

STAT 3690 Lecture 34

zhiyanggeezhou.github.io

Zhiyang Zhou (zhiyang.zhou@umanitoba.ca)

Apr 22, 2022

Clustering

- Problem: given observations $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ group the observations into K populations
 - Unknown K
 - Unsupervised: no label/training data
 - Why
 - Summarize a representation of the full data set
 - Exploration for structure of the data
 - Checking the validity of pre-existing group assignments
 - Assistance for prediction: sometimes clustering prior to prediction
 - Clustering $C : \mathbb{Z}^+ \rightarrow \mathbb{Z}^+$
 - $C(i) = k$: assign \mathbf{x}_i to group k
-

K -means

- Within-cluster scatter

$$W(K) = \frac{1}{2} \sum_{k=1}^K \frac{1}{n_k} \sum_{i:C(i)=k} \sum_{j:C(j)=k} \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 = \sum_{k=1}^K \sum_{i:C(i)=k} \|\mathbf{x}_i - \bar{\mathbf{x}}_k\|_2^2$$

- $\|\mathbf{x}_i - \mathbf{x}_j\|_2$: the Euclidean distance between \mathbf{x}_i and \mathbf{x}_j
 - $\bar{\mathbf{x}}_k = \frac{1}{n_k} \sum_{i:C(i)=k} \mathbf{x}_i$
 - Smaller $W(K)$ is better
 - To minimize the within-cluster scatter
$$\min_C W(K) = \min_{C, \mathbf{c}_1, \dots, \mathbf{c}_K} \sum_{k=1}^K \sum_{i:C(i)=k} \|\mathbf{x}_i - \mathbf{c}_k\|_2^2$$
 - Implementation:
 1. Specify K and start with an initial guess for $\mathbf{c}_1, \dots, \mathbf{c}_K$, then repeat
 - a. Labeling each point based the closest center: for each i , put \mathbf{x}_i to the k th cluster such that \mathbf{c}_k is closest to \mathbf{x}_i
 - b. Replacing each center by the average of points in its cluster: for each k , take $\mathbf{c}_k = \bar{\mathbf{x}}_k$
 2. Terminate when $W(K)$ doesn't change
-

- Comments

- Always converge
- No guarantee to lead to the smallest W
- Depend on K and initial cluster centers
 - * Typically run K -means multiple times and pick up the result with the smallest W
- Determination of K
 - Between-cluster variation

$$B(K) = \sum_{k=1}^K n_k \|\bar{\mathbf{x}}_k - \bar{\mathbf{x}}\|_2^2$$

- * $\bar{\mathbf{x}} = n^{-1} \sum_{i=1}^n \mathbf{x}_i$
- CH index (Caliński & Harabasz (1974), *Communications in Statistics*, 3:1–27)

$$\text{CH}(K) = \frac{B(K)/(K-1)}{W(K)/(n-K)}$$

- To choose K by maximizing $\text{CH}(K)$, i.e.,

$$\hat{K} = \arg \max_{K \in \{2, \dots, K_{\max}\}} \text{CH}(K)$$

- Example (iris)

```
options(digits = 4)
x = iris[, !(names(iris) %in% c('Species'))]
y = (iris$Species == unique(iris$Species)[1]) +
  2*(iris$Species == unique(iris$Species)[2]) +
  3*(iris$Species == unique(iris$Species)[3])
decomp = prcomp(x)
s = 2
PCscores = decomp$x[,1:s]
K = 3; cols = c("red", "darkgreen", "blue", "pink", "purple")
set.seed(3690); km = kmeans(PCscores, centers=K, nstart=100, algorithm="Lloyd", iter.max = 100)
# cluster plot with centers
plot(PCscores, col=cols[km$cluster]); points(km$centers, pch=19, cex=2, col=cols)
# comparison with true groups
par(mfrow=c(1,2)); plot(PCscores, col=cols[km$cluster], main="K-means"); plot(PCscores, col=cols[y], main="True groups")

# determine K
Ks = 2:20
Ws = numeric(length(Ks))
Bs = numeric(length(Ks))
CHs = numeric(length(Ks))

for(l in 1:length(Ks)){
  set.seed(3690); km = kmeans(PCscores, centers=Ks[l], nstart=25, algorithm="Lloyd", iter.max = 100)
  Ws[l] = km$tot.withinss
  Bs[l] = sum(km$size * rowSums(sweep(km$centers, 2, colMeans(PCscores))^2))
  CHs[l] = (Bs[l]/(Ks[l]-1))/(Ws[l]/(nrow(PCscores)-Ks[l]))
}
plot(Ks, CHs,
     type="b", pch = 19,
     xlab="Number of clusters K",
     ylab="CH index")
```

- An application to image compression/color quantization
 - Basic idea: compress images by reducing the color palette of an image to K colors

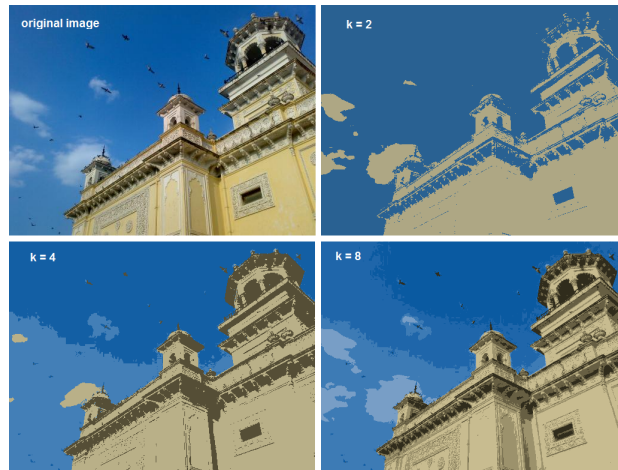


Figure 1: Image compression with K -means clustering (<http://opencvpython.blogspot.com/2012/12/k-means-clustering-2-working-with-scipy.html>)