Clustering 3: Hierarchical clustering (continued); choosing the number of clusters

> Ryan Tibshirani Data Mining: 36-462/36-662

> > January 31 2013

Optional reading: ISL 10.3, ESL 14.3

Even more linkages

Last time we learned about hierarchical agglomerative clustering, basic idea is to repeatedly merge two most similar groups, as measured by the linkage

Three linkages: single, complete, average linkage. Properties:

- Single and complete linkage can have problems with chaining and crowding, respectively, but average linkage doesn't
- Cutting an average linkage tree provides no interpretation, but there is a nice interpretation for single, complete linkage trees
- Average linkage is sensitive to a monotone transformation of the dissimilarities d_{ij}, but single and complete linkage are not
- All three linkages produce dendrograms with no inversions

Actually, there are many more linkages out there, each having different properties. Today: we'll look at two more

Reminder: linkages

Our setup: given X_1, \ldots, X_n and pairwise dissimilarities d_{ij} . (E.g., think of $X_i \in \mathbb{R}^p$ and $d_{ij} = ||X_i - X_j||_2$)

Single linkage: measures the closest pair of points

$$d_{\mathsf{single}}(G,H) = \min_{i \in G, \, j \in H} d_{ij}$$

Complete linkage: measures the farthest pair of points

$$d_{\mathsf{complete}}(G, H) = \max_{i \in G, \, j \in H} d_{ij}$$

Average linkage: measures the average dissimilarity over all pairs

$$d_{\text{average}}(G, H) = \frac{1}{n_G \cdot n_H} \sum_{i \in G, j \in H} d_{ij}$$

Centroid linkage

Centroid linkage¹ is commonly used. Assume that $X_i \in \mathbb{R}^p$, and $d_{ij} = ||X_i - X_j||_2$. Let \bar{X}_G, \bar{X}_H denote group averages for G, H. Then: $d = (C, H) = ||\bar{X}_G, \bar{X}_H||_2$

 $d_{\mathsf{centroid}}(G,H) = \|\bar{X}_G - \bar{X}_H\|_2$

Example (dissimilarities d_{ij} are distances, groups are marked by colors): centroid linkage score $d_{centroid}(G, H)$ is the distance between the group centroids (i.e., group averages)



¹Eisen et al. (1998), "Cluster Analysis and Display of Genome-Wide Expression Patterns"

Centroid linkage is the standard in biology

Centroid linkage is simple: easy to understand, and easy to implement. Maybe for these reasons, it has become the standard for hierarchical clustering in biology





Chang, Howard Y. et al. (2005) Proc. Natl. Acad. Sci. USA 102, 3738-3743

Centroid linkage example

Here n = 60, $X_i \in \mathbb{R}^2$, $d_{ij} = ||X_i - X_j||_2$. Cutting the tree at some heights wouldn't make sense ... because the dendrogram has inversions! But we can, e.g., still look at ouptut with 3 clusters



Cut interpretation: there isn't one, even with no inversions

Shortcomings of centroid linkage

- Can produce dendrograms with inversions, which really messes up the visualization
- Even if were we lucky enough to have no inversions, still no interpretation for the clusters resulting from cutting the tree
- ► Answers change with a monotone transformation of the dissimilarity measure d_{ij} = ||X_i X_j||₂. E.g., changing to d_{ij} = ||X_i X_j||₂² would give a different clustering



Minimax linkage

Minimax linkage² is a newcomer. First define radius of a group of points G around X_i as $r(X_i, G) = \max_{j \in G} d_{ij}$. Then:

$$d_{\min(G,H)} = \min_{i \in G \cup H} r(X_i, G \cup H)$$

Example (dissimilarities d_{ij} are distances, groups marked by colors): minimax linkage score $d_{\min\max}(G, H)$ is the smallest radius encompassing all points in G and H. The center X_c is the black point



²Bien et al. (2011), "Hierarchical Clustering with Prototypes via Minimax Linkage"

Minimax linkage example

Same data s before. Cutting the tree at $h=2.5 \ {\rm gives} \ {\rm clustering}$ assignments marked by the colors



Cut interpretation: each point X_i belongs to a cluster whose center X_c satisfies $d_{ic} \leq 2.5$

Properties of minimax linkage

- ► Cutting a minimax tree at a height h a nice interpretation: each point is ≤ h in dissimilarity to the center of its cluster. (This is related to a famous set cover problem)
- Produces dendrograms with no inversions
- Unchanged by monotone transformation of dissimilarities d_{ij}
- Produces clusters whose centers are chosen among the data points themselves. Remember that, depending on the application, this can be a very important property. (Hence minimax clustering is the analogy to *K*-medoids in the world of hierarchical clustering)

Example: Olivetti faces dataset



(From Bien et al. (2011))



(From Bien et al. (2011))

Centroid and minimax linkage in R

The function hclust in the base package performs hierarchical agglomerative clustering with centroid linkage (as well as many other linkages)

E.g.,

```
d = dist(x)
tree.cent = hclust(d, method="centroid")
plot(tree.cent)
```

The function protoclust in the package protoclust implements hierarchical agglomerative clustering with minimax linkage

Linkages summary

Linkage	No inversions?	Unchanged with monotone transformation?	Cut interpretation?	Notes
Single	\checkmark	\checkmark	\checkmark	chaining
Complete	\checkmark	\checkmark	\checkmark	crowding
Average	\checkmark	×	×	
Centroid	×	×	×	simple
Minimax	\checkmark	\checkmark	\checkmark	centers are data points

Note: this doesn't tell us what "best linkage" is

What's missing here: a detailed empirical comparison of how they perform. On top of this, remember that choosing a linkage can be very situation dependent

Designing a clever radio system (e.g., Pandora)

Suppose we have a bunch of songs, and dissimilarity scores between each pair. We're building a clever radio system—a user is going to give us an initial song, and a measure of how "risky" he is going to be, i.e., maximal tolerable dissimilarity between suggested songs



How could we use hierarchical clustering, and with what linkage?

Placing cell phone towers

Suppose we are helping to place cell phone towers on top of some buildings throughout the city. The cell phone company is looking to build a small number of towers, such that no building is further than half a mile from a tower



How could we use hierarchical clustering, and with what linkage?

How many clusters?

Sometimes, using K-means, K-medoids, or hierarchical clustering, we might have no problem specifying the number of clusters K ahead of time, e.g.,

- \blacktriangleright Segmenting a client database into K clusters for K salesman
- Compressing an image using vector quantization, where K controls the compression rate

Other times, K is implicitly defined by cutting a hierarchical clustering tree at a given height, e.g., designing a clever radio system or placing cell phone towers

But in most exploratory applications, the number of clusters K is unknown. So we are left asking the question: what is the "right" value of K?

This is a hard problem

Determining the number of clusters is a hard problem!

Why is it hard?

Determining the number of clusters is a hard task for humans to perform (unless the data are low-dimensional). Not only that, it's just as hard to explain what it is we're looking for. Usually, statistical learning is successful when at least one of these is possible

Why is it important?

- ► E.g., it might mean a big difference scientifically if we were convinced that there were K = 2 subtypes of breast cancer vs. K = 3 subtypes
- One of the (larger) goals of data mining/statistical learning is automatic inference; choosing K is certainly part of this

Reminder: within-cluster variation

We're going to focus on $K\mbox{-means},$ but most ideas will carry over to other settings

Recall: given the number of clusters K, the K-means algorithm approximately minimizes the within-cluster variation:

$$W = \sum_{k=1}^{K} \sum_{C(i)=k} \|X_i - \bar{X}_k\|_2^2$$

over clustering assignments C, where \bar{X}_k is the average of points in group k, $\bar{X}_k=\frac{1}{n_k}\sum_{C(i)=k}X_i$

Clearly a lower value of W is better. So why not just run K-means for a bunch of different values of K, and choose the value of K that gives the smallest W(K)?

That's not going to work

Problem: within-cluster variation just keeps decreasing

Example: $n = 250, p = 2, K = 1, \dots 10$



20

Between-cluster variation

Within-cluster variation measures how tightly grouped the clusters are. As we increase the number of clusters K, this just keeps going down. What are we missing?

Between-cluster variation measures how spread apart the groups are from each other:

$$B = \sum_{k=1}^{K} n_K \|\bar{X}_k - \bar{X}\|_2^2$$

where as before \bar{X}_k is the average of points in group k, and \bar{X} is the overall average, i.e.

$$\bar{X}_k = \frac{1}{n_k} \sum_{C(i)=k} X_i$$
 and $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$

Example: between-cluster variation

Example: n = 100, p = 2, K = 2



Still not going to work

Bigger B is better, can we use it to choose K? Problem: betweencluster variation just keeps increasing

Running example: n = 250, p = 2, $K = 1, \dots 10$



CH index

Ideally we'd like our clustering assignments C to simultaneously have a small W and a large B

This is the idea behind the CH index.³ For clustering assignments coming from K clusters, we record CH score:

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

To choose K, just pick some maximum number of clusters to be considered K_{max} (e.g., K = 20), and choose the value of K with the largest score CH(K), i.e.,

$$\hat{K} = \operatorname*{argmax}_{K \in \{2, \dots, K_{\mathsf{max}}\}} \mathsf{CH}(K)$$

³Calinski and Harabasz (1974), "A dendrite method for cluster analysis"

Example: CH index

Running example: n = 250, p = 2, K = 2, ... 10.



We would choose K = 4 clusters, which seems reasonable

General problem: the CH index is not defined for K = 1. We could never choose just one cluster (the null model)!

Gap statistic

It's true that W(K) keeps dropping, but how much it drops at any one K should be informative

The gap statistic⁴ is based on this idea. We compare the observed within-cluster variation W(K) to $W_{\text{unif}}(K)$, the within-cluster variation we'd see if we instead had points distributed uniformly (over an encapsulating box). The gap for K clusters is defined as

$$\mathsf{Gap}(K) = \log W_{\mathsf{unif}}(K) - \log W(K)$$

The quantity $\log W_{\text{unif}}(K)$ is computed by simulation: we average the log within-cluster variation over, say, 20 simulated uniform data sets. We also compute the standard error of s(K) of $\log W_{\text{unif}}(K)$ over the simulations. Then we choose K by $\hat{K} = \min \left\{ K \in \{1, \ldots K_{\max}\} : \operatorname{Gap}(K) \ge \operatorname{Gap}(K+1) - s(K+1) \right\}$

 $^{^{4}\}text{Tibshirani}$ et al. (2001), "Estimating the number of clusters in a data set via the gap statistic"

Example: gap statistic

Running example: n = 250, p = 2, K = 1, ... 10



We would choose K = 3 clusters, which is also reasonable

The gap statistic does especially well when the data fall into one cluster. (Why? Hint: think about the null distribution that it uses)

CH index and gap statistic in R

The CH index can be computed using the kmeans function in the base distribution, which returns both the within-cluster variation and the between-cluster variation (Homework 2)

```
E.g.,
k = 5
km = kmeans(x, k, alg="Lloyd")
names(km)
# Now use some of these return items to compute ch
```

The gap statistic is implemented by the function gap in the package lga, and by the function gap in the package SAGx. (Beware: these functions are poorly documented ... it's unclear what clustering method they're using)

Once again, it really is a hard problem



(Taken from George Cassella's CMU talk on January 16 2011)

an 7 (Confidence set)	
6 has highest posterior probability	
The posterior clearly supports three groups	
Optimal number of clusters is three	

(From George Cassella's CMU talk on January 16 2011)

Recap: more linkages, and determining K

Centroid linkage is commonly used in biology. It measures the distance between group averages, and is simple to understand and to implement. But it also has some drawbacks (inversions!)

Minimax linkage is a little more complex. It asks the question: "which point's furthest point is closest?", and defines the answer as the cluster center. This could be useful for some applications

Determining the number of clusters is both a hard and important problem. We can't simply try to find K that gives the smallest achieved within-class variation. We defined between-cluster variation, and saw we also can't choose K to just maximize this

Two methods for choosing K: the CH index, which looks at a ratio of between to within, and the gap statistic, which is based on the difference between within-class variation for our data and what we'd see from uniform data

Next time: principal components analysis

Finding interesting directions in our data set



(From ESL page 67)