

WEEK-8⇒ K-Nearest Neighbors

- ↳ KNN is a non-parametric method for classification.
- ↳ lazy learning algorithm. → all computation is deferred until classification.
- ↳ instance based learning algorithm where the function is approx locally.
- ↳ simplest
- ↳ no explicit training.
- ↳ used when there are non linear decision boundaries b/w classes.
- ↳ when the amount of data large.

→ Input Features

- ↳ both quantitative and qualitative

→ Outputs

- ↳ categorical

↳ which are classes of data.

KNN explains a categorical value with majority votes of nearest neighbors.

→ Assumptions:

- ↳ since it is non-parametric, there are no assumptions about the underlying data distribution.
- ↳ Select the parameter k based on the data.
- ↳ Requires a distance metric to define proximity b/w any two data points.
eg: Euclidian dist, Mahalanobis distance, Hamming distance.

→ Algorithm:

there are 4 steps:

- ① Compute the distance metric b/w the test data point and all the labeled data points.
- ② Order the labeled data points in the rising order of the distance metric.
- ③ Select the top k labeled data points and look at the class labels.
- ④ Find the class label that the majority of these k

Labeled data points have and assign it to the test data point.

→ things to consider

- ① parameter selection
- ② presence of noise
- ③ feature selection and scaling
- ④ Curse of dimensionality

→ Parameter Selection:

↳ The best choice of K depends on the data
↳ Larger values of K reduce the effect of noise on classification but makes the decision boundaries b/w classes less distinct.

↳ Smaller values of K tends to be affected by the noise, with clear separation b/w classes

2 Feature Selection:

- ↳ It is important to remove irrelevant features
- ↳ When the number of features is too large, and suspected to be highly redundant, feature extraction is required.
- ↳ If the features are carefully chosen then it is expected that the classification will be better.

3 K-means nearest neighbours implementation in R:

```
setwd("Path of the directory with data files")
```

```
rm(list=ls()) # clear the environment
```

```
library(caret) # for confusion matrix etc
```

```
library(class) # for Knn, not new choices
```

```
(unstable numbers x) function in
```

```
ServiceTrain <- read.csv("ServiceTrainData.csv")
```

```
ServiceTest <- read.csv("ServiceTestData.csv")
```

```
# this function is writing file without
```

```
function of m 02 [33 ... 2] -> 2792 (n>k)
```

```
(2.2 w) sample to make results consistent etc
```

Knn (train, test, cl, k=1)

- ↳ no. of neighbours do we consider
- factor of true classification
- and at size of clusters
- number of training set divided by number of training
- 1. fit with actual solution and contrast with K-means
- method of minimizing all distances between

K-means clustering

K-means clustering

A technique to partition N observations in k clusters ($k \leq N$) in which each observation belongs to cluster with nearest mean.

- ↳ one of the simplest unsupervised algorithms
 - ↳ works well for all distance metrics where mean is defined (ex Euclidian distance)

Given N observations (x_1, x_2, \dots, x_N) , k means clustering will partition n observations into k ($k \leq N$) sets $S = \{S_1, \dots, S_k\}$ so as to minimize the within cluster sum of squares (WCSS).

$$\arg \min \sum_{i=1}^k \sum_{x \in S_i} \|x - \mu_i\|^2$$

where μ_i is the mean of points in S_i

Determining the number of clusters

Elbow Method: looks at percentage of variance explained as a function of number of clusters

↳ The point where marginal decrease plateaus is an indicator of the optimal number of clusters.

→ Disadvantages:

- ↳ algorithm could converge to a local minima, therefore role of initial position is very important.
- ↳ if the clusters are not spherical, then K-means can fail to identify the correct number of clusters

IMPLEMENTATION IN R:

→

tripDetails = read.csv("tripDetails.csv", row.names=1, na.strings="")

- K-means clustering in R can be applied on data using 'kmeans()' function

object = kmeans(x, centers, iter.max = 10, nstart = 1)

↳ ↓ ↓ ↓
numeric either the max no. of
matrix of no. of iterations
data or an centers / allowed, if centers is
object that can set of a number,
be coerced to initial how many sets
such a matrix centers. should be
chosen.

tripCluster = kmeans(tripDetails, 3)

→ elbow method

Method to calculate optimal K.

```
k.max <- 10
wss <- rep(NA, k.max)
nclust <- list()
for (i in 1:k.max) {
  driveClasses <- kmeans(tripDetails, i)
  wss[i] <- driveClasses$tot.withinss
  nclust[[i]] <- driveClasses$size
}
plot(1:k.max, wss, type = "b", pch = 19,
  xlab = "Number of clusters K",
  ylab = "Total within-clusters sum of
  squares : Trips")
```