

Problem Set 4

Seth Harrison

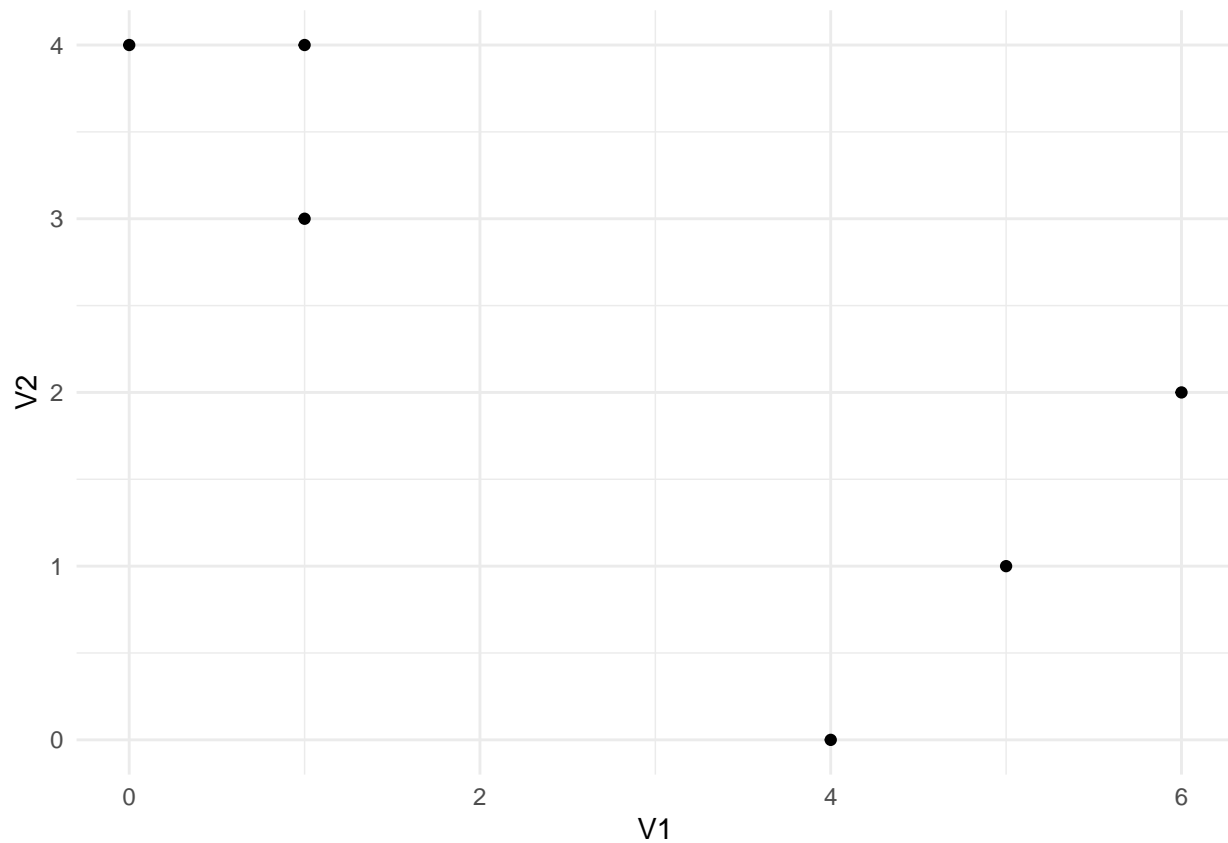
3/2/2020

```
library(tidyverse)
library(seriation)
library(knitr)
library(mixtools)
library(plotGMM)
library(clValid)
```

1. (5 points) Plot the observations.

```
data <- cbind(c(1, 1, 0, 5, 6, 4), c(4, 3, 4, 1, 2, 0)) %>%
  as.data.frame()

data%>%
  ggplot(aes(x = V1, y = V2)) +
  geom_point() +
  theme_minimal()
```



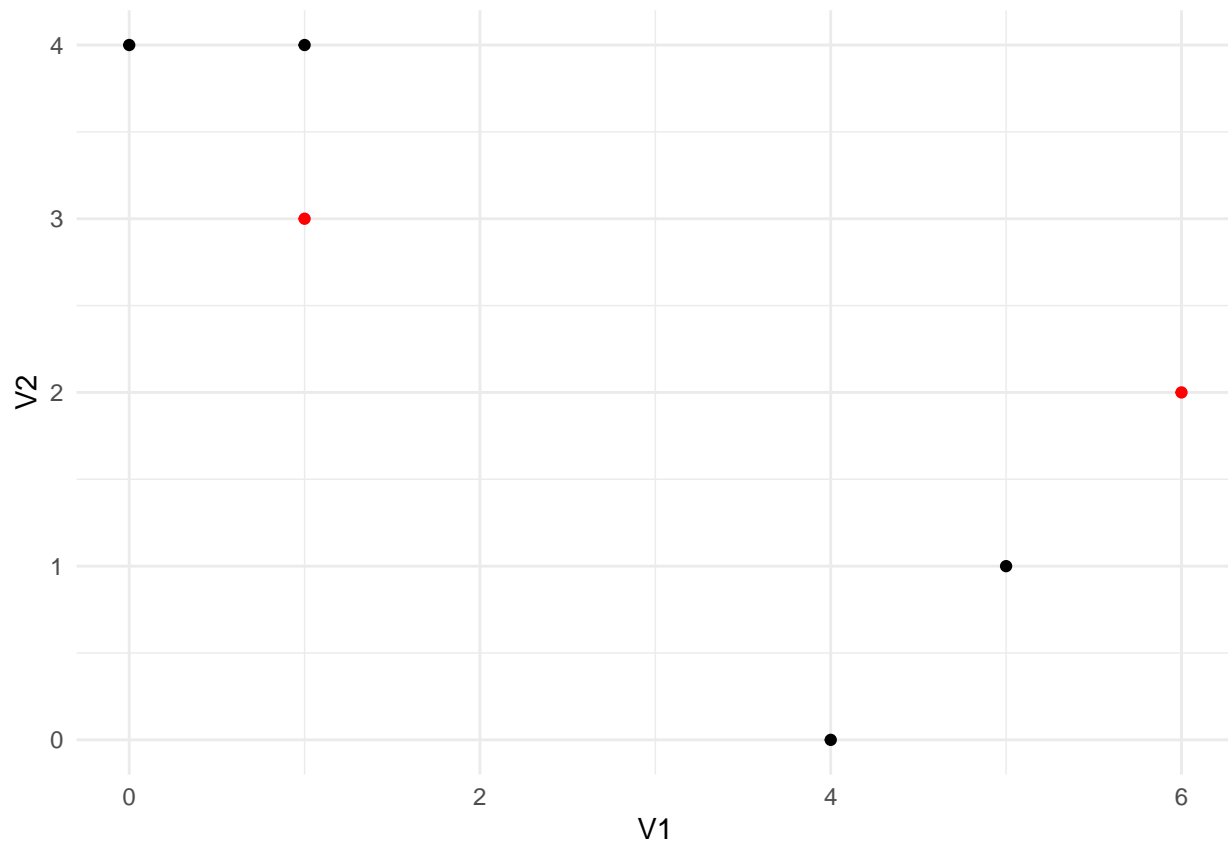
2. (5 points) Randomly assign a cluster label to each observation. Report the cluster labels for each observation and plot the results with a different color for each cluster (remember to set your seed first).

```
set.seed(1)

random <- sample(2, nrow(data), replace = TRUE)
random
```

```
## [1] 1 2 1 1 2 1
```

```
data %>%
  ggplot(aes(x = V1, y = V2)) +
  geom_point(color = random) +
  theme_minimal()
```

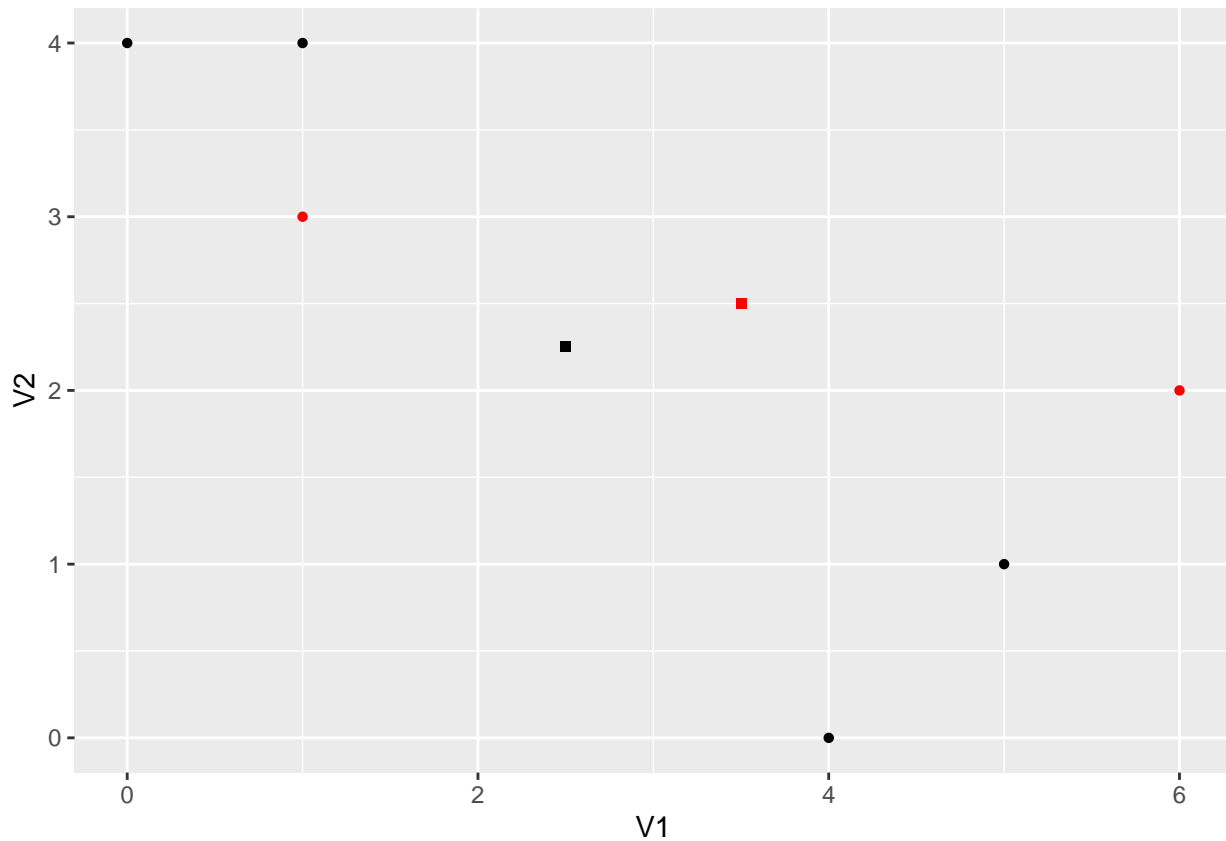


3. (10 points) Compute the centroid for each cluster.

```
centroid <- c(mean(data[random == 1, 1]),
              mean(data[random == 1, 2]))
centroid2 <- c(mean(data[random == 2, 1]),
               mean(data[random == 2, 2]))

data11 <- data %>%
  rbind(centroid, centroid2) %>%
  mutate(centroid_01 = if_else(V1 == 2.5 | V1 == 3.5, 15, 16)) %>%
  mutate(class = c(1,2,1,1,2,1,1,2))

data11 %>%
  ggplot(aes(x = V1, y = V2)) +
  geom_point(shape = data11$centroid_01,
            color = data11$class)
```



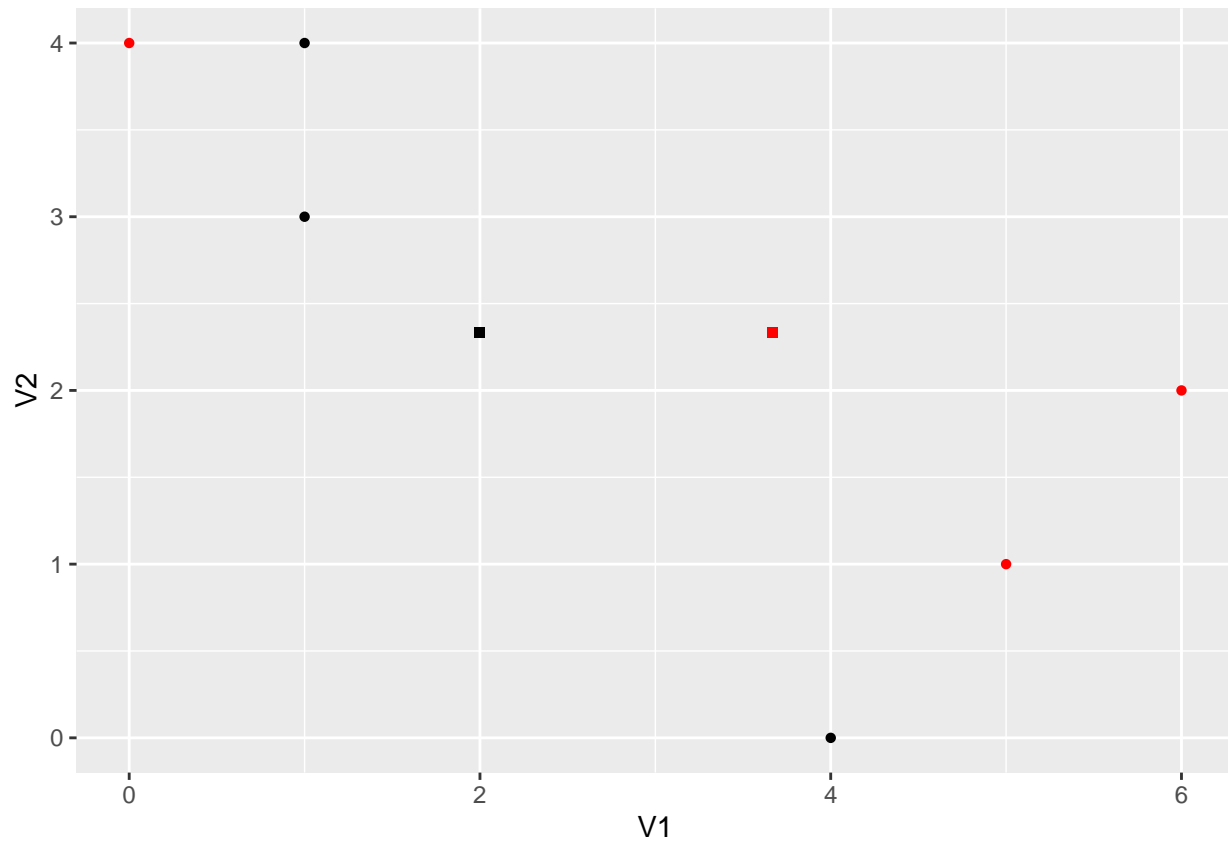
4. (10 points) Assign each observation to the centroid to which it is closest, in terms of Euclidean distance. Report the cluster labels for each observation.

```
notrandom <- c(1,1,2,2,2,1)

centroid3 <- c(mean(data[notrandom == 1, 1]),
               mean(data[notrandom == 1, 2]))
centroid4 <- c(mean(data[notrandom == 2, 1]),
               mean(data[notrandom == 2, 2]))

datai2 <- data %>%
  rbind(centroid3, centroid4) %>%
  mutate(centroid_01 = if_else(V1 == 2 | V1 == 11/3, 15, 16)) %>%
  mutate(class = c(1,1,2,2,2,1,1,2))

datai2 %>%
  ggplot(aes(x = V1, y = V2)) +
  geom_point(shape = datai2$centroid_01,
            color = datai2$class)
```



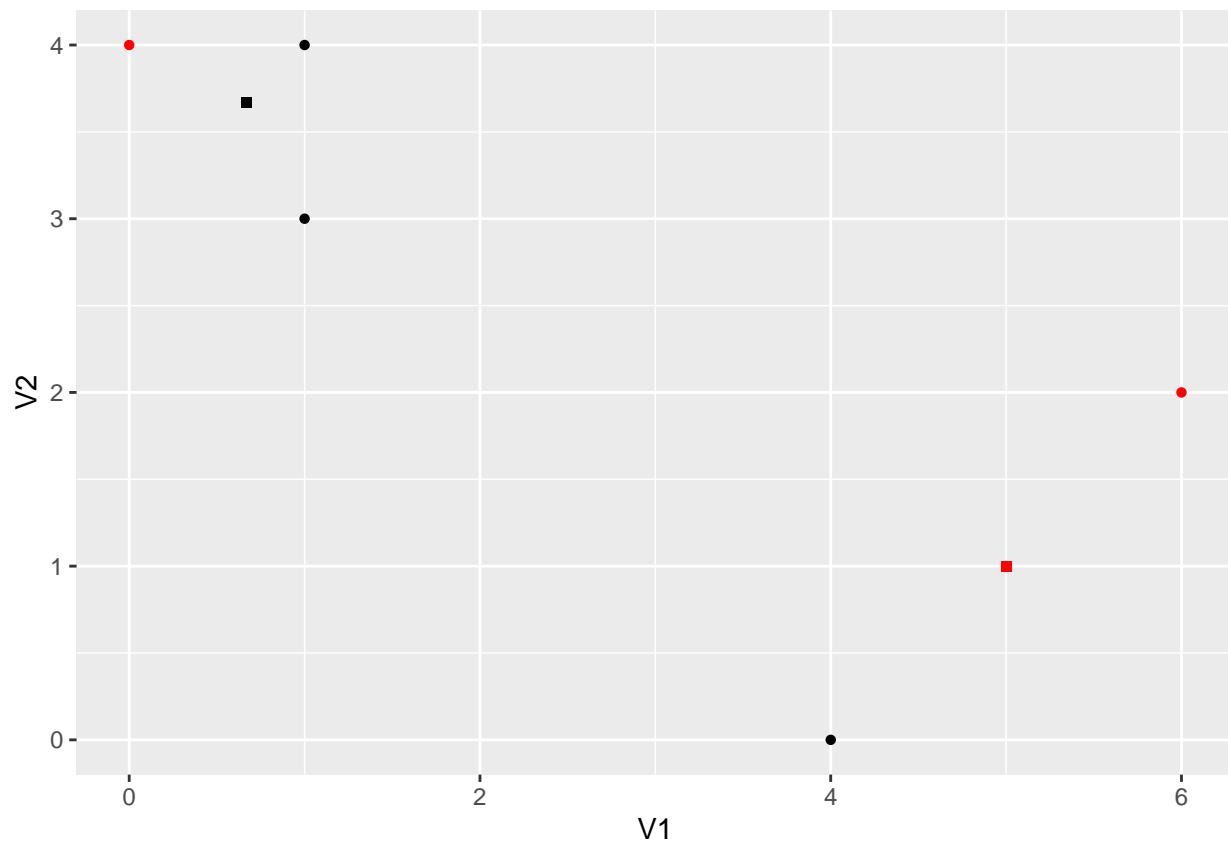
5. (5 points) Repeat (3) and (4) until the answers/clusters stop changing.

```
stable <- c(1,1,1,2,2,2)

centroid5 <- c(mean(data[stable == 1, 1]),
               mean(data[stable == 1, 2]))
centroid6 <- c(mean(data[stable == 2, 1]),
               mean(data[stable == 2, 2]))

datai3 <- data %>%
  rbind(centroid5, centroid6) %>%
  mutate(centroid_01 = if_else(V1 == 2/3 | V1 == 5, 15, 16)) %>%
  mutate(class = c(1,1,1,2,2,2,1,2))

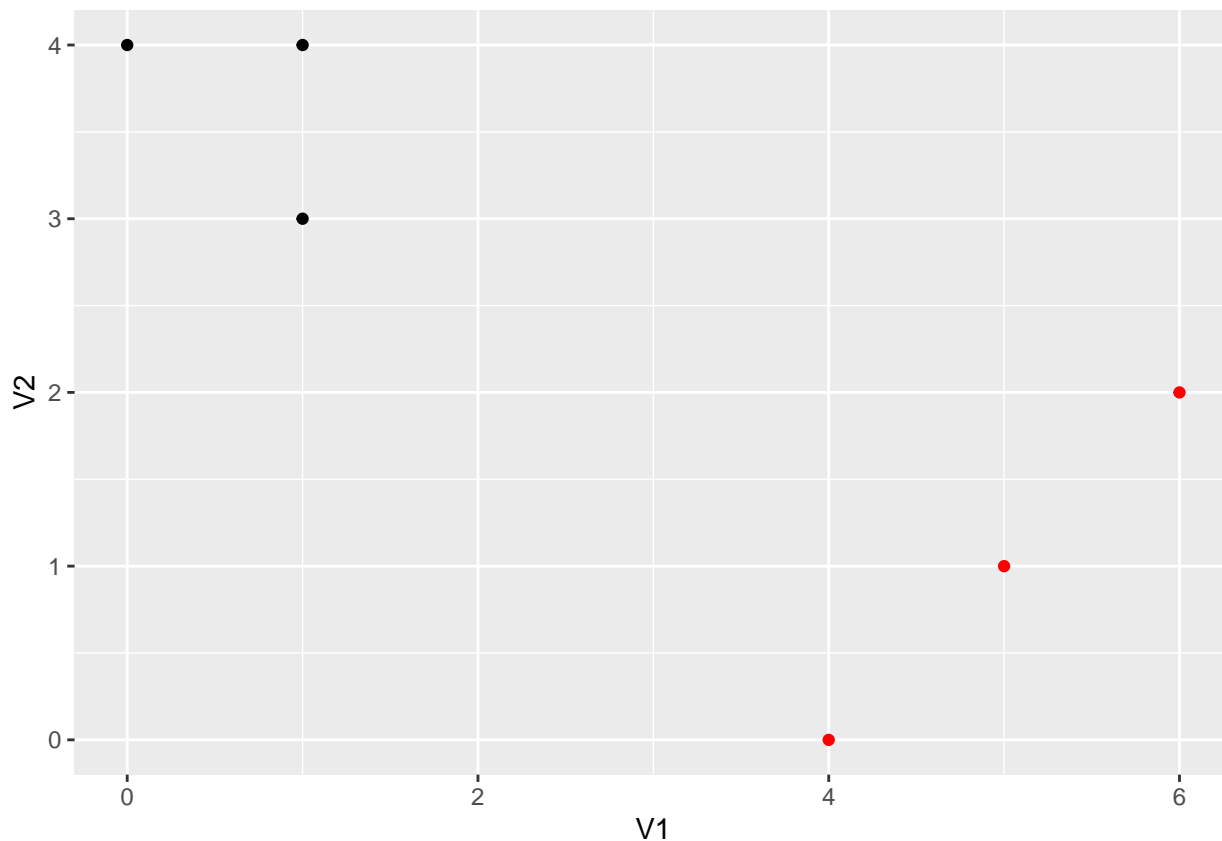
datai3 %>%
  ggplot(aes(x = V1, y = V2)) +
  geom_point(shape = datai2$centroid_01,
            color = datai2$class)
```



6. (10 points) Reproduce the original plot from (1), but this time color the observations according to the clusters labels you obtained by iterating the cluster centroid calculation and assignments.

```
data4 <- cbind(data,stable)

data4 %>%
  ggplot(aes(x = V1, y = V2)) +
  geom_point(color = stable)
```



Clustering State Legislative Professionalism

1. Load the state legislative professionalism data. See the codebook (or above) for further reference.

```
legprof <- load("~/myrepo/Machine Learning/Problem-Set-4/Data and Codebook/legprof-components.v1.0.RData")
legprof <- x
```

2. (5 points) Munge the data:

- select only the continuous features that should capture a state legislature's level of "professionalism" (session length (total and regular), salary, and expenditures);
- restrict the data to only include the 2009/10 legislative session for consistency;
- omit all missing values;
- standardize the input features;
- and anything else you think necessary to get this subset of data into workable form (hint: consider storing the state names as a separate object to be used in plotting later)

```
subset <- legprof %>%
  filter(sessid == "2009/10") %>%
  select(state,
         t_length,
```

```

        slength,
        salary_real,
        expend)

subset.scale <- scale(subset[, -c(1)])%>%
  na.omit()

states <- subset %>%
  select(state) %>%
  filter(state!= "Wisconsin")

```

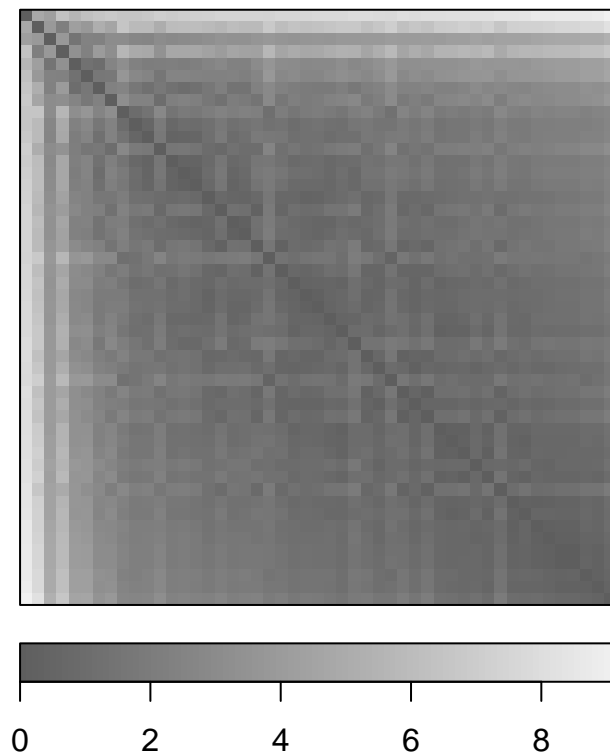
3. (5 points) Diagnose clusterability in any way you'd prefer (e.g., sparse sampling, ODI, etc.); display the results and discuss the likelihood that natural, non-random structure exist in these data. Hint: We didn't cover how to do this R in class, but consider `dissplot()` from the `seriation` package, the `factoextra` package, and others for calculating, presenting, and exploring the clusterability of some feature space.

```

ODI <- dist(subset.scale, method = "euclidean")

dissplot(ODI)

```

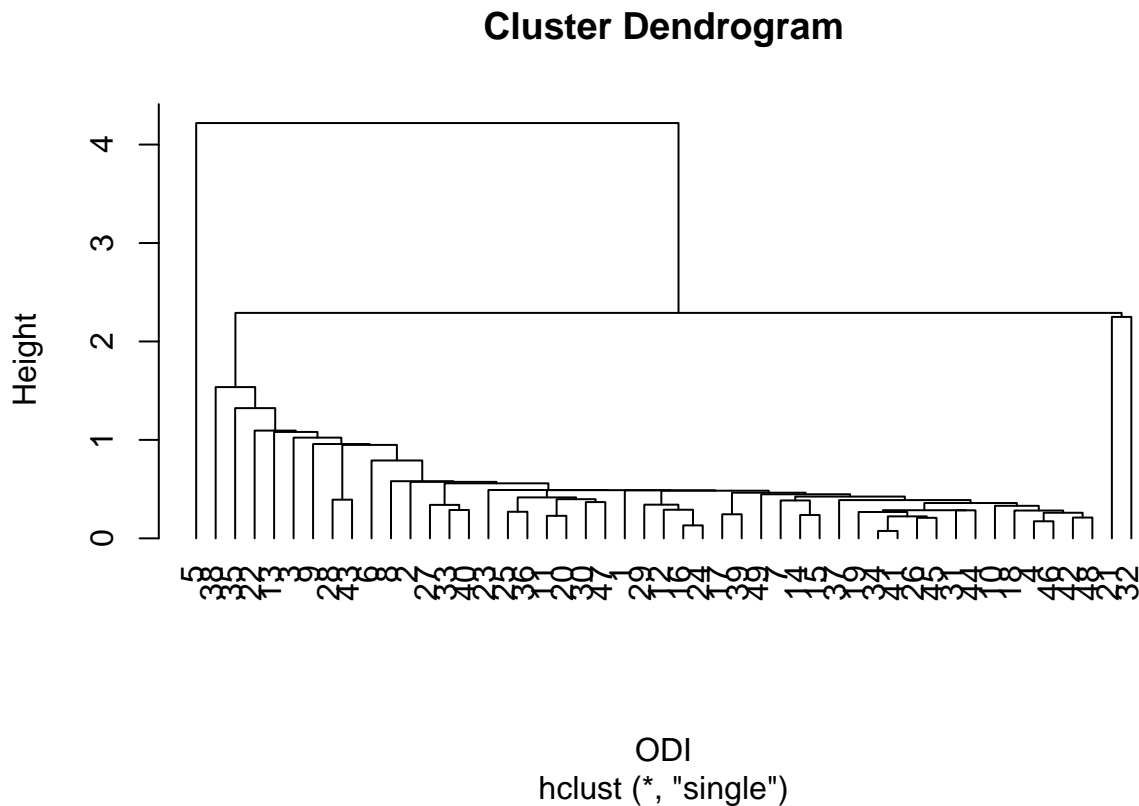


The data do not look particularly clusterable. I would like to see darker blocks along the diagonal of the matrix.

4. (5 points) Fit an agglomerative hierarchical clustering algorithm using any linkage method you prefer, to these data and present the results. Give a quick, high level summary of the output and general patterns.

```
AHC <- hclust(ODI, method = "single")
```

```
plot(AHC, hang = -1)
```



At the bottom, dendrogram shows several pairs of states with similarly professionalized legislatures. For example, Tennessee (42) is similar to West Virginia (48) and Arkansas (4) is similar to Virginia (46). At the top, the dendrogram shows that California (5) is unique compared to the other 48 (one is omitted) states in terms of professionalized legislatures.

5. (5 points) Fit a k-means algorithm to these data and present the results. Give a quick, high level summary of the output and general patterns. Initialize the algorithm at k=2, and then check this assumption in the validation questions below.

```
set.seed(1)

kmeans2 <- kmeans(subset.scale,
                  centers = 2,
                  nstart = 15)

t <- as.table(kmeans2$cluster)

t <- data.frame(t)
```

```
t$Var1 <- NULL

cbind(states,t) %>%
  kable()
```

state	Freq
Alabama	1
Alaska	1
Arizona	1
Arkansas	1
California	2
Colorado	1
Connecticut	1
Delaware	1
Florida	1
Georgia	1
Hawaii	1
Idaho	1
Illinois	1
Indiana	1
Iowa	1
Kansas	1
Kentucky	1
Louisiana	1
Maine	1
Maryland	1
Massachusetts	2
Michigan	2
Minnesota	1
Mississippi	1
Missouri	1
Montana	1
Nebraska	1
Nevada	1
New Hampshire	1
New Jersey	1
New Mexico	1
New York	2
North Carolina	1
North Dakota	1
Ohio	2
Oklahoma	1
Oregon	1
Pennsylvania	2
Rhode Island	1
South Carolina	1
South Dakota	1
Tennessee	1
Texas	1
Utah	1
Vermont	1
Virginia	1
Washington	1

state	Freq
West Virginia	1
Wyoming	1

Forty-three states were placed in the same cluster. Roughly, it seems that the six states placed in the second cluster tend to be more populous and all have major urban centers.

6. (5 points) Fit a Gaussian mixture model via the EM algorithm to these data and present the results. Give a quick, high level summary of the output and general patterns. Initialize the algorithm at $k = 2$, and then check this assumption in the validation questions below.

```
set.seed(1)

gmm1 <- normalmixEM(subset.scale, k = 2)

## number of iterations= 37

cbind(gmm1$mu, gmm1$sigma, gmm1$lambda) %>%
  as.data.frame() %>%
  rename(mu = V1) %>%
  rename(sigma = V2) %>%
  rename(lamda = V3) %>%
  kable()
```

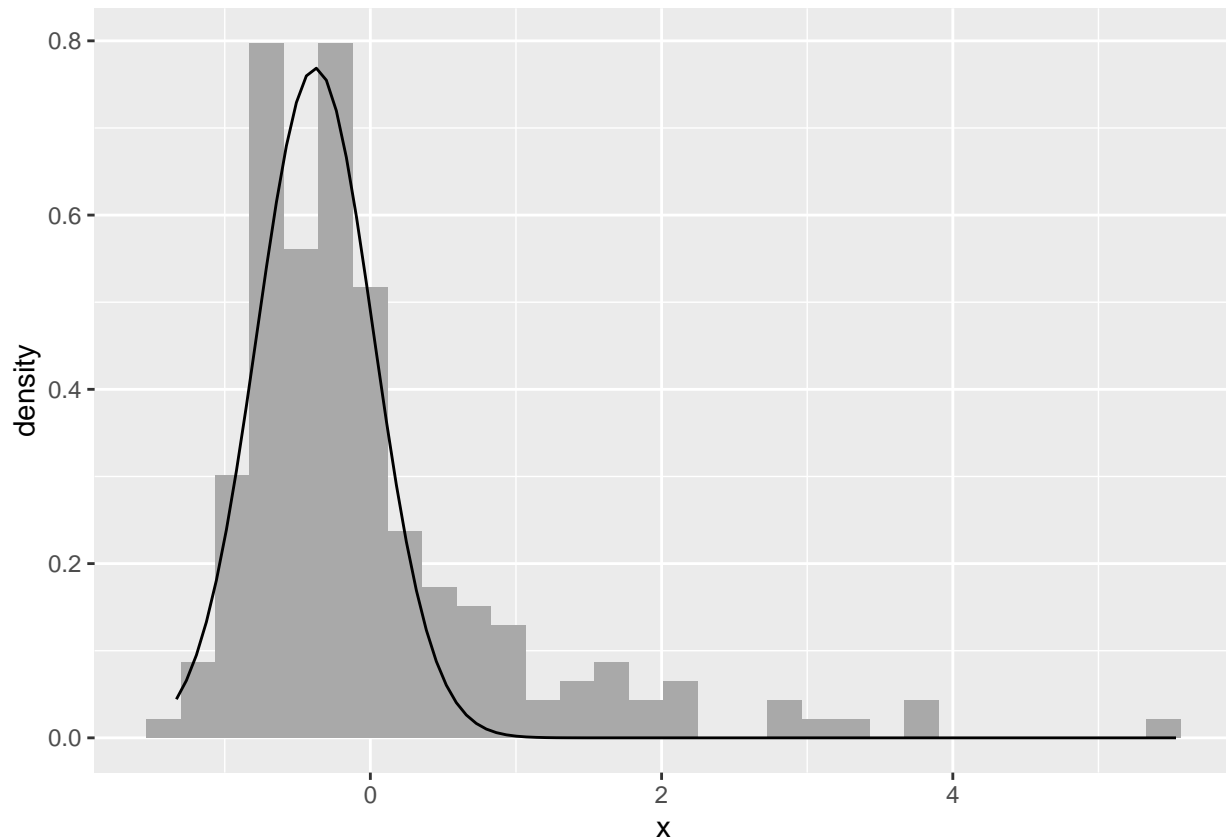
mu	sigma	lamda
-0.3784524	0.3991052	0.7690297
1.2378121	1.3143492	0.2309703

WORDS

7. (15 points) Compare output of all in visually useful, simple ways (e.g., present the dendrogram, plot by state cluster assignment across two features like salary and expenditures, etc.). There should be several plots of comparison and output.

```
ggplot(data.frame(x = gmm1$x)) +
  geom_histogram(aes(x, ..density..), fill = "darkgray") +
  stat_function(geom = "line", fun = plot_mix_comps,
    args = list(gmm1$mu[1][1], gmm1$sigma[1][1], lam = gmm1$lambda[1]),
    color = "black")
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



8. (5 points) Select a single validation strategy (e.g., compactness via $\min(\text{WSS})$, average silhouette width, etc.), and calculate for all three algorithms. Display and compare your results for all three algorithms you fit (hierarchical, k-means, GMM). Hint: Here again, we didn't cover this in R in class, but think about using the `clValid` package, though there are many other packages and ways to validate cluster patterns across iterations.

9. (10 points) Discuss the validation output, e.g., “What can you take away from the fit?”, “Which approach is optimal? And optimal at what value of k ?”, “What are reasons you could imagine selecting a technically “sub-optimal” clustering method, regardless of the validation statistics?”

Internal validation shows that AHC with 2 clusters is optimal. GMM allows for probabilistic assessments because, as a soft partitioning method, it allows clusters to overlap. AHC is also computationally expensive and requires that k and n be small.