

Clustering

Tony Ni, Antonella Basso, Jose Lopez

6/21/2020

Libraries

Reading in Data

```
setwd("~/harvard-summer-biostats")
df <- read_csv("data/wide_illinois.csv") #read in data
```

```
## Parsed with column specification:
## cols(
##   .default = col_double(),
##   well_id = col_character(),
##   site = col_character(),
##   disposal_area = col_character(),
##   type = col_character(),
##   gradient = col_character()
## )

## See spec(...) for full column specifications.
```

```
df1 <- na.omit(df) #get rid of na's
```

Hierarchical Methods

If we want to look for cereal groups via hierarchical clustering, we need to construct a distance matrix. Distances are constructed with the *dist* function, and we need to choose whether we compute them on scaled or unscaled variables (standardize or not).

```
df.dist <- dist(df1[, -c(1:5)])
```

Now we look at how hierarchical clustering is applied. The relevant function is *hclust*.

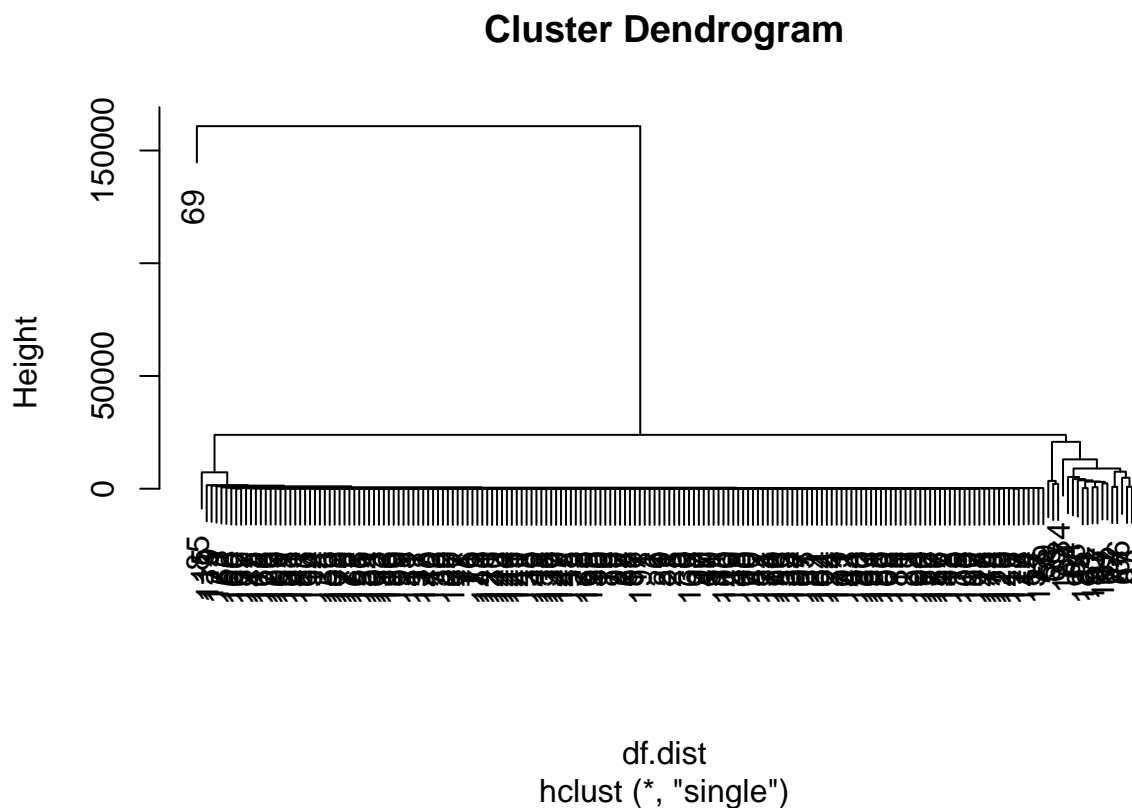
```
hcsingle <- hclust(df.dist, method = "single")
list(hcsingle) # reminds you of properties of the solution, if desired
```

```
## [[1]]
##
## Call:
```

```
## hclust(d = df.dist, method = "single")
##
## Cluster method   : single
## Distance        : euclidean
## Number of objects: 191
```

This creates the solution, and we can look at the dendrogram as:

```
plot(hcsingle)
```



The options for hclust in terms of linkages are provided in the help under options for method. The following options are listed: “ward.D”, “ward.D2”, “single”, “complete”, “average”, “mcquitty”, “median” or “centroid”.

In order to obtain cluster labels, we need to *cut* our dendrograms.

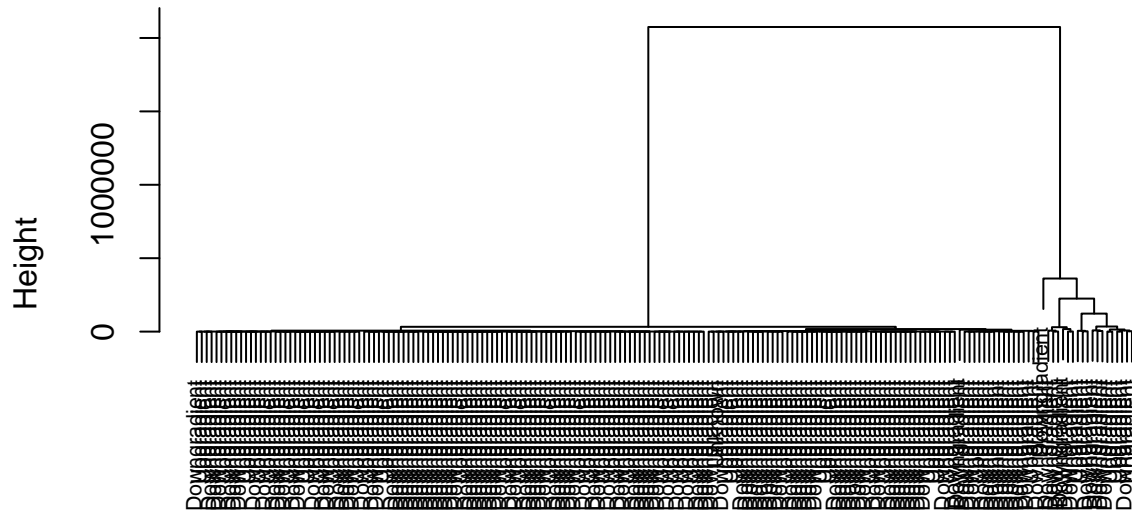
```
singleSol <- (cutree(hcsingle, k = 3)) #cluster labels are numeric, k= # clusters
summary(as.factor(singleSol)) #as factor to get table
```

```
##   1   2   3
## 172 18   1
```

To learn more details about the clusters we found:

```
hward <- hclust(df.dist, method = "ward.D")
plot(hward, labels = df1$gradient, cex = 0.7) #cex adjusts size of label
```

Cluster Dendrogram



df.dist
hclust (*, "ward.D")

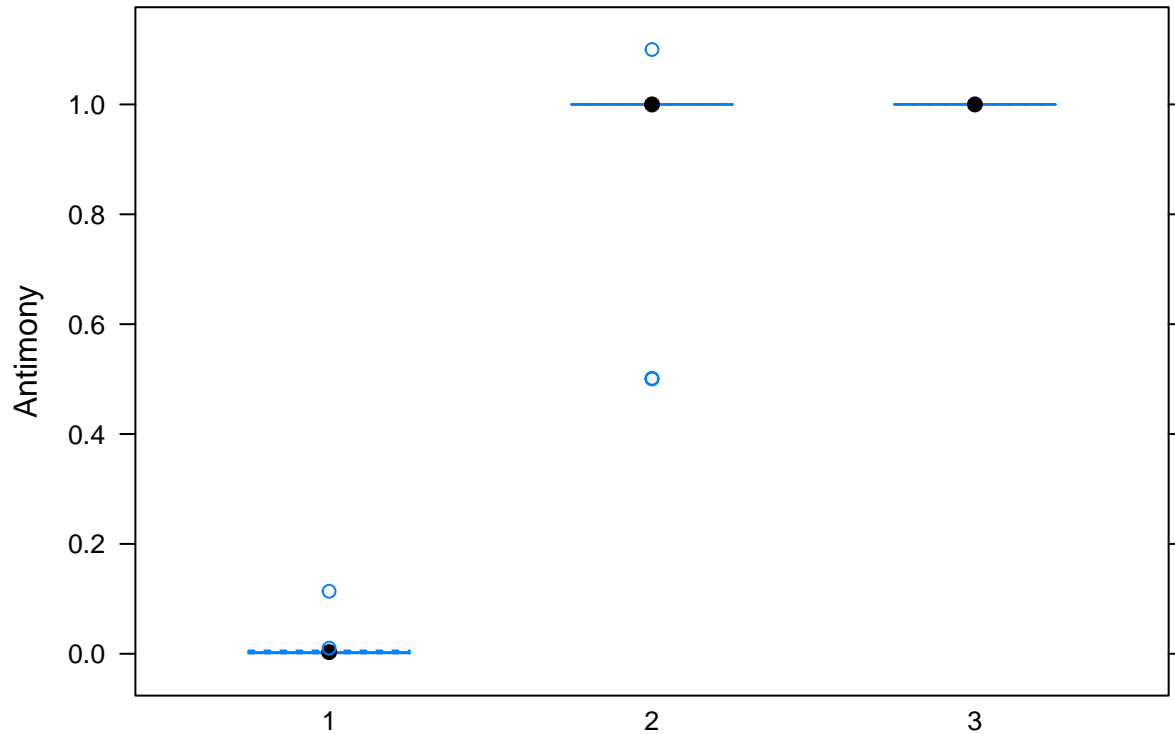
```
wardSol <- (cutree(hward, k = 3)) #cluster labels are numeric, k= # clusters
summary(as.factor(wardSol)) #as factor to get table
```

```
##   1   2   3
## 172  18   1
```

```
favstats(Antimony ~ wardSol, data = df1) #can choose any variable
```

```
##   wardSol      min    Q1 median    Q3      max      mean      sd    n
## 1       1 0.0007625 0.001  0.003 0.003 0.1137778 0.002903086 0.00858458 172
## 2       2 0.5005000 1.000  1.000 1.000 1.1000000 0.922361111 0.19538836  18
## 3       3 1.0000000 1.000  1.000 1.000 1.0000000 1.000000000      NA    1
##   missing
## 1       0
## 2       0
## 3       0
```

```
bwplot(Antimony ~ as.factor(wardSol), data = df1)
```



Our cluster sizes are extremely uneven... Our first cluster has 172 wells and the second has 18 and third has 1.

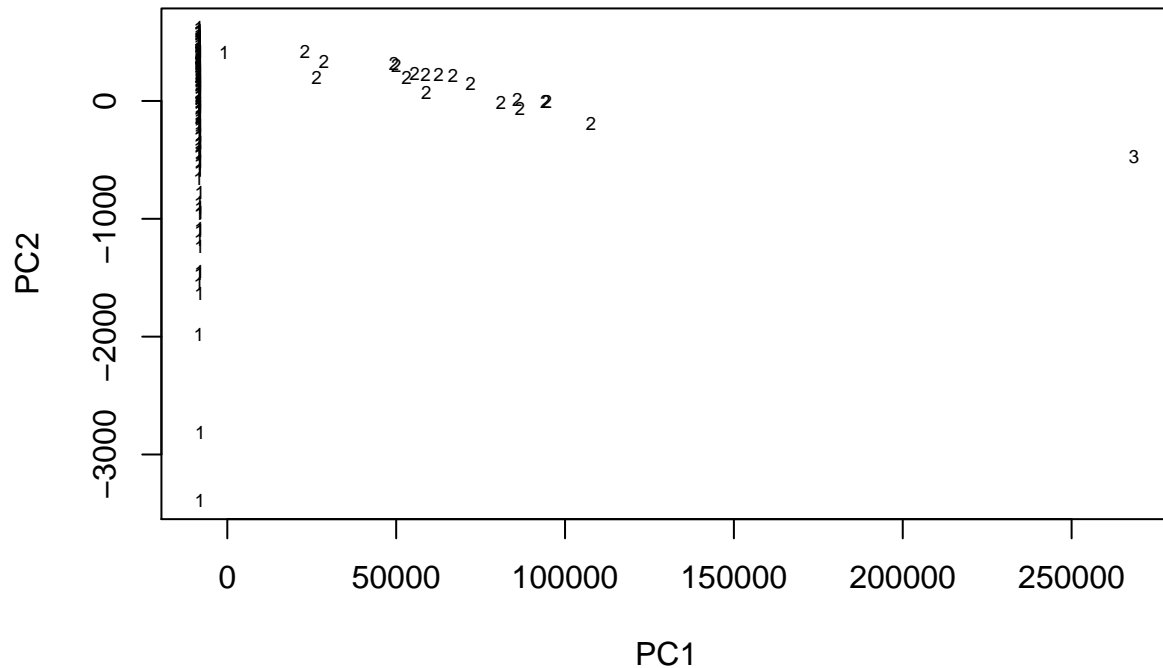
```
favstats(#chemical here ~ wardSol, data = df1)
bwplot(#chemical here ~ as.factor(wardSol), data = df1)
```

We can view the solution in the PC space (say 2-D) to see how well-separated the clusters are in that space. Because we used an unstandardized distance, we will run the PCA on the covariance matrix.

```
dfPCA <- princomp(df1[, -c(1:5)], cor = FALSE)

plot(dfPCA$scores[, 1:2], type = "n", xlab = "PC1", ylab = "PC2", main = "Ward's cluster solution") #bl
text(dfPCA$scores[, 1:2], labels = wardSol, cex = 0.6) #add the text
```

Ward's cluster solution



K-means Methods

For k-means, we don't need to compute the distance matrix ourselves. We feed the function the data set to operate on:

```
Ksol1 <- kmeans(scale(df1[, -c(1:5)]), centers = 3) #centers is the # of clusters desired
list(Ksol1) #so you can see what it gives you
```

```
## [[1]]
## K-means clustering with 3 clusters of sizes 145, 30, 16
##
## Cluster means:
##      Antimony      Arsenic      Barium  Beryllium      Boron      Cadmium      Calcium
## 1 -0.2874583 -0.2220554 -0.2411740 -0.2787967 -0.2896719 -0.2790253 -0.2544263
## 2 -0.3258770 -0.2579389 -0.2732023 -0.3141615 -0.2584159 -0.3120744 -0.2780177
## 3  3.2161104  2.4960121  2.6978935  3.1156481  3.1096817  3.1138065  2.8270214
##      Chloride      Chromium      Cobalt      Fluoride      Lead      Lithium      Mercury
## 1 -0.1476006 -0.1344754 -0.1548444 -0.2023032 -0.1277174 -0.2553406 -0.2872786
## 2  0.9239355 -0.1424800 -0.1477899  1.0504718 -0.1315048 -0.2683626 -0.3260804
## 3 -0.3947489  1.4858331  1.6803833 -0.1362618  1.4040106  2.8172043  3.2148635
##      Molybdenum      pH Radium 226+228      Selenium      Sulfate      Thallium
## 1 -0.230025 -0.09670553 -0.04453230 -0.279969 -0.2526068 -0.2891987
## 2 -0.236946  0.15741430  0.03779929 -0.317099  1.5046815 -0.3100661
## 3  2.528875  0.58124207  0.33270030  3.131780 -0.5320290  3.2022371
```

```
## Total Dissolved Solids
## 1 -0.3060546
## 2 1.7276146
## 3 -0.4656573
##
## Clustering vector:
## [1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1 1
## [38] 1 1 1 2 1 1 1 1 1 1 1 1 1 1 2 1 2 2 1 2 1 1 1 3 3 1 3 1 1 3 3 3 1 1 1 1 1
## [75] 1 1 3 1 1 1 1 1 1 3 3 1 3 3 1 1 1 1 1 1 1 3 1 1 1 1 1 1 1 1 1 1 1 1 2 1 1
## [112] 2 1 2 2 1 2 2 2 2 1 1 2 1 1 1 2 1 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## [149] 1 1 1 1 2 1 1 1 1 1 1 2 2 2 2 1 1 2 1 1 1 2 2 2 1 2 1 2 1 2 1 1 1 1 3 3 3
## [186] 1 3 1 1 1 1
##
## Within cluster sum of squares by cluster:
## [1] 482.6601 478.4676 792.9099
## (between_SS / total_SS = 56.0 %)
##
## Available components:
##
## [1] "cluster" "centers" "totss" "withinss" "tot.withinss"
## [6] "betweenss" "size" "iter" "ifault"
```

The list option provides us with lots of information. We can pull out the cluster means as:

```
Ksol1$centers
```

```
## Antimony Arsenic Barium Beryllium Boron Cadmium Calcium
## 1 -0.2874583 -0.2220554 -0.2411740 -0.2787967 -0.2896719 -0.2790253 -0.2544263
## 2 -0.3258770 -0.2579389 -0.2732023 -0.3141615 -0.2584159 -0.3120744 -0.2780177
## 3 3.2161104 2.4960121 2.6978935 3.1156481 3.1096817 3.1138065 2.8270214
## Chloride Chromium Cobalt Fluoride Lead Lithium Mercury
## 1 -0.1476006 -0.1344754 -0.1548444 -0.2023032 -0.1277174 -0.2553406 -0.2872786
## 2 0.9239355 -0.1424800 -0.1477899 1.0504718 -0.1315048 -0.2683626 -0.3260804
## 3 -0.3947489 1.4858331 1.6803833 -0.1362618 1.4040106 2.8172043 3.2148635
## Molybdenum pH Radium 226+228 Selenium Sulfate Thallium
## 1 -0.230025 -0.09670553 -0.04453230 -0.279969 -0.2526068 -0.2891987
## 2 -0.236946 0.15741430 0.03779929 -0.317099 1.5046815 -0.3100661
## 3 2.528875 0.58124207 0.33270030 3.131780 -0.5320290 3.2022371
## Total Dissolved Solids
## 1 -0.3060546
## 2 1.7276146
## 3 -0.4656573
```

We can also get the clustering vector (with the cluster labels) as:

```
Ksol1$cluster
```

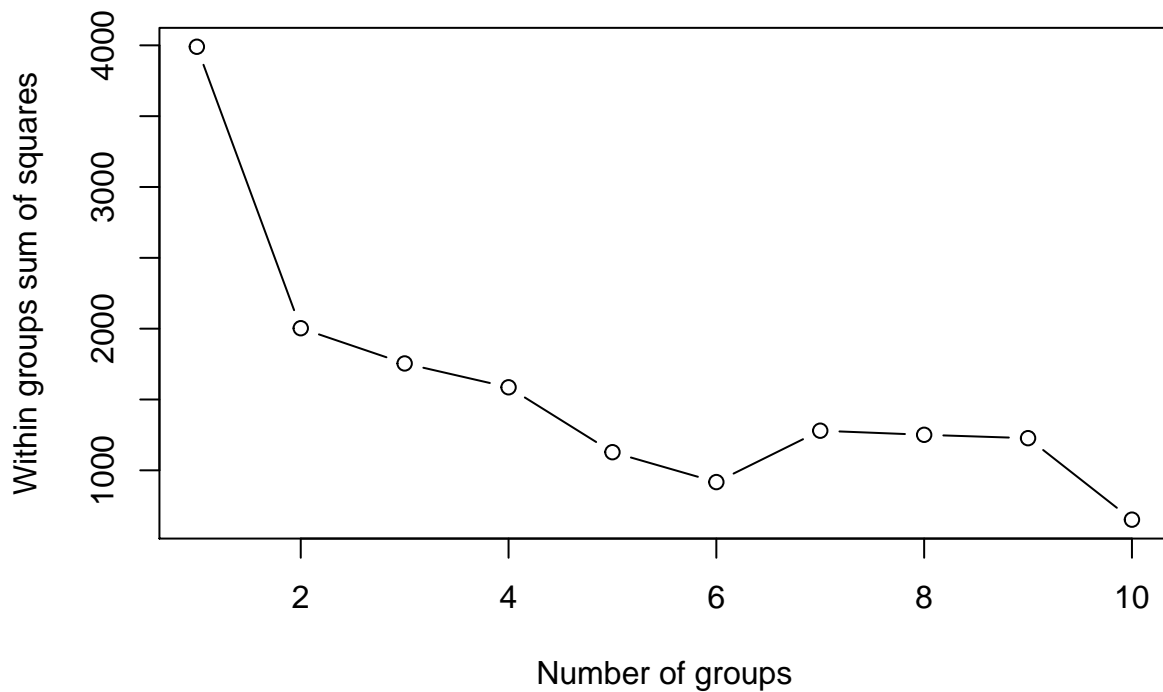
```
## [1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1 1
## [38] 1 1 1 2 1 1 1 1 1 1 1 1 1 1 2 1 2 2 1 2 1 1 1 3 3 1 3 1 1 3 3 3 1 1 1 1 1
## [75] 1 1 3 1 1 1 1 1 1 3 3 1 3 3 1 1 1 1 1 1 1 3 1 1 1 1 1 1 1 1 1 1 1 2 1 1
## [112] 2 1 2 2 1 2 2 2 2 1 1 2 1 1 1 2 1 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## [149] 1 1 1 1 2 1 1 1 1 1 1 2 2 2 2 1 1 2 1 1 1 2 2 2 1 2 1 2 1 2 1 1 1 3 3 3
## [186] 1 3 1 1 1 1
```

In order to determine if we have chosen a “good” value of the number of clusters, we can look at the within cluster sum of squares for this solution and a few other options for k, the number of clusters. This runs the solution from 1 to 10 clusters and pulls the within group sum of squares from each.

```
n <- nrow(df1) #number of observations

wss <- rep(0, 10) #creates 10 copies of 0 to create an empty vector
for(i in 1:10){
  wss[i] <- sum(kmeans(scale(df1[, -c(1:5)]), centers = i)$withinss)
}

plot(1:10, wss, type = "b", xlab = "Number of groups", ylab = "Within groups sum of squares")
```



We look for elbows in the plot - here there are elbows at 2 and 8, maybe these values will be good to use?

With two clusters, we should see if there is any relationship with sites...

```
tally(Ksol1$cluster ~ type, data = df1, format = "count")
```

```
##           type
## Ksol1$cluster  L  M  SI
##           1  32  5 108
##           2   0  0  30
##           3  16  0   0
```

There does seem to be something of interest here... The first cluster has 30 wells that are “SI,” the second

cluster has much more of a mix with 32 “L”, 5 “M”, and 108 “SI” wells. The third cluster has 16 wells that are “L”.

We can compare clustering solutions with similar tables. How do the K-means and Ward’s solutions overlap?

```
tally(Ksol1$cluster ~ wardSol, data = df1, format = "count")
```

```
##           wardSol
## Ksol1$cluster  1  2  3
##           1 142  3  0
##           2  30  0  0
##           3   0 15  1
```

Our solution with k-means puts many wells into cluster 2 which the ward’s solution put into cluster 1...