



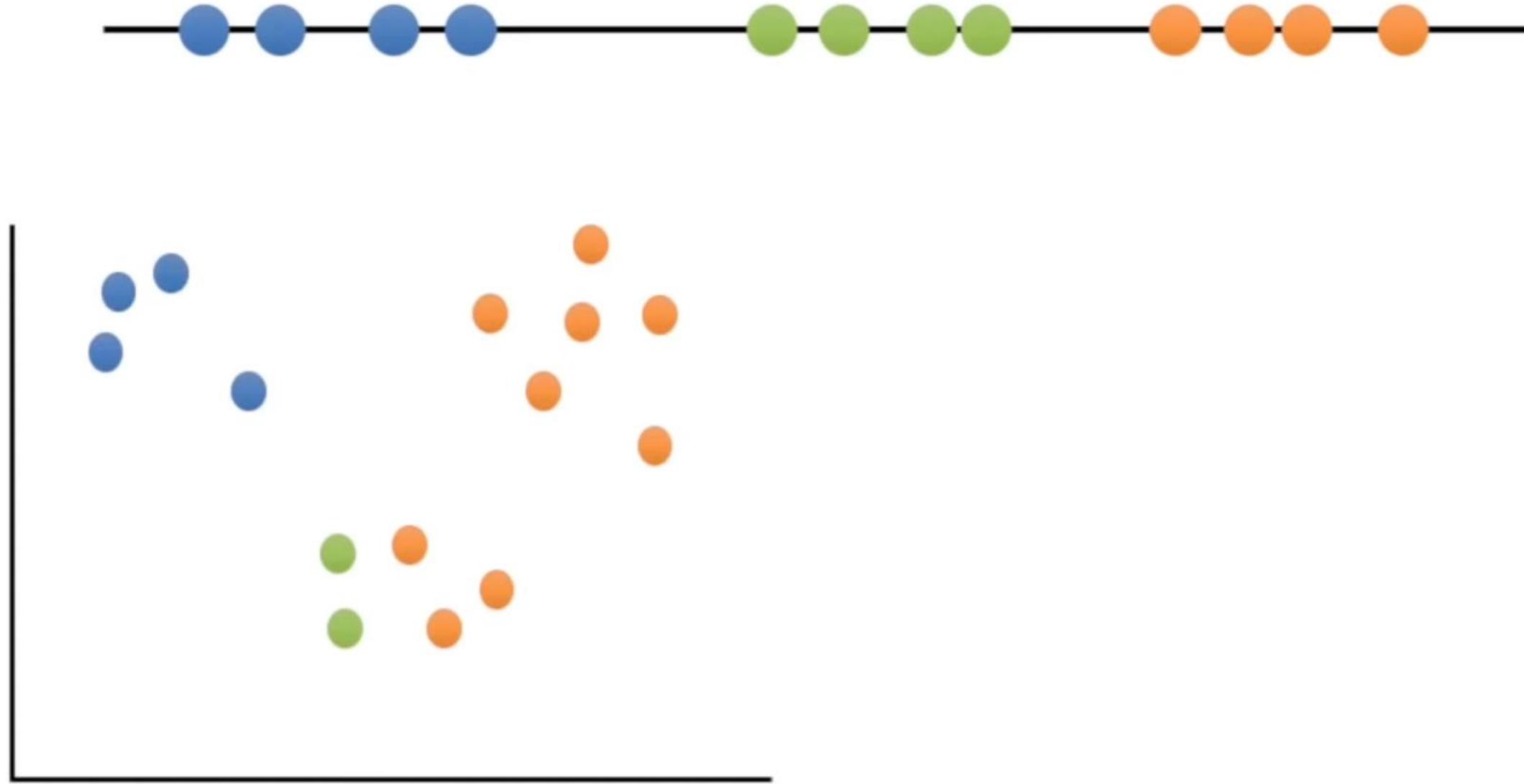
Machine Learning

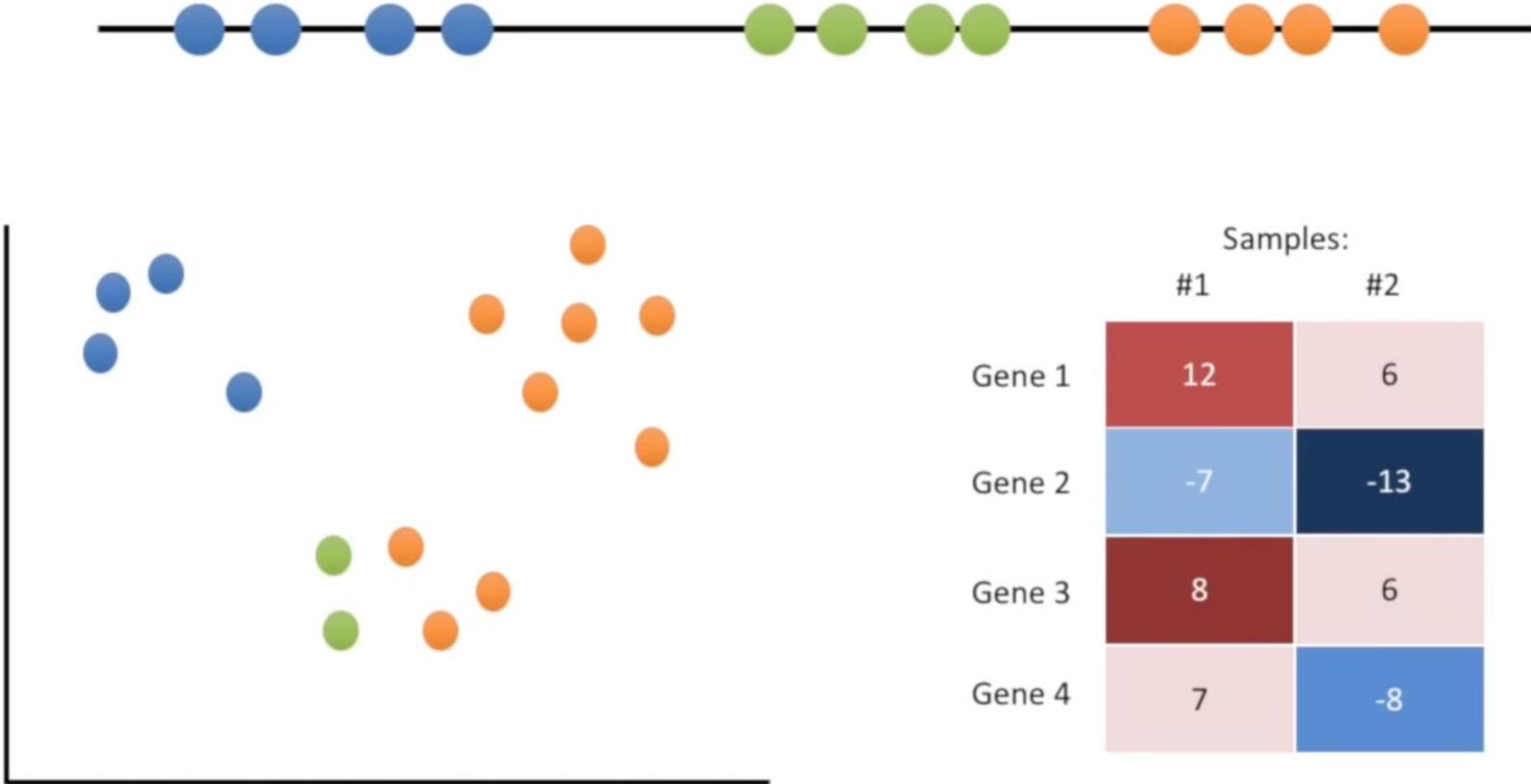
K-Means

Phd. César Astudillo | Facultad de Ingeniería

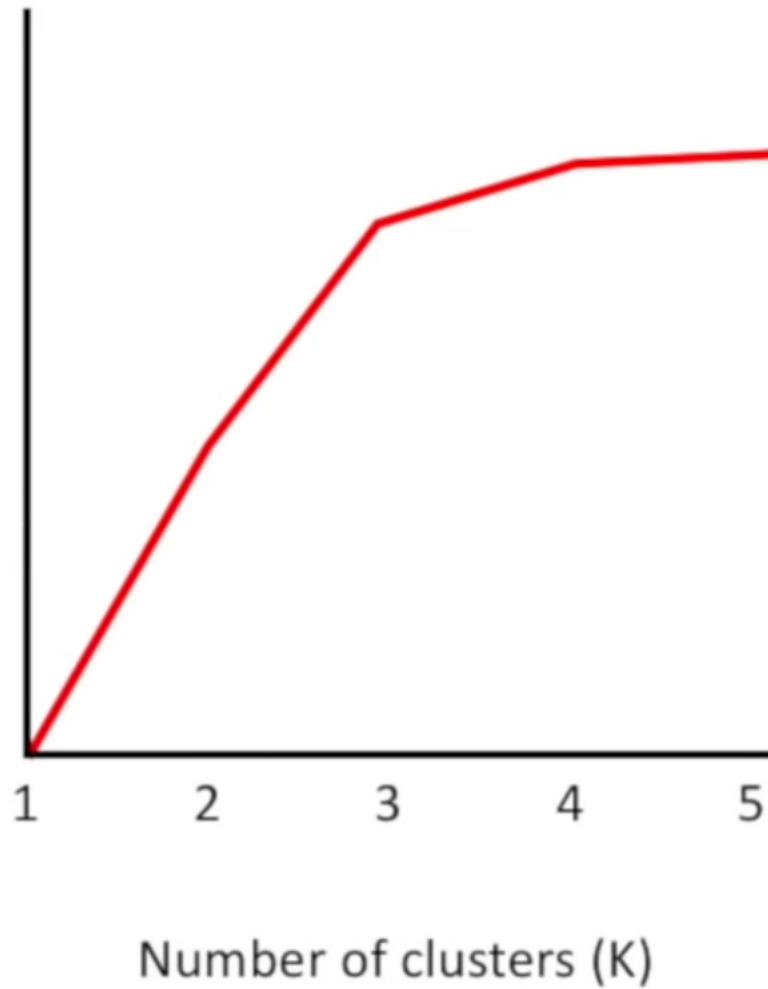
Introduction







Reduction is
Variation



K Means in one Dimension

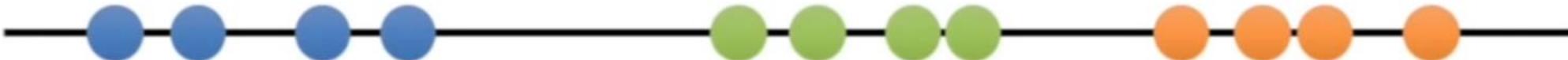
Imagine you had some data that you could plot on a line, and you knew you needed to put it into 3 clusters. Maybe they are measurements from 3 different types of tumors or other cell types.



In this case the data make three, relatively obvious, clusters.

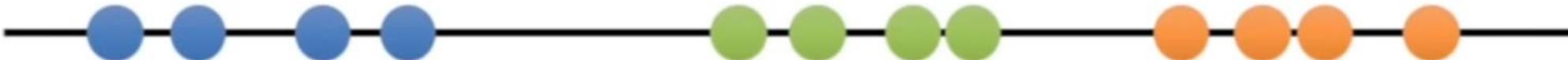


In this case the data make three, relatively obvious, clusters.



But, rather than rely on our eye, let's see if we can
get a computer to identify the same 3 clusters.

In this case the data make three, relatively obvious, clusters.



But, rather than rely on our eye, let's see if we can get a computer to identify the same 3 clusters.

To do this, we'll use K-means clustering.



Step 1: Select the number of clusters you want to identify in your data. This is the “K” in “K-means clustering”.



Step 1: Select the number of clusters you want to identify in your data. This is the “K” in “K-means clustering”.

In this case, we'll select K=3. That is to say, we want to identify 3 clusters.



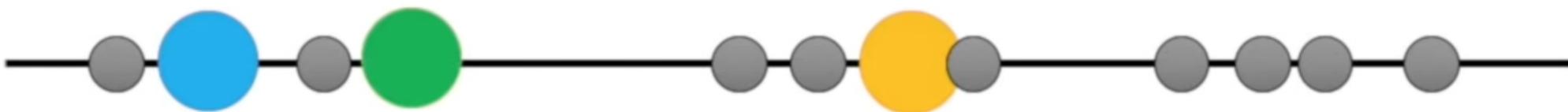
Step 1: Select the number of clusters you want to identify in your data. This is the “K” in “K-means clustering”.

In this case, we'll select K=3. That is to say, we want to identify 3 clusters.



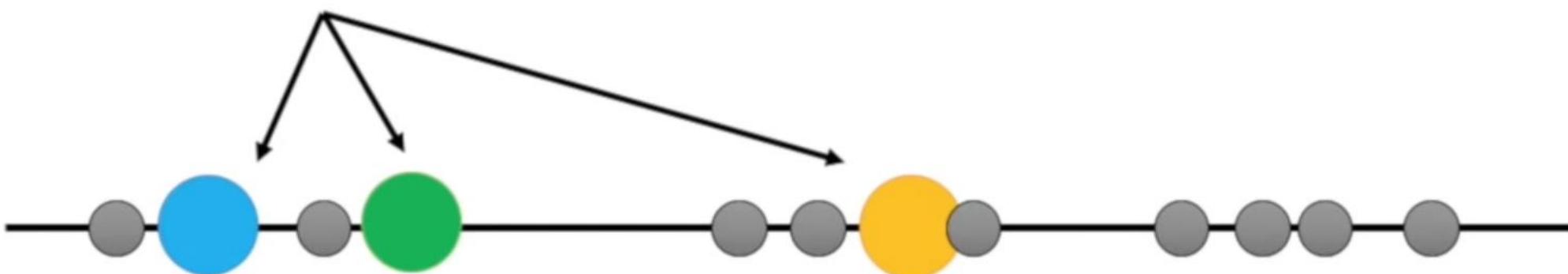
There is a fancier way to select a value for “K”, but we'll talk about that later.

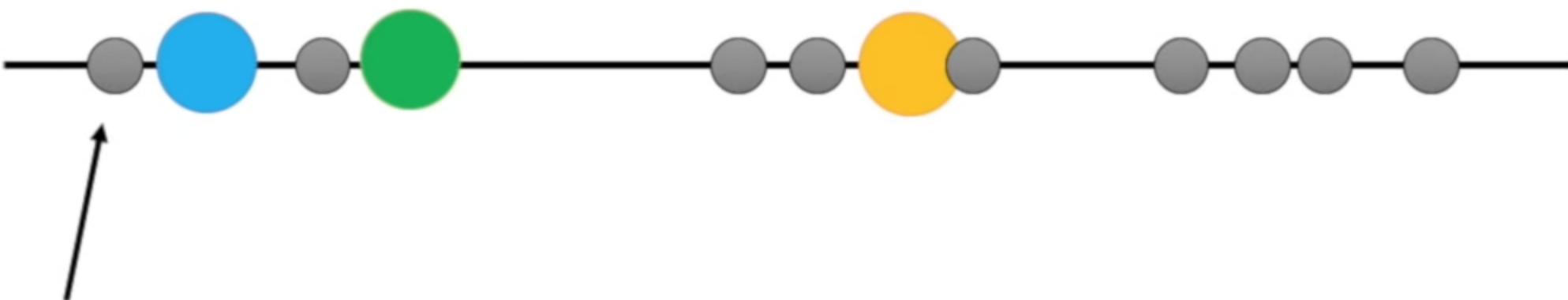
Step 2: Randomly select 3 distinct data points.



Step 2: Randomly select 3 distinct data points.

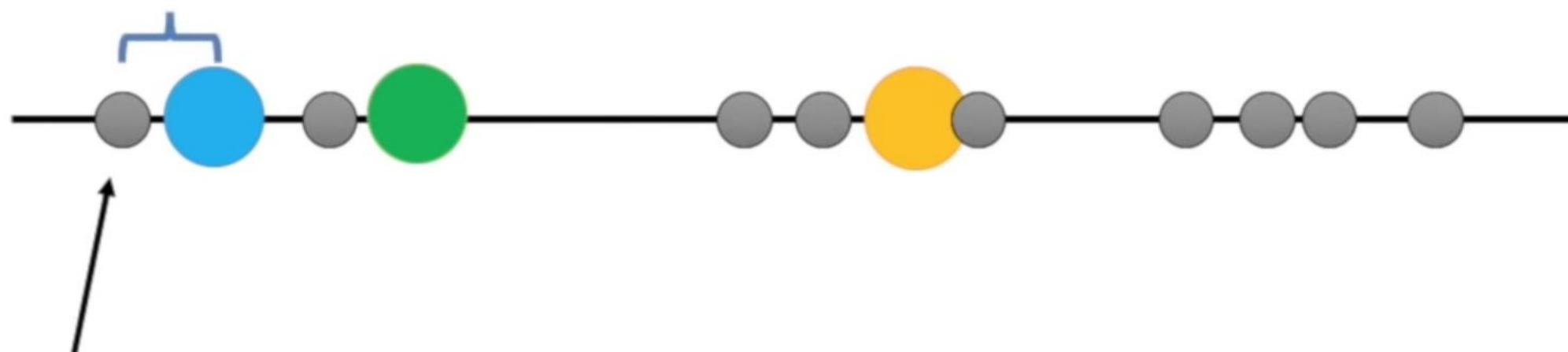
These are the initial clusters.





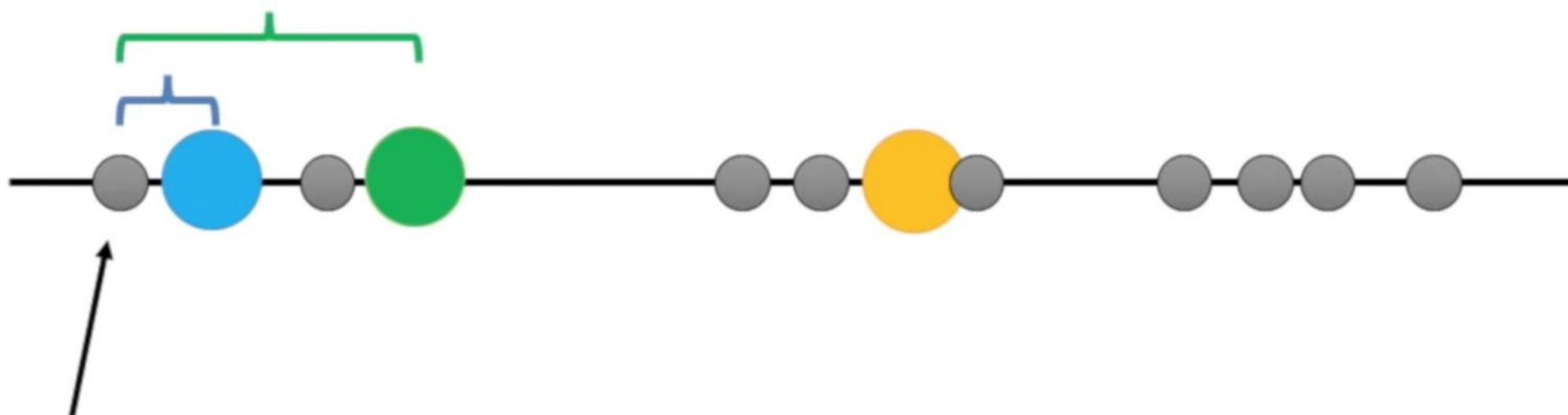
Step 3: Measure the distance
between the 1st point and the three
initial clusters.

Distance from the 1st
point to the **blue**
cluster



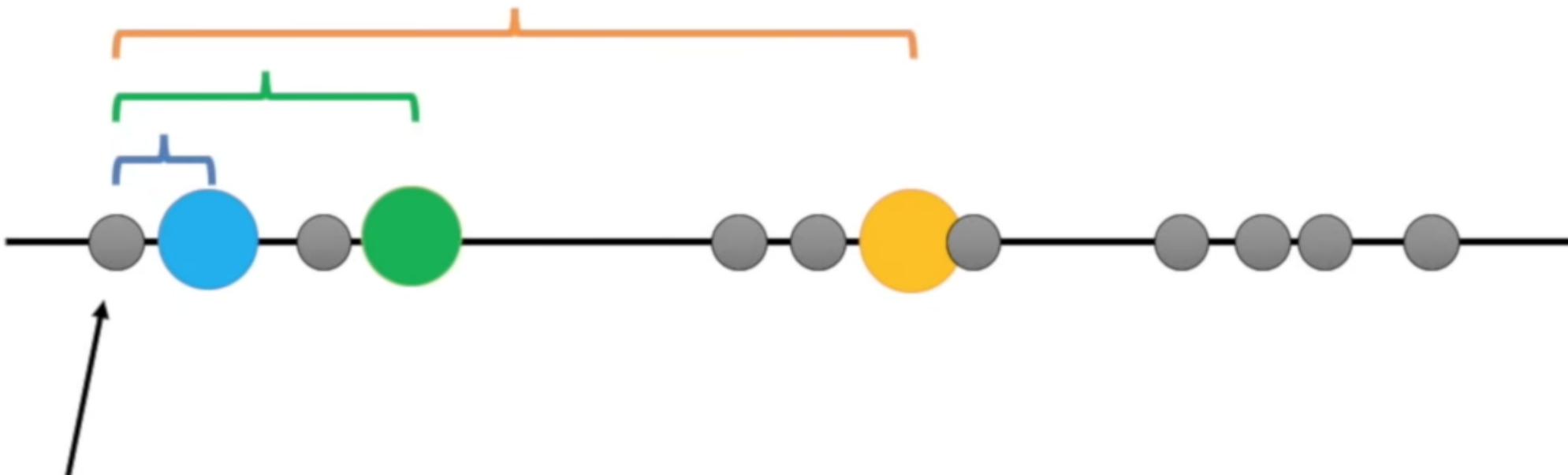
Step 3: Measure the distance
between the 1st point and the three
initial clusters.

Distance from the 1st
point to the **green**
cluster

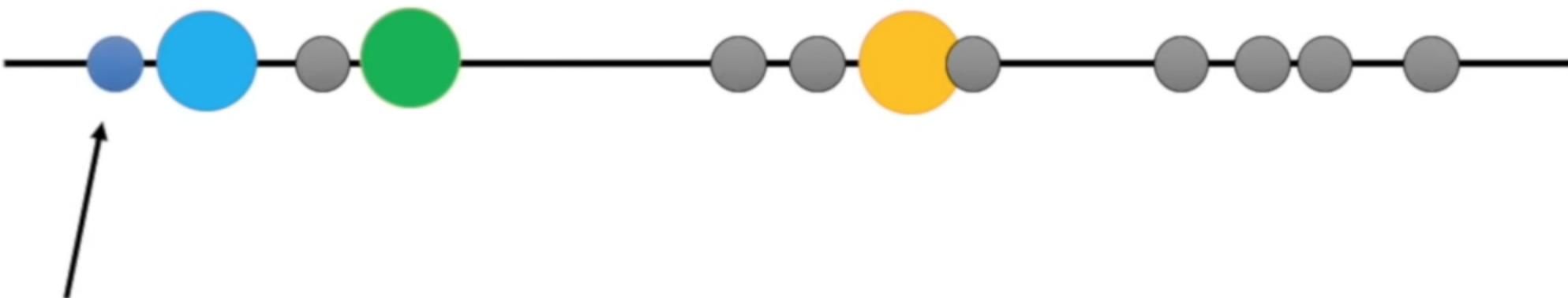


Step 3: Measure the distance
between the 1st point and the three
initial clusters.

Distance from the 1st
point to the **orange**
cluster

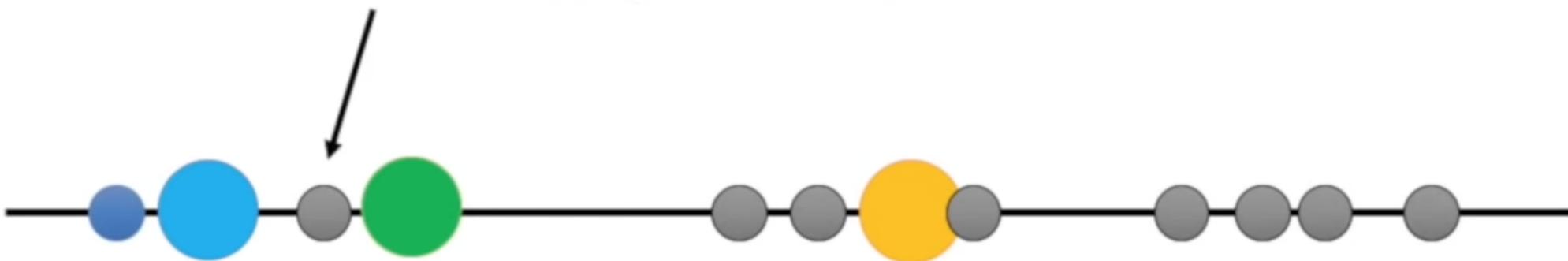


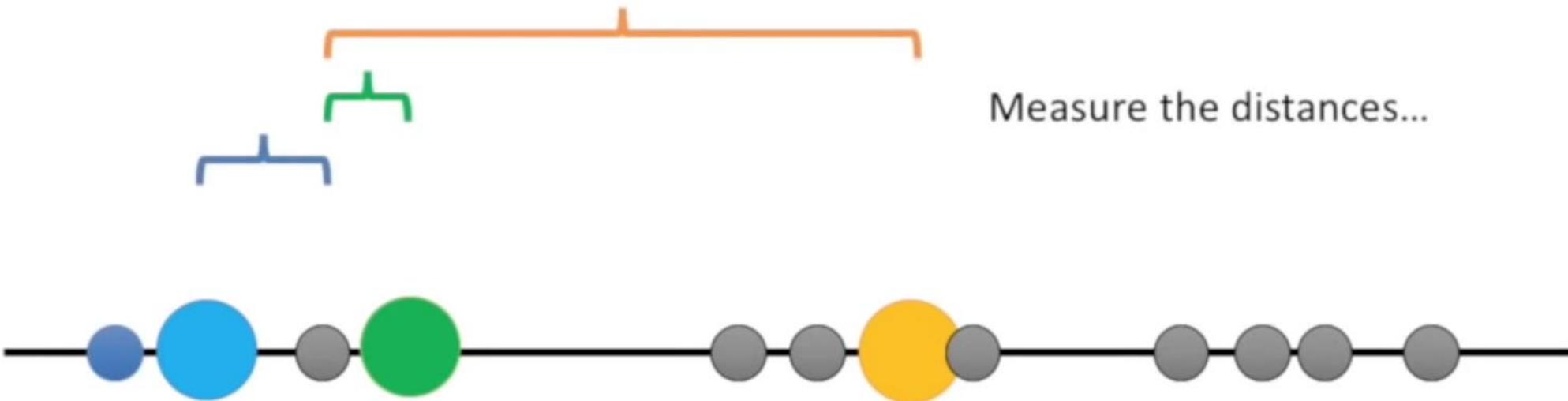
Step 3: Measure the distance
between the 1st point and the three
initial clusters.

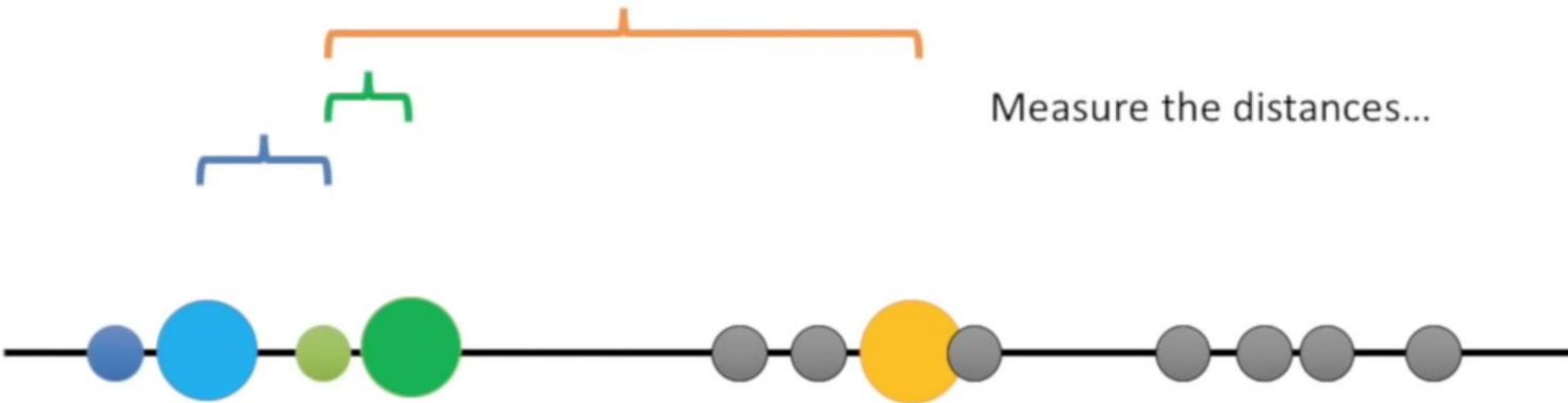


Step 4: Assign the 1st point to the nearest cluster. In this case, the nearest cluster is the **blue** cluster.

Now do the same thing for the next point.

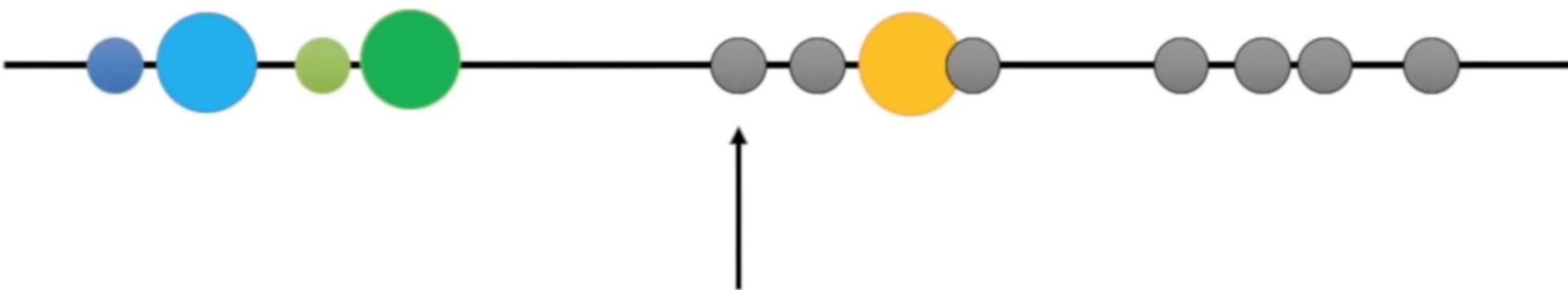




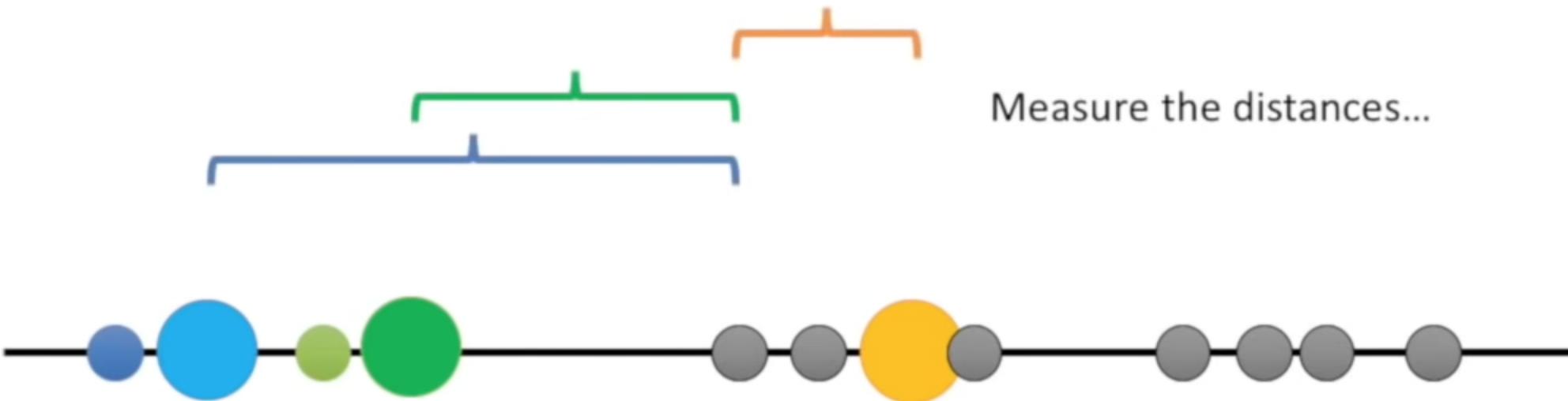


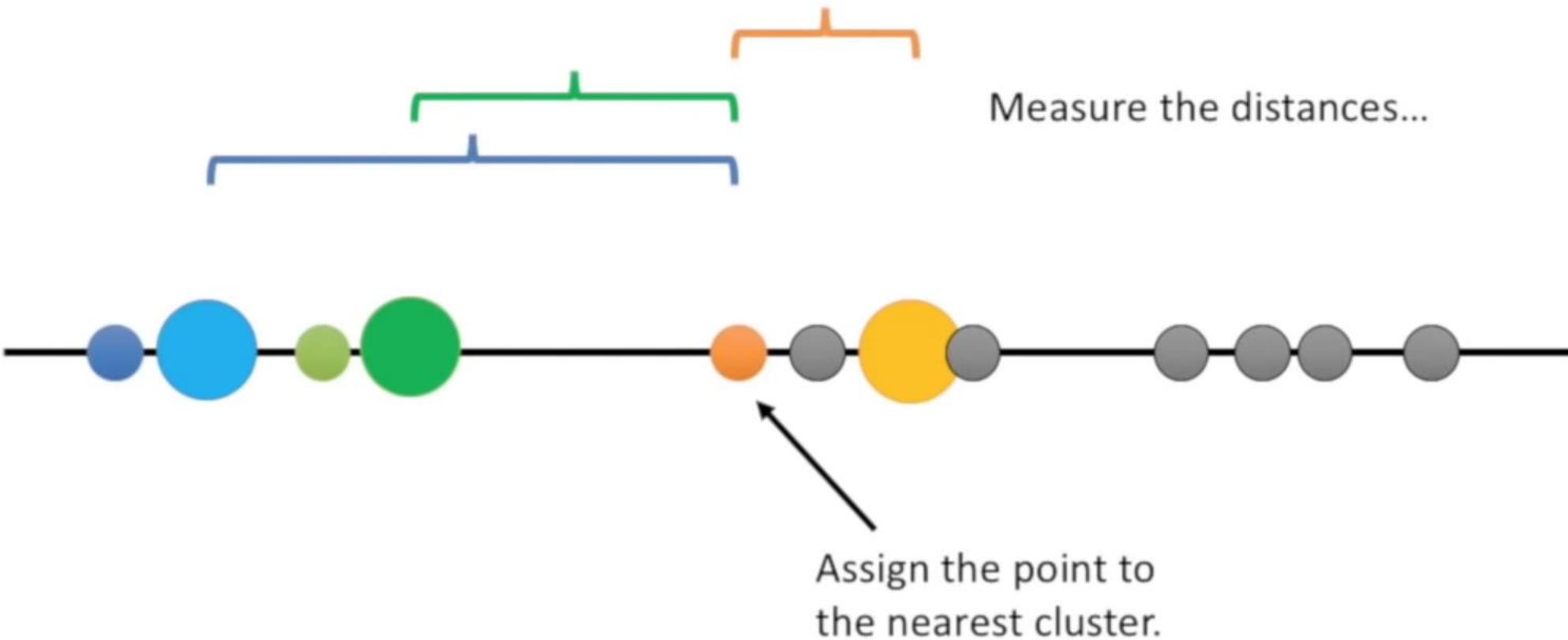
Assign the point
to the nearest
cluster.

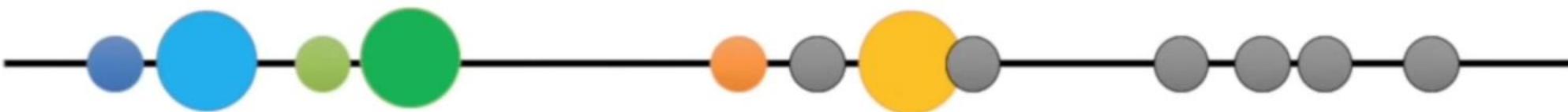
Measure the distances...



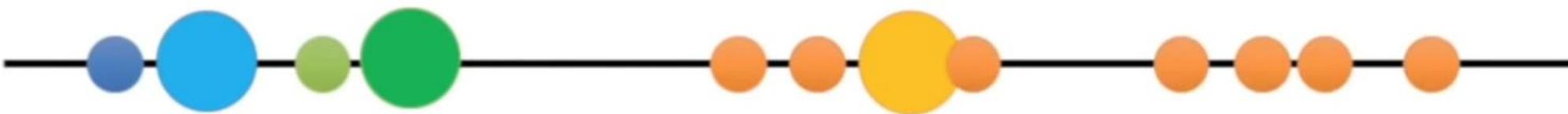
Now figure out which cluster the 3rd point belongs to.





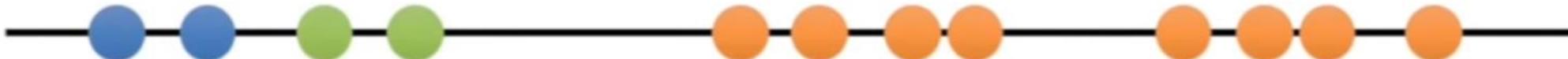


The rest of these points are closest to the **orange** cluster, so they'll go in that one, too.

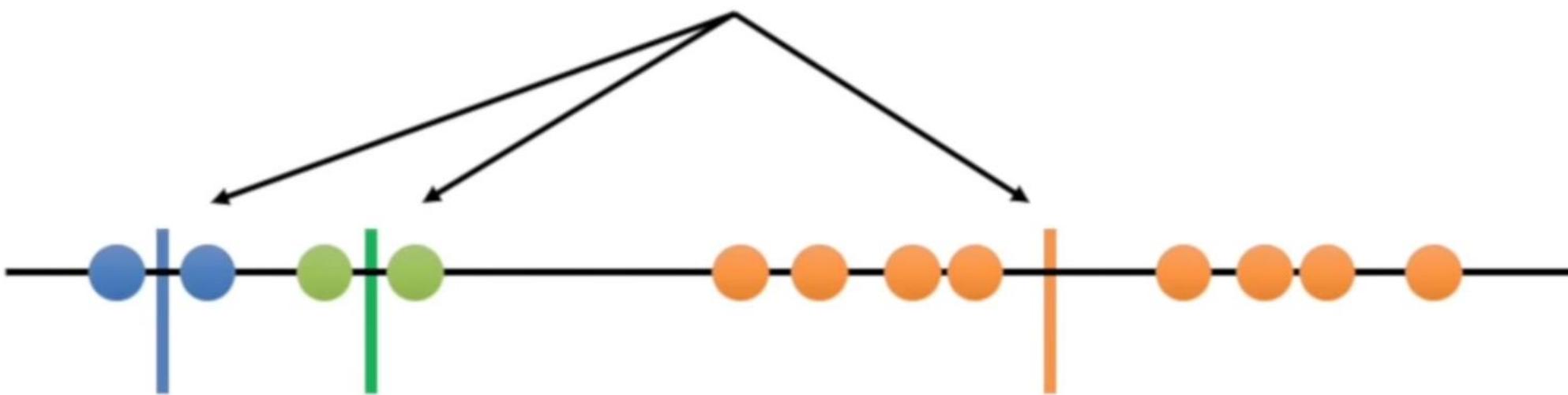


The rest of these points are
closest to the **orange** cluster,
so they'll go in that one, too.

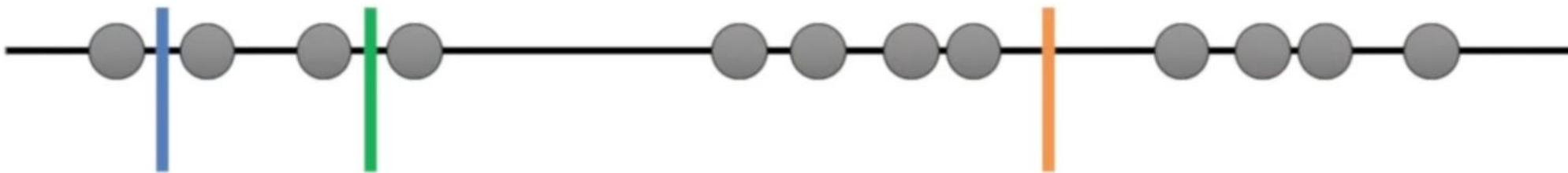
Now that all of the points are
in clusters, we go on to...



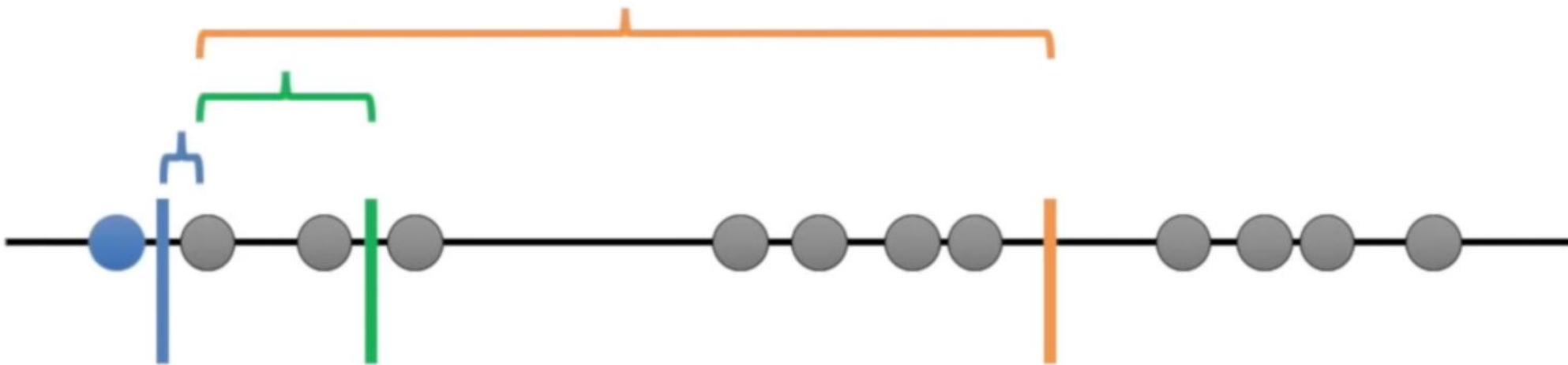
Step 5: calculate the mean of each cluster.



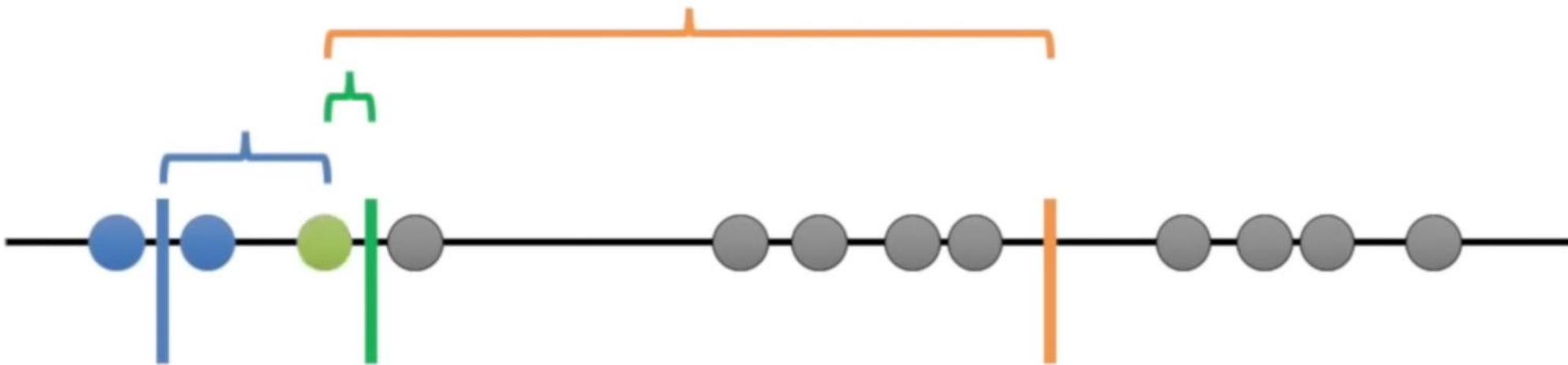
Then we repeat what we just
did (measure and cluster)
using the mean values.



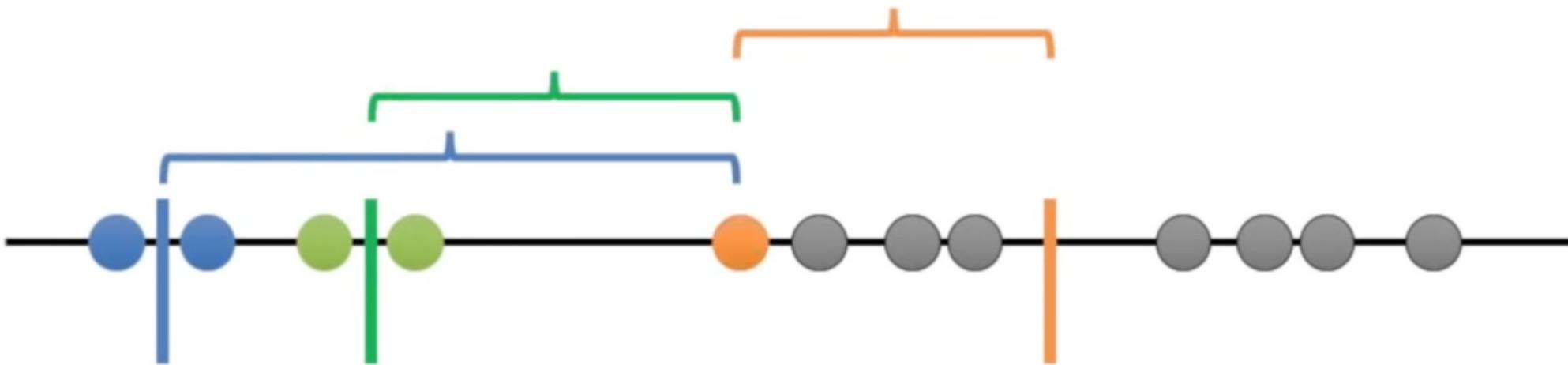
Then we repeat what we just
did (measure and cluster)
using the mean values.



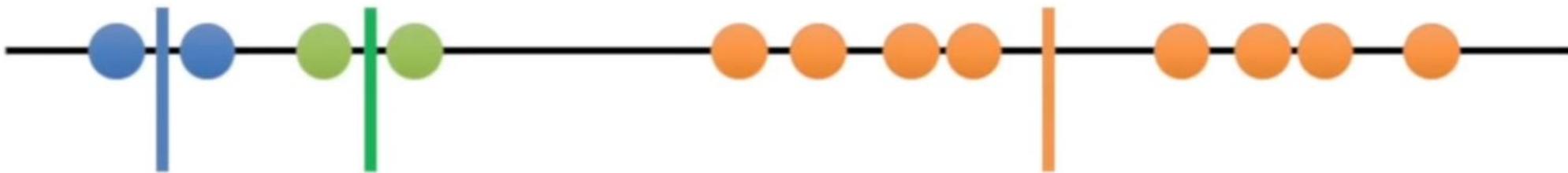
Then we repeat what we just
did (measure and cluster)
using the mean values.



Then we repeat what we just
did (measure and cluster)
using the mean values.



Since the clustering did not change at all during the last iteration, we're done...

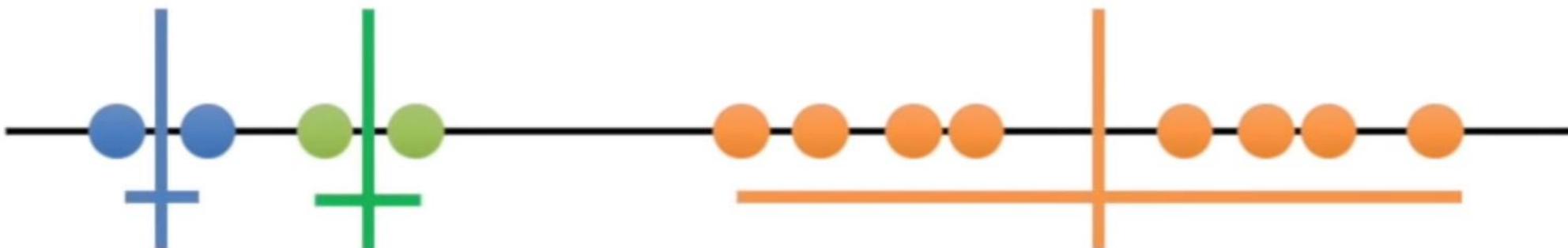




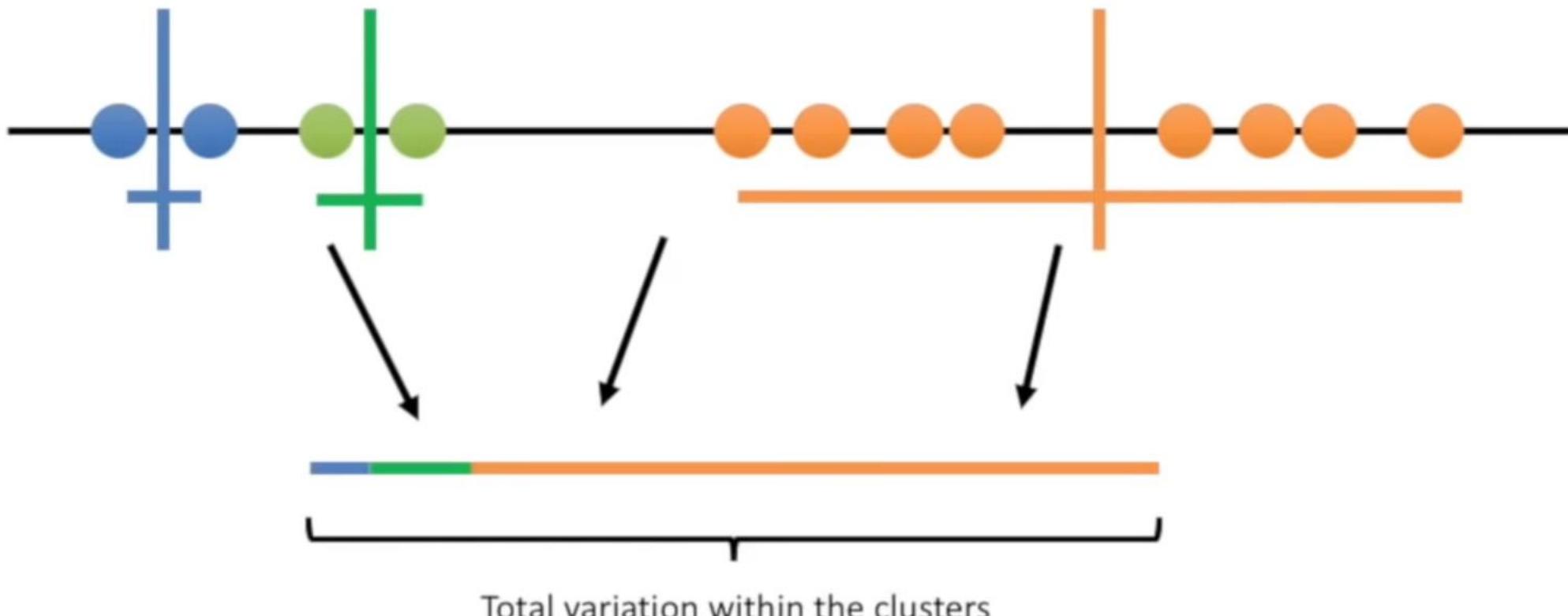
The K-means clustering is pretty terrible compared to what
we did by eye.



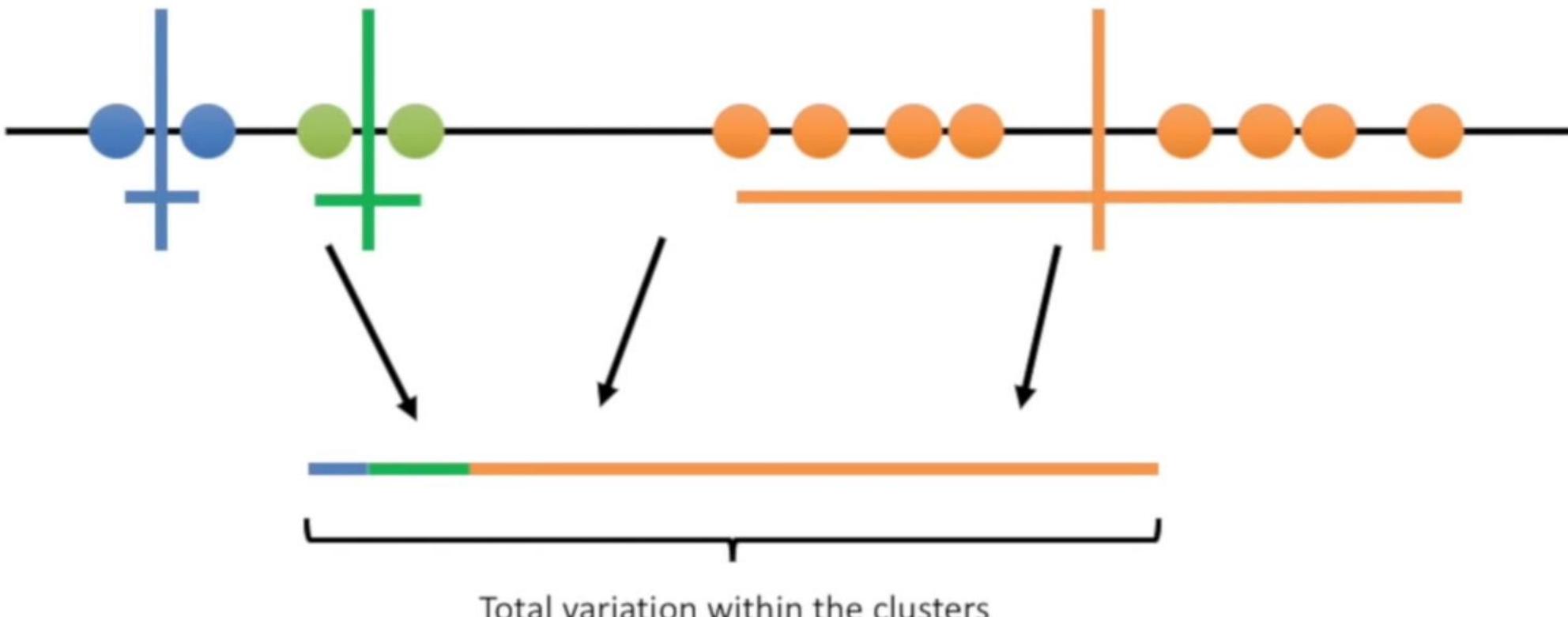
We can assess the quality of the clustering by adding up the variation within each cluster.



We can assess the quality of the clustering by adding up the variation within each cluster.

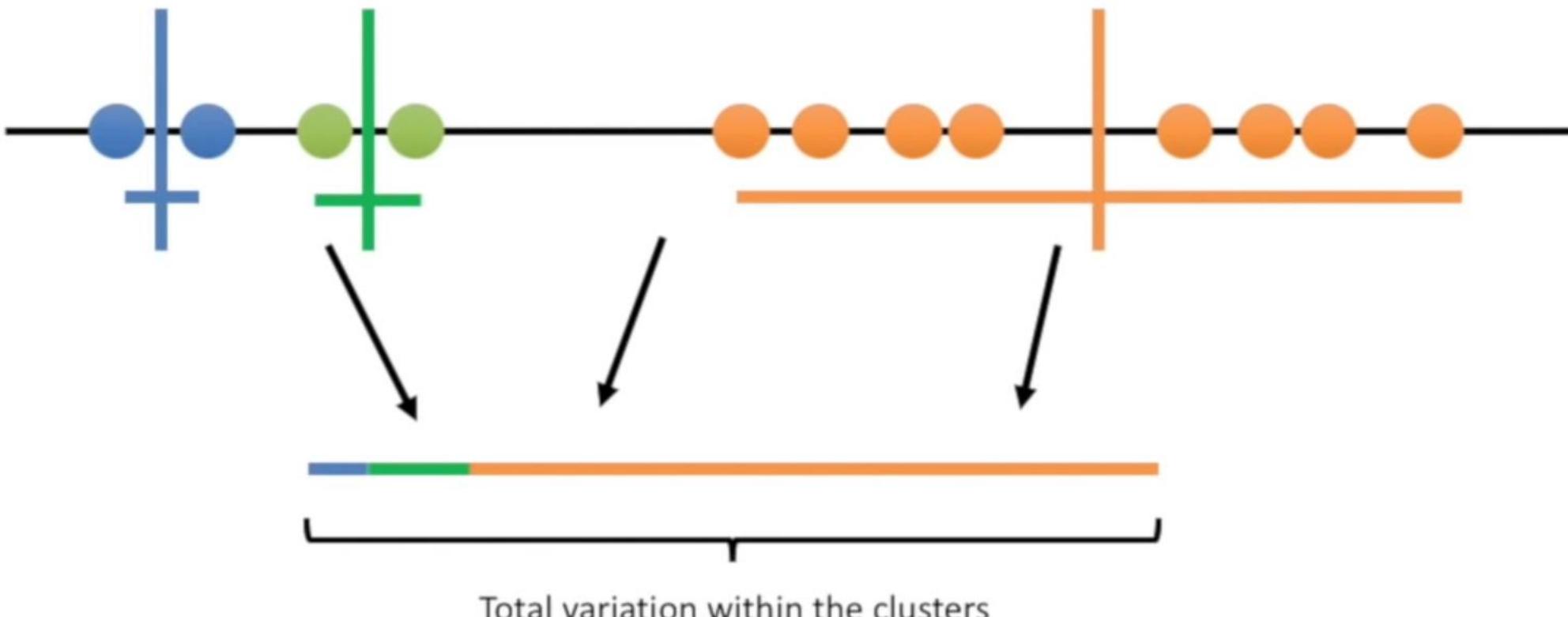


We can assess the quality of the clustering by adding up the variation within each cluster.



Since K-means clustering can't "see" the best clustering, its only option is to keep track of these clusters, and their total variance, and do the whole thing over again with different starting points.

We can assess the quality of the clustering by adding up the variation within each cluster.



Since K-means clustering can't "see" the best clustering, its only option is to keep track of these clusters, and their total variance, and do the whole thing over again with different starting points.

So, here we are again, back at the beginning.



K-means clustering picks 3 initial clusters...



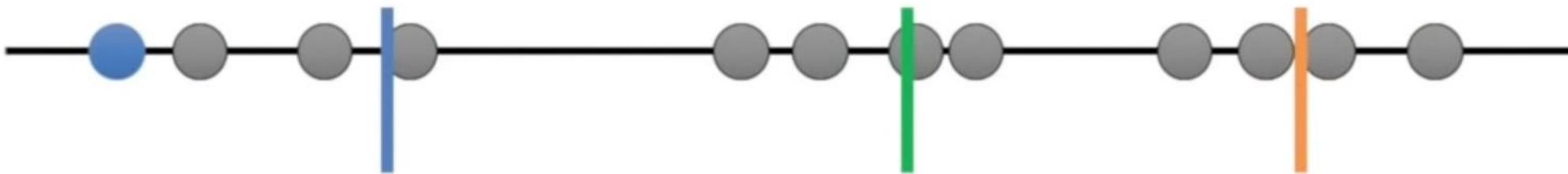
...and then clusters all the remaining points, calculates the mean of each cluster and then reclusters based on the new means. It repeats until the clusters no longer change.



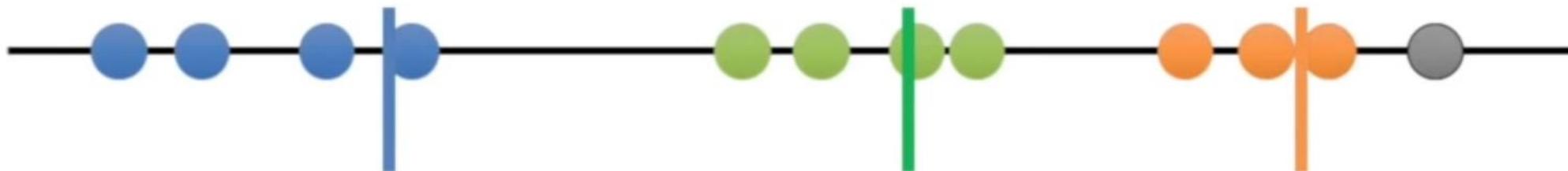
...and then clusters all the remaining points, calculates the mean of each cluster and then reclusters based on the new means. It repeats until the clusters no longer change.



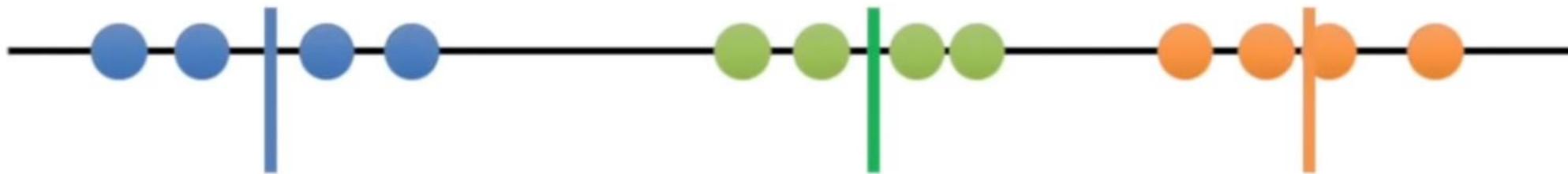
...and then clusters all the remaining points, calculates the mean of each cluster and then reclusters based on the new means. It repeats until the clusters no longer change.



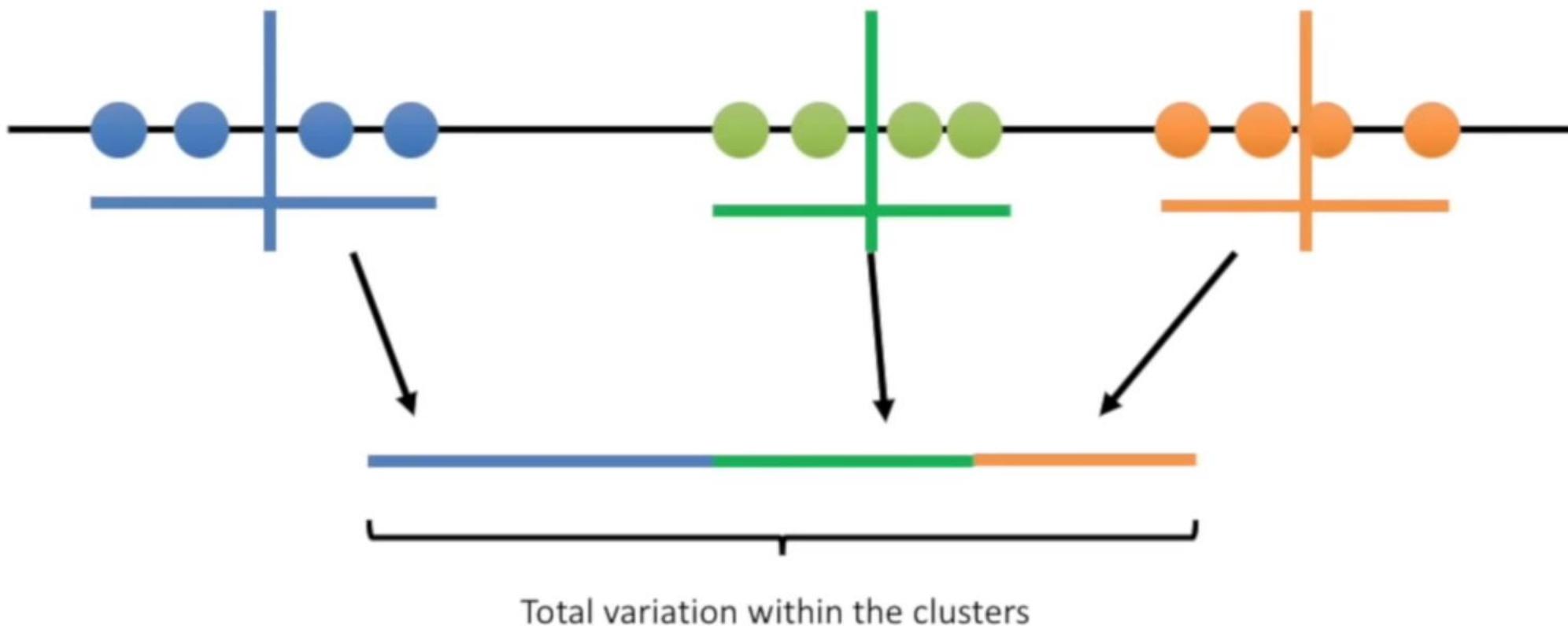
...and then clusters all the remaining points, calculates the mean of each cluster and then reclusters based on the new means. It repeats until the clusters no longer change.



...and then clusters all the remaining points, calculates the mean of each cluster and then reclusters based on the new means. It repeats until the clusters no longer change.



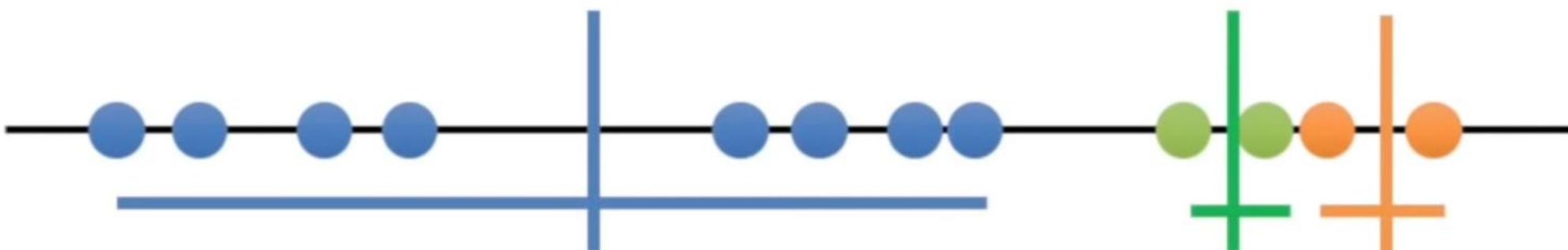
Now that the data are clustered, we sum the variation within each cluster.



And then do it all again...



At this point, K-means clustering knows that *the 2nd clustering is the best clustering so far*. But it doesn't know if it's *the best overall*, so it will do a few more clusters (it does as many as you tell it to do) and then come back and return that one if it is still the best.



1st cluster attempt:



2nd cluster attempt:

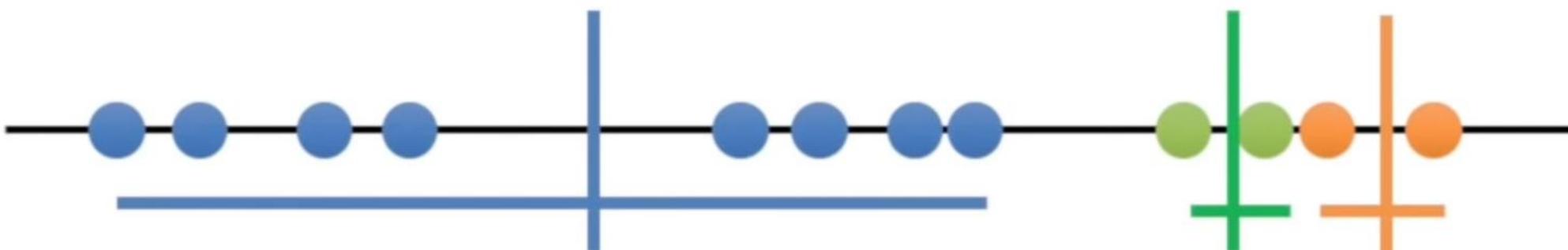


The winner!!

3rd cluster attempt:



At this point, K-means clustering knows that *the 2nd clustering is the best clustering so far*. But it doesn't know if it's *the best overall*, so it will do a few more clusters (it does as many as you tell it to do) and then come back and return that one if it is still the best.



1st cluster attempt:



2nd cluster attempt:



The winner!!

3rd cluster attempt:

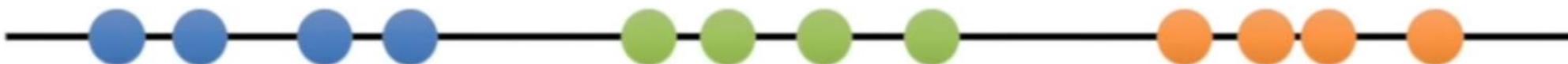


Determining a suitable value for K

Question: How do you figure out what value to use for “K”?



Question: How do you figure out what value to use for “K”?



With this data, it's obvious that we should set K to 3, but
other times it is not so clear.

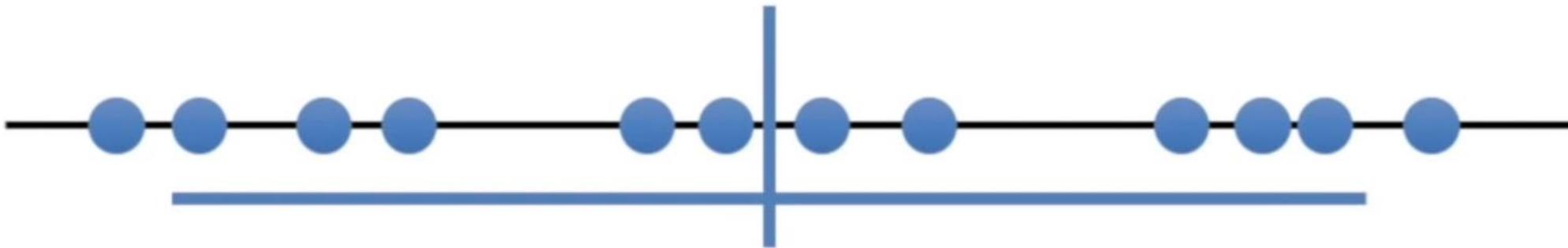
One way to decide is to just try different values for K.



Start with K = 1



Start with $K = 1$

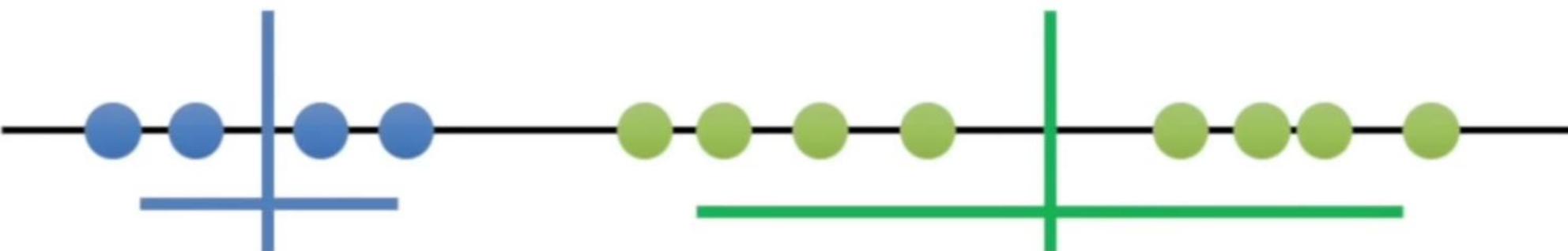


$K = 1$ is the worst case scenario. We can quantify its “badness” with the total variation.

Now try K = 2

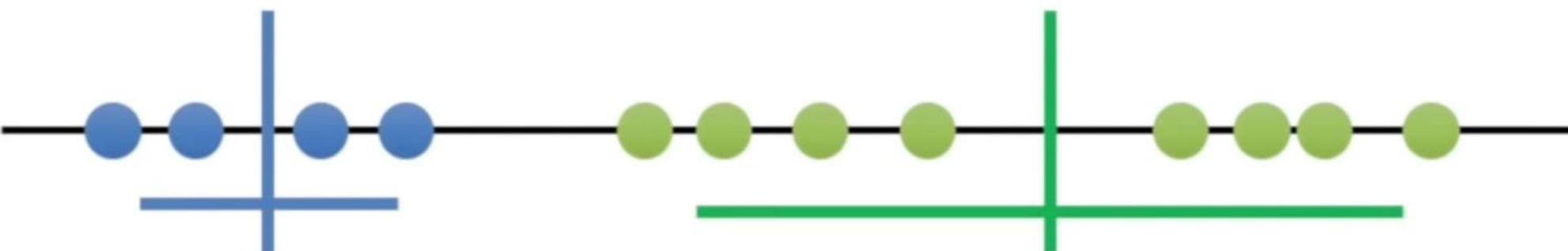


Now try $K = 2$



$K = 2$ is better, and we can quantify how much better by comparing the total variation within the 2 clusters to $K = 1$

Now try $K = 2$



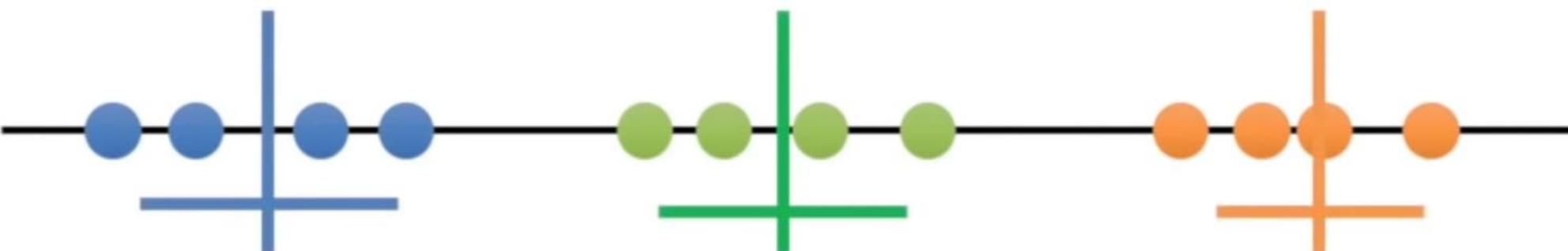
$K = 2$ is better, and we can quantify how much better by comparing the total variation within the 2 clusters to $K = 1$



Now try K = 3

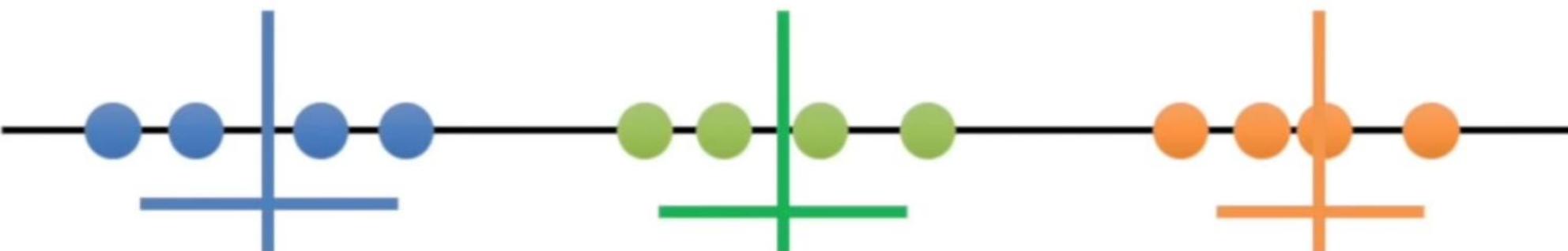


Now try $K = 3$



$K = 3$ is even better! We can quantify how much better by comparing the total variation within the 3 clusters to $K = 2$

Now try $K = 3$



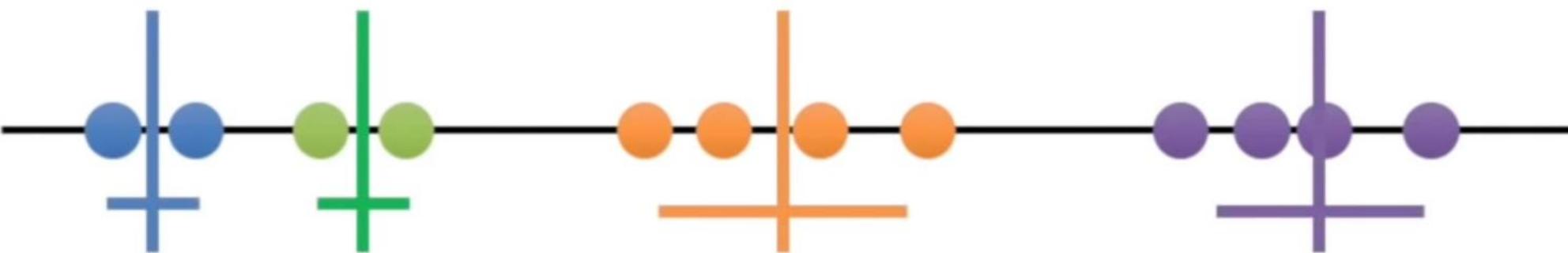
$K = 3$ is even better! We can quantify how much better by comparing the total variation within the 3 clusters to $K = 2$



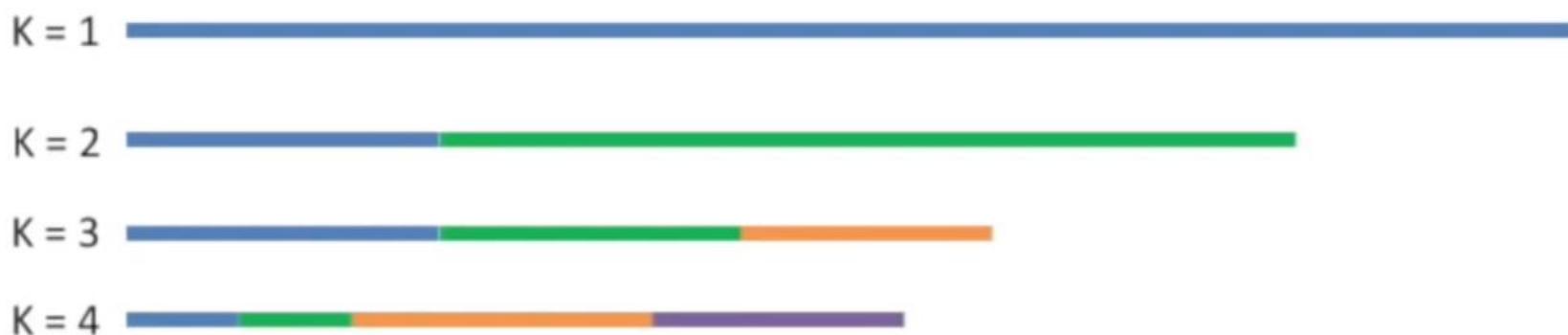
Now try K = 4



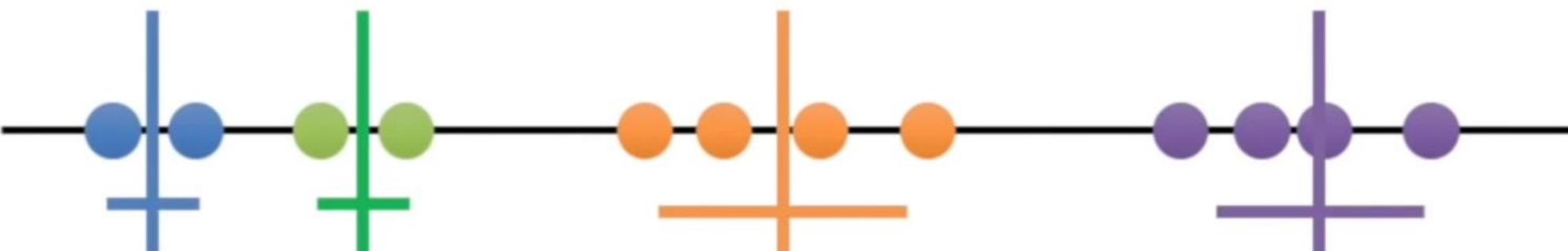
Now try $K = 4$



The total variation within each cluster is less than when $K=3$



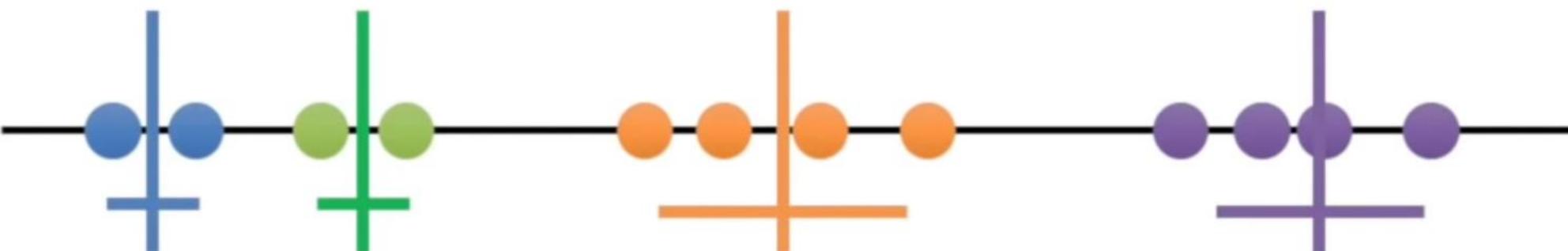
Now try K = 4



The total variation within each cluster is less than when K=3

Each time we add a new cluster, the total variation within each cluster is smaller than before. And when there is only one point per cluster, the variation = 0.

Now try K = 4

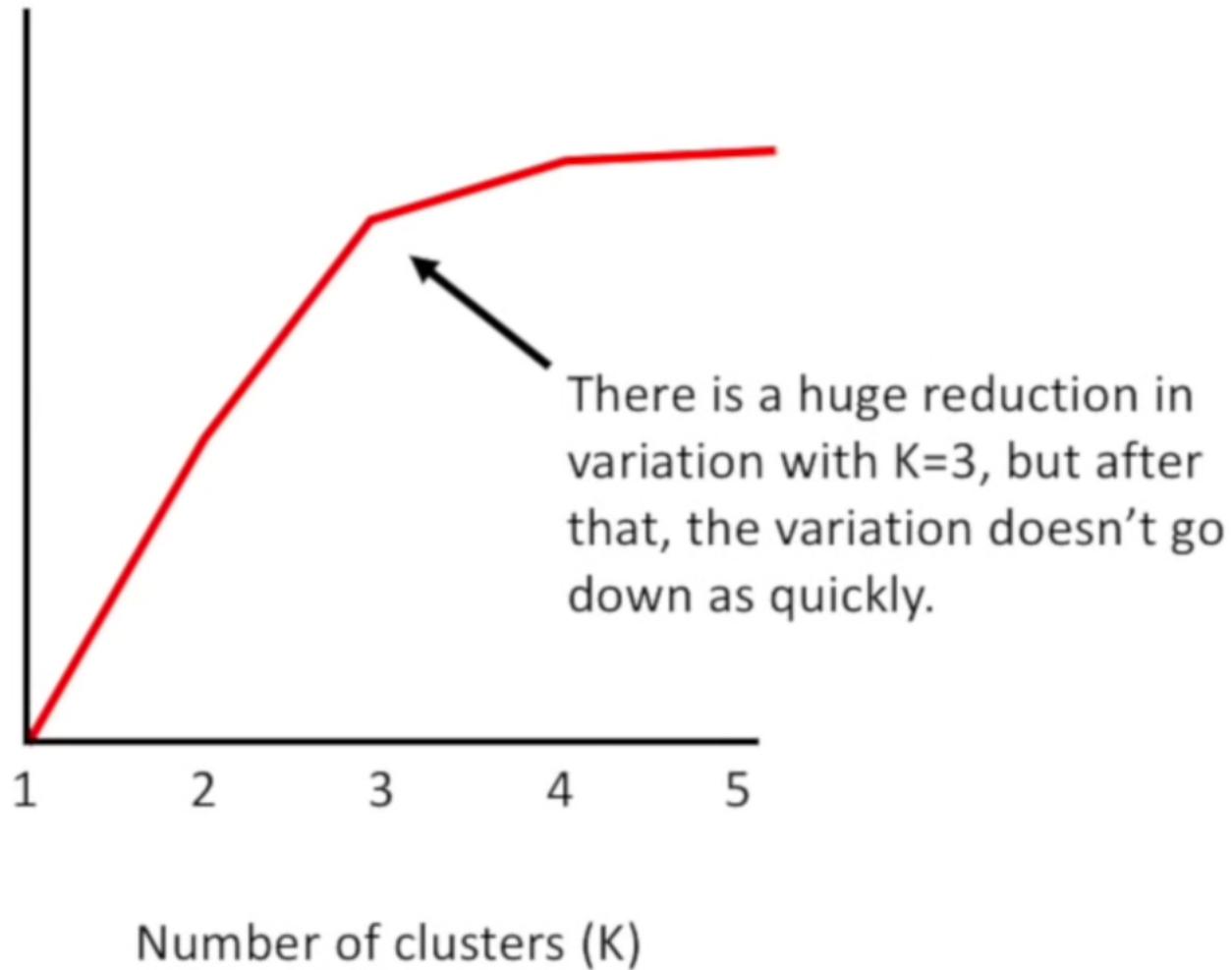


The total variation within each cluster is less than when K=3

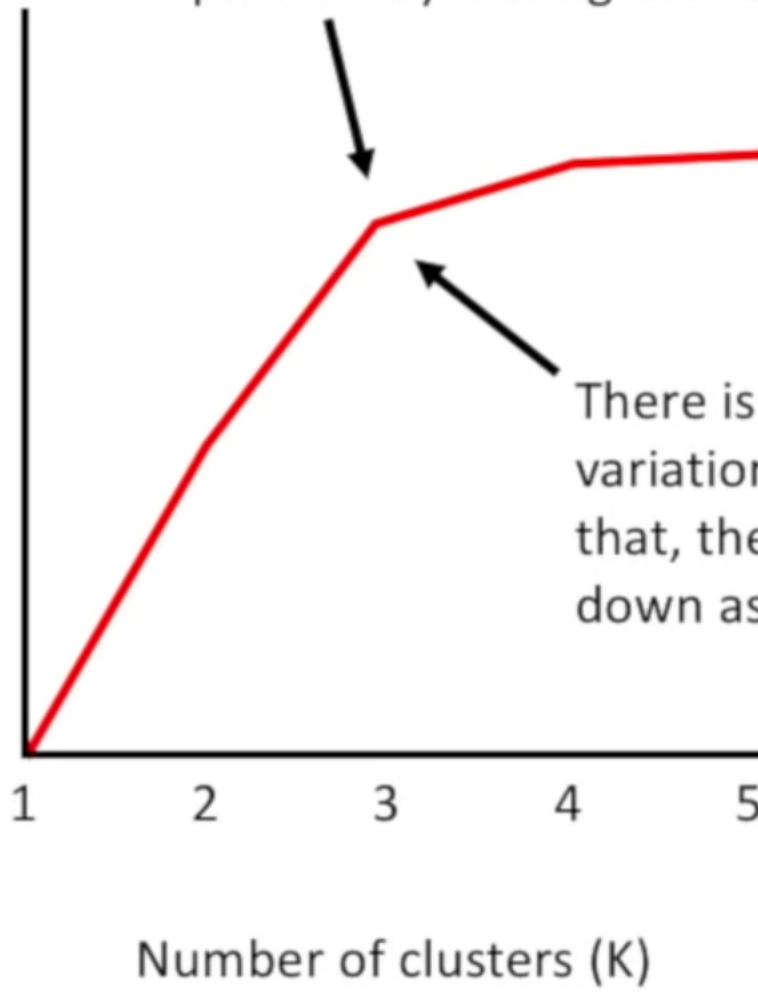
Each time we add a new cluster, the total variation within each cluster is smaller than before. And when there is only one point per cluster, the variation = 0.

However, if we plot the reduction in variance per value for K...

Reduction is Variation



Reduction is Variation

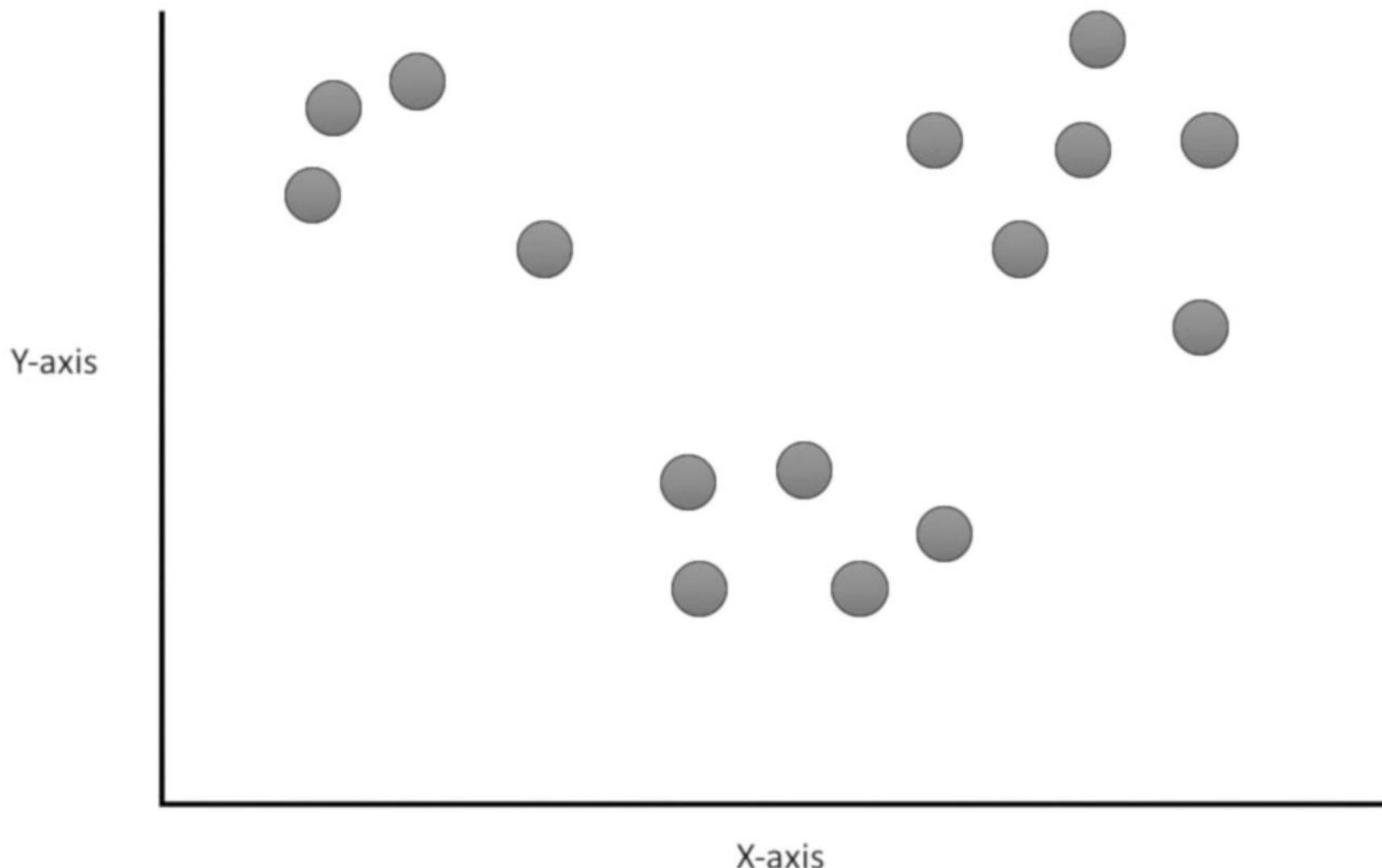


This is called an “elbow plot”, and you can pick “K” by finding the “elbow” in the plot

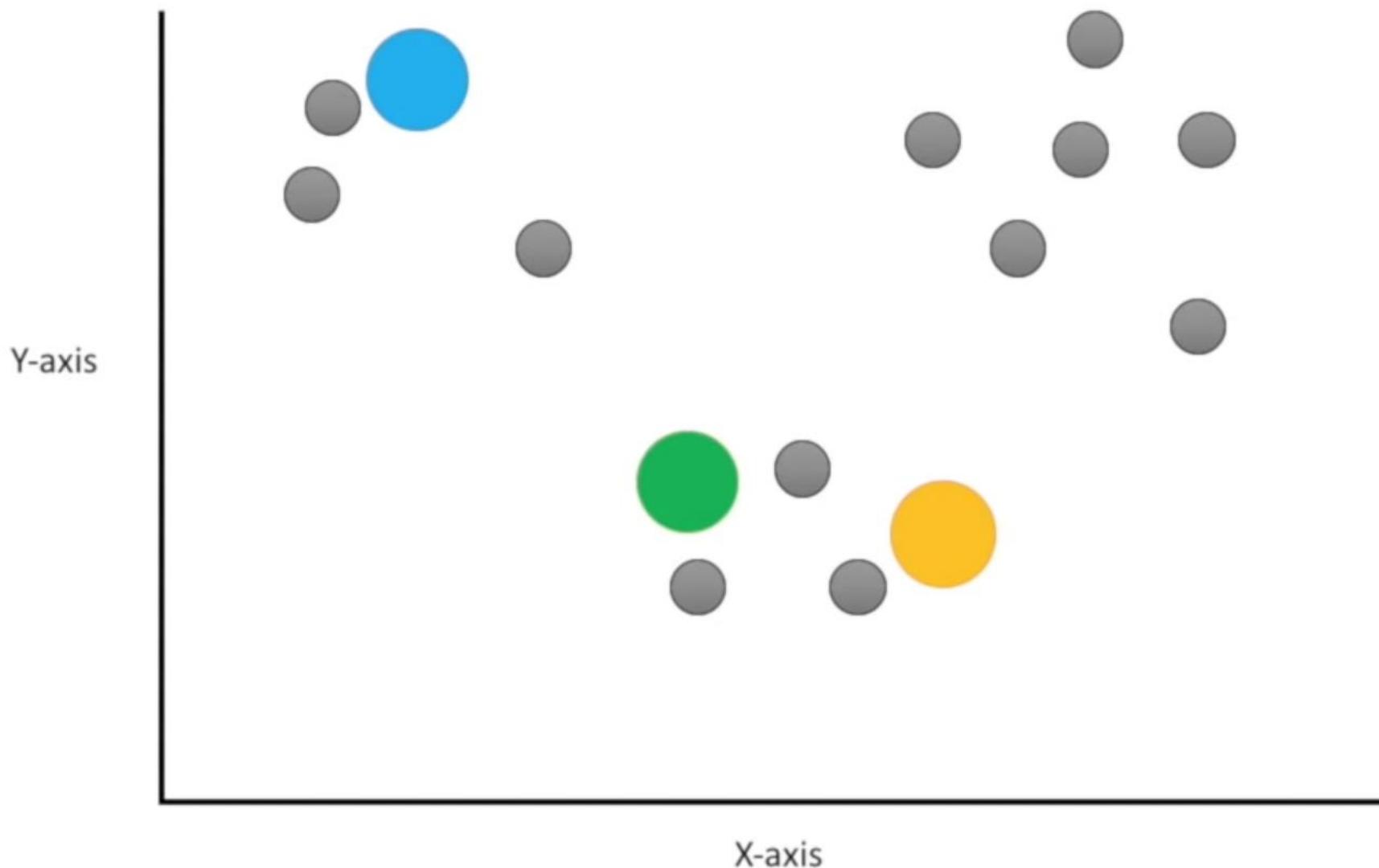
There is a huge reduction in variation with $K=3$, but after that, the variation doesn't go down as quickly.

K Means in two dimensions

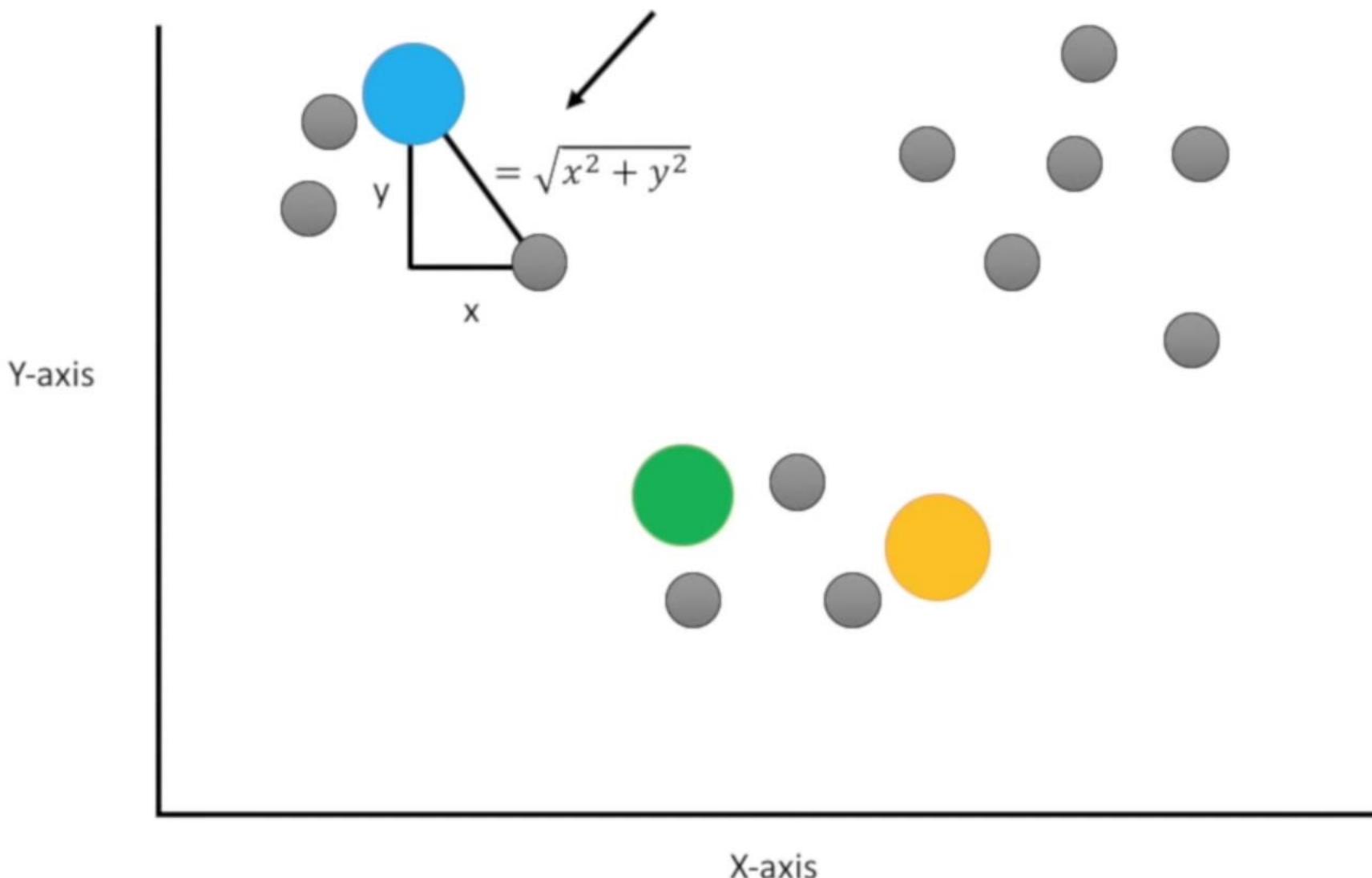
Question: What if our data isn't plotted on a number line?



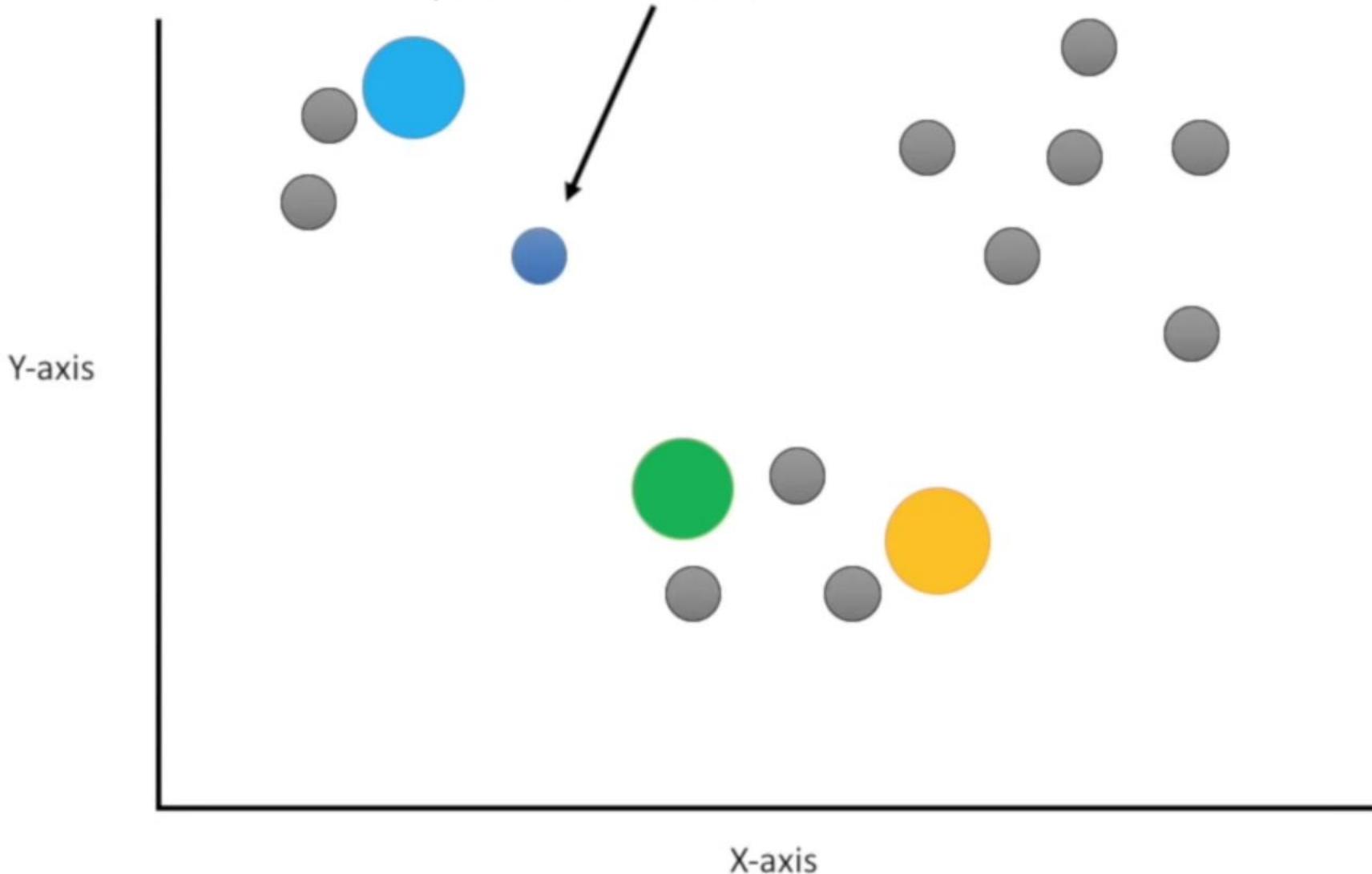
Just like before, you pick three random points...



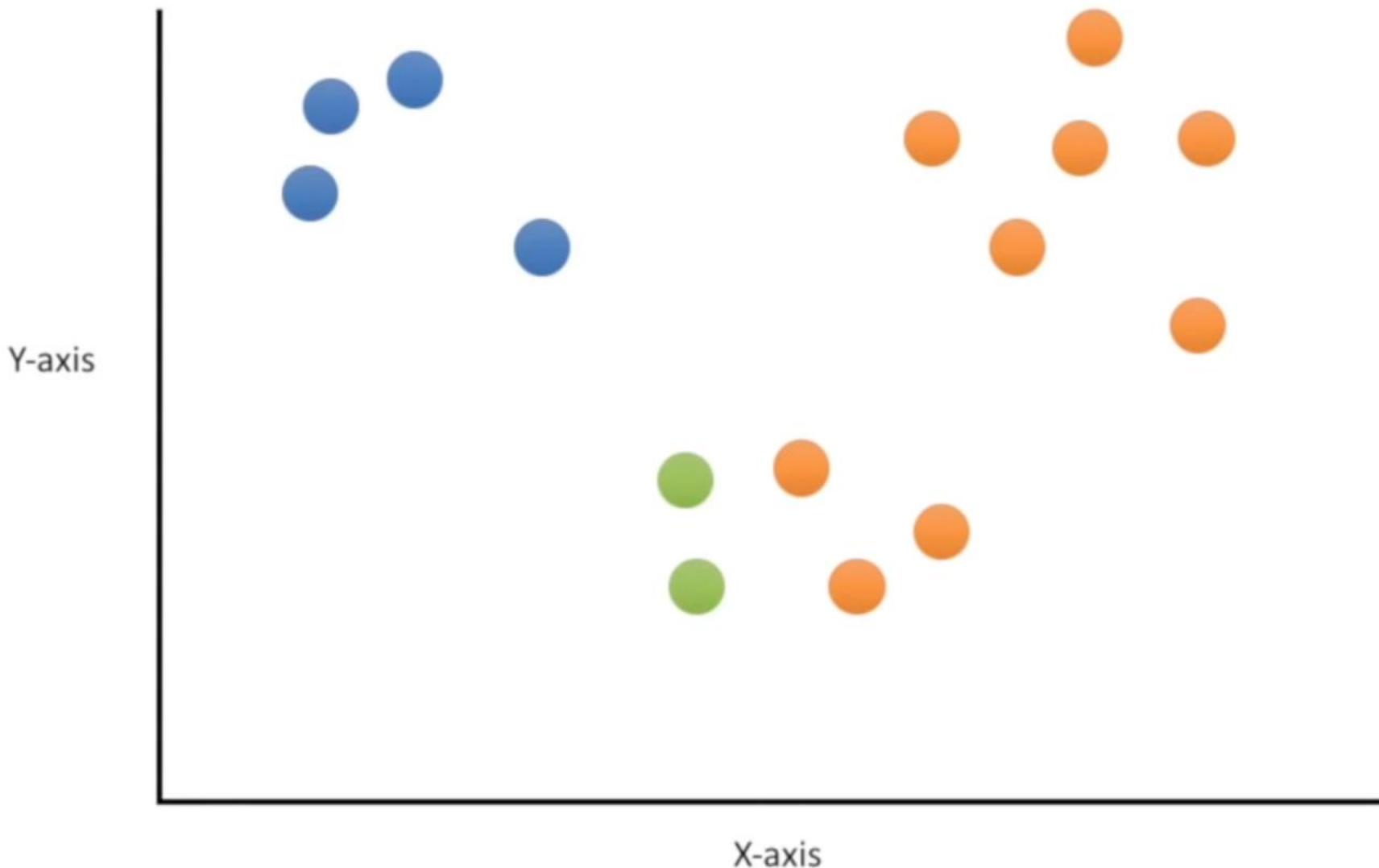
And we use the Euclidean distance. In 2 dimensions, the Euclidian distance is the same thing as the Pythagorean theorem.



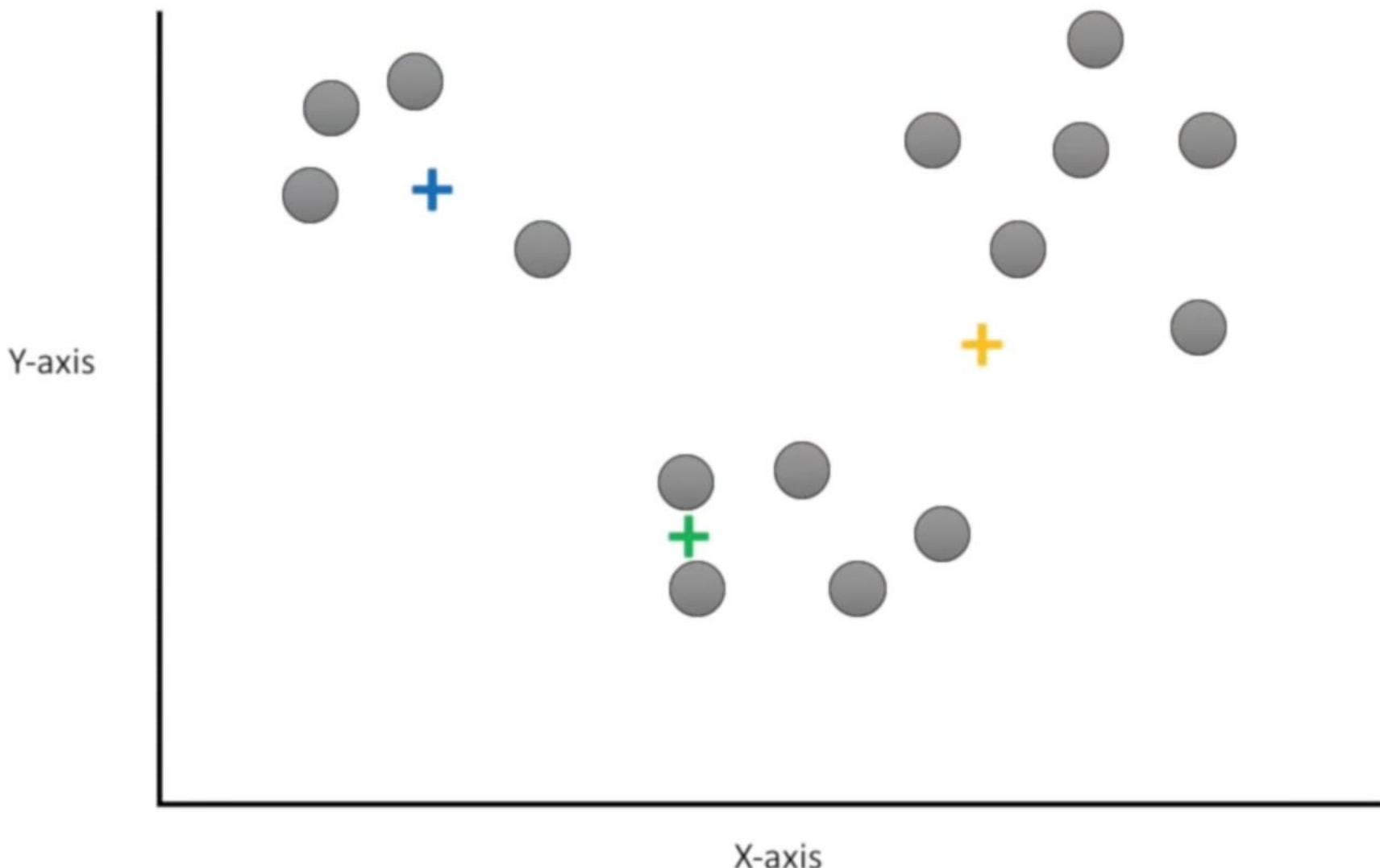
Then, just like before, we assign the point to the nearest cluster.



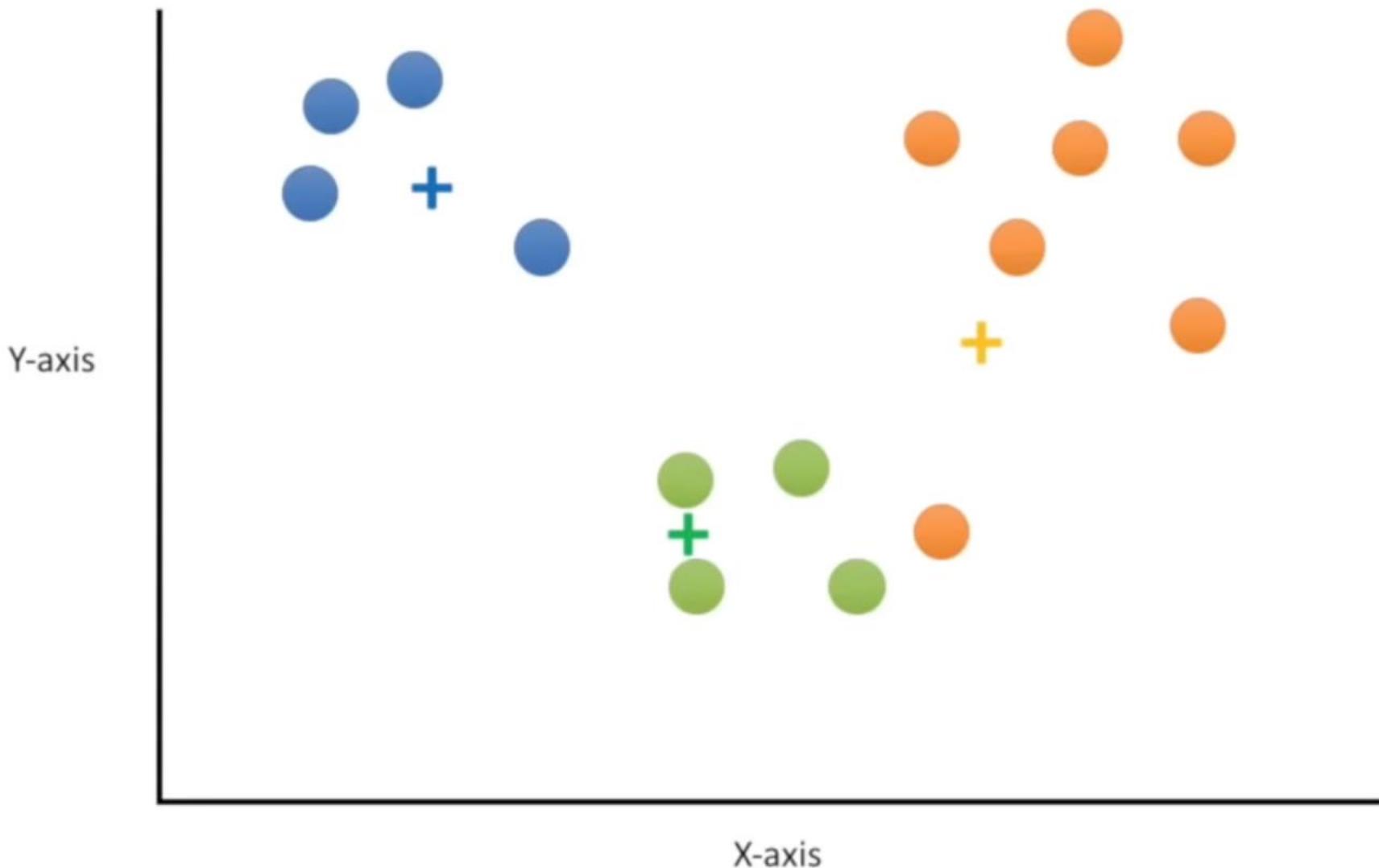
And, just like before, we then calculate the center of each cluster and recluster...



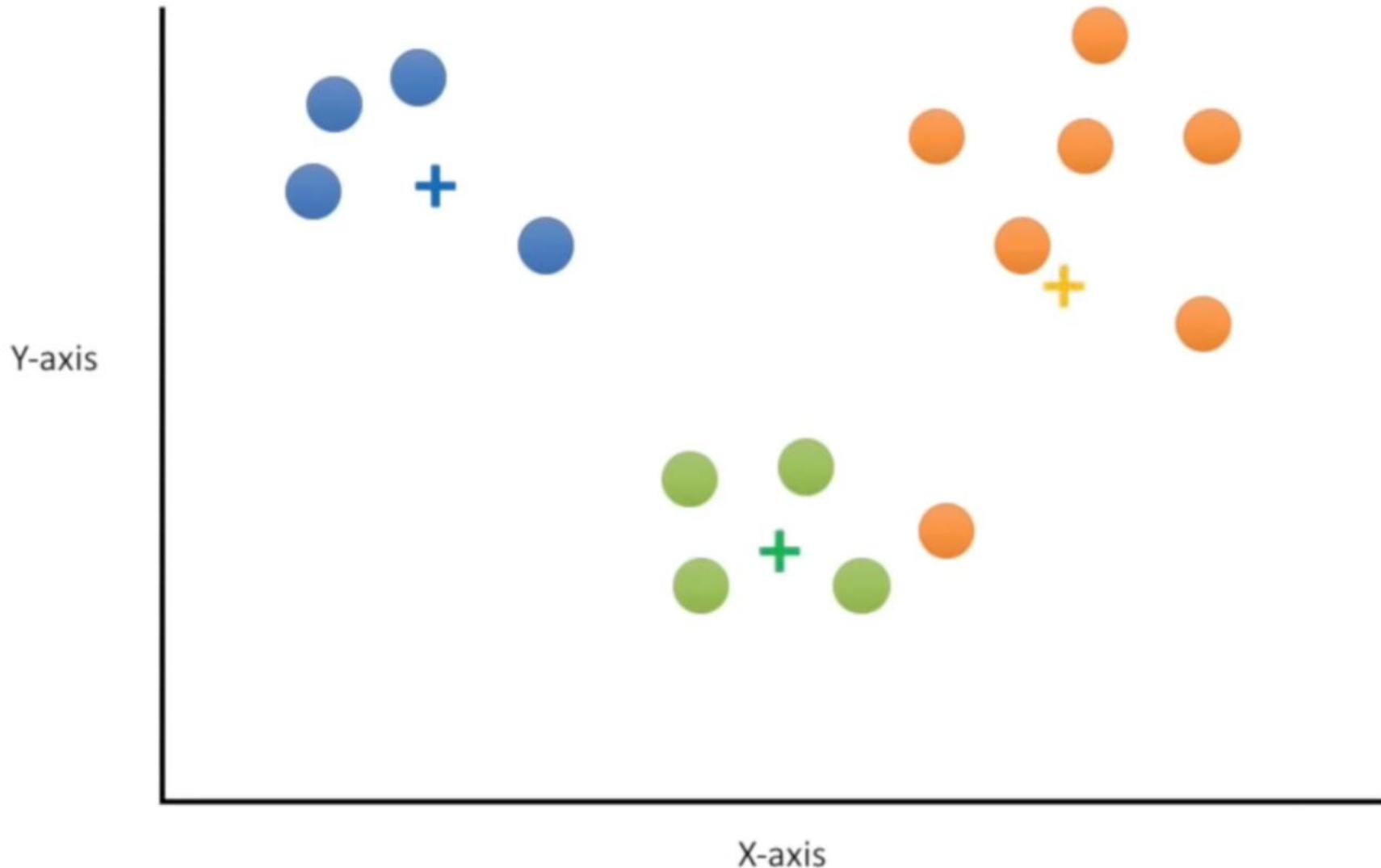
And, just like before, we then calculate the center of each cluster and recluster...



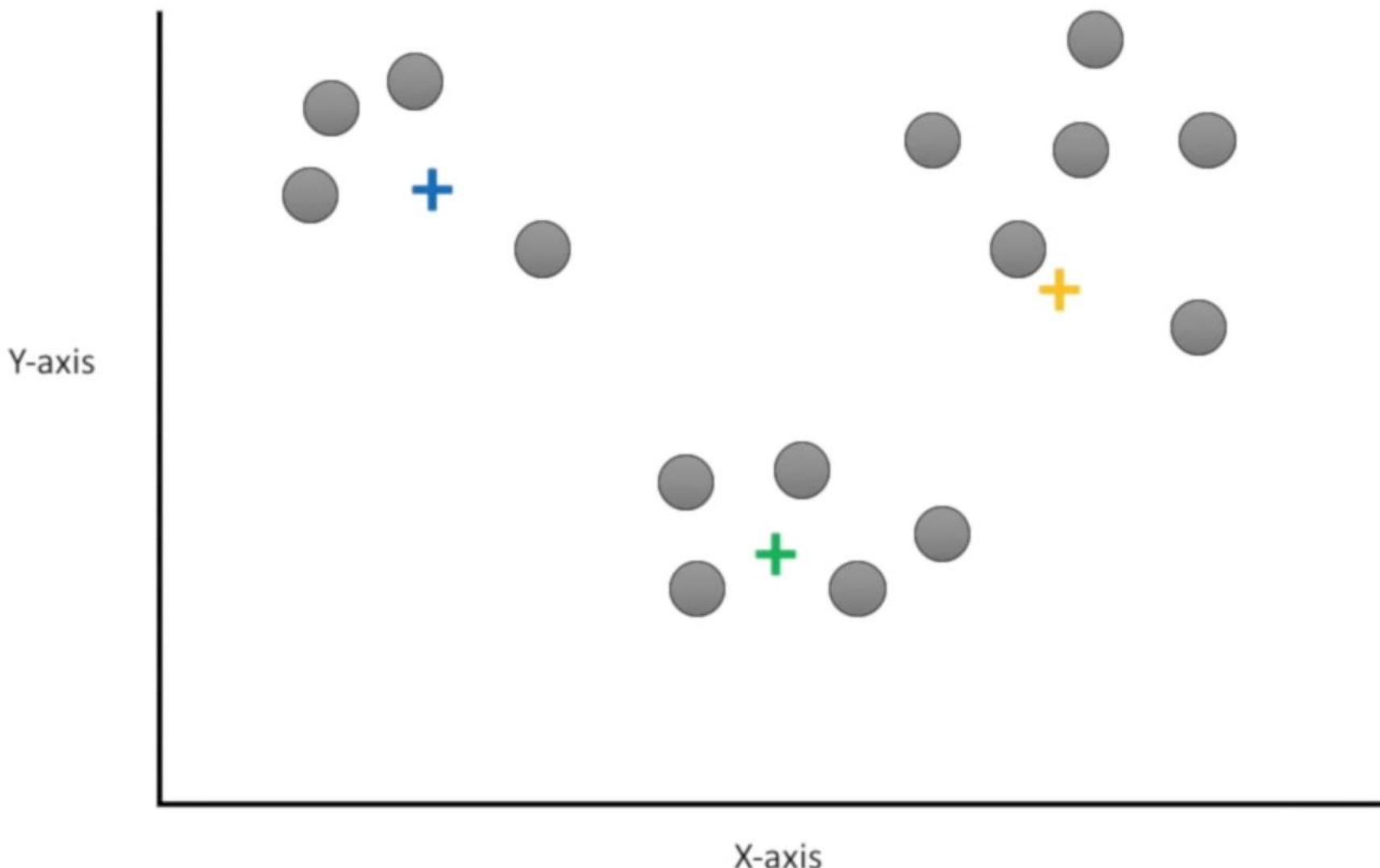
And, just like before, we then calculate the center of each cluster and recluster...



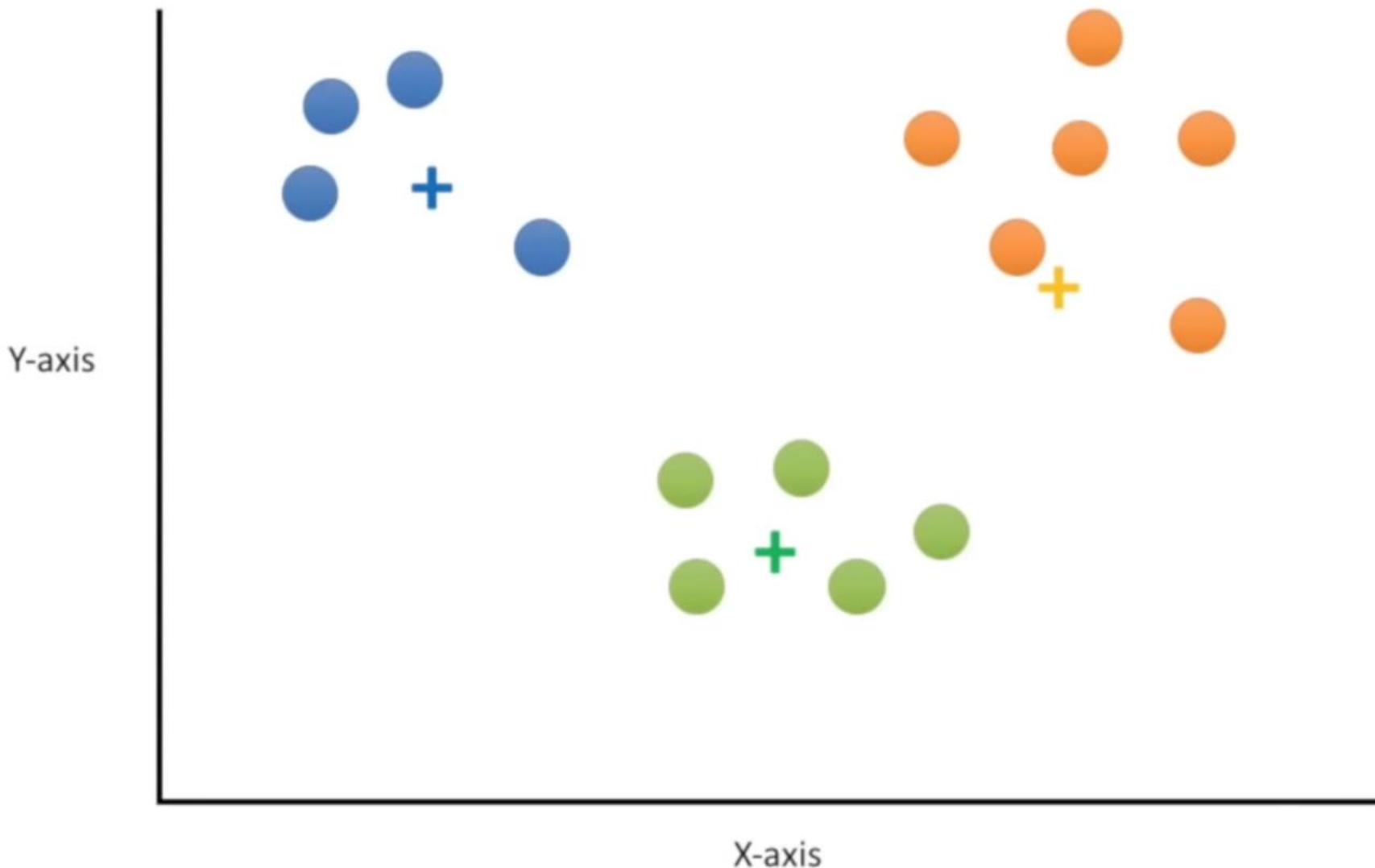
And, just like before, we then calculate the center of each cluster and recluster...



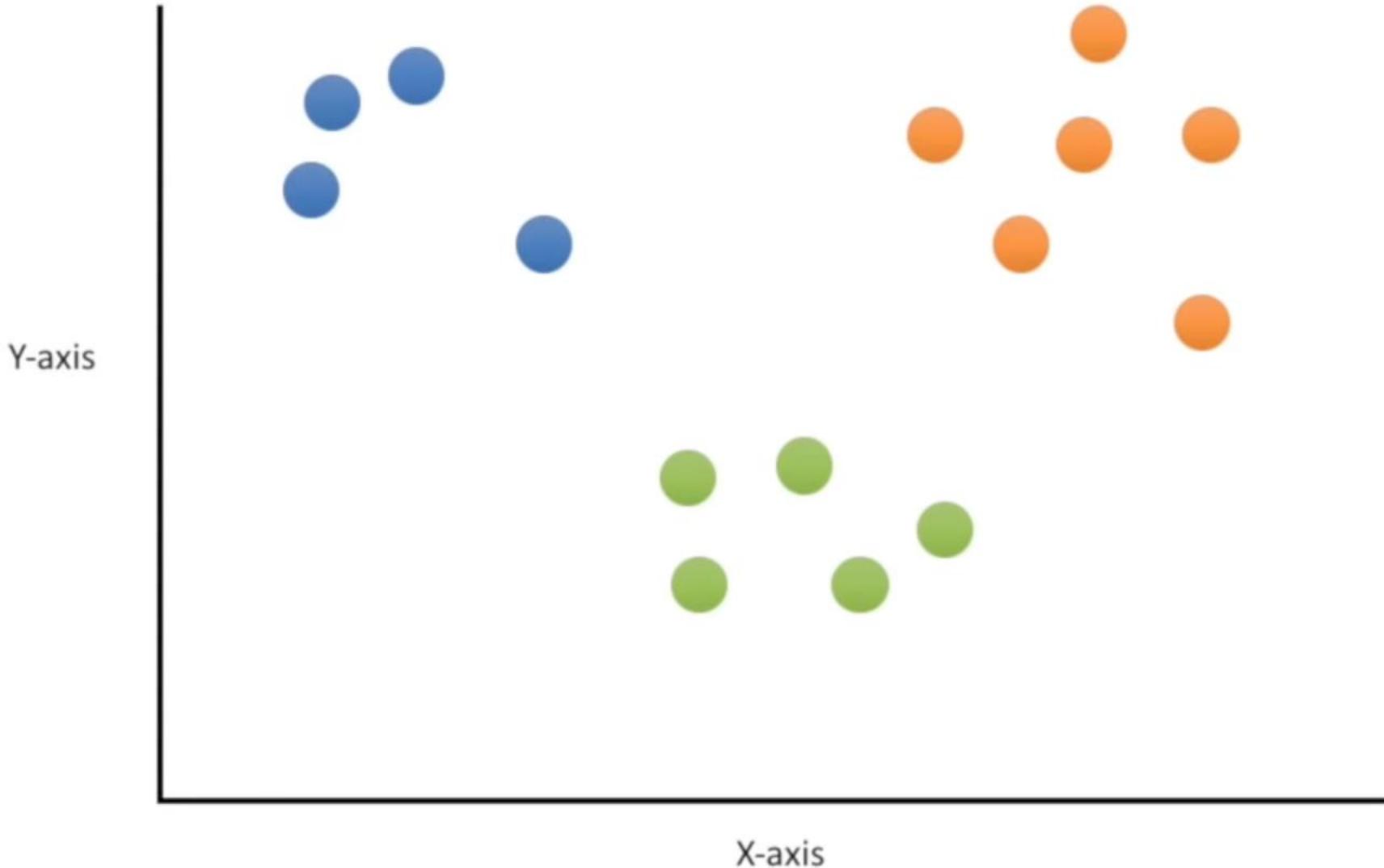
And, just like before, we then calculate the center of each cluster and recluster...



And, just like before, we then calculate the center of each cluster and recluster...



Bam? Although this looks good, the computer doesn't know that until it does the clustering a few more times.

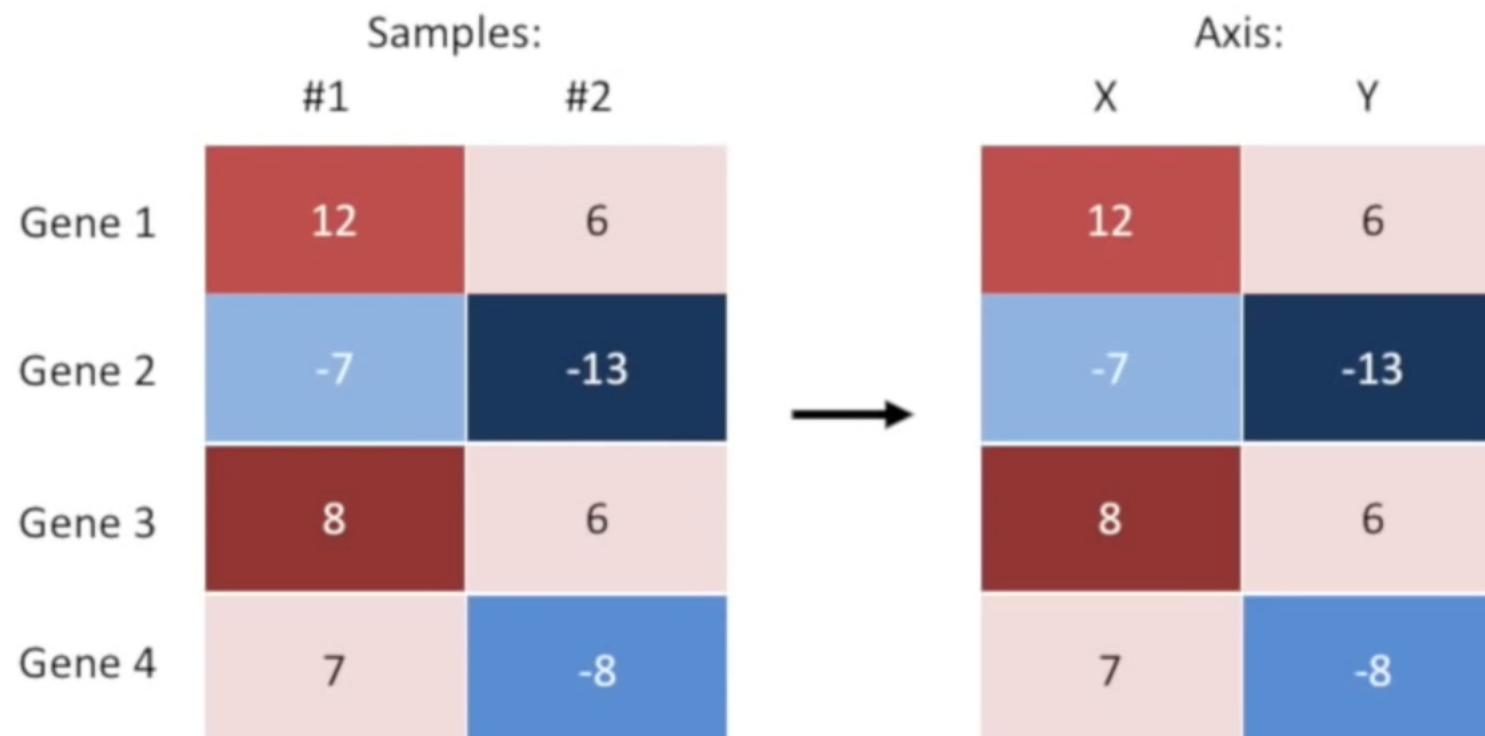


K Means for Heatmaps

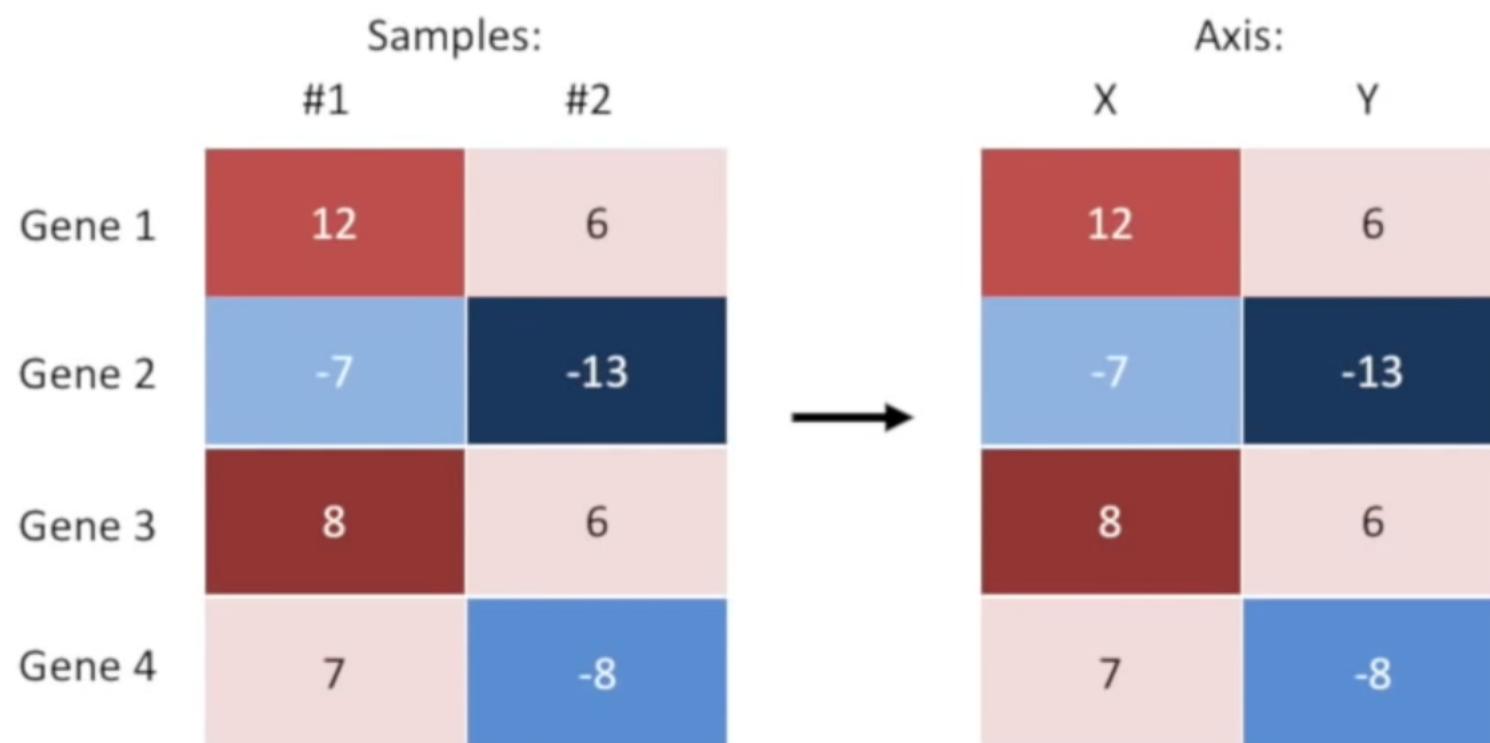
Question: What if my data is a heatmap?

Samples:		
	#1	#2
Gene 1	12	6
Gene 2	-7	-13
Gene 3	8	6
Gene 4	7	-8

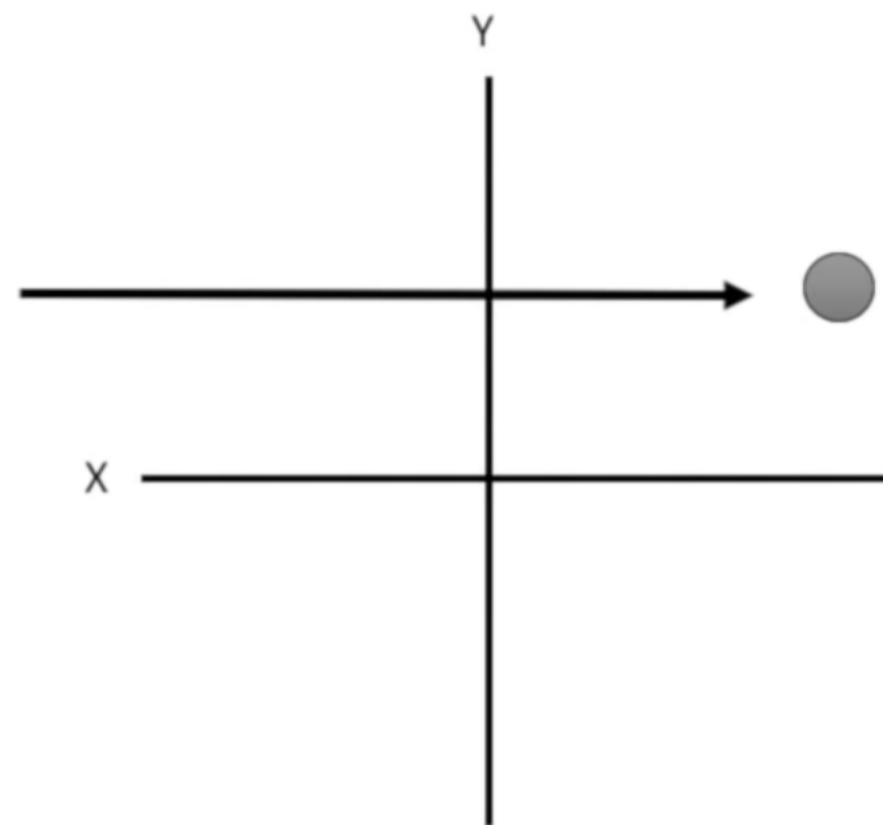
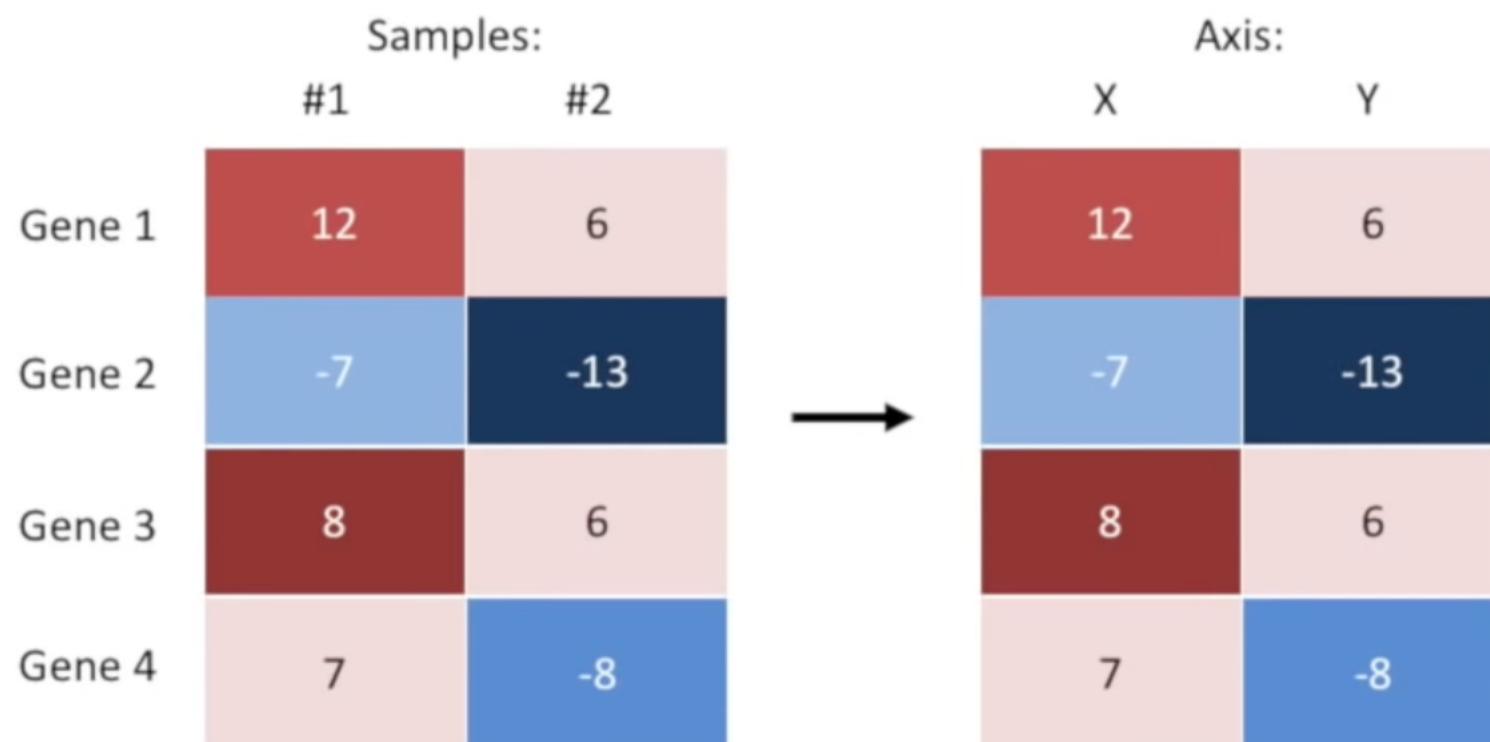
Well, if we just have 2 samples, we can rename them “X” and “Y”



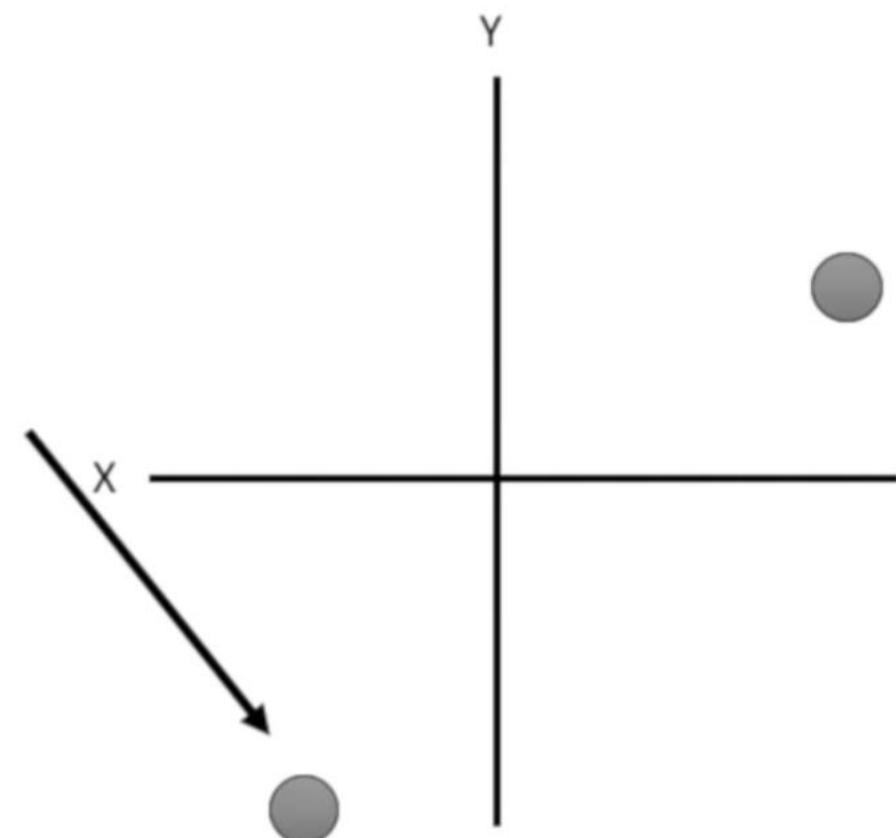
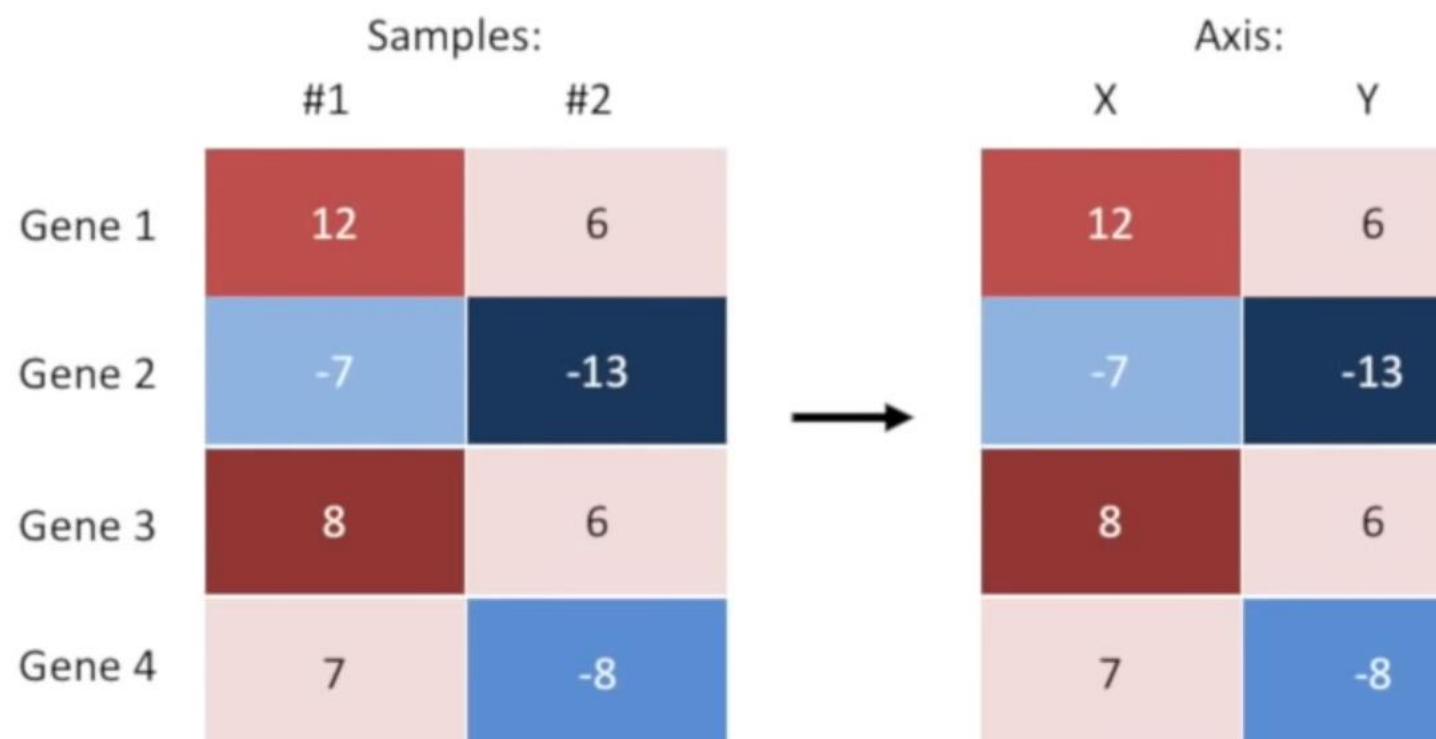
And we then plot the data in an X/Y graph



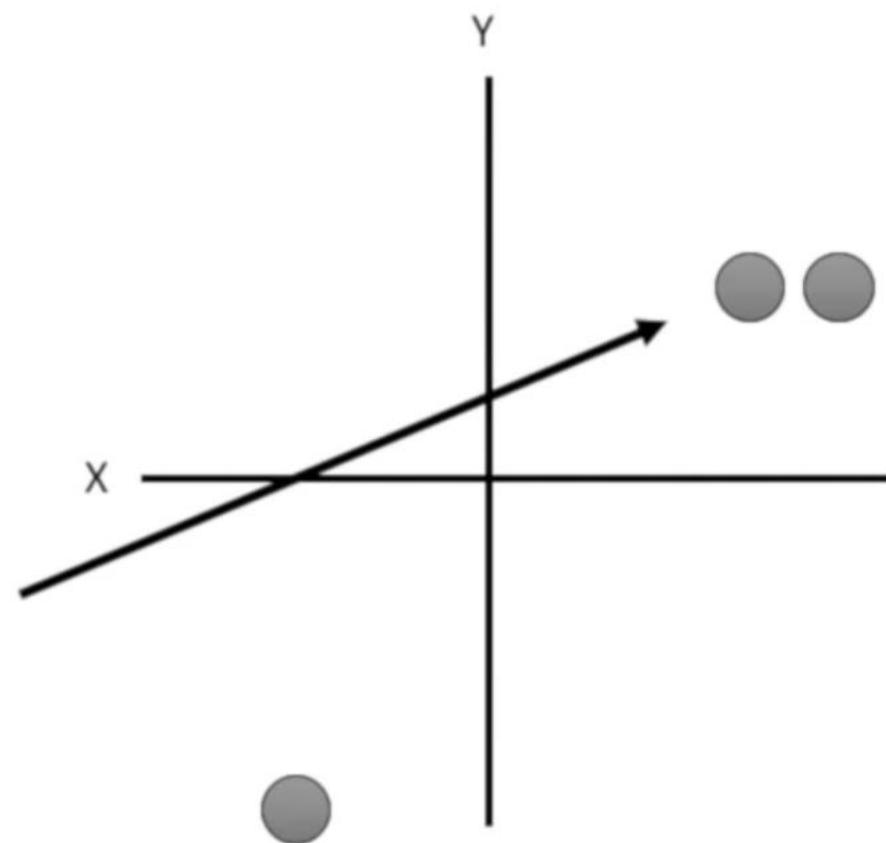
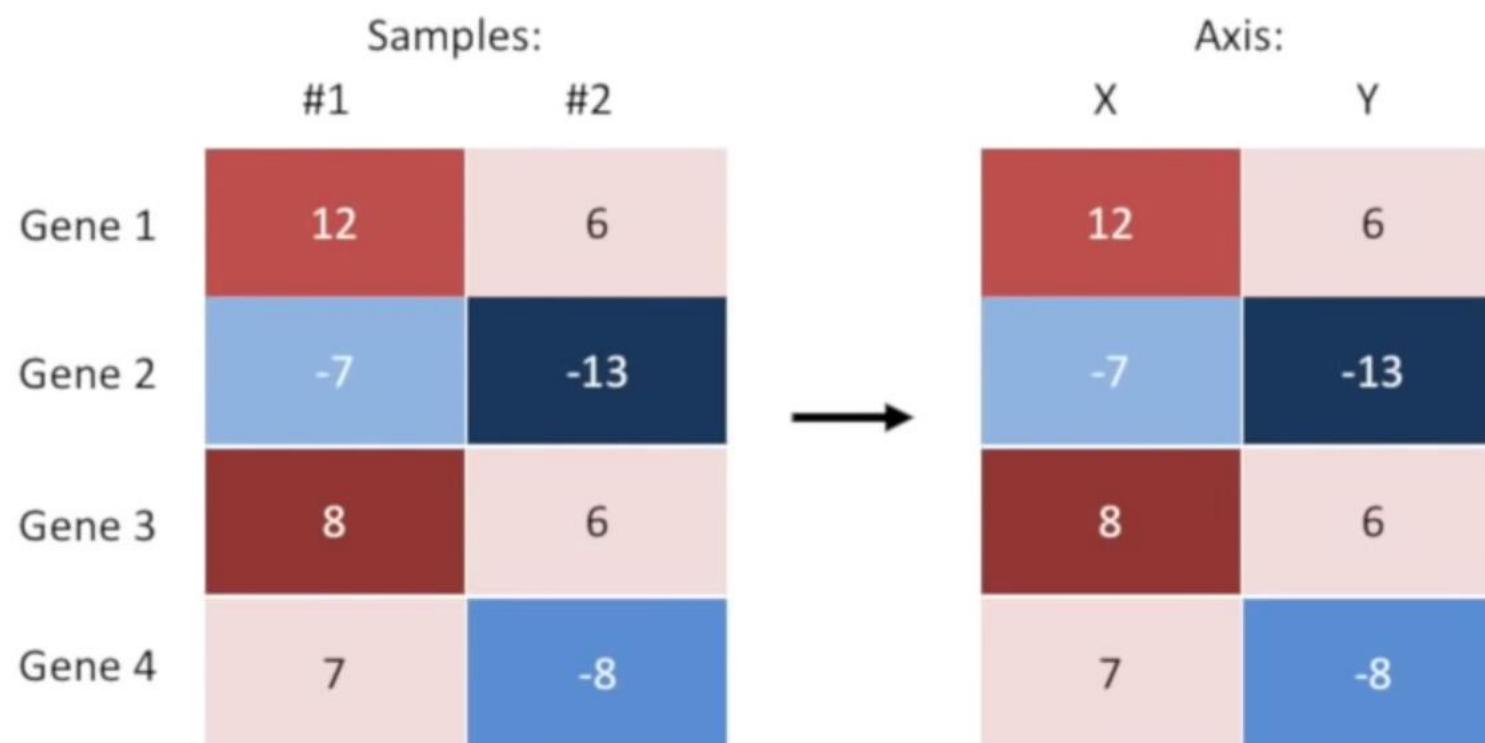
And we then plot the data in an X/Y graph



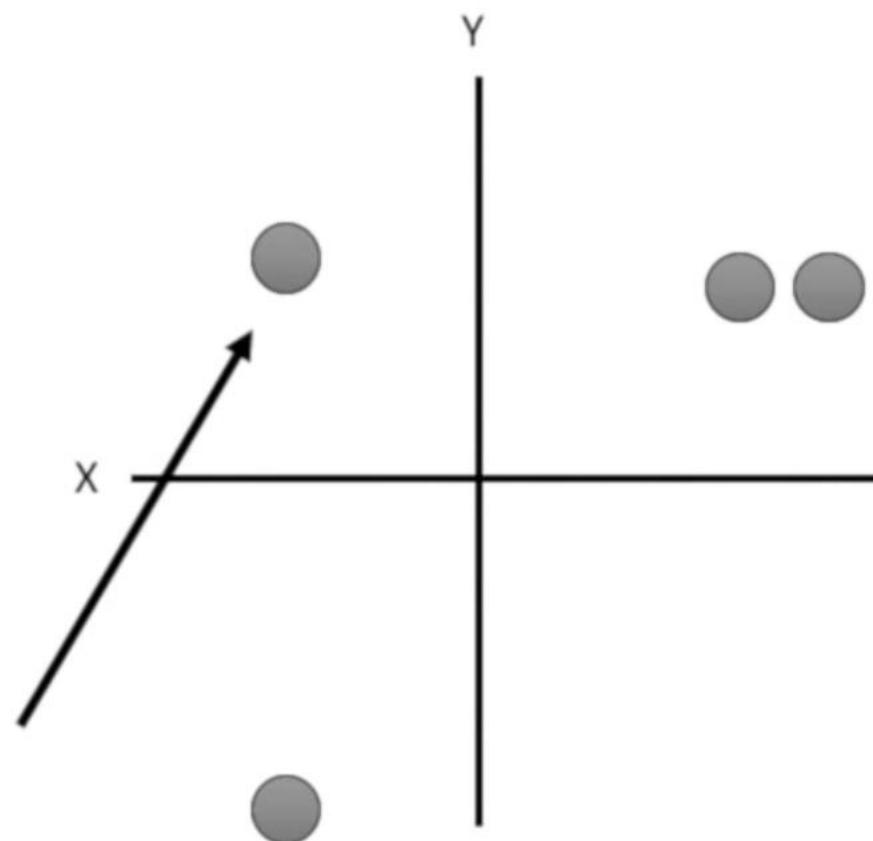
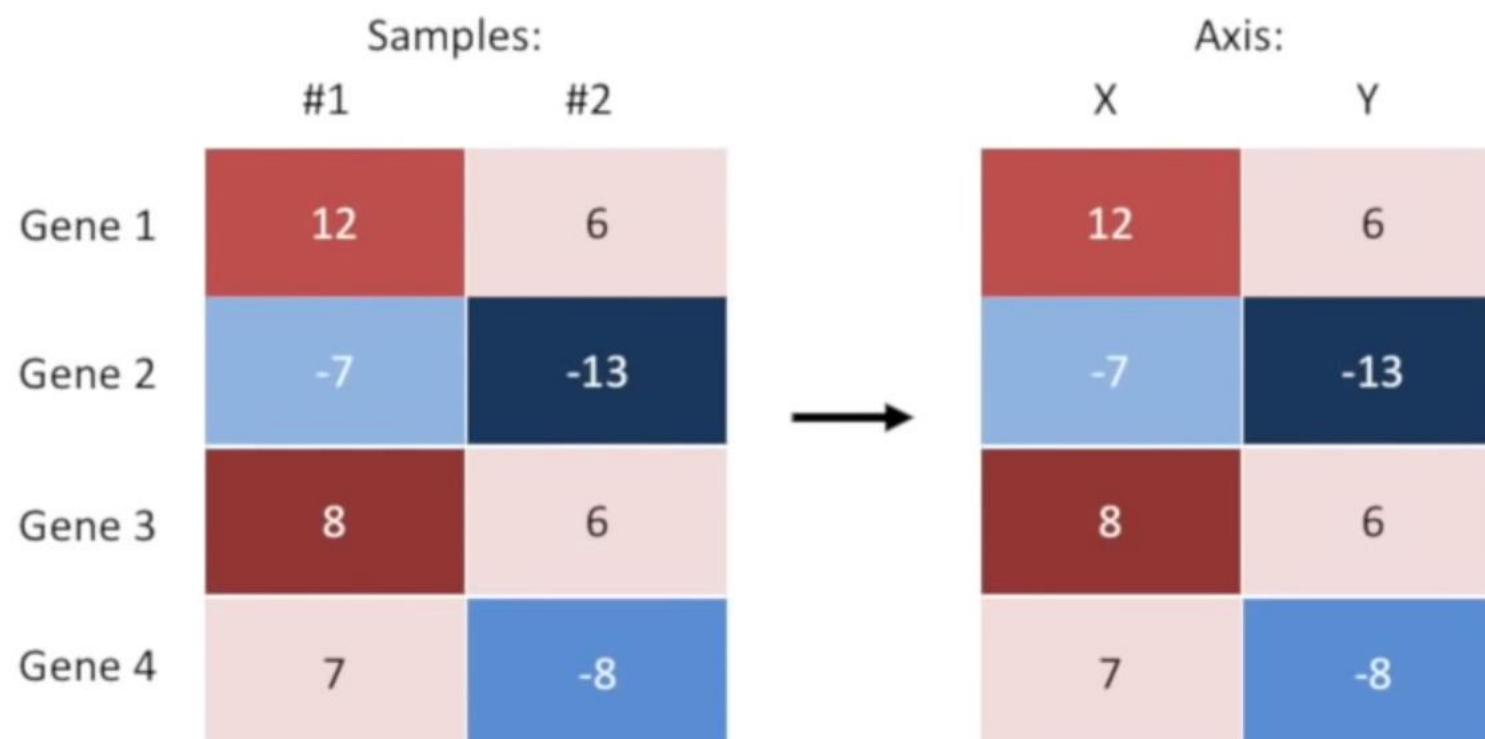
And we then plot the data in an X/Y graph



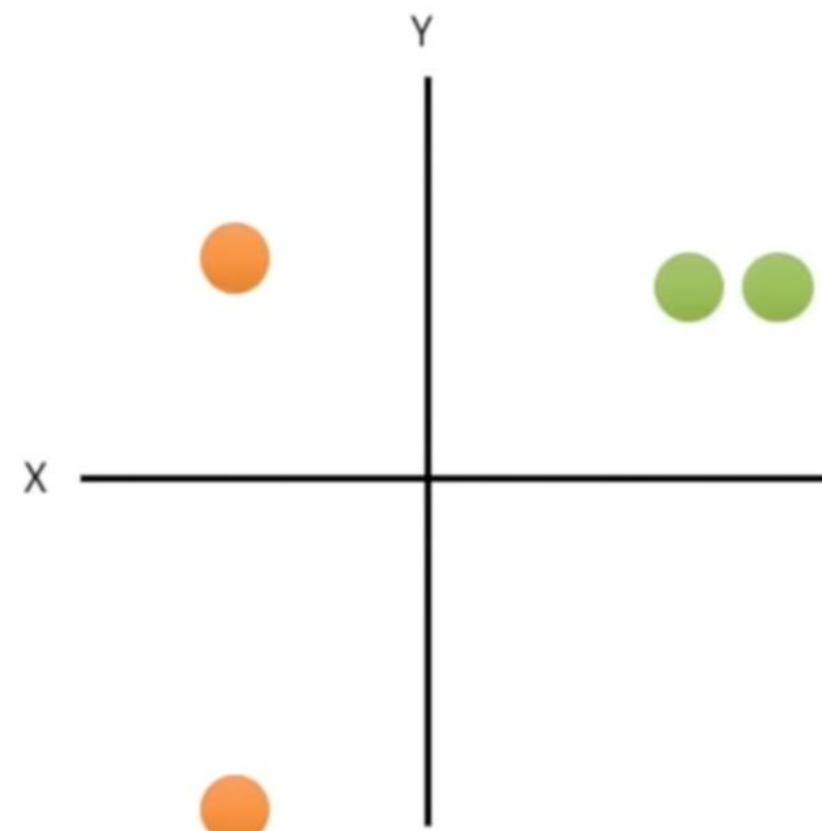
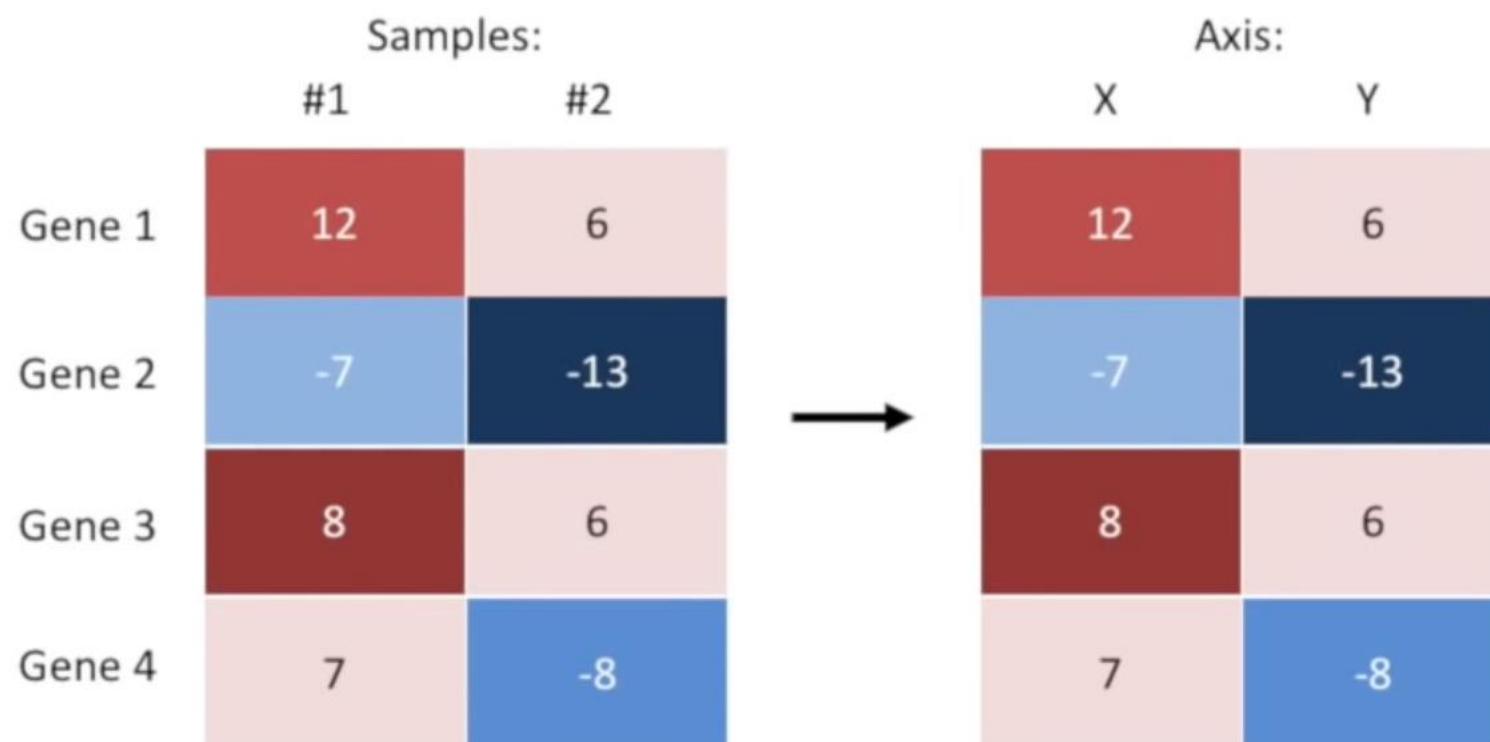
And we then plot the data in an X/Y graph

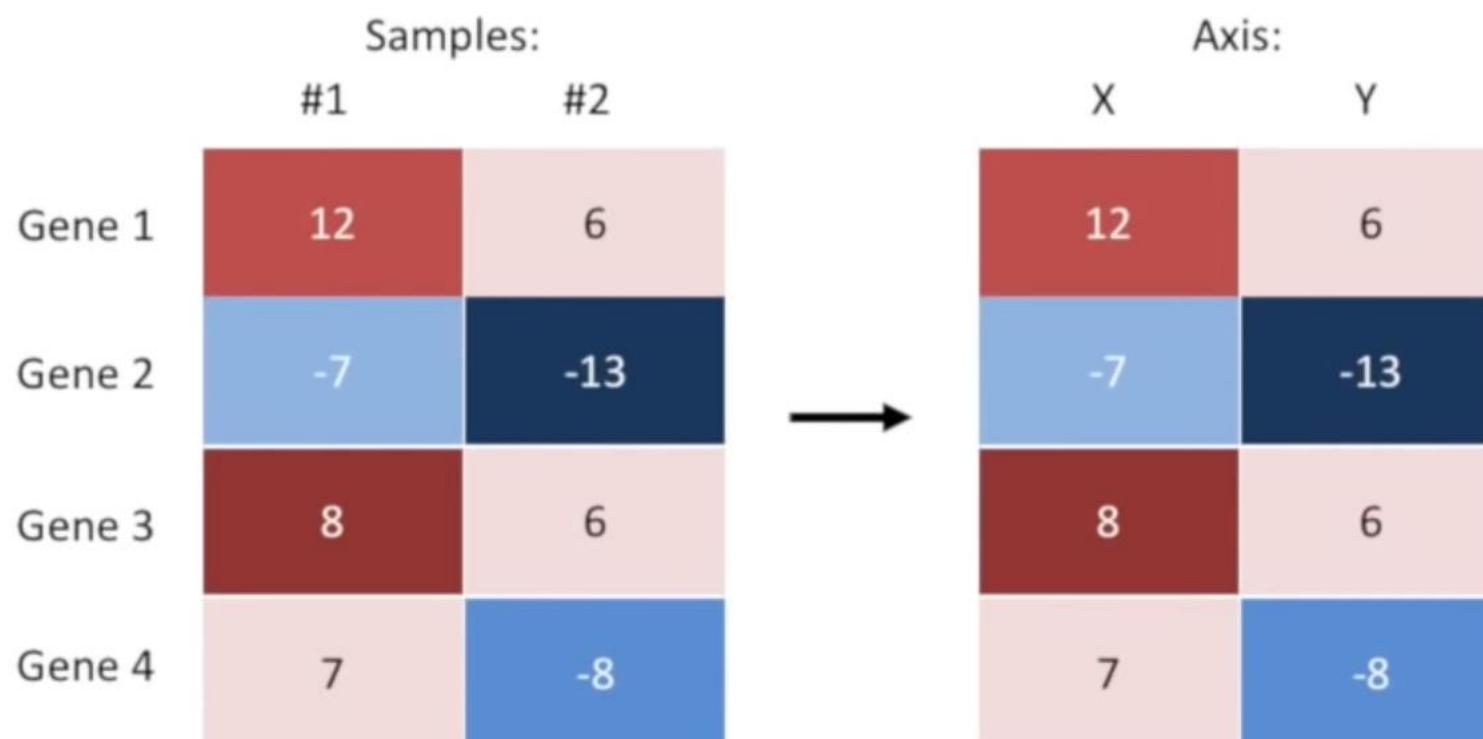


And we then plot the data in an X/Y graph

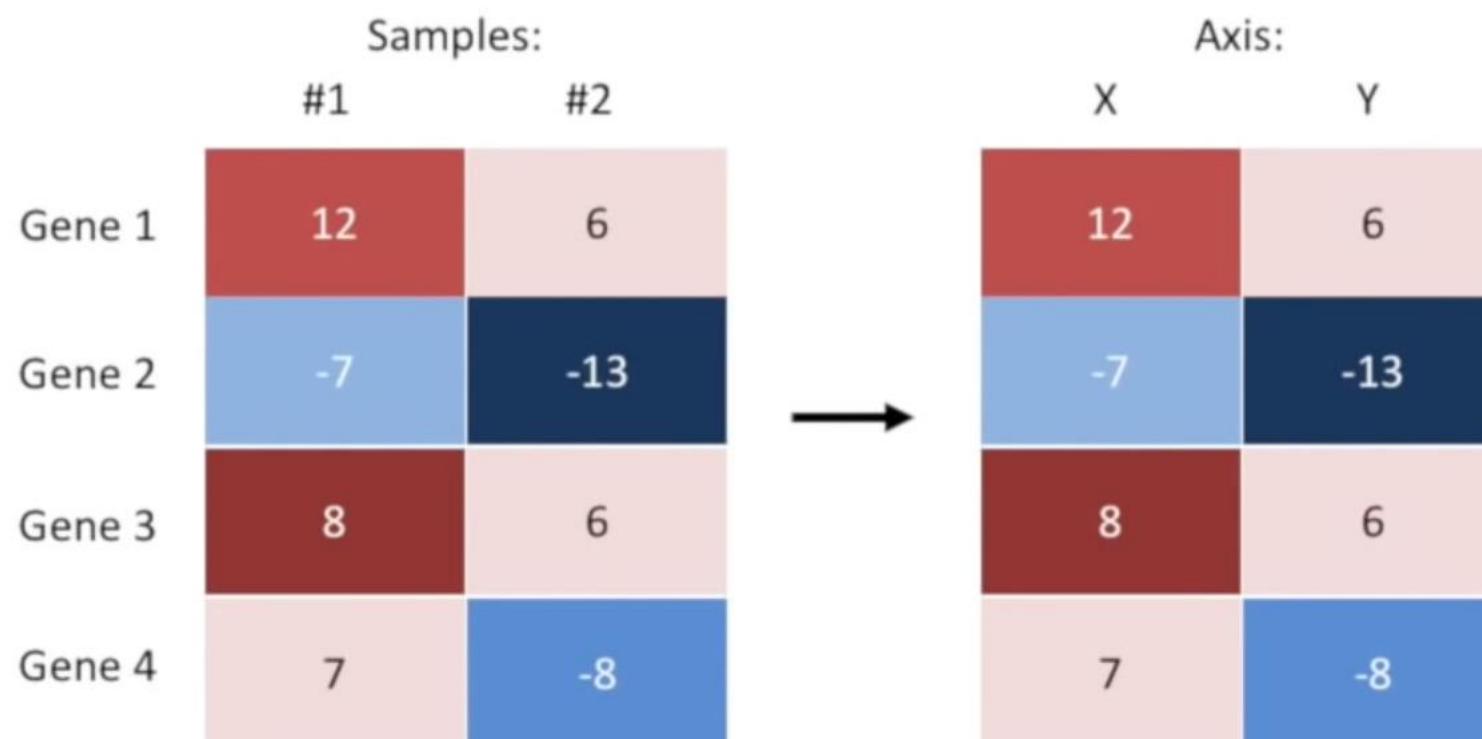


Then we can cluster just like before!





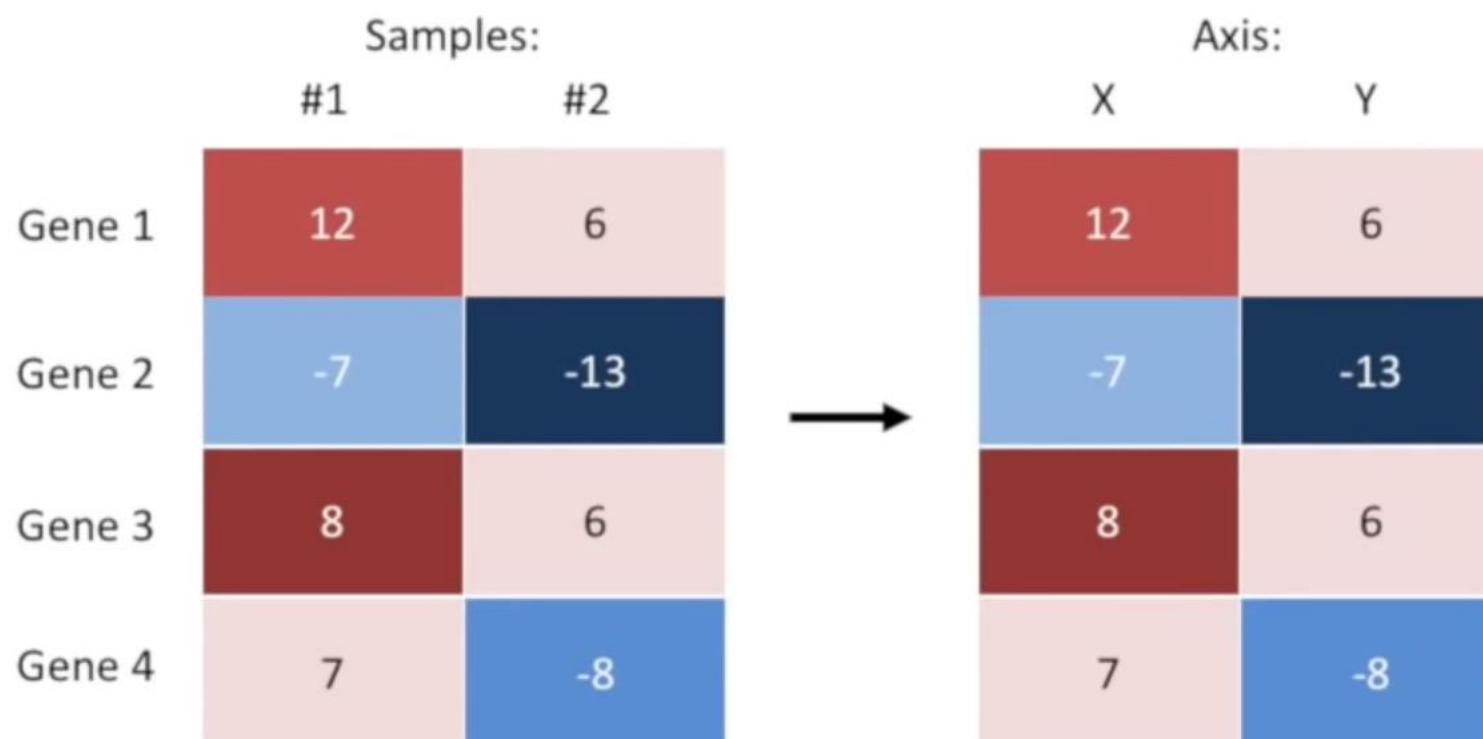
Note: We don't actually need to plot the data in order to cluster it. We just need to calculate the distances between things.



Note: We don't actually need to plot the data in order to cluster it. We just need to calculate the distances between things.

When we have 2 samples, or 2 axes, the Euclidean distance is:

$$\sqrt{x^2 + y^2}$$



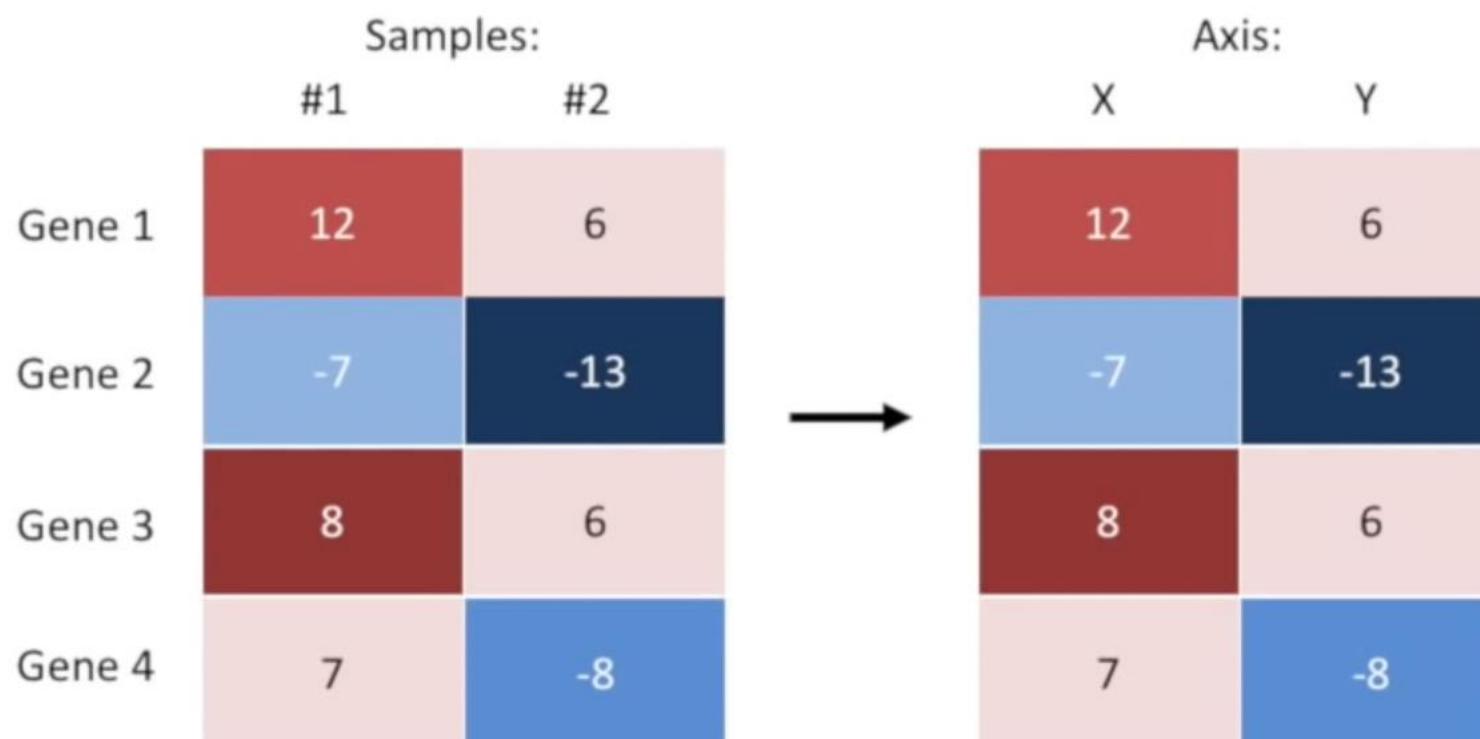
Note: We don't actually need to plot the data in order to cluster it. We just need to calculate the distances between things.

When we have 2 samples, or 2 axes, the Euclidean distance is:

$$\sqrt{x^2 + y^2}$$

When we have 3 samples, or 3 axes, the Euclidean distance is:

$$\sqrt{x^2 + y^2 + z^2}$$



Note: We don't actually need to plot the data in order to cluster it. We just need to calculate the distances between things.

When we have 2 samples, or 2 axes, the Euclidean distance is:

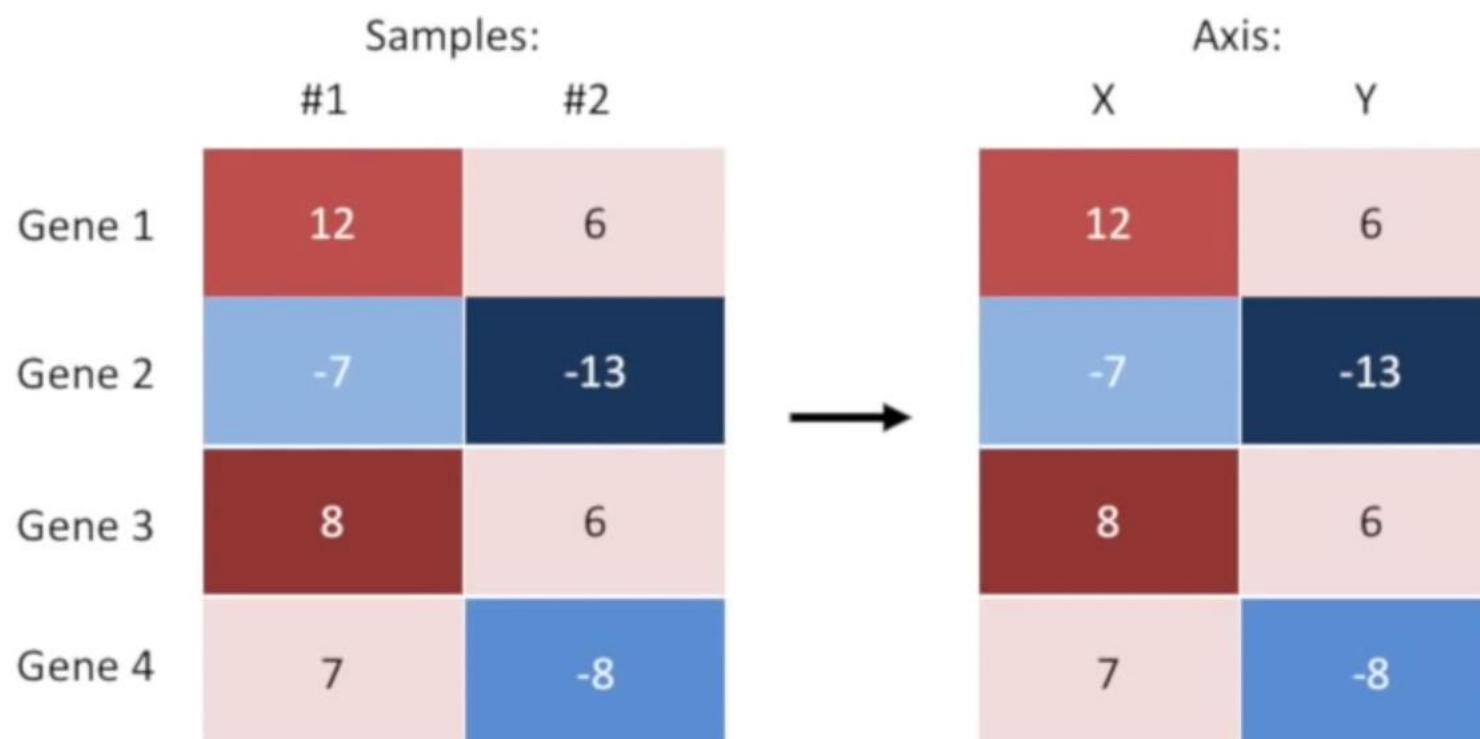
$$\sqrt{x^2 + y^2}$$

When we have 3 samples, or 3 axes, the Euclidean distance is:

$$\sqrt{x^2 + y^2 + z^2}$$

When we have 4 samples, or 4 axes, the Euclidean distance is:

$$\sqrt{x^2 + y^2 + z^2 + a^2}$$



Note: We don't actually need to plot the data in order to cluster it. We just need to calculate the distances between things.

When we have 2 samples, or 2 axes, the Euclidean distance is:

$$\sqrt{x^2 + y^2}$$

When we have 3 samples, or 3 axes, the Euclidean distance is:

$$\sqrt{x^2 + y^2 + z^2}$$

When we have 4 samples, or 4 axes, the Euclidean distance is:

$$\sqrt{x^2 + y^2 + z^2 + a^2}$$

etc. etc. etc.



Machine Learning

K-Means

Phd. César Astudillo | Facultad de Ingeniería