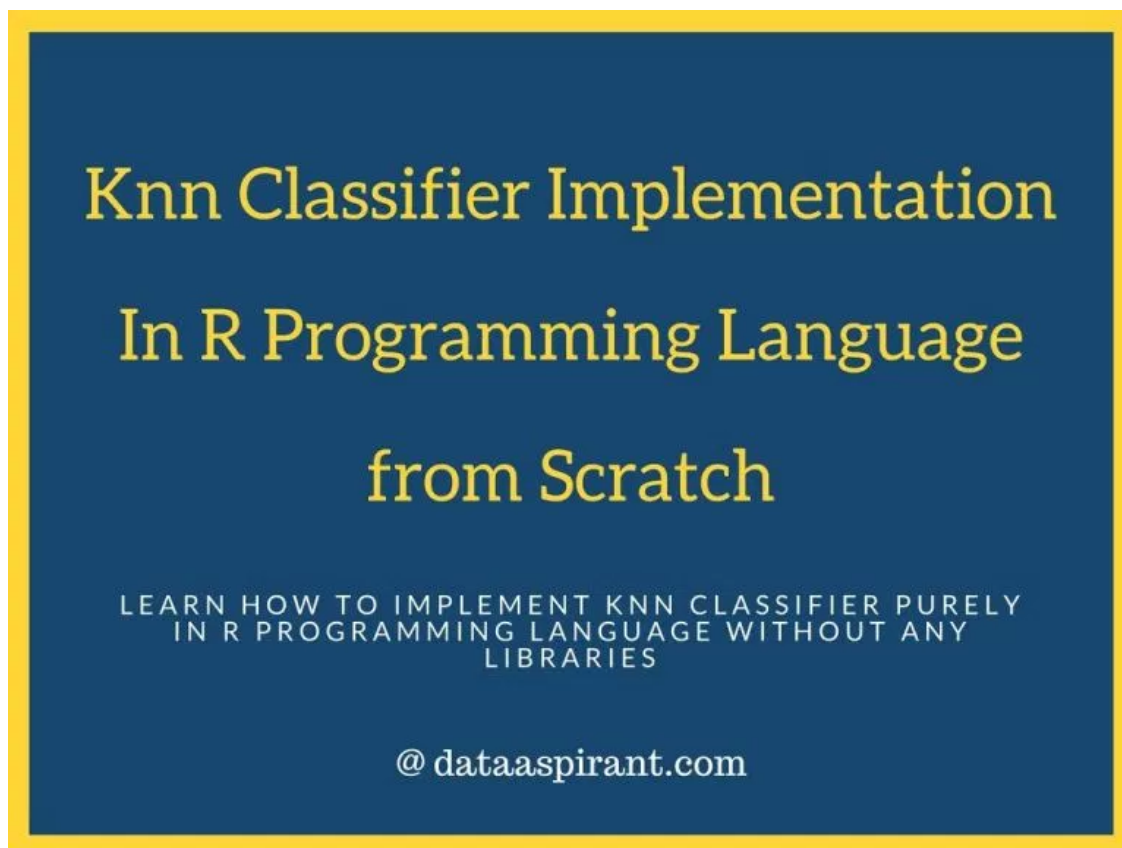


KNN R, K-NEAREST NEIGHBOR CLASSIFIER IMPLEMENTATION IN R PROGRAMMING FROM SCRATCH

📅 January 2, 2017 👤 Rahul Saxena 💬 3 Comments 📂 Data Science, Machine Learning, R



knn in implementation r from scratch

K-Nearest neighbor algorithm implement in R Programming from scratch

In the introduction to [k-nearest-neighbor algorithm](#) article, we have learned the core concepts of the knn algorithm. Also learned about the applications using knn algorithm to solve the real world problems.

In this post, we will be implementing K-Nearest Neighbor Algorithm on a dummy data set using R programming language from scratch. Along the way, we will implement a prediction model to predict classes for data.

Knn Implementation in R

Why we need to implement knn algorithm from scratch in R Programming Language

Implementation of K-Nearest Neighbor algorithm in R language from scratch will help us to apply the concepts of Knn algorithm. As we are going to implement each and every component of the knn algorithm and the other components like how to use the datasets and find the accuracy of our implemented model etc.

Problem Set

We will use a sample dataset extracted from ionosphere database by John Hopkins University. We have converted the database into a small dataset so as to simplify the learning curve for our readers.

Our objective is to program a **Knn classifier in R programming language without using any machine learning package**. We have two classes "g"(good) or "b"(bad), it is the response of radar from the ionosphere. The classifier could be capable of predicting "g" or "b" class for new records from training data.

Ionosphere Dataset Description

This dummy dataset consists of 6 attributes and 30 records. Out of these 5 attributes are continuous variables with values ranging from -1 to +1 i.e, [-1,+1]. Last(6th) attribute is a categorical variable with values as "g"(good) or "b"(bad) according to the definition summarized above. This is a binary classification task.

K-Nearest Neighbor Algorithm Pseudocode

Let (X_i, C_i) where $i = 1, 2, \dots, n$ be data points. X_i denotes feature values & C_i denotes labels for X_i for each i .

Assuming the number of classes as 'c'

$C_i \in \{1, 2, 3, \dots, c\}$ for all values of i

Let x be a point for which label is not known, and we would like to find the label class using k-nearest neighbor algorithms.

Procedure:

1. Calculate " $d(x, x_i)$ " $i = 1, 2, \dots, n$; where **d** denotes the [Euclidean distance](#) between the points.
2. Arrange the calculated **n** Euclidean distances in non-decreasing order.
3. Let **k** be a +ve integer, take the first **k** distances from this sorted list.
4. Find those **k**-points corresponding to these **k**-distances.
5. Let k_i denotes the number of points belonging to the i^{th} class among **k** points i.e. $k \geq 0$
6. If $k_i > k_j \forall i \neq j$ then put x in class i .

Let's use the above pseudocode for implementing the knn algorithm in R Language.

Prerequisites:

1. Basic programming experience is required
2. [Install](#) R-Studio on your system.

K-Nearest neighbor algorithm implement in R Language from scratch

We are going to follow the below workflow for implementing the knn algorithm in R:

1. Getting Data
2. Train & Test Data Split
3. Euclidean Distance Calculation

4. KNN prediction function
5. Accuracy calculation

Let's get our hands dirty and start the coding stuff.

Getting Data in R

For any programmatic implementation on the dataset, we first need to import it. Using `read.csv()`, we are importing dataset into `knn.df` dataframe. Since dataset has no header so, we are using `header = FALSE`. `sep` parameter is to define the literal which separates values our document. `knn.df` is a dataframe. A dataframe is a table or 2-D array, in which each column contains measurements on one variable, and each row contains one record.

```
1 knn.df <- read.csv('i_data_sample_30.csv', header = FALSE, sep = ',')
```

For checking dimensions of the dataset, we can call `dim()` method and be passing data frame as a parameter.

```
1 dim(knn.df)
2 [1] 30 6
```

It shows that the data frame consists of 30 records and 6 columns.

To check summary of our dataset, we can use `summary()` method.

```
1 summary(knn.df)
2
3      V1      V2      V3      V4
4 Min.   :-1.0000 Min.   :-1.0000 Min.   :-1.0000 Min.   :-1.0000 Min.
5 1st Qu.: -0.1079 1st Qu.: 0.0000 1st Qu.: 0.0000 1st Qu.: 0.0000 1st Qu.
6 Median : 0.0000 Median : 0.8099 Median : 0.4059 Median : 0.5891 Median
7 Mean   :-0.0137 Mean   : 0.5611 Mean   : 0.3246 Mean   : 0.2769 Mean
8 3rd Qu.: 0.1336 3rd Qu.: 1.0000 3rd Qu.: 0.9730 3rd Qu.: 0.8269 3rd Qu.
9 Max.   : 1.0000 Max.   : 1.0000 Max.   : 1.0000 Max.   : 1.0000 Max.
10 V6
11 b:15
12 g:15
```

It shows that records with bad class and good class are 15 each.

Train & Test Data split in R

Before Train & Test data split, we need to distribute it randomly. In R, we can use `sample()` method. It helps to randomize all the records of dataframe.

Please use `set.seed(2)`, `seed()` method is used to produce reproducible results.

70
Shares

58

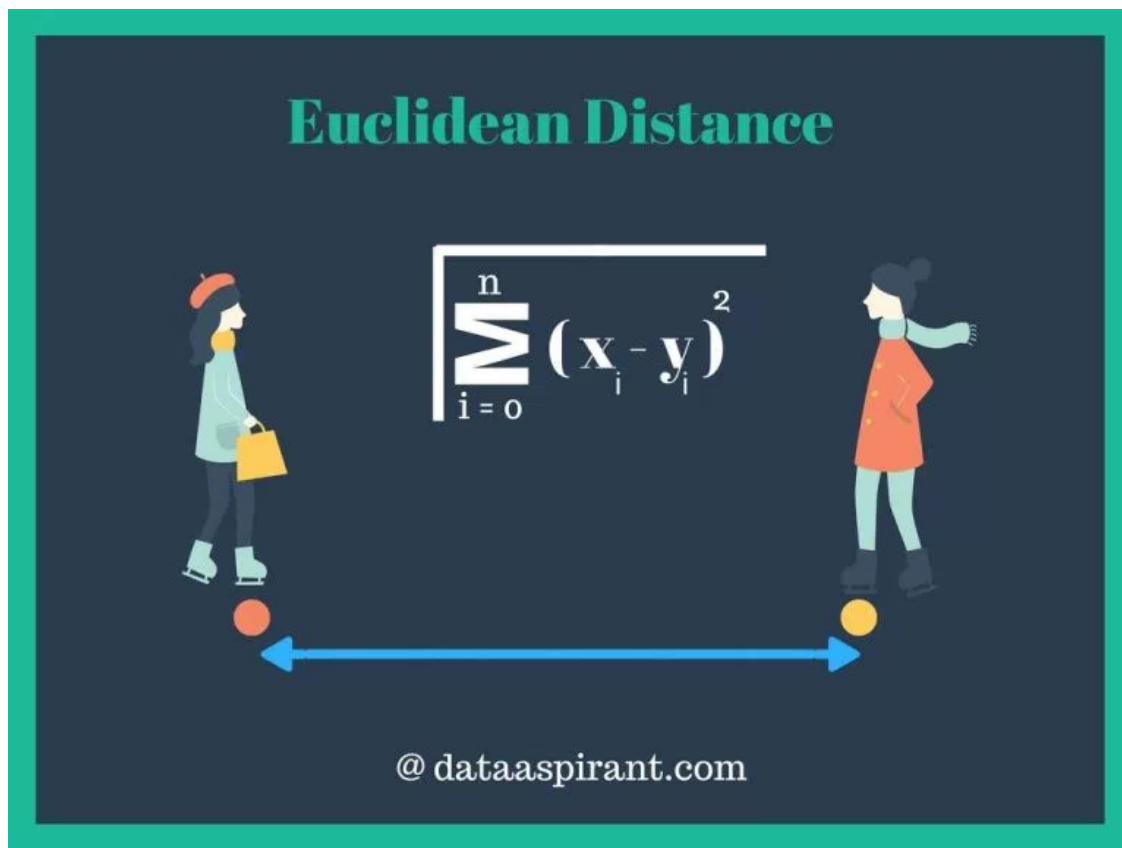
12

In the next line we are passing sample() method inside dataframe. This is to randomize all 30 records of knn.df. Now, we are ready for a split. For dividing train, test data we are splitting them in 70:30 ratio i.e., 70% of data will be considered as train set & 30% as the test set.

```
1 set.seed(2)
2 knn.df<- knn.df[sample(nrow(knn.df)),]
3 train.df <- knn.df[1:as.integer(0.7*30),]
4 test.df  <- knn.df[as.integer(0.7*30 +1):30,]
```

Euclidean Distance Calculation in R

Below snippet consists of a function defined in R to calculate Euclidean distance between 2 points a & b. The formula of Euclidean distance is:



Euclidean Distance

```
1 euclideanDist <- function(a, b){
2   d = 0
3   for(i in c(1:(length(a)-1) ))
4   {
5     d = d + (a[[i]]-b[[i]])^2
6   }
7   d = sqrt(d)
8   return(d)
9 }
```

KNN prediction function in R

This function is the core part of this tutorial. We are writing a function `knn_predict`. It takes 3 arguments: test data, train data & value of K. It loops over all the records of test data and train data. It returns the predicted class labels of test data.

```

1 knn_predict <- function(test_data, train_data, k_value){
2   pred <- c() #empty pred vector
3   #LOOP-1
4   for(i in c(1:nrow(test_data))){ #looping over each record of test data
5     eu_dist = c() #eu_dist & eu_char empty vector
6     eu_char = c()
7     good = 0 #good & bad variable initialization with 0 value
8     bad = 0
9
10    #LOOP-2-looping over train data
11    for(j in c(1:nrow(train_data))){
12
13      #adding euclidean distance b/w test data point and train data to eu_dist
14      eu_dist <- c(eu_dist, euclideanDist(test_data[i,], train_data[j,]))
15
16      #adding class variable of training data in eu_char
17      eu_char <- c(eu_char, as.character(train_data[j,][[6]]))
18    }
19
20    eu <- data.frame(eu_char, eu_dist) #eu dataframe created with eu_char & eu_
21
22    eu <- eu[order(eu$eu_dist),] #sorting eu dataframe to gettop K neighb
23    eu <- eu[1:k_value,] #eu dataframe with top K neighbors
24
25    #Loop 3: loops over eu and counts classes of neighbors.
26    for(k in c(1:nrow(eu))){
27      if(as.character(eu[k,"eu_char"]) == "g"){
28        good = good + 1
29      }
30      else
31        bad = bad + 1
32    }
33
34    # Compares the no. of neighbors with class label good or bad
35    if(good > bad){ #if majority of neighbors are good then put "g" in
36
37      pred <- c(pred, "g")
38    }
39    else if(good < bad){
40      #if majority of neighbors are bad then put "b" in pred vecto
41      pred <- c(pred, "b")
42    }
43  }
44 }
45 return(pred) #return pred vector
46 }
```

It returns a vector with predicted classes of test dataset. These predictions can be used to calculate accuracy metric.

Accuracy Calculation in R

The accuracy metric calculates the ratio of the number of correctly predicted class labels to the total number of predicted labels.

```
1 accuracy <- function(test_data){  
2   correct = 0  
3   for(i in c(1:nrow(test_data))){  
4     if(test_data[i,6] == test_data[i,7]){  
5       correct = correct+1  
6     }  
7   }  
8   accu = correct/nrow(test_data) * 100  
9   return(accu)  
10 }
```

KNN Algorithm accuracy print: In this code snippet we are joining all our functions. We are calling the knn_predict function with train and test dataframes that we split earlier and K value as 5.

We are appending the prediction vector as the 7th column in our test dataframe and then using accuracy() method we are printing accuracy of our KNN model.

```
1 K = 5  
2 predictions <- knn_predict(test.df, train.df, K) #calling knn_predict()  
3  
4 test.df[,7] <- predictions #Adding predictions in test data as 7th column  
5 print(accuracy(test.df))
```

Script Output:

```
1 Accuracy of our KNN model is  
2 77.77778
```

It prints accuracy of our knn model. Here our accuracy is **77.78%**. That's pretty good 😊 for our randomly selected dummy dataset.

You can download the RMD file of this code from our GitHub repository.

Finally, we have implemented our KNN model in R programming without using any specific R packages. Hope you enjoyed learning it.

Related Articles To Read

- “ Introduction to knn algorithm
- “ Knn algorithm implementation purely in python without any machine learning libraries

- “ Cancer tumor detection with knn sklearn

Follow us:

[FACEBOOK](#) | [QUORA](#) | [TWITTER](#) | [GOOGLE+](#) | [LINKEDIN](#) | [REDDIT](#) | [FLIPBOARD](#) | [MEDIUM](#) | [GITHUB](#)

I hope you like this post. If you have any questions, then feel free to comment below. If you want me to write on one particular topic, then do tell it to me in the comments below.

Related Courses:

Do check out [unlimited data science courses](#)

Title of the course	Course Link	Course Link
R Programming A-Z: R For Data Science With Real Exercises!	R Programming A-Z: R For Data Science With Real Exercises!	<ul style="list-style-type: none"> • This course is a step-by-step tutorial that will help you build or already move on forward • After every lesson, you will learn a concept and apply it to the best of your learning example • In summary, this course is designed for all levels of

have no
program
statistic
you will
in this c

R Programming: Advanced Analytics In R For Data Science

R Programming: Advanced
Analytics In R For Data
Science

- Perform
Prepara
locate n
datafram
- Apply th
Analysis
replace
records
- Work w
and sub
for repl
- Use lapp
apply()
working
vectors.
- Use lapp
apply()
working
vectors.

Data Mining with R: Go from Beginner to Advanced!

Data Mining with R: Go
from Beginner to
Advanced!

- Use R sc
data im
export,
explora
visualiza
data an

includin

a compi

of data

operatio

- Apply th
included
cases ar
using re
scripts t
unique
and dat
problem

- Effective
number
contem
mining
techniq
demand
includin
Decision
classific
regressi
(CART);
forests;
logistic
and (4) '
cluster ;
techniq

Share this:



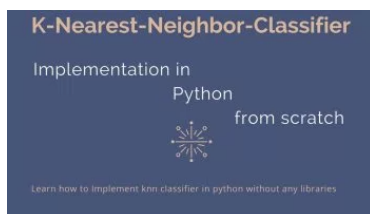
12



58



Related



K-nearest neighbor
algorithm implementation
in Python from scratch

December 27, 2016

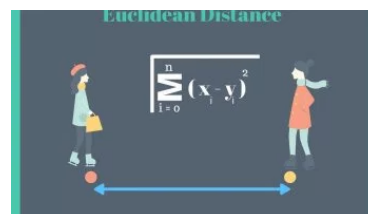
In "Machine Learning"



KNN R, K-Nearest Neighbor
implementation in R using
caret package

January 9, 2017

In "Machine Learning"



Knn sklearn, K-Nearest
Neighbor implementation
with scikit learn

December 30, 2016

In "Data Science"

📌 Classification 📌 k-nearest neighbor 📌 R Programming