

APSTA-GE 2352

Statistical Computing: Lecture 4

Klinton Kanopka

New York University



NYU Grey Art Gallery

RIVER CENTER

WASHINGTON PL



Table of Contents

1. Statistical Computing - Week 4

1. Table of Contents

2. Announcements

3. Total PS Points Needed by Desired Grade

4. Check-In

2. Motivating Problem

1. Making Groups from Data

3. Tools

1. Measuring Distances

2. `while` Loops

4. Back to Clustering

1. The k -Means Algorithm

2. Generalizing our k -Means Implementation

5. Wrap Up

1. Recap

2. Final Thoughts

Announcements

- PS1 Grades released
 - Correlation between PS0 and PS1 grades was low ($r \approx 0.28$)
 - Mean score was 7.7pts higher on PS1 than PS0
 - Seems fine to me
- Answer keys for both PS0 and PS1 are posted in the Week 1 materials
- PS2 is due next week before class

Total PS Points Needed by Desired Grade

Grade	PS Points
A	752
A-	707
B+	680
B	645

Note: this does not account for the final exam!

Check-In

- PollEv.com/klintkanopka

Motivating Problem

Making Groups from Data

- Often called "clustering"
- Big idea:
 - Observations in your dataset belong to *latent* groups
 - *Latent* means unobservable
 - The nature and number of groups is unobservable
 - The assignment of each point to a group is unobservable
 - We also need to figure all of this out at once
 - Conceptually, the clustering task takes observations that are "close together" and puts them into groups
 - Typically there are three ways to do clustering:
 - *Iterative*: Pick a number of clusters and rearrange them until you get your "best fit"
 - *Agglomerative*: Build up clusters by sticking nearby points together
 - *Divisive*: Start with one big group and repeatedly cut it into pieces

Making Groups from Data

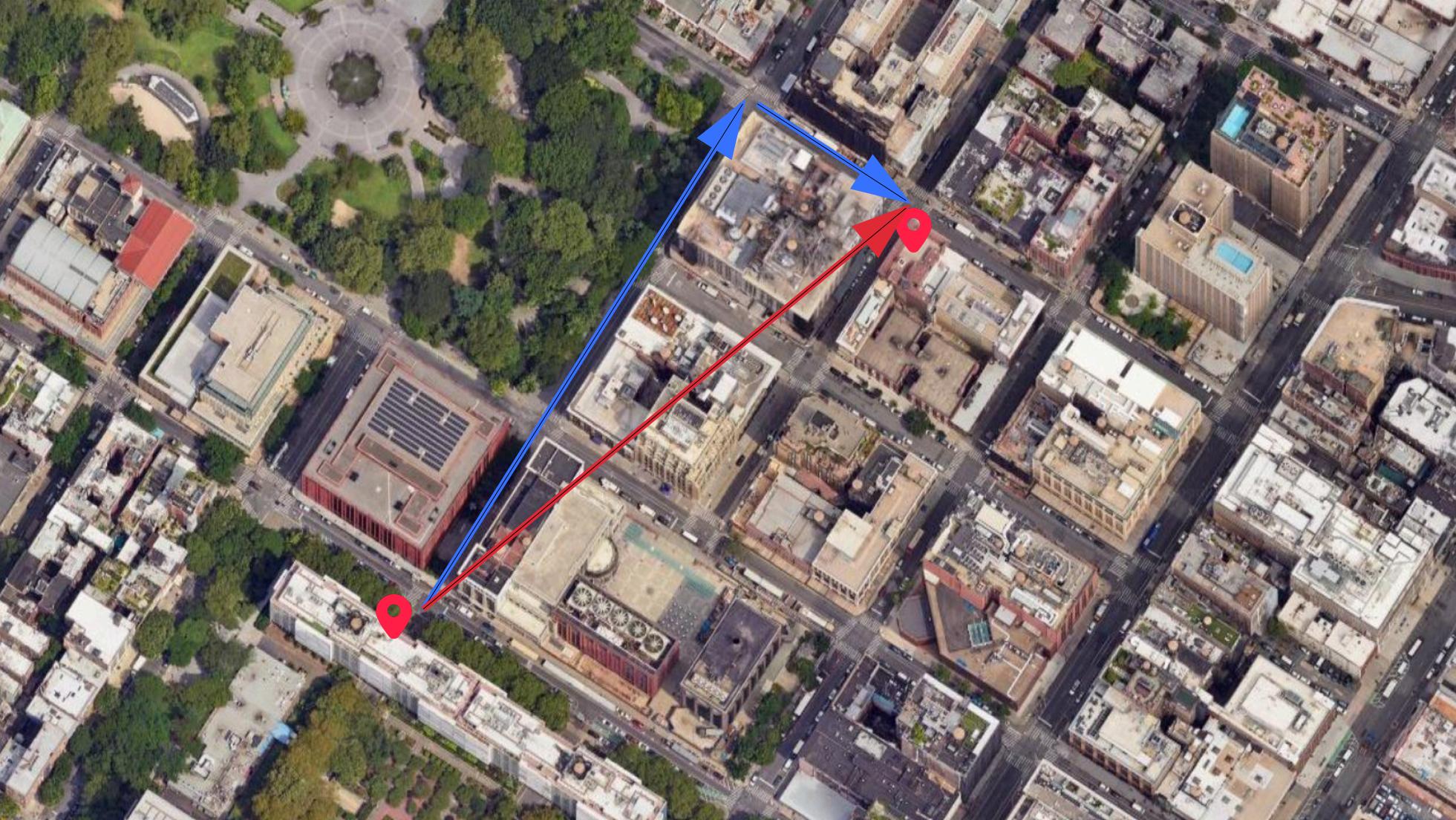
- This is different from building a model to predict known group membership and classifying new observations
 - This type of problem is called "supervised learning"
 - You can train a model to predict a known outcome and check how well it does
- Clustering is, on the other hand, an "unsupervised learning" problem
 - You don't know that there really are groups
 - If there are groups, you don't know how many groups there should be
 - You don't even know that the groups you find are the "right" groups
- This makes clustering really tricky task; you can't check your work!

Tools

Measuring Distances

Measuring Distances

- If you want to talk about how "close together" two things are, you need some notion of distance
- Distances are *dissimilarity* metrics
 - Larger magnitudes mean less similar (or farther apart) objects
 - Assign a value of 0 to identical objects
- Contrast these with *similarity* metrics
 - Larger magnitudes mean more similar (or closer together) objects
 - Assign maximal values (sometimes 1) to identical objects



Measuring Distances

- Most generically, we often think about *Euclidean distance*
 - Sometimes called "as the crow flies"
 - It's straight line distance

$$d(p, q) = \sqrt{(x_p - x_q)^2 + (y_p - y_q)^2}$$

$$d(p, q) = \sqrt{\sum_i (p_i - q_i)^2}$$

Measuring Distances

- Depending on the situation, something else might make more sense!
- Another spatial distance metric is *Manhattan distance*
- It's what you might expect from the name

$$d(p, q) = |x_p - x_q| + |y_p - y_q|$$

$$d(p, q) = \sum_i |p_i - q_i|$$

while Loops

while Loops

- Recall: loops allow us to specify a chunk of code to be executed repeatedly
 - `for` loops execute a pre-specified number of times
- What if you don't know how many times a loop will take, but you know when you want it to stop?
 - This is the job of the `while` loop
 - When possible, prefer `for` loops to `while` loops!

while Loops

- In R, a while loop has a few main components
 - The call: `while`
 - The condition, specified in `()` after the call and contain:
 - A logical statement that returns a single value
 - The code, wrapped in `{ }`:
 - At the beginning of the loop, the condition is checked
 - If it evaluates to `TRUE`, all of the code in the loop is executed
 - If it evaluates to `FALSE`, the code in the loop is skipped
 - After this is done, the condition is checked again to decide if the loop should be repeated
 - This continues until the condition evaluates to `FALSE`

Example `while` Loops 1 and 2

- What will be the result of executing these loops?

Example `while` Loops 1 and 2

- What will be the result of executing these loops?

```
1 i <- 0
2
3 while (i <= 5){
4     print(i)
5 }
6
7 # [1] 0
8 # [1] 0
9 # [1] 0
10 # [1] 0
11 # [1] 0
12
13 # this will not end...
14
15 while (i < 5){
16     print(i)
17 }
```

Example `while` Loops 3 and 4

- What will be the result of executing these loops?

```
1 i <- 0
2
3 while (i <= 5){
4   print(i)
5   i <- i + 1
6 }
7
8 # [1] 0
9 # [1] 1
10 # [1] 2
11 # [1] 3
12 # [1] 4
13 # [1] 5
14
15 while (i < 5){
16   print(i)
17   i <- i + 1
18 }
19
20 # nothing happens?
```

Example `while` Loop 4 (for real)

```
1 i <- 0
2
3 while (i < 5){
4   print(i)
5   i <- i + 1
6 }
7
8 # [1] 0
9 # [1] 1
10 # [1] 2
11 # [1] 3
12 # [1] 4
13
14 # all done!
```

Example while Loop 5

- You can check for convergence

```
1 last <- Inf
2 eps <- 1e-3
3 current <- 10
4
5 while (abs(last - current) > eps){
6   last <- current
7   current <- sqrt(current)
8   print(current)
9 }
10
11 # [1] 3.162278
12 # [1] 1.778279
13 # [1] 1.333521
14 # [1] 1.154782
15 # [1] 1.074608
16 # [1] 1.036633
17 # [1] 1.018152
18 # [1] 1.009035
19 # [1] 1.004507
20 # [1] 1.002251
21 # [1] 1.001125
22 # [1] 1.000562
```

Example while Loop 6

- You can combine both methods to set a maximum number of iterations while looking for convergence

```
1 last <- Inf
2 eps <- 1e-3
3 max_iter <- 3
4 current <- 10
5 i <- 0
6
7 while (abs(last - current) > eps && i < max_iter){
8   last <- current
9   current <- sqrt(current)
10  print(current)
11  i <- i + 1
12 }
13
14 # [1] 3.162278
15 # [1] 1.778279
16 # [1] 1.333521
```

Back to Clustering

The k -Means Algorithm

Iterative Clustering using k -Means

- We specify some number of clusters, k
- Clusters are described by *centroids*, or the center point
- We want to find locations of k centroids that, when we assign our observed points to them, minimize the *within-cluster sum of squares* (or variance)
- We call it k -Means because the location of each centroid is the mean of the coordinates of the points assigned to it!

Iterative Clustering using k -Means

- We follow a two step process:
 1. Assign each point to a cluster by finding the nearest centroid
 2. Move each centroid to the middle of the points assigned to it
- Do this over and over until *convergence*
 - An algorithm *converges* when additional iterations fail to improve the solution
 - For k -Means, convergence is when the centroids stop moving
 - What tool that we have used is a good fit for this task?
 - `while()` loops!

The k -Means Algorithm

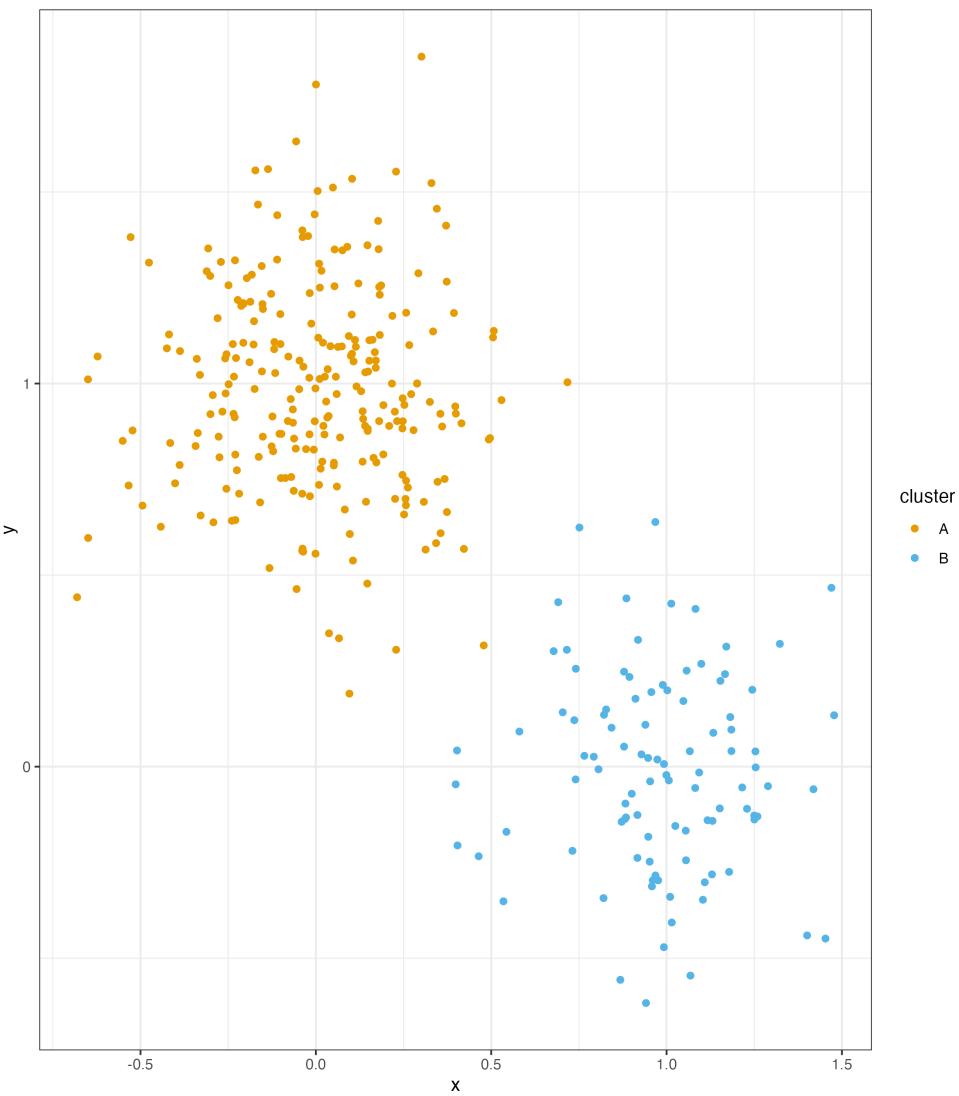
1. Select a number of clusters, k , and initialize the center of each cluster to a different point in space
2. Assign each observation from the dataset its nearest centroid using the Euclidean distance
3. Move the center of each cluster to the mean value of the coordinates of the points assigned to it
4. Repeat steps 2&3 until the centers stop moving

Generating Clustered Data

```
1 set.seed(8675309)
2 d <- data.frame(x = c(rnorm(250, 0, 0.25),
3                      rnorm(100, 1, 0.25)),
4                  y = c(rnorm(250, 1, 0.25),
5                      rnorm(100, 0, 0.25)),
6                  cluster = c(rep('A', 250),
7                             rep('B', 100)))
```

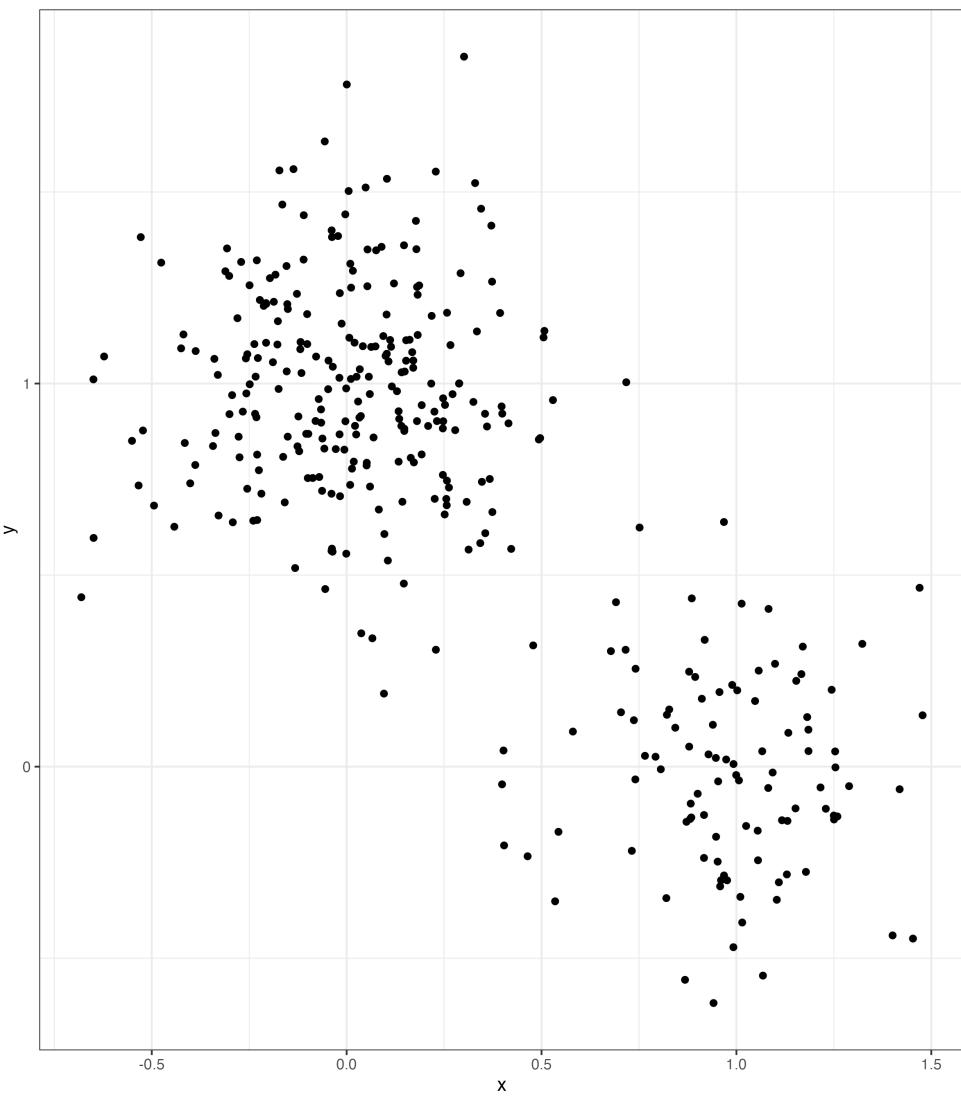
Visualizing Clustered Data

```
1 ggplot(d, aes(x=x, y=y, color=cluster)) +  
2   geom_point() +  
3   scale_color_okabeito() +  
4   theme_bw()
```



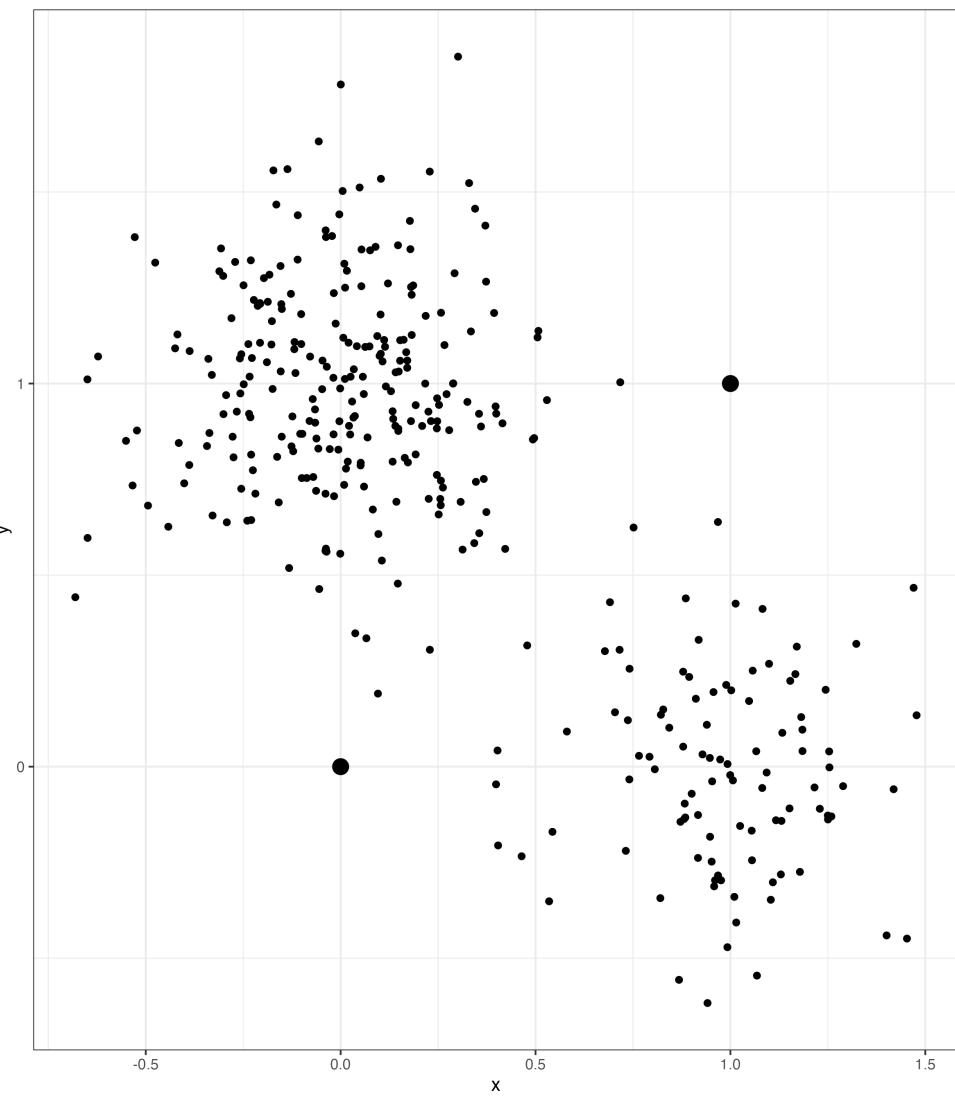
Visualizing Clustered Data

```
1 ggplot(d, aes(x=x, y=y)) +  
2   geom_point() +  
3   theme_bw()
```



Generate and Plot Starting Centroids

```
1 centroids <- data.frame(x = c(0,1),  
2                           y = c(0,1))  
3  
4  
5 ggplot(d, aes(x=x, y=y)) +  
6   geom_point() +  
7   geom_point(data=centroids, aes(x=x, y=y),  
8               color="black", size=4) +  
9   theme_bw()
```

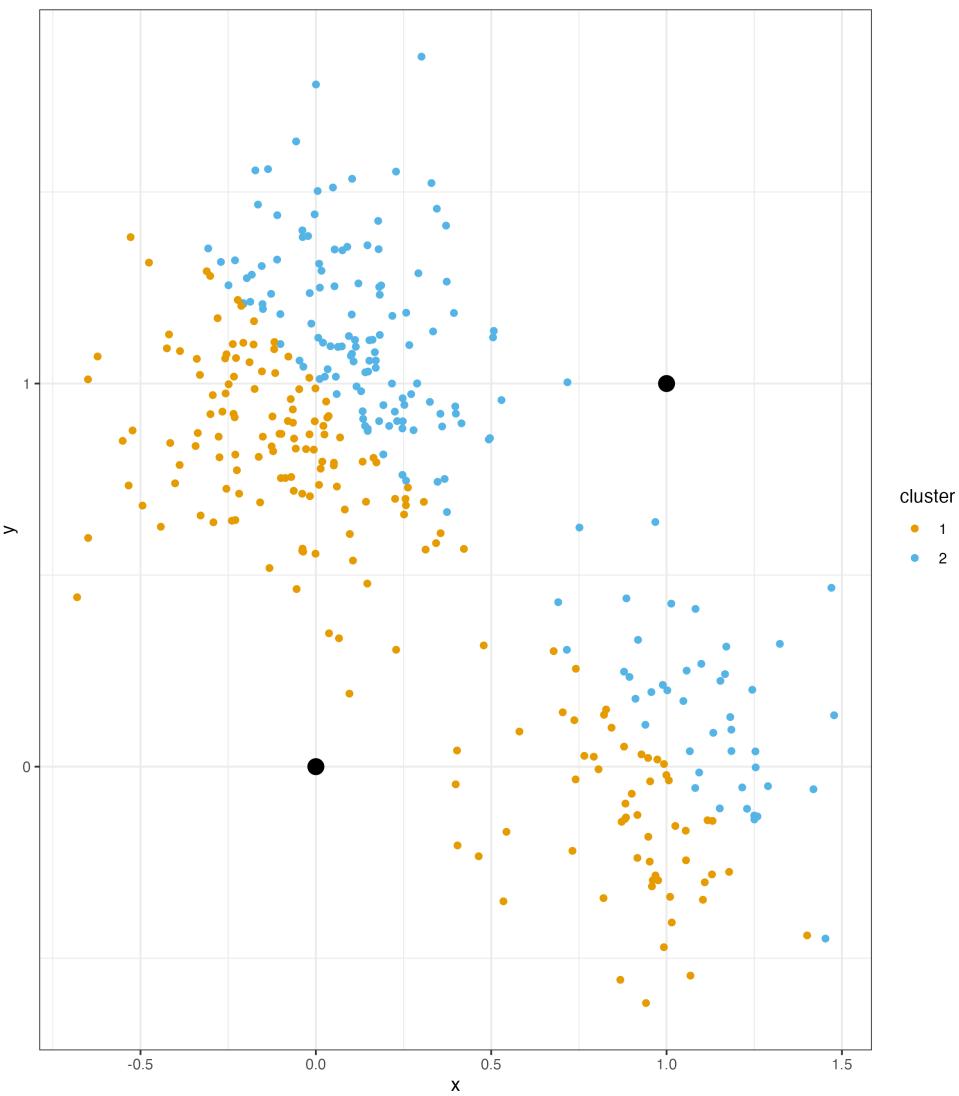


Compute Distances from Each Point to a Centroid

```
1 GetDistances <- function(data, centroid){  
2   # TODO: Compute distance from each point in data  
3   #   to a single centroid  
4  
5   distances <- sqrt((data[['x']] - centroid[['x']])^2 +  
6                     (data[['y']] - centroid[['y']])^2)  
7  
8   return(distances)  
9 }  
10  
11 GetDistances(d, centroids[2,])[1:5]  
12  
13 # [1] 1.2752070 0.8253962 1.1943523 0.5114553 0.7404843
```

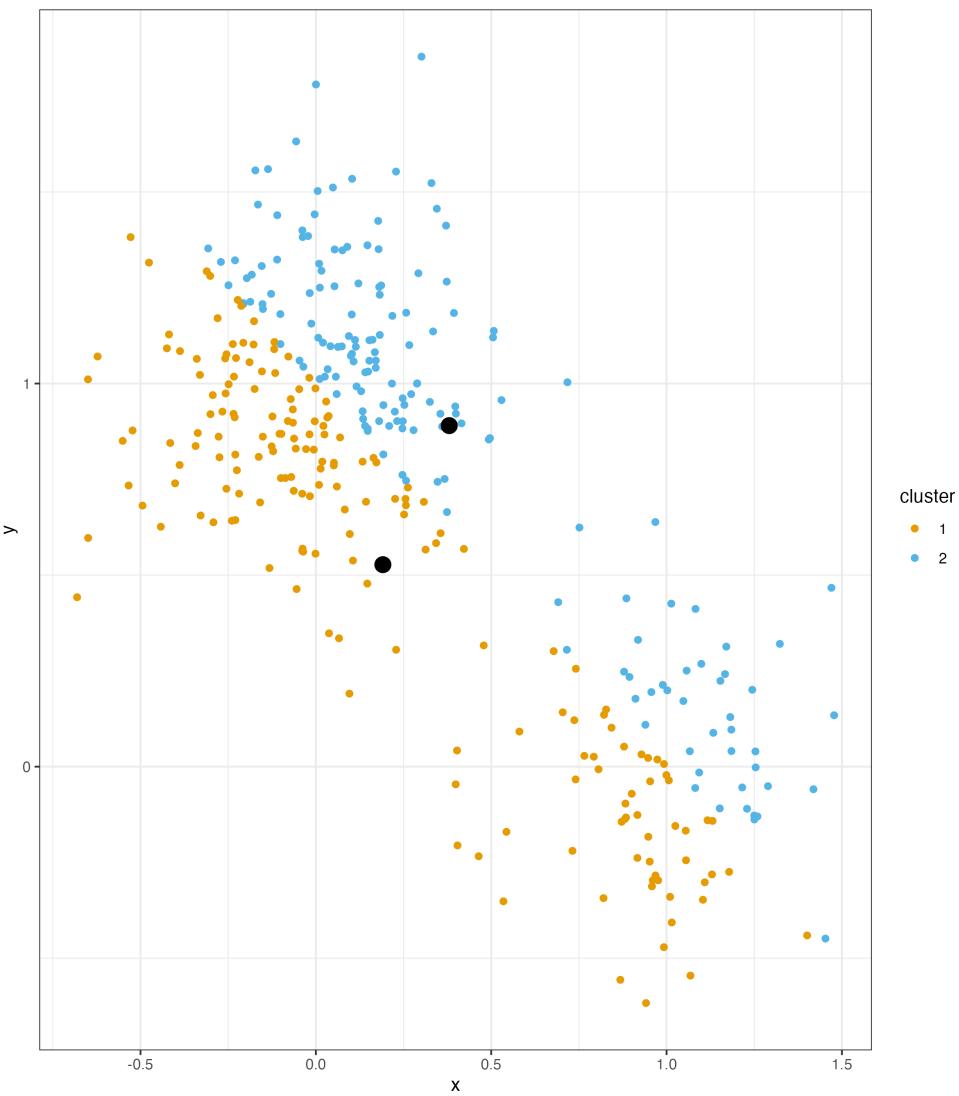
Assign Each Point to the Nearest Centroid

```
1 d$assignment <-  
2   ifelse(GetDistances(d, centroids[1,]) <=  
3     GetDistances(d, centroids[2,]),  
4     1, 2)  
5  
6 ggplot(d,  
7         aes(x=x,  
8                 y=y,  
9                 color=as.factor(assignment))) +  
10    geom_point() +  
11    geom_point(data=centroids,  
12                  aes(x=x, y=y),  
13                  color="black",  
14                  size=4) +  
15    labs(color='cluster') +  
16    scale_color_okabeito() +  
17    theme_bw()
```



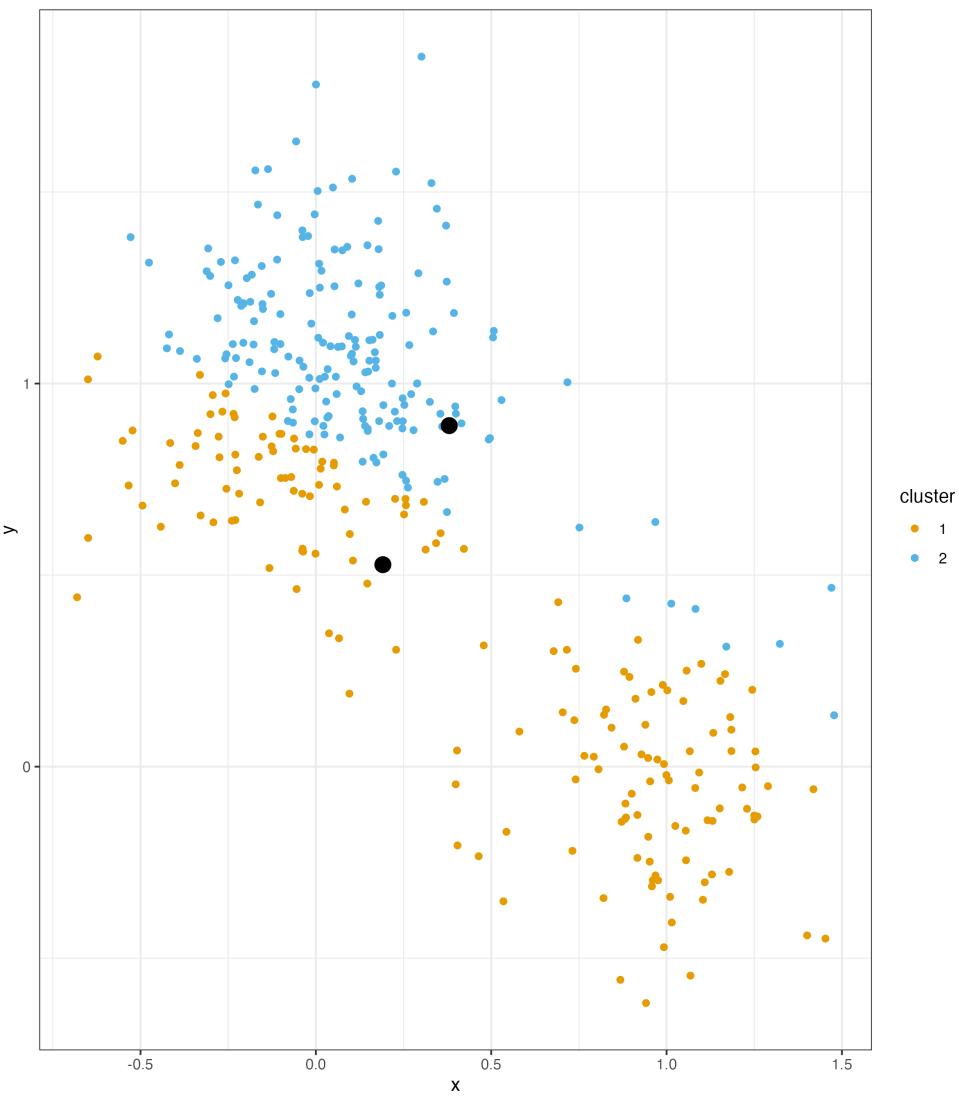
Move Centroids to the Mean of their Cluster

```
1  centroids <-  
2    data.frame(x = c(mean(d$x[d$assignment == 1]  
3                      mean(d$x[d$assignment == 2]  
4                      y = c(mean(d$y[d$assignment == 1]  
5                      mean(d$y[d$assignment == 2]  
6  
7  ggplot(d,  
8          aes(x=x,  
9                  y=y,  
10                 color=as.factor(assignment))) +  
11  geom_point() +  
12  geom_point(data=centroids,  
13                 aes(x=x, y=y),  
14                 color="black",  
15                 size=4) +  
16  labs(color='cluster') +  
17  scale_color_okabeito() +  
18  theme_bw()
```



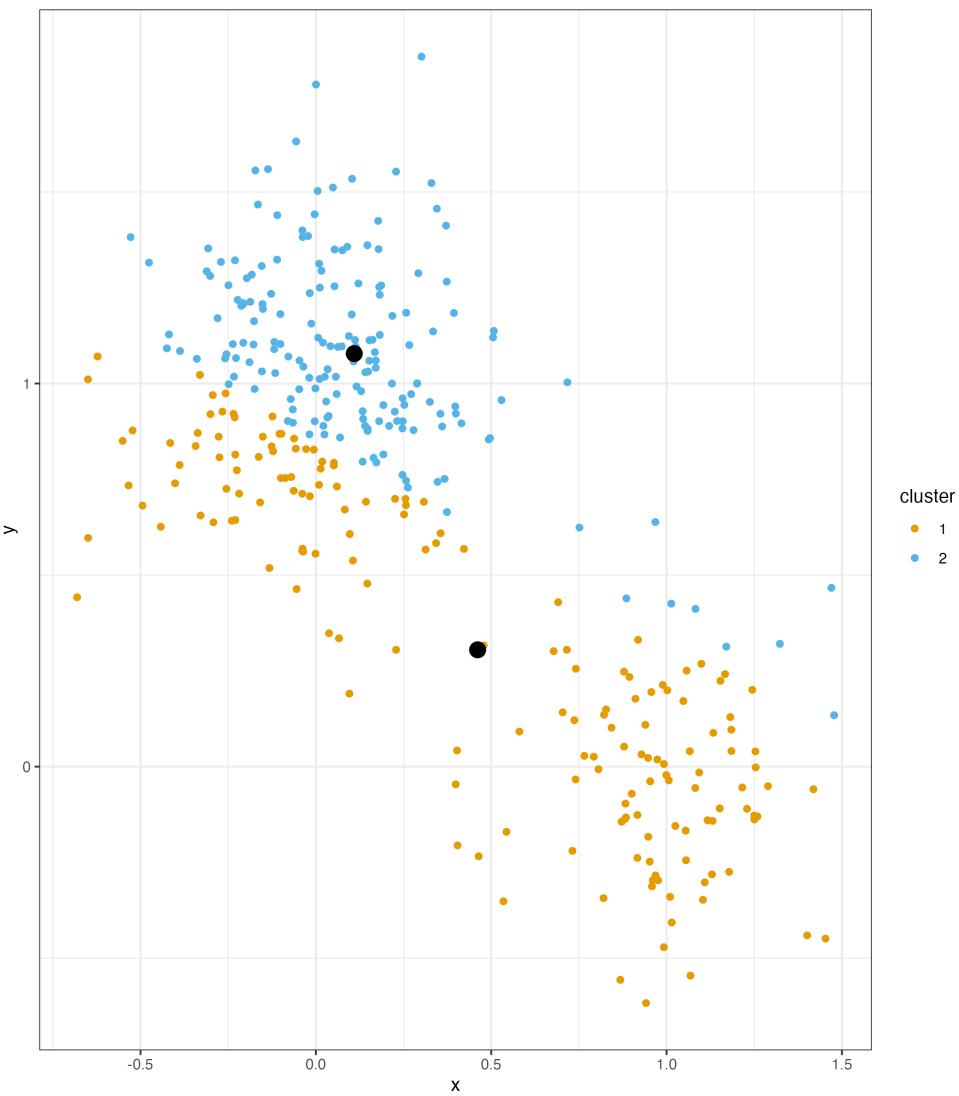
Second Iteration: Assign Points

```
1 d$assignment <-  
2   ifelse(GetDistances(d, centroids[1,]) <=  
3         GetDistances(d, centroids[2,]),  
4         1, 2)  
5  
6 ggplot(d,  
7         aes(x=x,  
8                 y=y,  
9                 color=as.factor(assignment))) +  
10    geom_point() +  
11    geom_point(data=centroids,  
12                  aes(x=x, y=y),  
13                  color="black",  
14                  size=4) +  
15    labs(color='cluster') +  
16    scale_color_okabeito() +  
17    theme_bw()
```



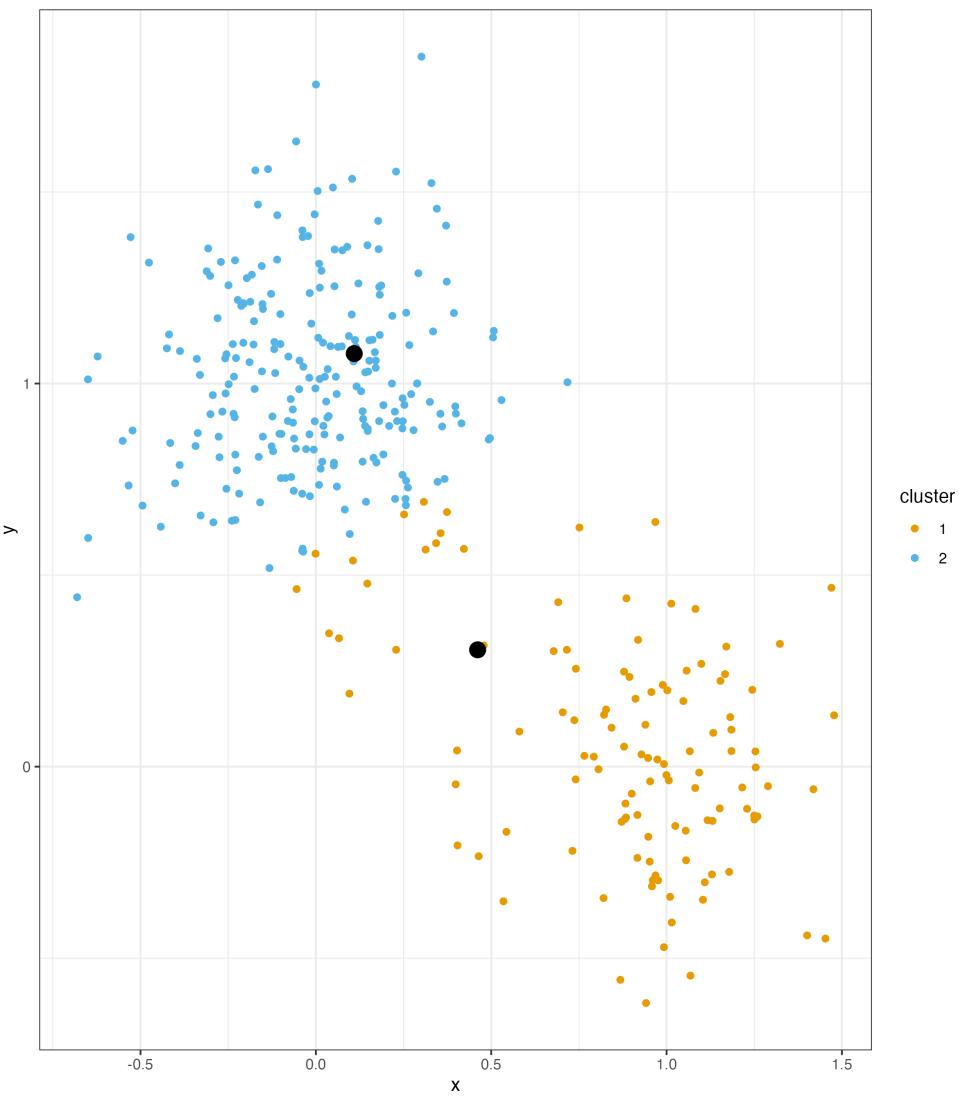
Second Iteration: Move Centroids

```
1 centroids <-  
2   data.frame(x = c(mean(d$x[d$assignment == 1]  
3                   mean(d$x[d$assignment == 2]  
4                   y = c(mean(d$y[d$assignment == 1]  
5                   mean(d$y[d$assignment == 2]  
6  
7 ggplot(d,  
8         aes(x=x,  
9                 y=y,  
10                color=as.factor(assignment))) +  
11   geom_point() +  
12   geom_point(data=centroids,  
13                 aes(x=x, y=y),  
14                 color="black",  
15                 size=4) +  
16   labs(color='cluster') +  
17   scale_color_okabeito() +  
18   theme_bw()
```



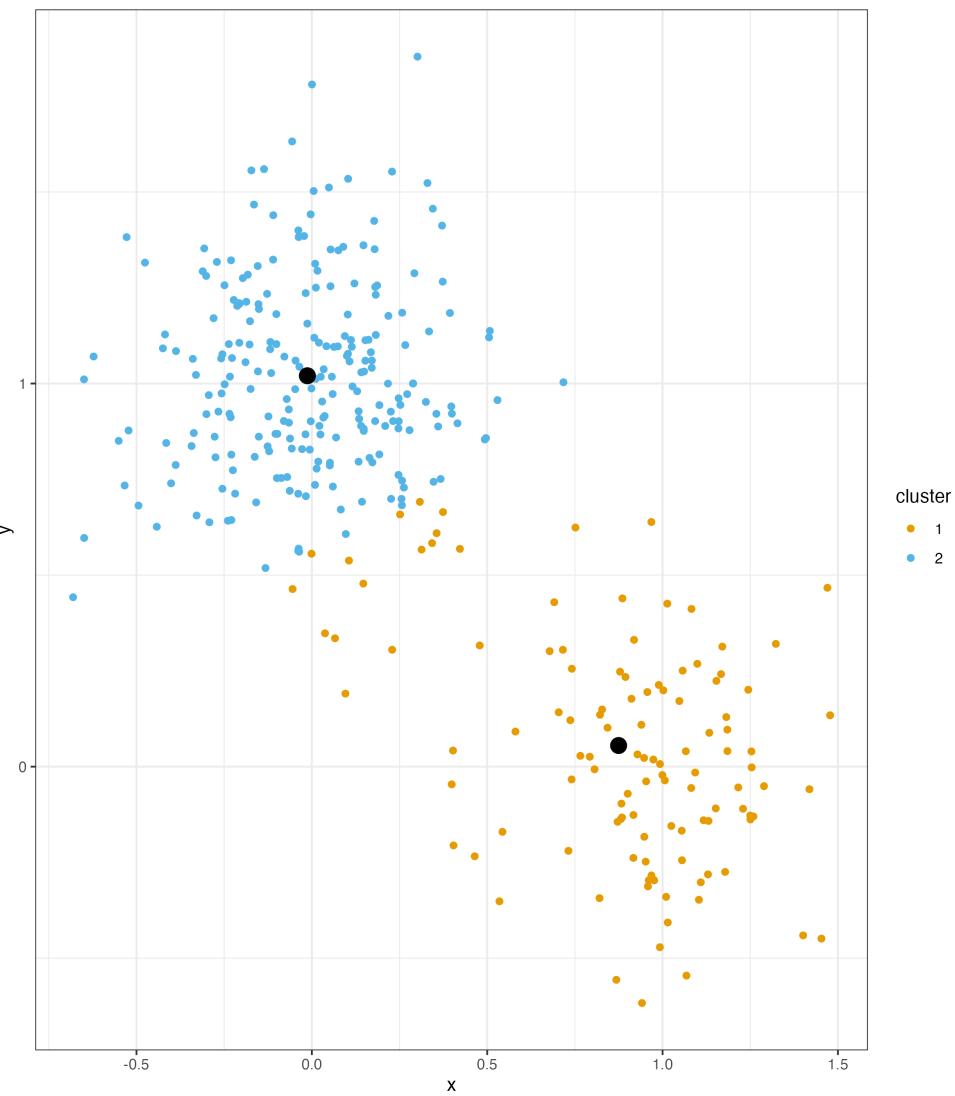
Third Iteration: Assign Points

```
1 d$assignment <-  
2   ifelse(GetDistances(d, centroids[1,]) <=  
3         GetDistances(d, centroids[2,]),  
4         1, 2)  
5  
6 ggplot(d,  
7         aes(x=x,  
8                 y=y,  
9                 color=as.factor(assignment))) +  
10    geom_point() +  
11    geom_point(data=centroids,  
12                  aes(x=x, y=y),  
13                  color="black",  
14                  size=4) +  
15    labs(color='cluster') +  
16    scale_color_okabeito() +  
17    theme_bw()
```



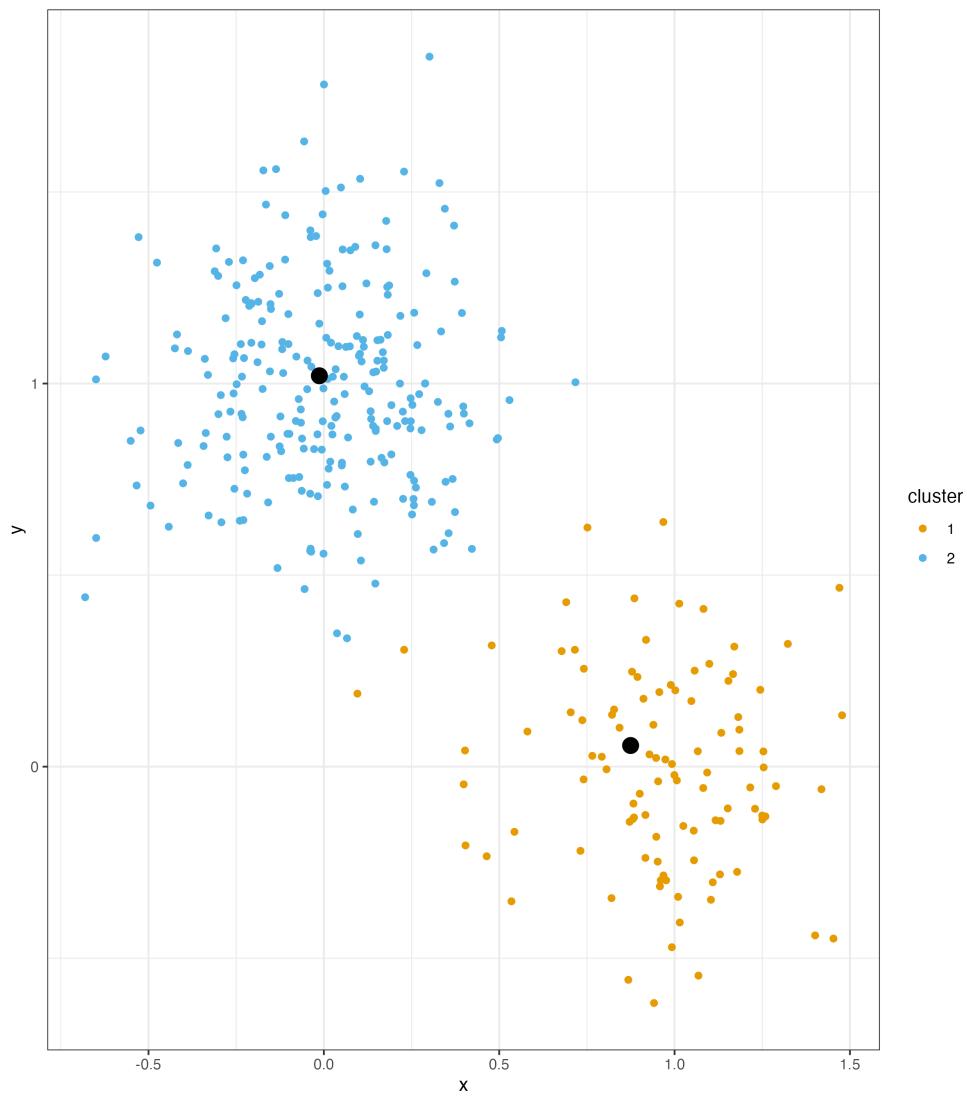
Third Iteration: Move Centroids

```
1 centroids <-  
2   data.frame(x = c(mean(d$x[d$assignment == 1]  
3                   mean(d$x[d$assignment == 2]  
4                   y = c(mean(d$y[d$assignment == 1]  
5                   mean(d$y[d$assignment == 2]  
6  
7 ggplot(d,  
8         aes(x=x,  
9                 y=y,  
10                color=as.factor(assignment))) +  
11   geom_point() +  
12   geom_point(data=centroids,  
13                 aes(x=x, y=y),  
14                 color="black",  
15                 size=4) +  
16   labs(color='cluster') +  
17   scale_color_okabeito() +  
18   theme_bw()
```



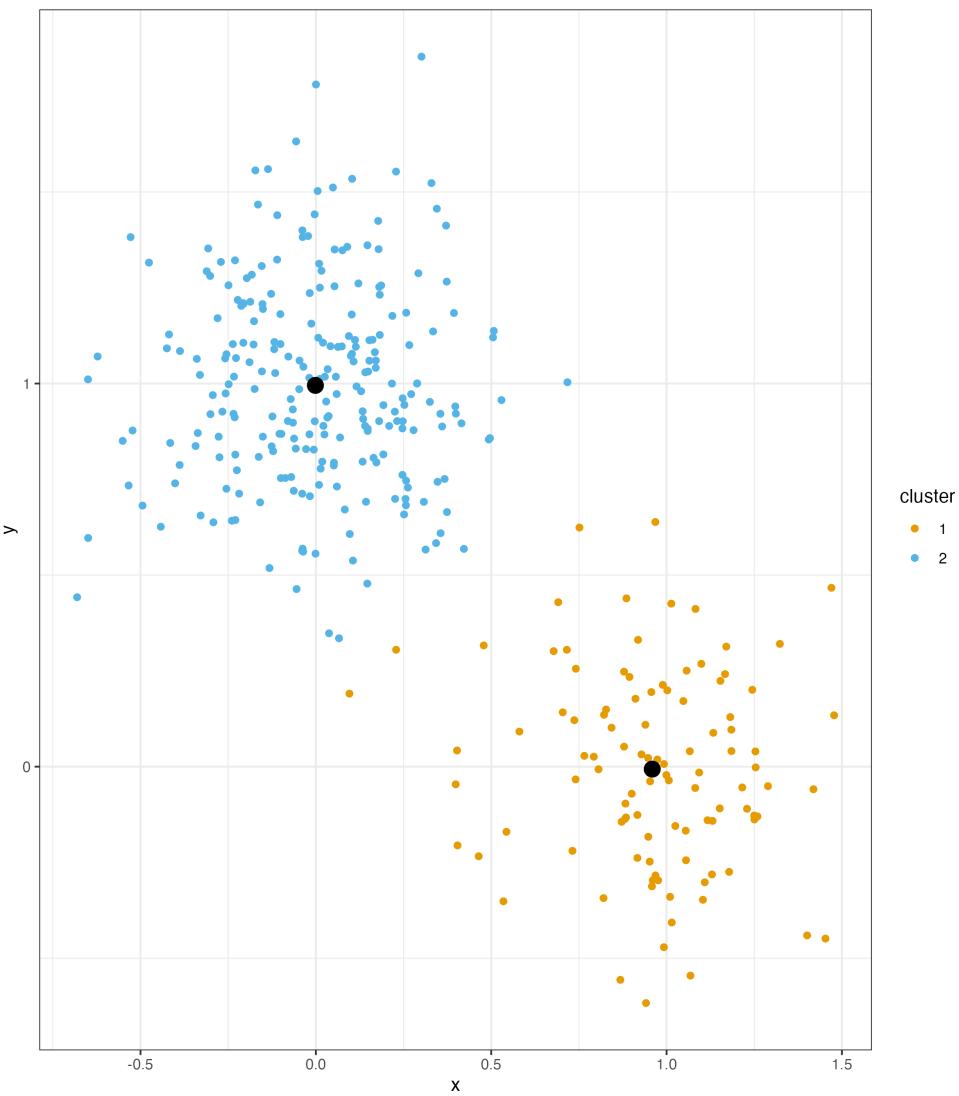
Fourth Iteration: Assign Points

```
1 d$assignment <-  
2   ifelse(GetDistances(d, centroids[1,]) <=  
3         GetDistances(d, centroids[2,]),  
4         1, 2)  
5  
6 ggplot(d,  
7         aes(x=x,  
8                 y=y,  
9                 color=as.factor(assignment))) +  
10    geom_point() +  
11    geom_point(data=centroids,  
12                  aes(x=x, y=y),  
13                  color="black",  
14                  size=4) +  
15    labs(color='cluster') +  
16    scale_color_okabeito() +  
17    theme_bw()
```



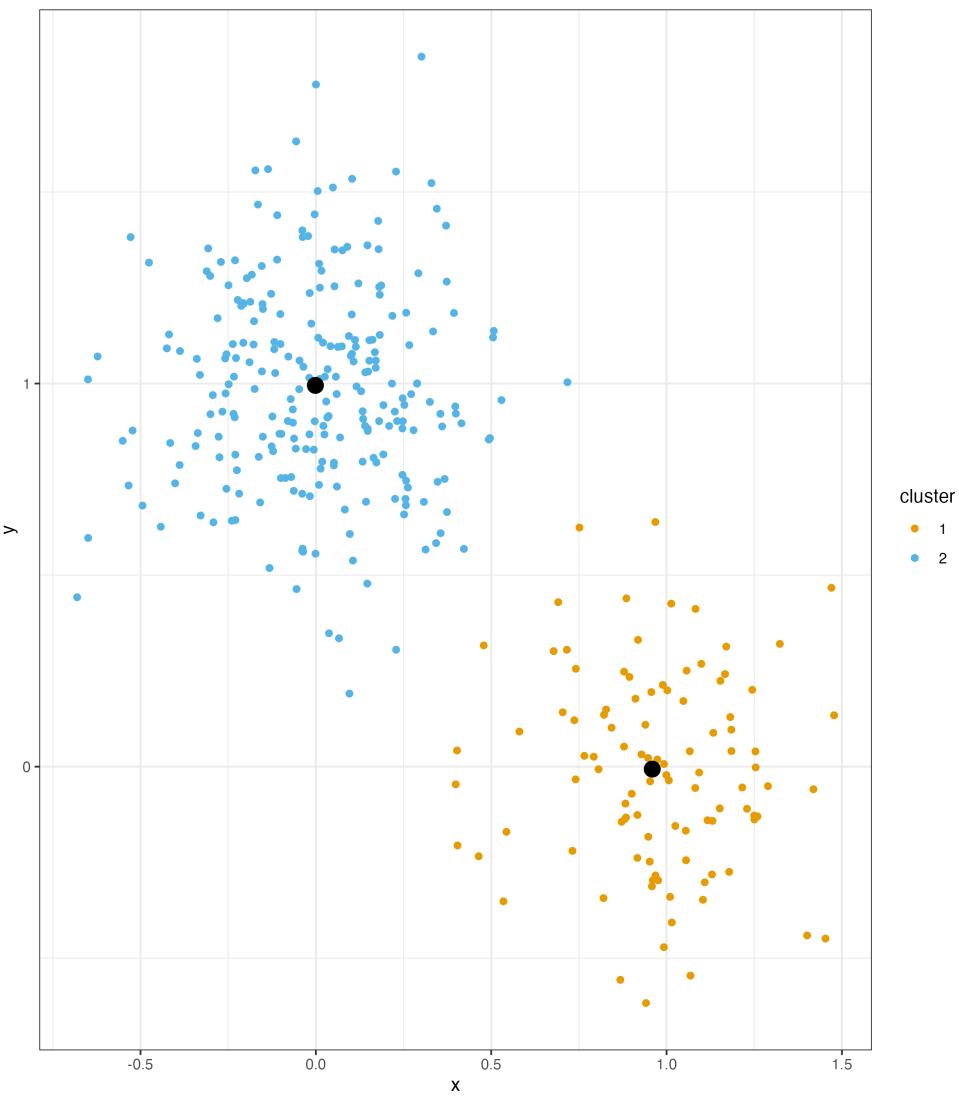
Fourth Iteration: Move Centroids

```
1 centroids <-  
2   data.frame(x = c(mean(d$x[d$assignment == 1]  
3                   mean(d$x[d$assignment == 2]  
4                   y = c(mean(d$y[d$assignment == 1]  
5                   mean(d$y[d$assignment == 2]  
6  
7 ggplot(d,  
8         aes(x=x,  
9                 y=y,  
10                color=as.factor(assignment))) +  
11   geom_point() +  
12   geom_point(data=centroids,  
13                 aes(x=x, y=y),  
14                 color="black",  
15                 size=4) +  
16   labs(color='cluster') +  
17   scale_color_okabeito() +  
18   theme_bw()
```



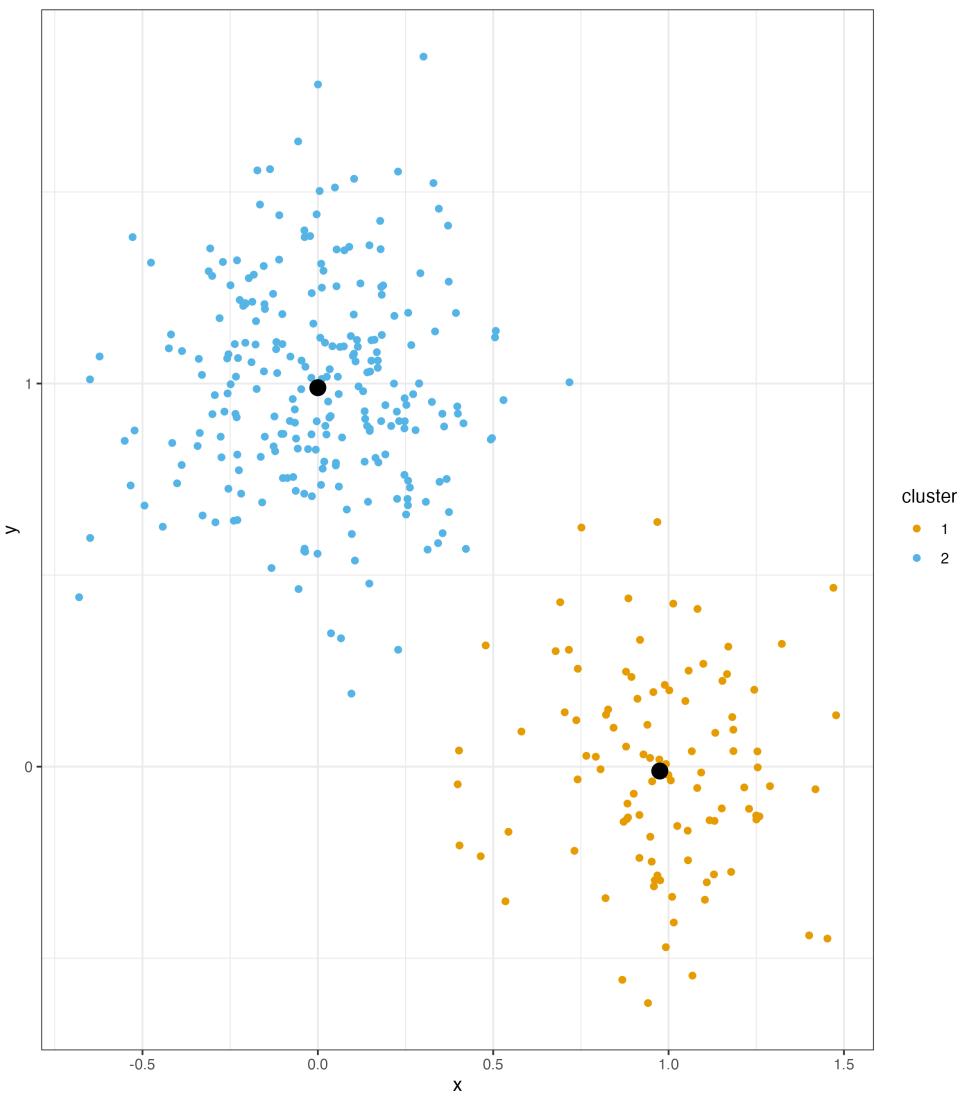
Fifth Iteration: Assign Points

```
1 d$assignment <-  
2   ifelse(GetDistances(d, centroids[1,]) <=  
3         GetDistances(d, centroids[2,]),  
4         1, 2)  
5  
6 ggplot(d,  
7         aes(x=x,  
8                 y=y,  
9                 color=as.factor(assignment))) +  
10    geom_point() +  
11    geom_point(data=centroids,  
12                  aes(x=x, y=y),  
13                  color="black",  
14                  size=4) +  
15    labs(color='cluster') +  
16    scale_color_okabeito() +  
17    theme_bw()
```



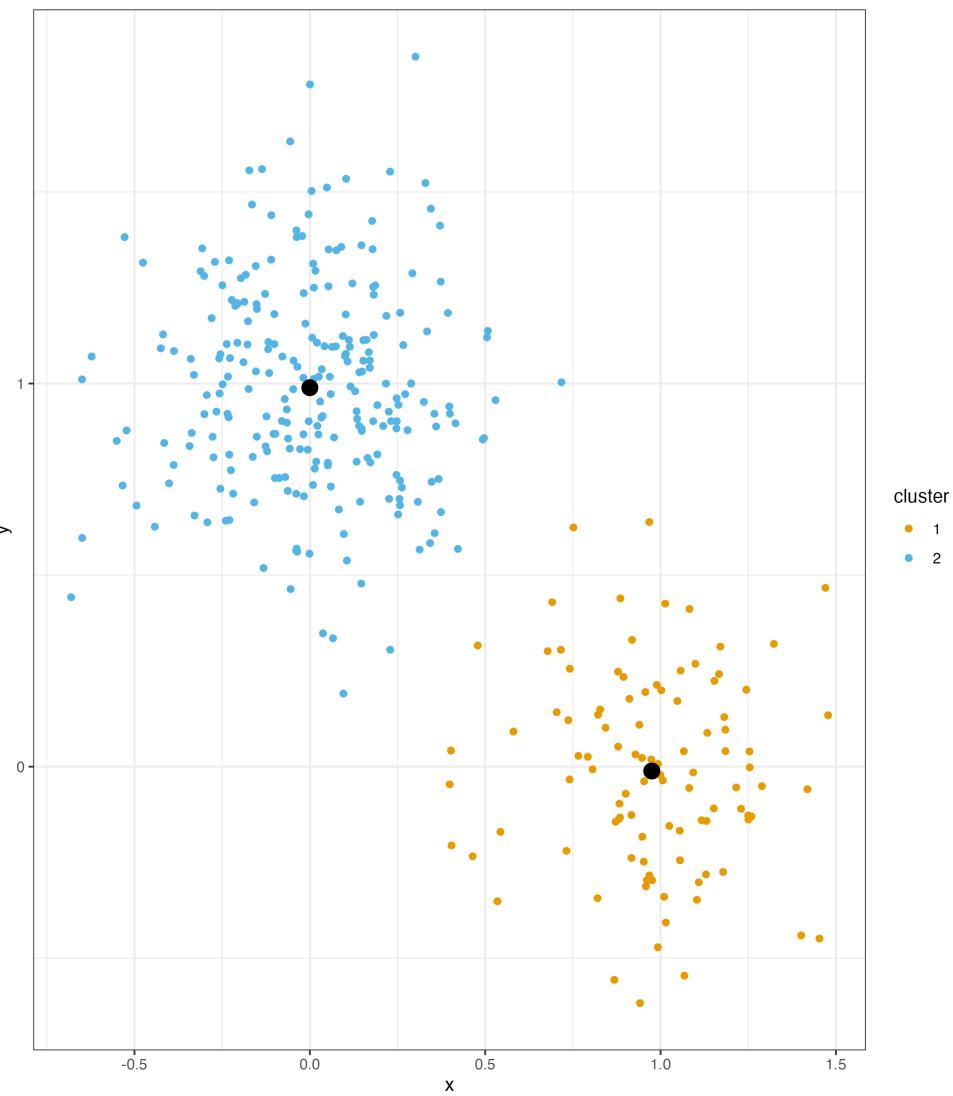
Fifth Iteration: Move Centroids

```
1  centroids <-  
2    data.frame(x = c(mean(d$x[d$assignment == 1]  
3                      mean(d$x[d$assignment == 2]  
4                      y = c(mean(d$y[d$assignment == 1]  
5                      mean(d$y[d$assignment == 2]  
6  
7  ggplot(d,  
8          aes(x=x,  
9                  y=y,  
10                 color=as.factor(assignment))) +  
11  geom_point() +  
12  geom_point(data=centroids,  
13                 aes(x=x, y=y),  
14                 color="black",  
15                 size=4) +  
16  labs(color='cluster') +  
17  scale_color_okabeito() +  
18  theme_bw()
```



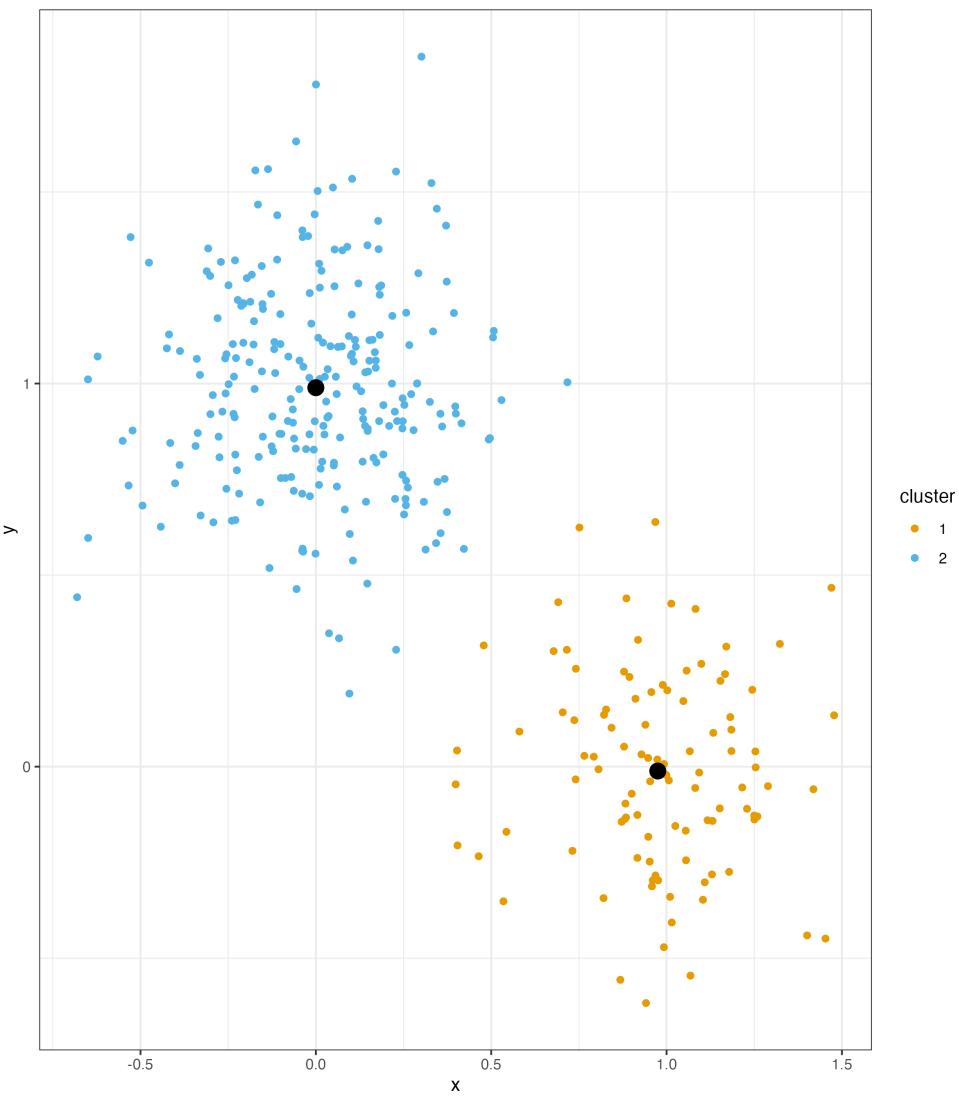
Sixth Iteration: Assign Points

```
1 d$assignment <-  
2   ifelse(GetDistances(d, centroids[1,]) <=  
3         GetDistances(d, centroids[2,]),  
4         1, 2)  
5  
6 ggplot(d,  
7         aes(x=x,  
8                 y=y,  
9                 color=as.factor(assignment))) +  
10    geom_point() +  
11    geom_point(data=centroids,  
12                  aes(x=x, y=y),  
13                  color="black",  
14                  size=4) +  
15    labs(color='cluster') +  
16    scale_color_okabeito() +  
17    theme_bw()
```



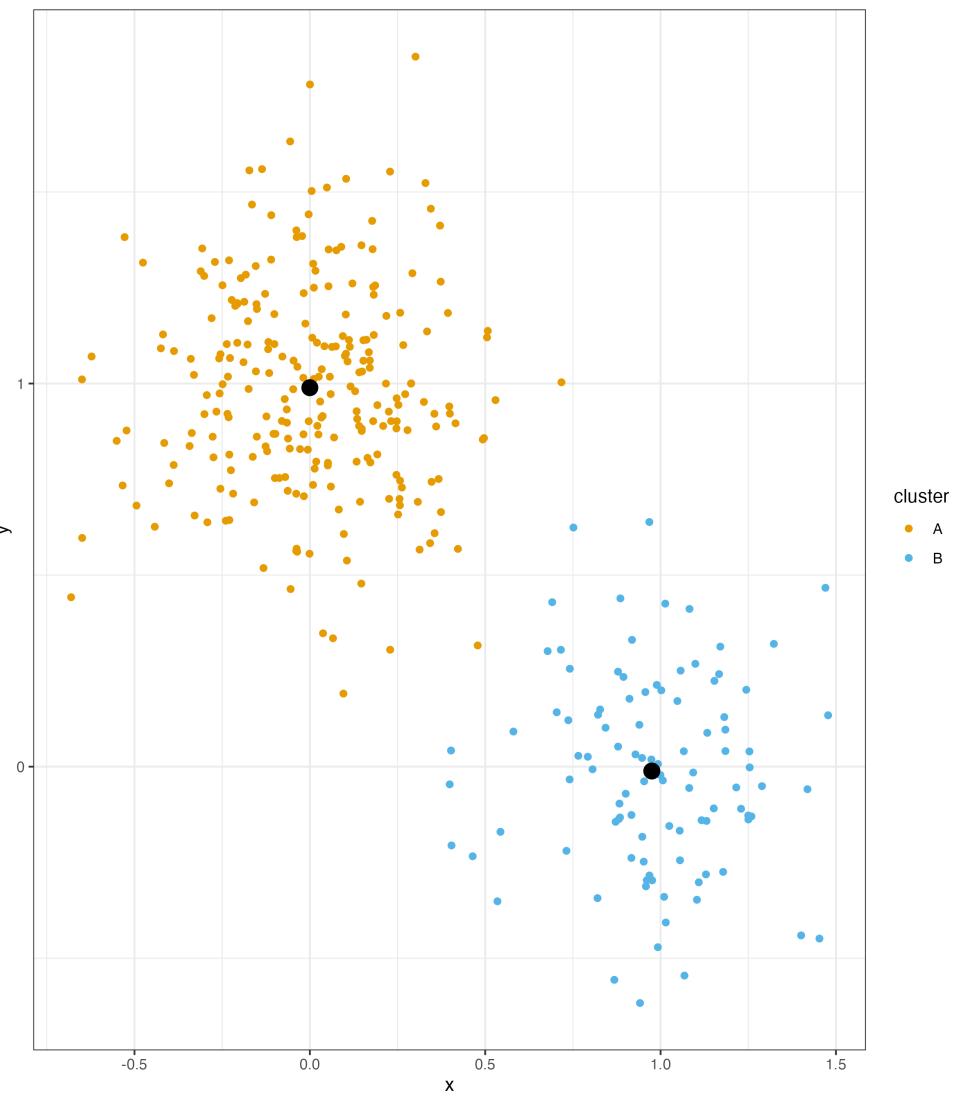
Sixth Iteration: Move Centroids

```
1  centroids <-  
2    data.frame(x = c(mean(d$x[d$assignment == 1]  
3                      mean(d$x[d$assignment == 2]  
4                      y = c(mean(d$y[d$assignment == 1]  
5                      mean(d$y[d$assignment == 2]  
6  
7  ggplot(d,  
8          aes(x=x,  
9                  y=y,  
10                 color=as.factor(assignment))) +  
11  geom_point() +  
12  geom_point(data=centroids,  
13                 aes(x=x, y=y),  
14                 color="black",  
15                 size=4) +  
16  labs(color='cluster') +  
17  scale_color_okabeito() +  
18  theme_bw()
```



"True" Assignments

```
1 ggplot(d, aes(x=x, y=y, color=cluster)) +  
2   geom_point() +  
3   geom_point(data=centroids,  
4               aes(x=x, y=y),  
5               color="black",  
6               size=4) +  
7   labs(color='cluster') +  
8   scale_color_okabeito() +  
9   theme_bw()
```



Generalizing our k -Means Implementation

What Do We Need to Do?

1. Initialize Centroids:

- *Output*: Dataframe of k centroids
- *Inputs*: The dataframe we are clustering and k

2. Get Cluster Assignments:

- *Output*: Vector of cluster assignments
- *Inputs*: Data, centroids, a distance function

3. Recalculate Centroids:

- *Output*: A dataframe of new centroid locations
- *Inputs*: Data, centroids, and cluster assignments

4. Evaluation Function:

- *Output*: Within-cluster sum of squares
- *Inputs*: Data, centroids, and cluster assignments

Initialize Centroids

```
1 InitializeCentroids <- function(data, k){  
2     # TODO: Return a dataframe of k starting  
3     # centroids, one centroid per row  
4     c <- sample(1:nrow(data), k)  
5     centroids <- data[c, ]  
6  
7     return(centroids)  
8 }
```

Get Cluster Assignments

```
1 GetAssignments <- function(data, centroids, distfun){  
2   # TODO: Return a vector of numeric cluster  
3   # assignments by matching each point with  
4   # its nearest centroid  
5  
6   dist <- matrix(rep(NA, nrow(data)*nrow(centroids)),  
7                   ncol=nrow(centroids))  
8  
9   for (i in 1:nrow(centroids)){  
10     dist[,i] <- distfun(data, centroids[i,])  
11   }  
12   assignment <- max.col(-dist)  
13  
14   return(assignment)  
15 }
```

Recalculate Centroids

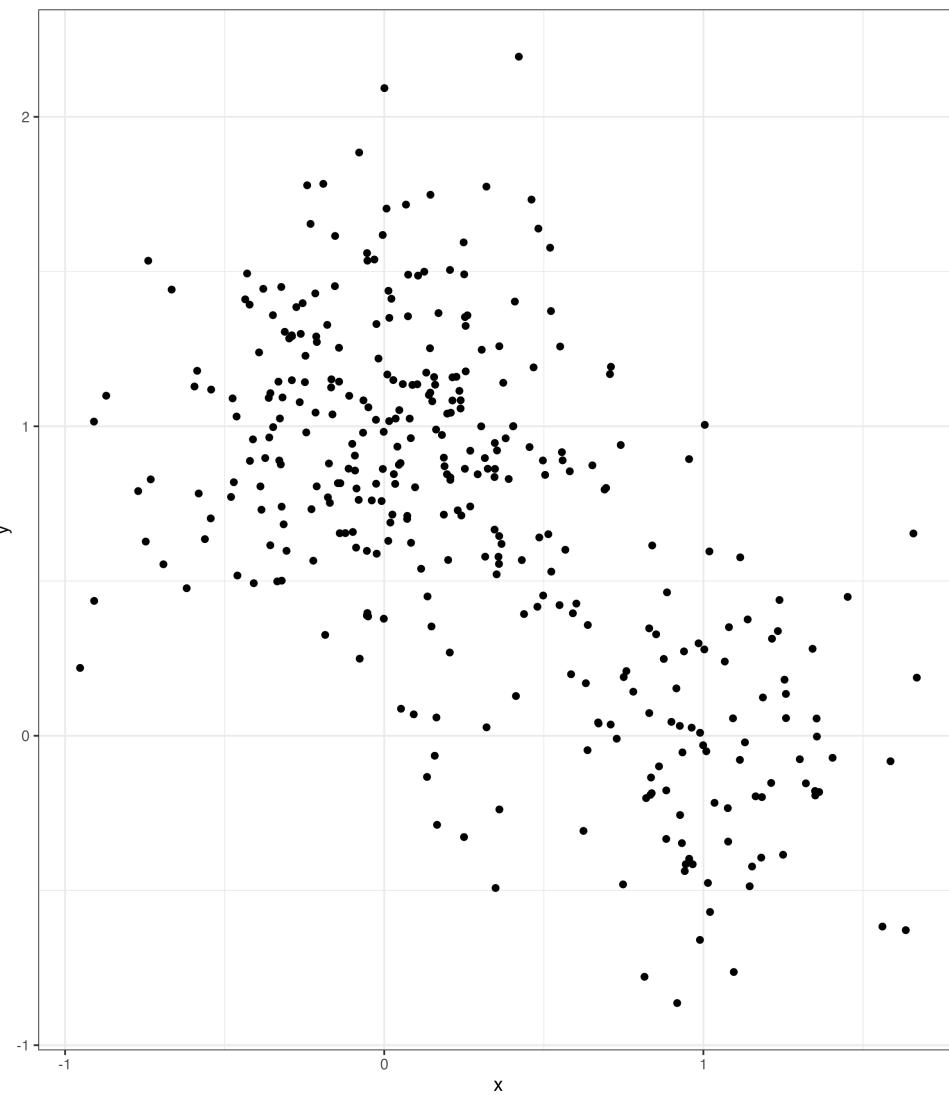
```
1 RecalculateCentroids <- function(data, centroids, assignment){  
2   # TODO: Return a dataframe of new centroids  
3   # calculated from the means of the points  
4   # in each cluster  
5  
6   for (i in 1:nrow(centroids)){  
7     centroids[i,] <- colMeans(data[assignment==i,])  
8   }  
9  
10  return(centroids)  
11 }
```

Evaluation Function

```
1 WithinClusterSumSquares <- function(data, centroids, assignment){  
2   # TODO: Return the total within-cluster sum of squares  
3  
4   WCSS <- 0  
5   for (i in 1:nrow(centroids)){  
6     WCSS <- WCSS + sum(GetDistances(data[assignment==i,], centroids[i,])^2)  
7   }  
8  
9   return(WCSS)  
10 }
```

Some New Data

```
1 set.seed(8675309)
2
3 # New data, same solution but larger variance
4 d <- data.frame(x = c(rnorm(250, 0, 0.35),
5                      rnorm(100, 1, 0.35)),
6                   y = c(rnorm(250, 1, 0.35),
7                      rnorm(100, 0, 0.35)))
8
9 # shuffle to destroy any "known" solution
10 d <- d[sample(1:nrow(d)),]
11
12 # looks worse!
13 ggplot(d, aes(x=x, y=y)) +
14   geom_point() +
15   theme_bw()
```

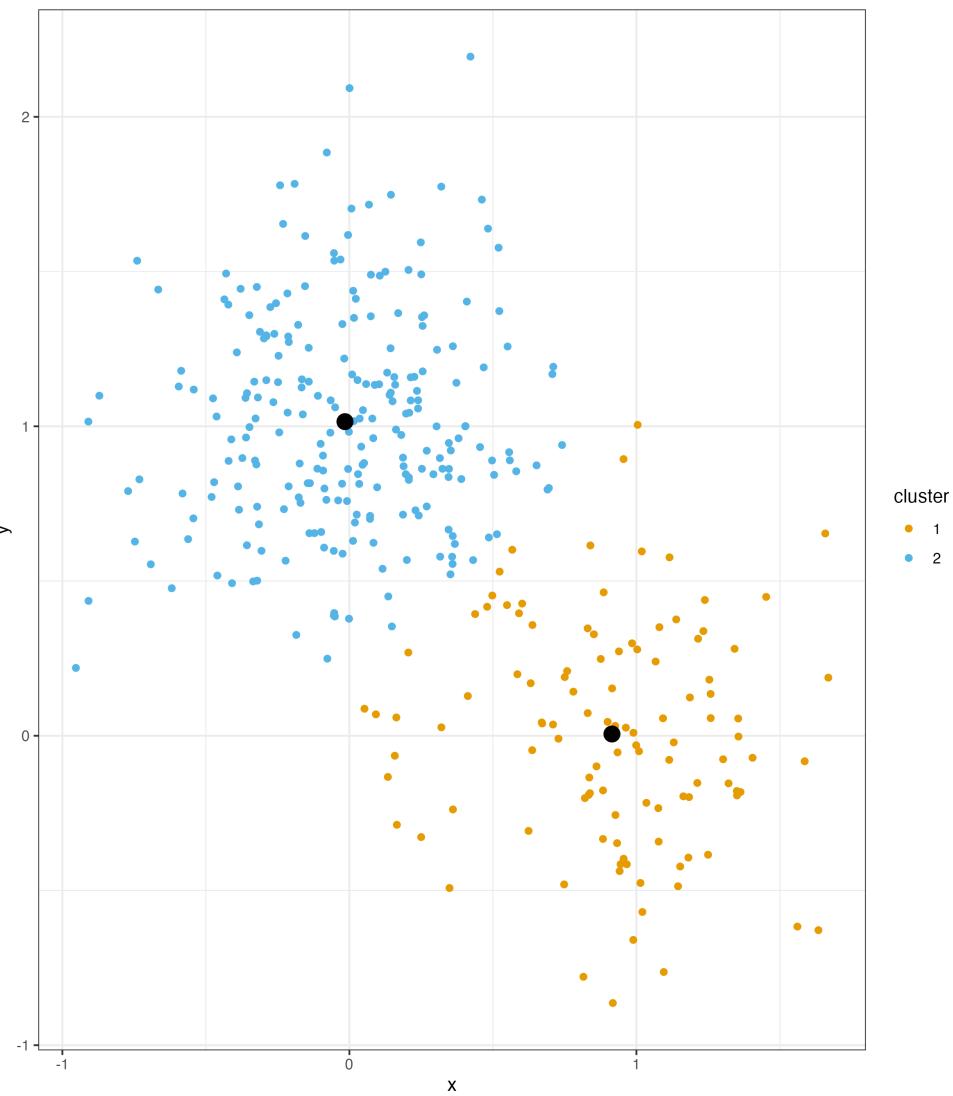


Putting it All Together

```
1 # TODO: initialize starting values
2 k <- 2
3 centroids <- InitializeCentroids(d, k)
4 old_centroids <- InitializeCentroids(d, k)
5
6 # TODO: run until convergence
7 while(!identical(old_centroids, centroids)){
8   old_centroids <- centroids
9   assignment <- GetAssignments(d, centroids, GetDistances)
10  centroids <- RecalculateCentroids(d, centroids, assignment)
11 }
12
13 # TODO: get WCSS for cluster solution
14 WithinClusterSumSquares(d, centroids, assignment)
```

Visualize the Solution: $k = 2$

```
1 ggplot(d, aes(x=x,
2                 y=y,
3                 color=as.factor(assignment))) +
4   geom_point() +
5   geom_point(data=centroids, aes(x=x, y=y),
6               color="black", size=4) +
7   labs(color='cluster') +
8   scale_color_okabeito() +
9   theme_bw()
```



What About $k = 3$?

```
1 # TODO: initialize starting values
2 k <- 3
3 centroids <- InitializeCentroids(d, k)
4 old_centroids <- InitializeCentroids(d, k)
5
6 # TODO: run until convergence
7 while(!identical(old_centroids, centroids)){
8   old_centroids <- centroids
9   assignment <- GetAssignments(d, centroids, GetDistances)
10  centroids <- RecalculateCentroids(d, centroids, assignment)
11 }
12
13 # TODO: get WCSS for cluster solution
14 WithinClusterSumSquares(d, centroids, assignment)
```

Visualize the Solution: $k = 3$

```
1 ggplot(d, aes(x=x,
2                 y=y,
3                 color=as.factor(assignment))) +
4   geom_point() +
5   geom_point(data=centroids, aes(x=x, y=y),
6               color="black", size=4) +
7   labs(color='cluster') +
8   scale_color_okabeito() +
9   theme_bw()
```



What About $k = 4$?

```
1 # TODO: initialize starting values
2 k <- 4
3 centroids <- InitializeCentroids(d, k)
4 old_centroids <- InitializeCentroids(d, k)
5
6 # TODO: run until convergence
7 while(!identical(old_centroids, centroids)){
8   old_centroids <- centroids
9   assignment <- GetAssignments(d, centroids, GetDistances)
10  centroids <- RecalculateCentroids(d, centroids, assignment)
11 }
12
13 # TODO: get WCSS for cluster solution
14 WithinClusterSumSquares(d, centroids, assignment)
```

Visualize the Solution: $k = 4$

```
1 ggplot(d, aes(x=x,
2                 y=y,
3                 color=as.factor(assignment))) +
4   geom_point() +
5   geom_point(data=centroids, aes(x=x, y=y),
6               color="black", size=4) +
7   labs(color='cluster') +
8   scale_color_okabeito() +
9   theme_bw()
```

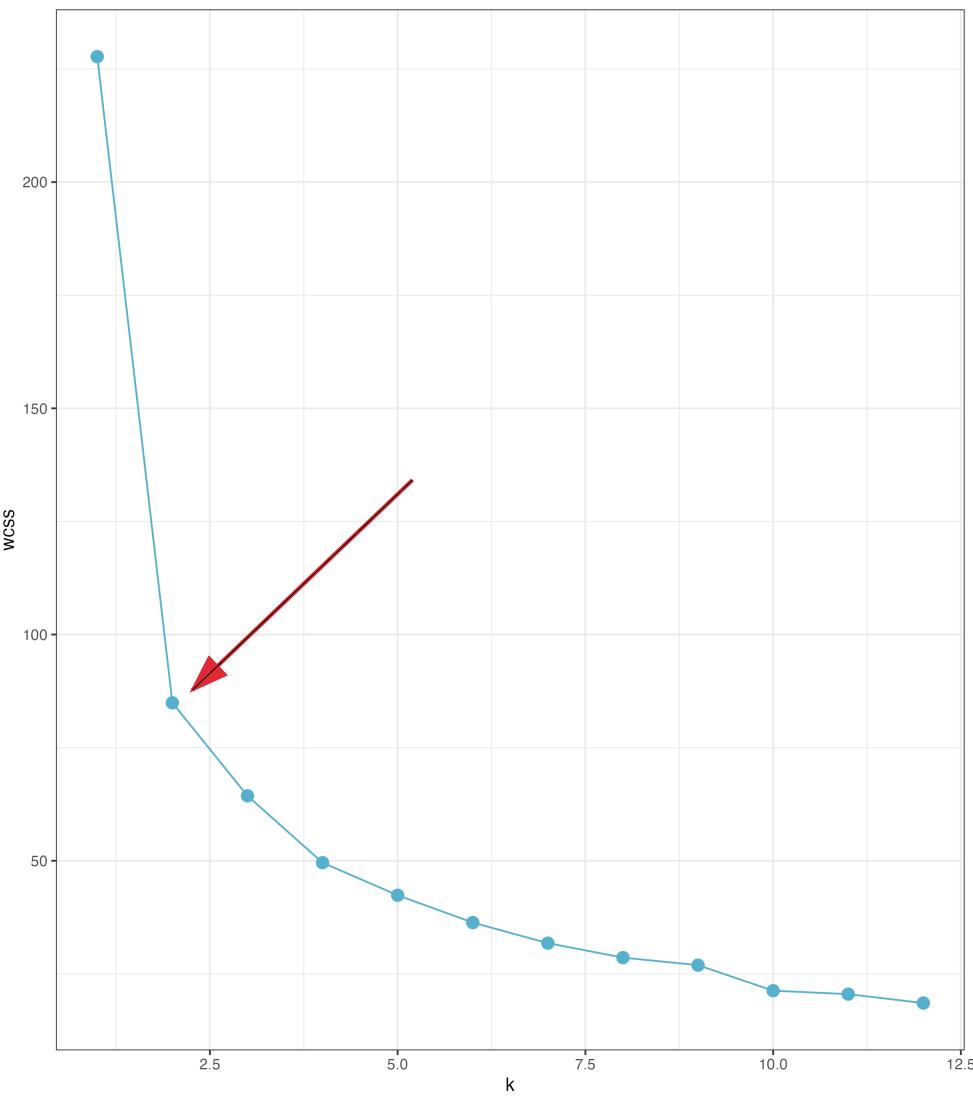


How Do You Decide What the Right Answer is?

- You'll never know the right answer!
- Look at some evaluation metric
 - Here we use within-cluster variance
 - Other ones exist (like silhouette scores)
- The problem is that adding an additional cluster always makes the evaluation metric go down
- Ask "when does adding another cluster stop making a big difference in my evaluation metric?"
 - This is a judgement call!
 - Often we look at "scree plot" and try and identify the "elbow"
 - Plot WCSS vs. k
 - Pick the "elbow," or the point where adding another cluster doesn't give much improvement

The Elbow Plot

```
1  ks <- 1:12
2  wcss <- vector('numeric', length(ks))
3  for (k in ks){
4    centroids <- InitializeCentroids(d, k)
5    old_centroids <- InitializeCentroids(d, k)
6    while(!identical(old_centroids, centroids)){
7      old_centroids <- centroids
8      assignment <- GetAssignments(d,
9                                centroids,
10                               GetDistances)
11      centroids <- RecalculateCentroids(d,
12                                            centroid
13                                            assignme
14  }
15  wcss[k] <- WithinClusterSumSquares(d,
16                                      centroids
17                                      assignmen
18}
19 cluster_results <- data.frame(k=ks, wcss=wcss)
20
21 ggplot(cluster_results, aes(x=k, y=wcss)) +
22   geom_point(size=3, color="#59B2D1") +
23   geom_line(color="#59B2D1") +
24   theme_bw()
```



Wrap Up

Recap

- Clustering has tons of applications in social science!
 - Putting things into groups is helpful, but sometimes we really don't know what the groups should be
 - If you want to learn more about clustering and similar unsupervised learning paradigms, take a look at my course "Modern Approaches in Measurement" happening in Spring Semester!
- Distance metrics give us a way to quantify how similar (or different) two observations are
 - The distance measure you choose dictates what features of the data are prioritized during this comparison
- `while` loops are our last major tool for control flow
 - Great for checking convergence conditions
 - If you can do a job with a `for` loop, however, prioritize using that

Final Thoughts

- PollEv.com/klintkanopka