

Classification K-Nearest-Neighbors

Chris Bentz

14/10/2023

Description

k-nearest neighbor analyses of the feature vectors per character string (loaded from NaLaFi/results/features.csv). The results are stored in NaLaFi/results/KNN. Note that the number of characters has to be chosen manually (via `num.char = “`”), likewise, the features to be included can be chosen in the lines defining `estimations.subset`” below. Also, subcorpora can be excluded via the `“selected”` object.

Load Packages

If the libraries are not installed yet, you need to install them using, for example, the command: `install.packages(“ggplot2”)`. For the Hrate package this is different, since it comes from github. The devtools library needs to be installed, and then the `install_github()` function is used.

```
library(ggplot2)
library(dplyr)
library(class)
library(gridExtra)
library(gmodels)
library(caret)
library(ggExtra)
library(ggpubr)
```

Load Data

Load data table with values per text file.

```
# load estimations from stringBase corpus
estimations.df <- read.csv("~/Github/NaLaFi/results/features.csv")
#head(features.csv)
```

Exclude subcorpora (if needed). Choose the subcorpus to be excluded via `“selected”`.

```
selected <- c("heraldics",
             "procune",
             "animal",
             "random",
             "morse",
             "shuffled",
             "weather",
             "natural"
             #"pycode")
```

```
)
estimations.df <- estimations.df[!(estimations.df$subcorpus %in% selected), ]
```

Split into separate files by length of chunks in characters.

```
# choose number of characters
num.char = 100
# subset data frame
estimations.df <- estimations.df[estimations.df$num.char == num.char, ]
nrow(estimations.df)
```

```
## [1] 1626
```

Select relevant columns of the data frame, i.e. the measures to be included in classification and the ‘corpus’ or ‘subcorpus’ column.

```
estimations.subset <- estimations.df[c("corpus", "subcorpus",
                                       "huni.chars",
                                       "hrate.chars",
                                       "ttr.chars",
                                       "rm.chars"
                                       )]
```

Remove NAs (whole row)

```
estimations.subset <- na.omit(estimations.subset)
```

Center and scale the data

```
estimations.scaled <- cbind(estimations.subset[1:2], scale(estimations.subset[3:ncol(estimations.subset),
nrow(estimations.scaled)
```

```
## [1] 1626
```

Create Training and Test Sets

```
# Generating seed
set.seed(1234)
# Randomly generating our training and test samples with a respective ratio of 2/3 and 1/3
datasample <- sample(2, nrow(estimations.scaled), replace = TRUE, prob = c(0.67, 0.33))
# Generate training set
estimations.training <- estimations.scaled[datasample == 1, 3:ncol(estimations.scaled)]
nrow(estimations.training)
```

```
## [1] 1114
```

```
# Generate test set
estimations.test <- estimations.scaled[datasample == 2, 3:ncol(estimations.scaled)]
nrow(estimations.test)
```

```
## [1] 512
```

Get training and test labels

```
# Generate training labels
training.labels <- estimations.scaled[datasample == 1, 1]
# Generate test labels
test.labels <- estimations.scaled[datasample == 2, 1]
```

Initialize data frame

```
knn.results <- data.frame(k = numeric(0), accuracy = numeric(0), precision = numeric(0),
                          recall = numeric(0), f1 = numeric(0))
```

Building knn classifier

```
# choose maximum number of neighbors n
n = 10
# run a loop over different numbers of neighbors up to n
for (k in 1:n){
  # knn estimation of labels
  predictions.knn <- knn(train = as.data.frame(estimations.training),
                        test = as.data.frame(estimations.test),
                        cl = training.labels, k = k)

  # model evaluation
  # creating a dataframe from known (true) test labels
  test.labels <- data.frame(test.labels)
  # combining predicted and known classes
  class.comparison <- data.frame(predictions.knn, test.labels)
  # giving appropriate column names
  names(class.comparison) <- c("predicted", "observed")
  # inspecting our results table
  head(class.comparison)
  # get confusion matrix
  cm <- confusionMatrix(data = class.comparison$predicted,
                        reference = as.factor(class.comparison$observed))

  print(cm)
  # get precision, recall, and f1 from the output list of confusionMatrix()
  accuracy <- cm$overall['Accuracy']
  f1 <- cm[["byClass"]][["F1"]]
  recall <- cm[["byClass"]][["Recall"]]
  precision <- cm[["byClass"]][["Precision"]]

  # prepare data frame with results
  local.results <- data.frame(k, accuracy, precision, recall, f1, row.names = NULL)
  local.results.rounded <- round(local.results, 2)
  # print(local.results.rounded)
  knn.results <- rbind(knn.results, local.results.rounded)
}
```

```

## Confusion Matrix and Statistics
##
##           Reference
## Prediction  non-writing writing
## non-writing      5         6
## writing          5        496
##
##           Accuracy : 0.9785
##           95% CI : (0.9619, 0.9892)
##       No Information Rate : 0.9805
##       P-Value [Acc > NIR] : 0.6979
##
##           Kappa : 0.4652
##
## Mcnemar's Test P-Value : 1.0000
##
##           Sensitivity : 0.500000
##           Specificity : 0.988048
##       Pos Pred Value : 0.454545
##       Neg Pred Value : 0.990020
##           Prevalence : 0.019531
##       Detection Rate : 0.009766
##       Detection Prevalence : 0.021484
##       Balanced Accuracy : 0.744024
##
##       'Positive' Class : non-writing
##
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  non-writing writing
## non-writing      4         6
## writing          6        496
##
##           Accuracy : 0.9766
##           95% CI : (0.9594, 0.9878)
##       No Information Rate : 0.9805
##       P-Value [Acc > NIR] : 0.7934
##
##           Kappa : 0.388
##
## Mcnemar's Test P-Value : 1.0000
##
##           Sensitivity : 0.400000
##           Specificity : 0.988048
##       Pos Pred Value : 0.400000
##       Neg Pred Value : 0.988048
##           Prevalence : 0.019531
##       Detection Rate : 0.007812
##       Detection Prevalence : 0.019531
##       Balanced Accuracy : 0.694024
##
##       'Positive' Class : non-writing
##

```

```

## Confusion Matrix and Statistics
##
##               Reference
## Prediction    non-writing writing
## non-writing         4         4
## writing             6       498
##
##               Accuracy : 0.9805
##               95% CI : (0.9644, 0.9906)
##       No Information Rate : 0.9805
##       P-Value [Acc > NIR] : 0.5830
##
##               Kappa : 0.4346
##
## Mcnemar's Test P-Value : 0.7518
##
##       Sensitivity : 0.400000
##       Specificity : 0.992032
##       Pos Pred Value : 0.500000
##       Neg Pred Value : 0.988095
##       Prevalence : 0.019531
##       Detection Rate : 0.007812
##       Detection Prevalence : 0.015625
##       Balanced Accuracy : 0.696016
##
##       'Positive' Class : non-writing
##
## Confusion Matrix and Statistics
##
##               Reference
## Prediction    non-writing writing
## non-writing         5         3
## writing             5       499
##
##               Accuracy : 0.9844
##               95% CI : (0.9694, 0.9932)
##       No Information Rate : 0.9805
##       P-Value [Acc > NIR] : 0.3306
##
##               Kappa : 0.5477
##
## Mcnemar's Test P-Value : 0.7237
##
##       Sensitivity : 0.500000
##       Specificity : 0.994024
##       Pos Pred Value : 0.625000
##       Neg Pred Value : 0.990079
##       Prevalence : 0.019531
##       Detection Rate : 0.009766
##       Detection Prevalence : 0.015625
##       Balanced Accuracy : 0.747012
##
##       'Positive' Class : non-writing
##

```

```

## Confusion Matrix and Statistics
##
##           Reference
## Prediction  non-writing writing
## non-writing      4         1
## writing           6        501
##
##           Accuracy : 0.9863
##           95% CI : (0.972, 0.9945)
##       No Information Rate : 0.9805
##       P-Value [Acc > NIR] : 0.2176
##
##           Kappa : 0.5272
##
## Mcnemar's Test P-Value : 0.1306
##
##           Sensitivity : 0.400000
##           Specificity : 0.998008
##       Pos Pred Value : 0.800000
##       Neg Pred Value : 0.988166
##           Prevalence : 0.019531
##       Detection Rate : 0.007812
##       Detection Prevalence : 0.009766
##       Balanced Accuracy : 0.699004
##
##       'Positive' Class : non-writing
##
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  non-writing writing
## non-writing      5         3
## writing           5        499
##
##           Accuracy : 0.9844
##           95% CI : (0.9694, 0.9932)
##       No Information Rate : 0.9805
##       P-Value [Acc > NIR] : 0.3306
##
##           Kappa : 0.5477
##
## Mcnemar's Test P-Value : 0.7237
##
##           Sensitivity : 0.500000
##           Specificity : 0.994024
##       Pos Pred Value : 0.625000
##       Neg Pred Value : 0.990079
##           Prevalence : 0.019531
##       Detection Rate : 0.009766
##       Detection Prevalence : 0.015625
##       Balanced Accuracy : 0.747012
##
##       'Positive' Class : non-writing
##

```

```

## Confusion Matrix and Statistics
##
##               Reference
## Prediction    non-writing writing
## non-writing         5         4
## writing             5        498
##
##               Accuracy : 0.9824
##               95% CI : (0.9669, 0.9919)
##               No Information Rate : 0.9805
##               P-Value [Acc > NIR] : 0.4567
##
##               Kappa : 0.5174
##
## Mcnemar's Test P-Value : 1.0000
##
##               Sensitivity : 0.500000
##               Specificity : 0.992032
##               Pos Pred Value : 0.555556
##               Neg Pred Value : 0.990060
##               Prevalence : 0.019531
##               Detection Rate : 0.009766
##               Detection Prevalence : 0.017578
##               Balanced Accuracy : 0.746016
##
##               'Positive' Class : non-writing
##
## Confusion Matrix and Statistics
##
##               Reference
## Prediction    non-writing writing
## non-writing         5         3
## writing             5        499
##
##               Accuracy : 0.9844
##               95% CI : (0.9694, 0.9932)
##               No Information Rate : 0.9805
##               P-Value [Acc > NIR] : 0.3306
##
##               Kappa : 0.5477
##
## Mcnemar's Test P-Value : 0.7237
##
##               Sensitivity : 0.500000
##               Specificity : 0.994024
##               Pos Pred Value : 0.625000
##               Neg Pred Value : 0.990079
##               Prevalence : 0.019531
##               Detection Rate : 0.009766
##               Detection Prevalence : 0.015625
##               Balanced Accuracy : 0.747012
##
##               'Positive' Class : non-writing
##

```

```

## Confusion Matrix and Statistics
##
##               Reference
## Prediction   non-writing writing
## non-writing         5         2
## writing             5        500
##
##               Accuracy : 0.9863
##               95% CI : (0.972, 0.9945)
##               No Information Rate : 0.9805
##               P-Value [Acc > NIR] : 0.2176
##
##               Kappa : 0.5815
##
## Mcnemar's Test P-Value : 0.4497
##
##               Sensitivity : 0.500000
##               Specificity : 0.996016
##               Pos Pred Value : 0.714286
##               Neg Pred Value : 0.990099
##               Prevalence : 0.019531
##               Detection Rate : 0.009766
##               Detection Prevalence : 0.013672
##               Balanced Accuracy : 0.748008
##
##               'Positive' Class : non-writing
##
## Confusion Matrix and Statistics
##
##               Reference
## Prediction   non-writing writing
## non-writing         4         3
## writing             6        499
##
##               Accuracy : 0.9824
##               95% CI : (0.9669, 0.9919)
##               No Information Rate : 0.9805
##               P-Value [Acc > NIR] : 0.4567
##
##               Kappa : 0.4619
##
## Mcnemar's Test P-Value : 0.5050
##
##               Sensitivity : 0.400000
##               Specificity : 0.994024
##               Pos Pred Value : 0.571429
##               Neg Pred Value : 0.988119
##               Prevalence : 0.019531
##               Detection Rate : 0.007812
##               Detection Prevalence : 0.013672
##               Balanced Accuracy : 0.697012
##
##               'Positive' Class : non-writing
##

```



```
print(knn.results)
```

```
##      k accuracy precision recall   f1
## 1    1      0.98      0.45    0.5 0.48
## 2    2      0.98      0.40    0.4 0.40
## 3    3      0.98      0.50    0.4 0.44
## 4    4      0.98      0.62    0.5 0.56
## 5    5      0.99      0.80    0.4 0.53
## 6    6      0.98      0.62    0.5 0.56
## 7    7      0.98      0.56    0.5 0.53
## 8    8      0.98      0.62    0.5 0.56
## 9    9      0.99      0.71    0.5 0.59
## 10 10      0.98      0.57    0.4 0.47
```

Write to file. Note that the file names have to be changed manually here.

```
write.csv(knn.results, file = paste("~/Github/NaLaFi/results/KNN/knn_results_justPyCode100",
                                     paste(num.char, ".csv", sep = ""),
                                     sep = ""), row.names = F)
```