2     0 74.97176 87.92922
## 3     0 73.78333 92.20325
## 4     0 66.40747 81.10617
## 5     0 69.07399 84.53739
## 6     0 72.23616 86.38403
```

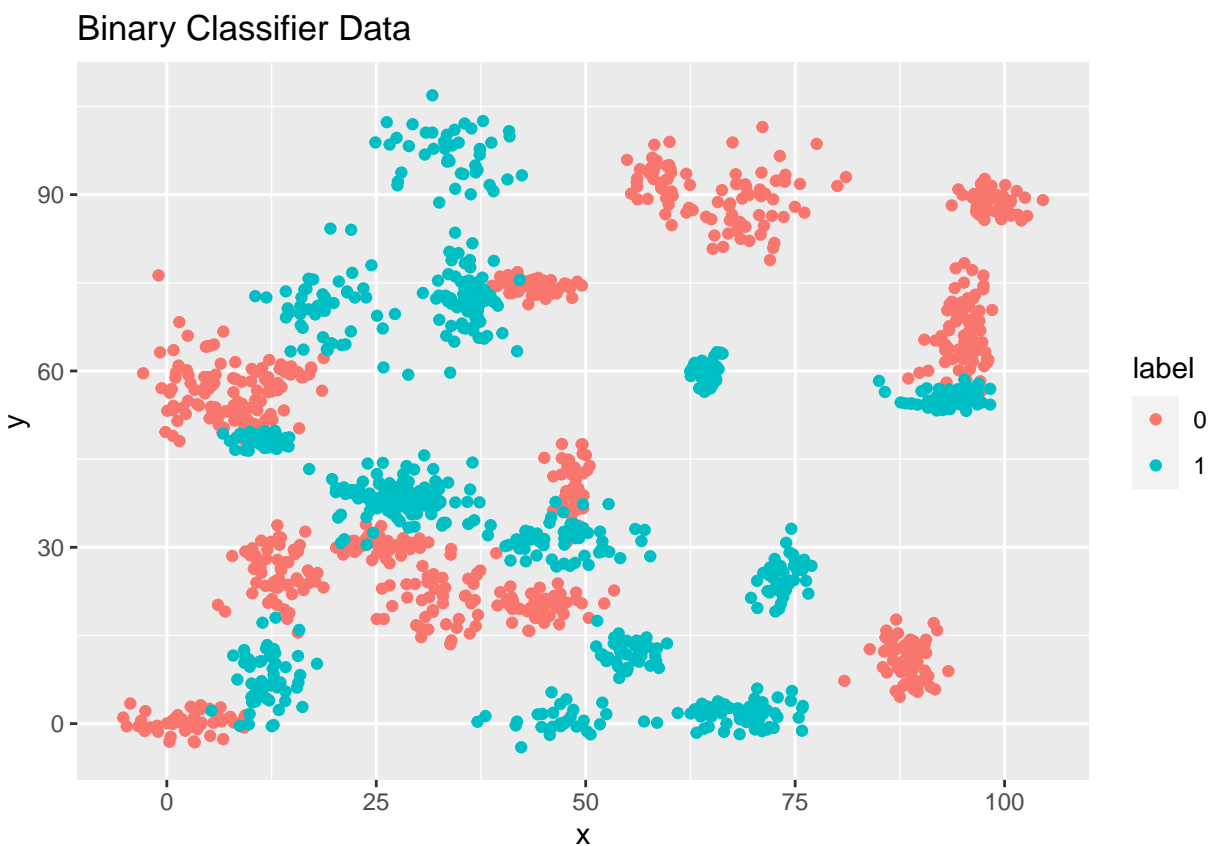
```
str(bi_classifier_df)
```

```
## 'data.frame':   1498 obs. of  3 variables:
##  $ label: int   0 0 0 0 0 0 0 0 0 0 ...
##  $ x    : num  70.9 75 73.8 66.4 69.1 ...
##  $ y    : num  83.2 87.9 92.2 81.1 84.5 ...
```

```
summary(bi_classifier_df)
```

```
##      label      x      y
## Min.   :0.000 Min.   : -5.20 Min.   : -4.019
## 1st Qu.:0.000 1st Qu.: 19.77 1st Qu.: 21.207
## Median :0.000 Median : 41.76 Median : 44.632
## Mean   :0.488 Mean   : 45.07 Mean   : 45.011
## 3rd Qu.:1.000 3rd Qu.: 66.39 3rd Qu.: 68.698
## Max.   :1.000 Max.   :104.58 Max.   :106.896
```

```
bi_classifier_df$label <- as.factor(bi_classifier_df$label)
ggplot(bi_classifier_df, aes(x=x, y=y, color=label)) + geom_point() + ggtitle('Binary Classifier Data')
```



```
tri_classifier_df <- read.csv("data/trinary-classifier-data.csv")
head(tri_classifier_df)
```

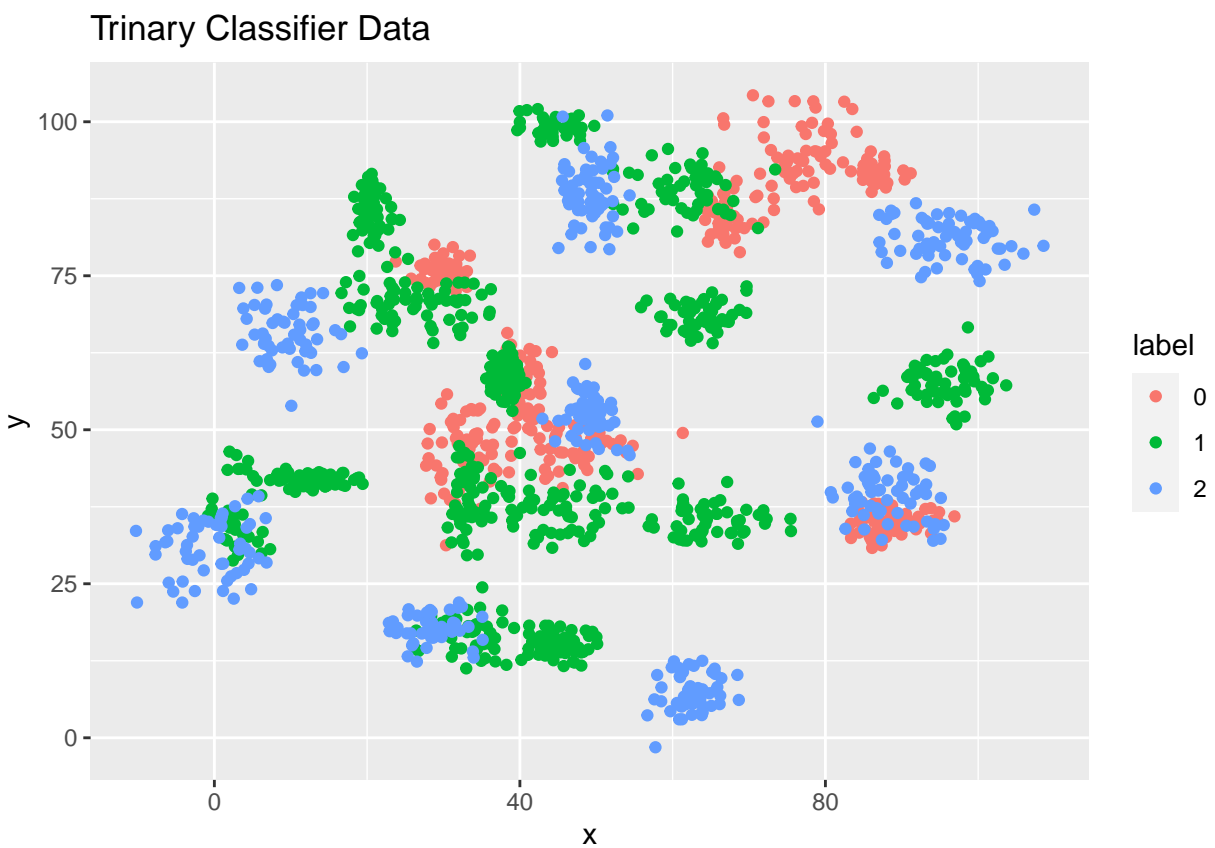
```
##  label      x      y
## 1     0 30.08387 39.63094
## 2     0 31.27613 51.77511
## 3     0 34.12138 49.27575
## 4     0 32.58222 41.23300
## 5     0 34.65069 45.47956
## 6     0 33.80513 44.24656
```

```
summary(tri_classifier_df)
```

```
##      label      x      y
##  Min.   :0.000  Min.   :-10.26  Min.    : -1.541
## 1st Qu.:0.000  1st Qu.: 31.15  1st Qu.: 35.906
## Median :1.000  Median : 45.59  Median : 55.073
## Mean   :1.037  Mean   : 48.86  Mean    : 55.282
## 3rd Qu.:2.000  3rd Qu.: 66.27  3rd Qu.: 77.403
## Max.   :2.000  Max.   :108.56  Max.    :104.293
```

```
tri_classifier_df$label <- as.factor(tri_classifier_df$label)
```

```
ggplot(tri_classifier_df, aes(x=x, y=y, color=label)) + geom_point() + ggtitle('Trinary Classifier Data')
```



```
head(tri_classifier_df)
```

```
##  label      x      y
## 1     0 30.08387 39.63094
## 2     0 31.27613 51.77511
## 3     0 34.12138 49.27575
## 4     0 32.58222 41.23300
## 5     0 34.65069 45.47956
## 6     0 33.80513 44.24656
```

```
str(tri_classifier_df)
```

```
## 'data.frame': 1568 obs. of 3 variables:
## $ label: Factor w/ 3 levels "0","1","2": 1 1 1 1 1 1 1 1 1 1 ...
## $ x : num 30.1 31.3 34.1 32.6 34.7 ...
## $ y : num 39.6 51.8 49.3 41.2 45.5 ...
```

```
summary(tri_classifier_df)
```

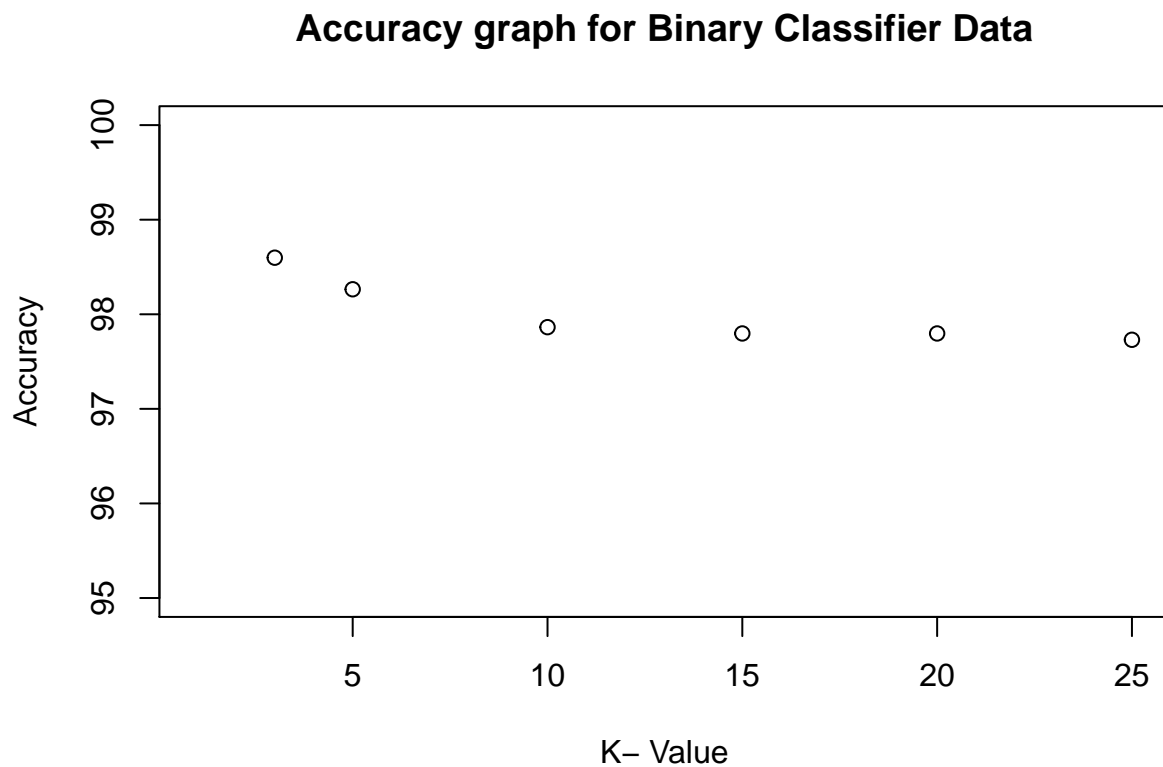
```
## label      x      y
## 0:394  Min.   :-10.26  Min.   : -1.541
## 1:722  1st Qu.: 31.15  1st Qu.: 35.906
## 2:452  Median : 45.59  Median : 55.073
##      Mean   : 48.86  Mean   : 55.282
##      3rd Qu.: 66.27  3rd Qu.: 77.403
##      Max.   :108.56  Max.   :104.293
```

##b. The k nearest neighbors algorithm categorizes an input value by looking at the labels for the k nearest points and assigning a category based on the most common label. In this problem, you will determine which points are nearest by calculating the Euclidean distance between two points. As a refresher, the Euclidean distance between two points: $\text{##p1}=(x_1, y_1)$ and $\text{##p2}=(x_2, y_2)$ is $\text{##d} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$. Fitting a model is when you use the input data to create a predictive model. There are various metrics you can use to determine how well your model fits the data. You will learn more about these metrics in later lessons. For this problem, you will focus on a single metric; accuracy. Accuracy is simply the percentage of how often the model predicts the correct result. If the model always predicts the correct result, it is 100% accurate. If the model always predicts the incorrect result, it is 0% accurate. Fit a k nearest neighbors model for each dataset for $k=3$, $k=5$, $k=10$, $k=15$, $k=20$, and $k=25$. Compute the accuracy of the resulting models for each value of k. Plot the results in a graph where the x-axis is the different values of k and the y-axis is the accuracy of the model.

```
set.seed(42)
bi_split<-sample.split(bi_classifier_df, SplitRatio=0.80)
tri_split<-sample.split(tri_classifier_df, SplitRatio=0.80)
bi_train <- subset(bi_classifier_df, bi_split=="TRUE")
bi_test  <- subset(bi_classifier_df, bi_split=="FALSE")
tri_train <- subset(tri_classifier_df, tri_split=="TRUE")
tri_test  <- subset(tri_classifier_df, tri_split=="FALSE")
list_of_k <- list(3,5,10,15,20,25)
accuracy_binary = 1
for (i in list_of_k) {
  knn_bi <- knn(train=bi_train, test=bi_test, cl=bi_train$label, k=i )
  accuracy_binary[i] <- 100 * sum(bi_test$label == knn_bi)/nrow(bi_test)
}
accuracy_binary
```

```
## [1] 1.00000 NA 98.59813 NA 98.26435 NA NA NA
## [9] NA 97.86382 NA NA NA NA 97.79706 NA
## [17] NA NA NA 97.79706 NA NA NA NA
## [25] 97.73031
```

```
plot(accuracy_binary, type="b", xlab="K- Value",ylab="Accuracy", ylim = c(95,100), main = "Accuracy graph")
```

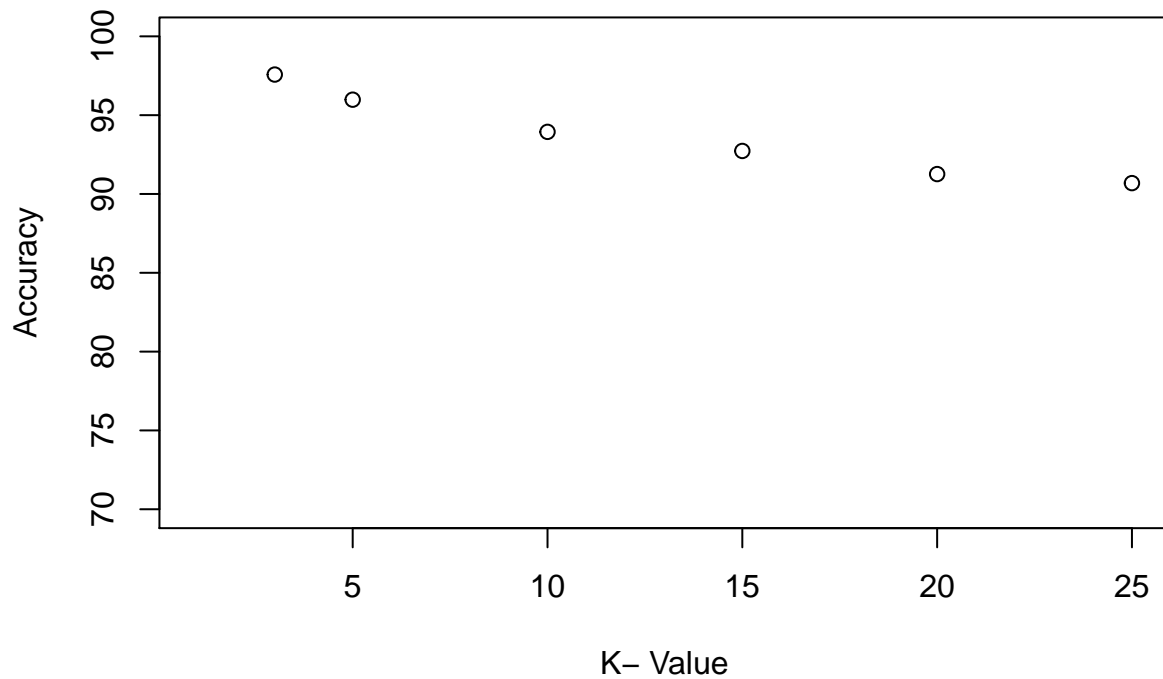


```
accuracy_trinary = 1
for (i in list_of_k) {
  knn_tri <- knn(train=tri_train, test=tri_test, cl=tri_train$label, k=i )
  accuracy_trinary[i] <- 100 * sum(tri_test$label == knn_tri)/nrow(tri_test)
}
accuracy_trinary
```

```
## [1] 1.00000      NA 97.57653      NA 95.98214      NA      NA      NA
## [9]      NA 93.94133      NA      NA      NA      NA 92.72959      NA
## [17]      NA      NA      NA 91.26276      NA      NA      NA      NA
## [25] 90.68878
```

```
plot(accuracy_trinary, type="b", xlab="K- Value",ylab="Accuracy", ylim = c(70,100), main = "Accuracy graph")
```

Accuracy graph for Trinary Classifier Data



##c. In later lessons, you will learn about linear classifiers. These algorithms work by defining a decision boundary that separates the different categories. ##Looking back at the plots of the data, do you think a linear classifier would work well on these datasets?

##Response Notes: No. As per scattered plot of the data is widely spread. Also value of K-Value & 'Accuracy' is dropping gradually. Linear classifier may be helpful because they will form a classification boundary based on the characteristics.