



САНКТ-ПЕТЕРБУРГСКИЙ  
ГОСУДАРСТВЕННЫЙ  
ЭКОНОМИЧЕСКИЙ  
УНИВЕРСИТЕТ

# Метод главных компонент (РСА)

Занятие 4

Глазунова Е.В.

# PCA

	Mouse 1	Mouse 2	Mouse 3	Mouse 4	Mouse 5	Mouse 6
Gene 1	10	11	8	3	1	2
Gene 2	6	4	5	3	2.8	1

# PCA

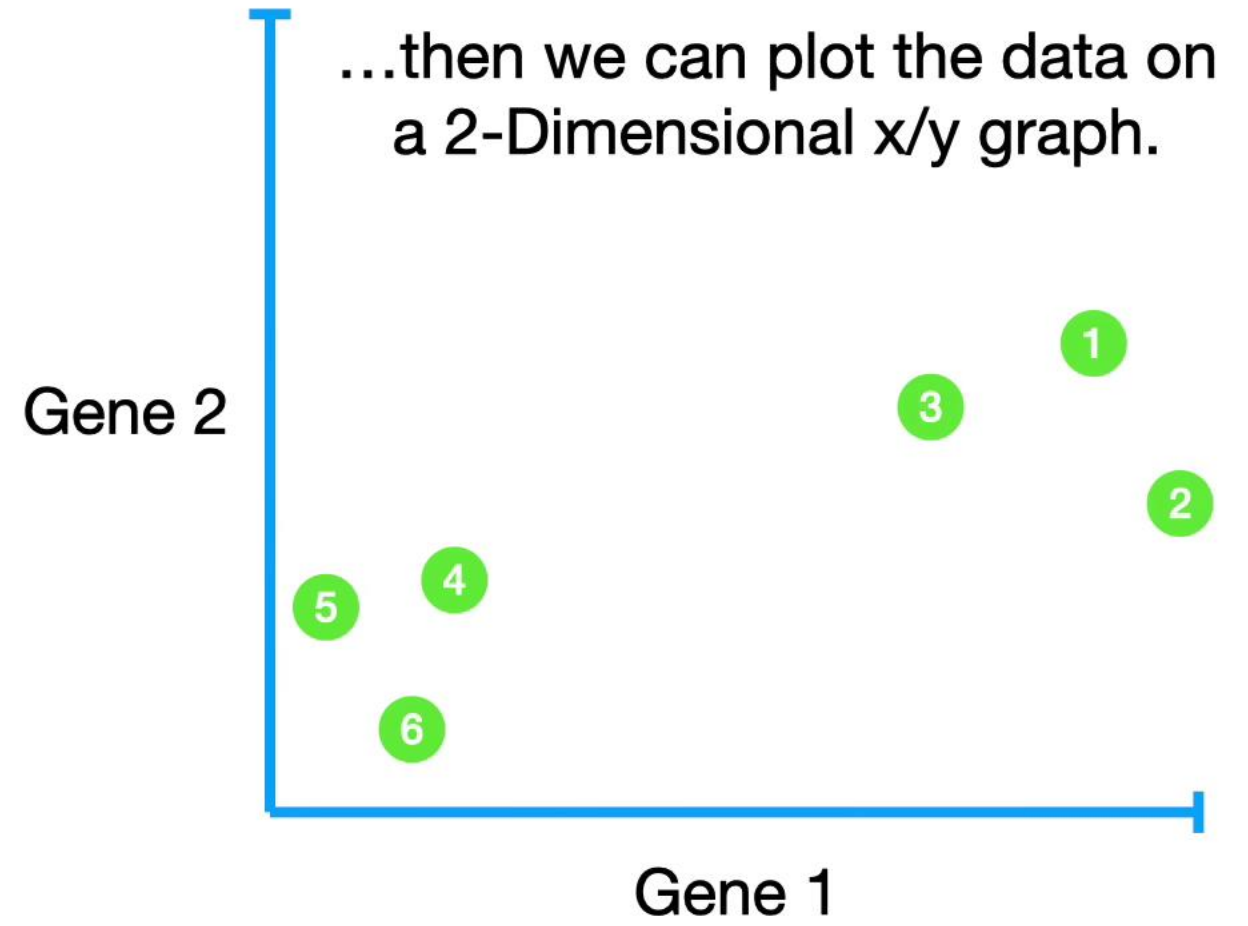
	Mouse 1	Mouse 2	Mouse 3	Mouse 4	Mouse 5	Mouse 6
Gene 1	10	11	8	3	1	2

If we only measure 1 gene,  
we can plot the data on a  
number line...



# PCA

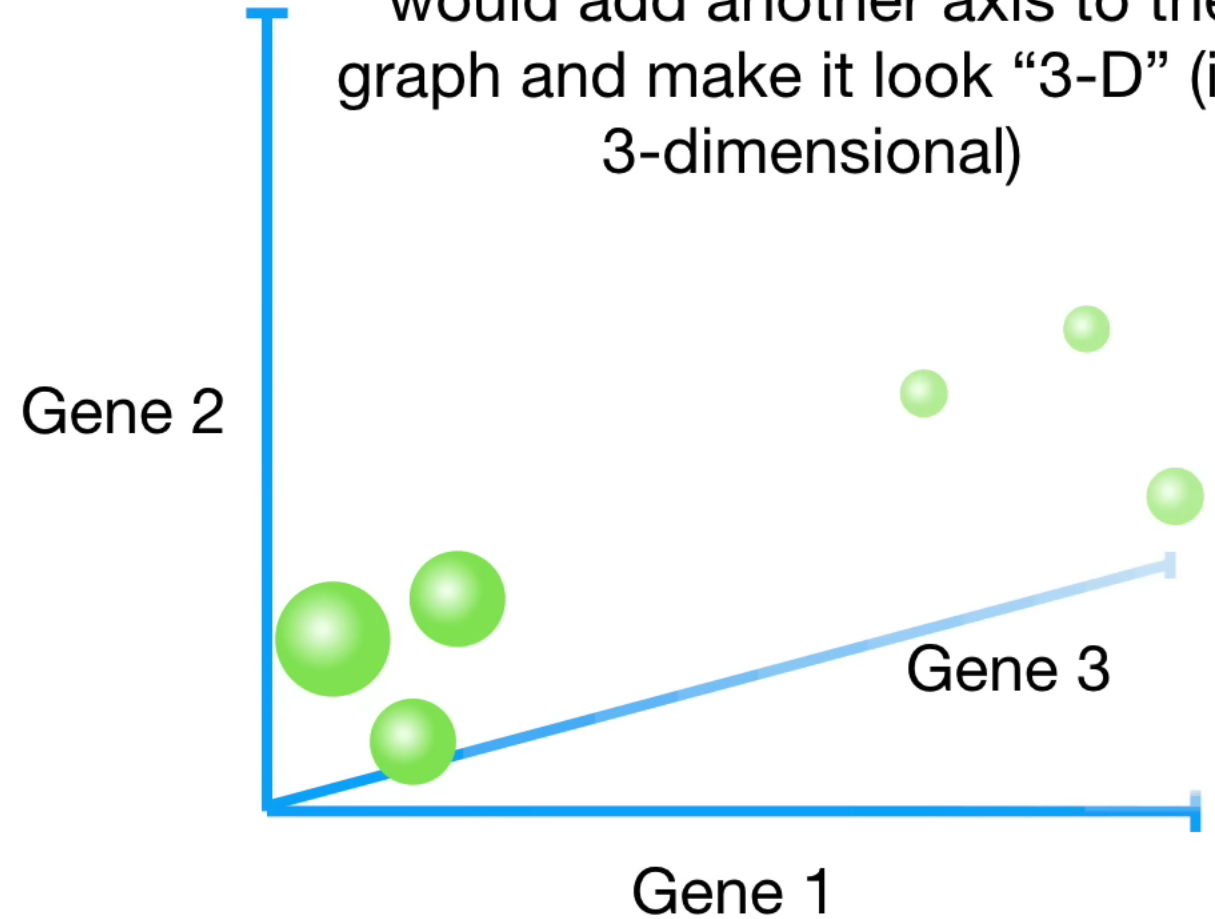
	Mouse 1	Mouse 2	Mouse 3	Mouse 4	Mouse 5	Mouse 6
Gene 1	10	11	8	3	1	2
Gene 2	6	4	5	3	2.8	1



# PCA

	Mouse 1	Mouse 2	Mouse 3	Mouse 4	Mouse 5	Mouse 6
Gene 1	10	11	8	3	2	1
Gene 2	6	4	5	3	2.8	1
Gene 3	12	9	10	2.5	1.3	2

If we measured 3 genes, we would add another axis to the graph and make it look “3-D” (i.e. 3-dimensional)



# PCA

	Mouse 1	Mouse 2	Mouse 3	Mouse 4	Mouse 5	Mouse 6
Gene 1	10	11	8	3	2	1
Gene 2	6	4	5	3	2.8	1
Gene 3	12	9	10	2.5	1.3	2
Gene 4	5	7	6	2	4	7

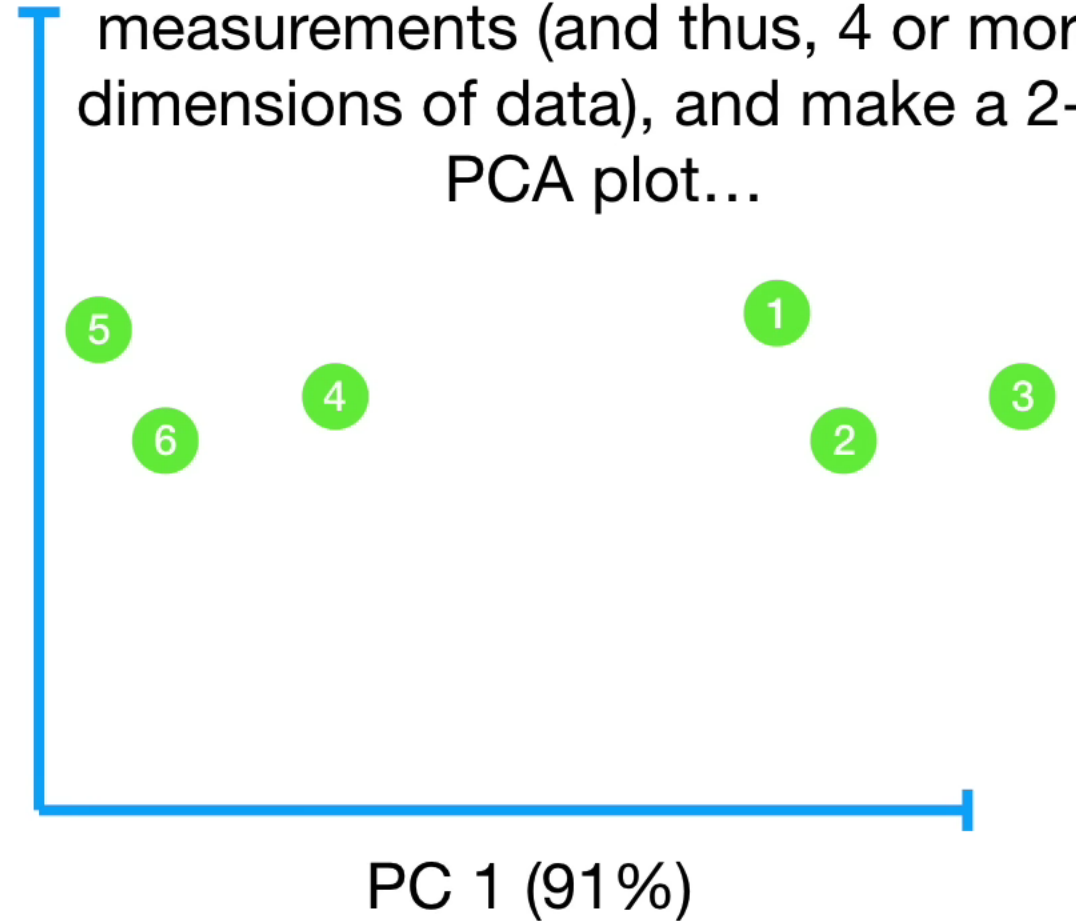
If we measured 4 genes,  
however, we can no longer  
plot the data - 4 genes require  
4 dimensions.

# PCA

	Mouse 1	Mouse 2	Mouse 3	Mouse 4	Mouse 5	Mouse 6
Gene 1	10	11	8	3	2	1
Gene 2	6	4	5	3	2.8	1
Gene 3	12	9	10	2.5	1.3	2
Gene 4	5	7	6	2	4	7

So we're going to talk about how PCA can take 4 or more gene measurements (and thus, 4 or more dimensions of data), and make a 2-D PCA plot...

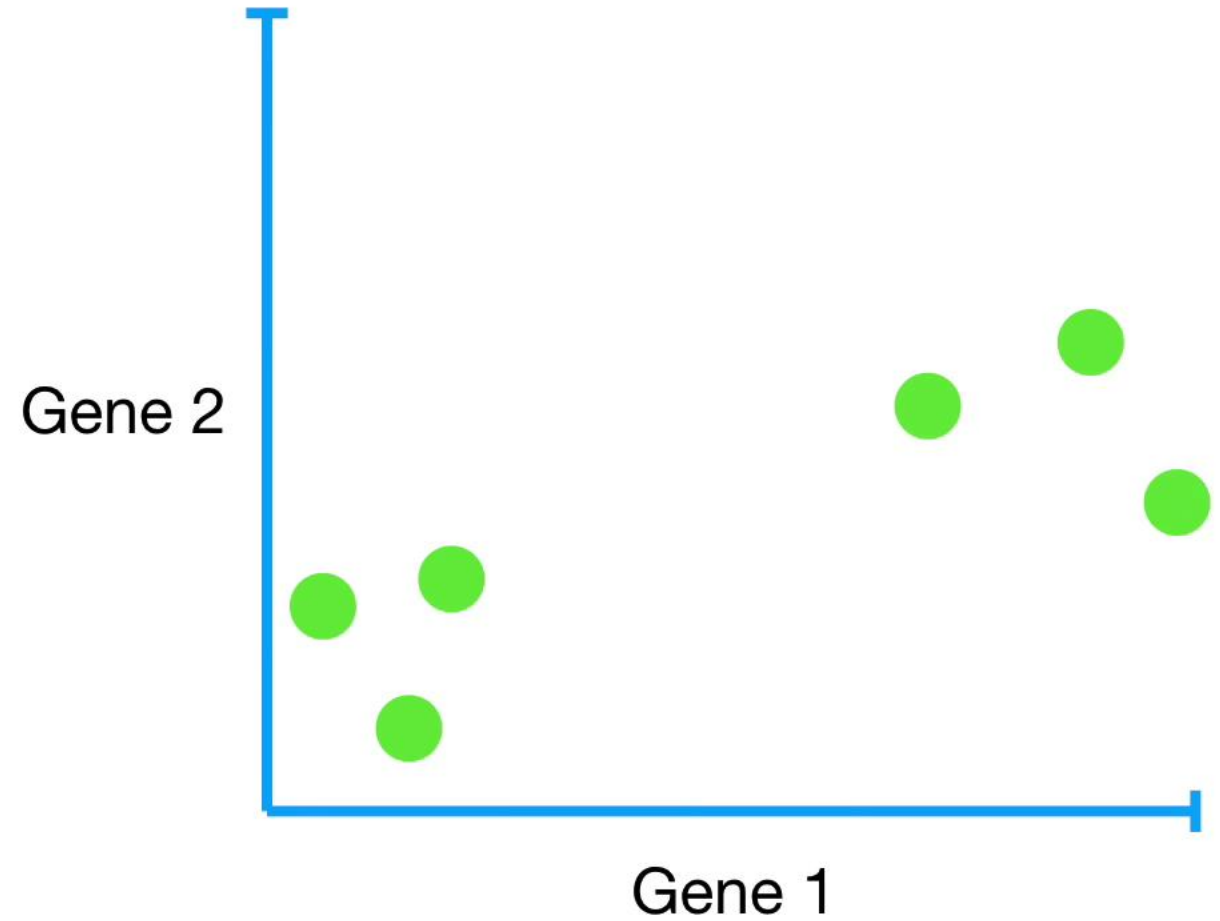
PC 2  
(4%)



# PCA

	Mouse 1	Mouse 2	Mouse 3	Mouse 4	Mouse 5	Mouse 6
Gene 1	10	11	8	3	2	1
Gene 2	6	4	5	3	2.8	1

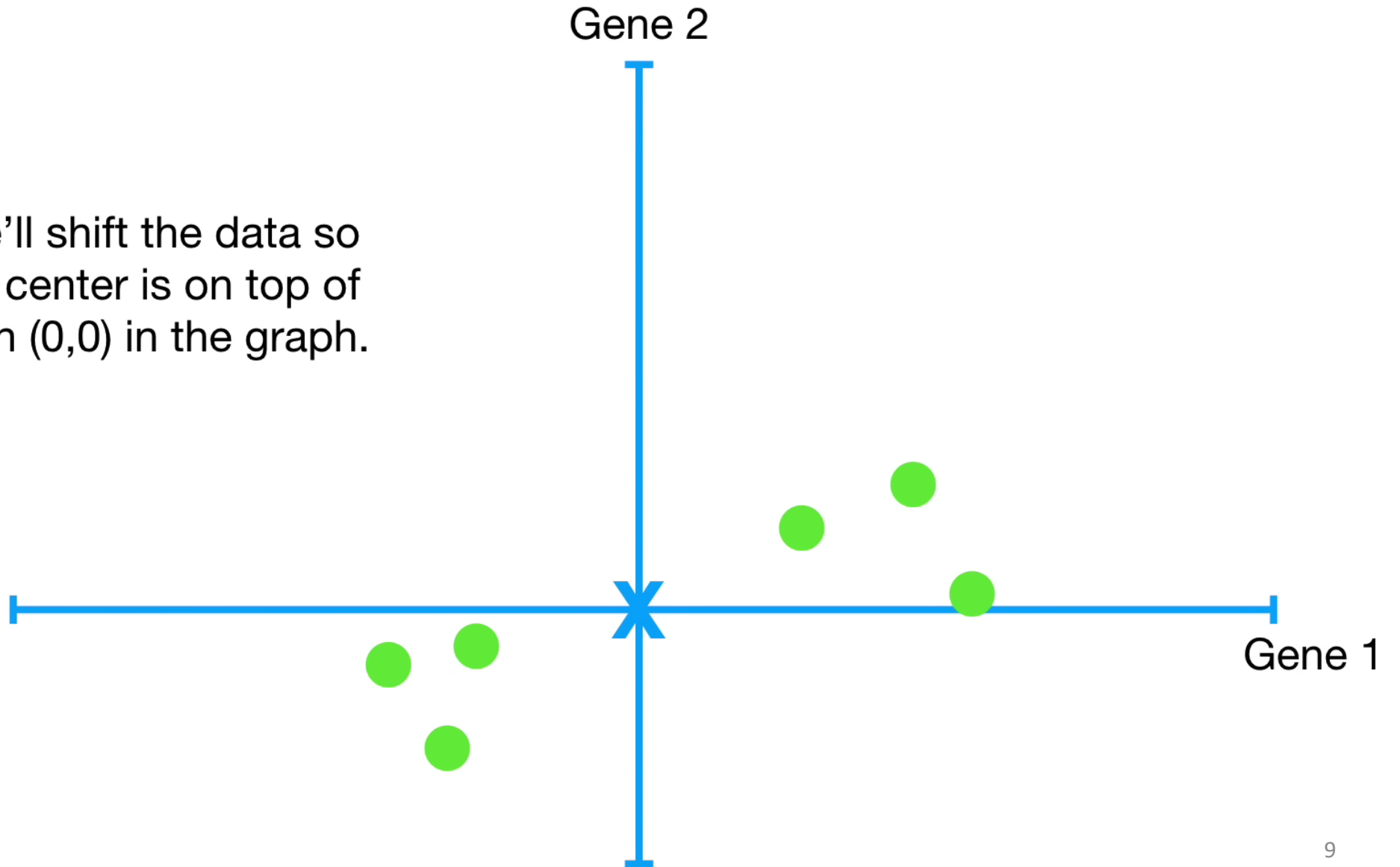
We'll start by plotting the data...





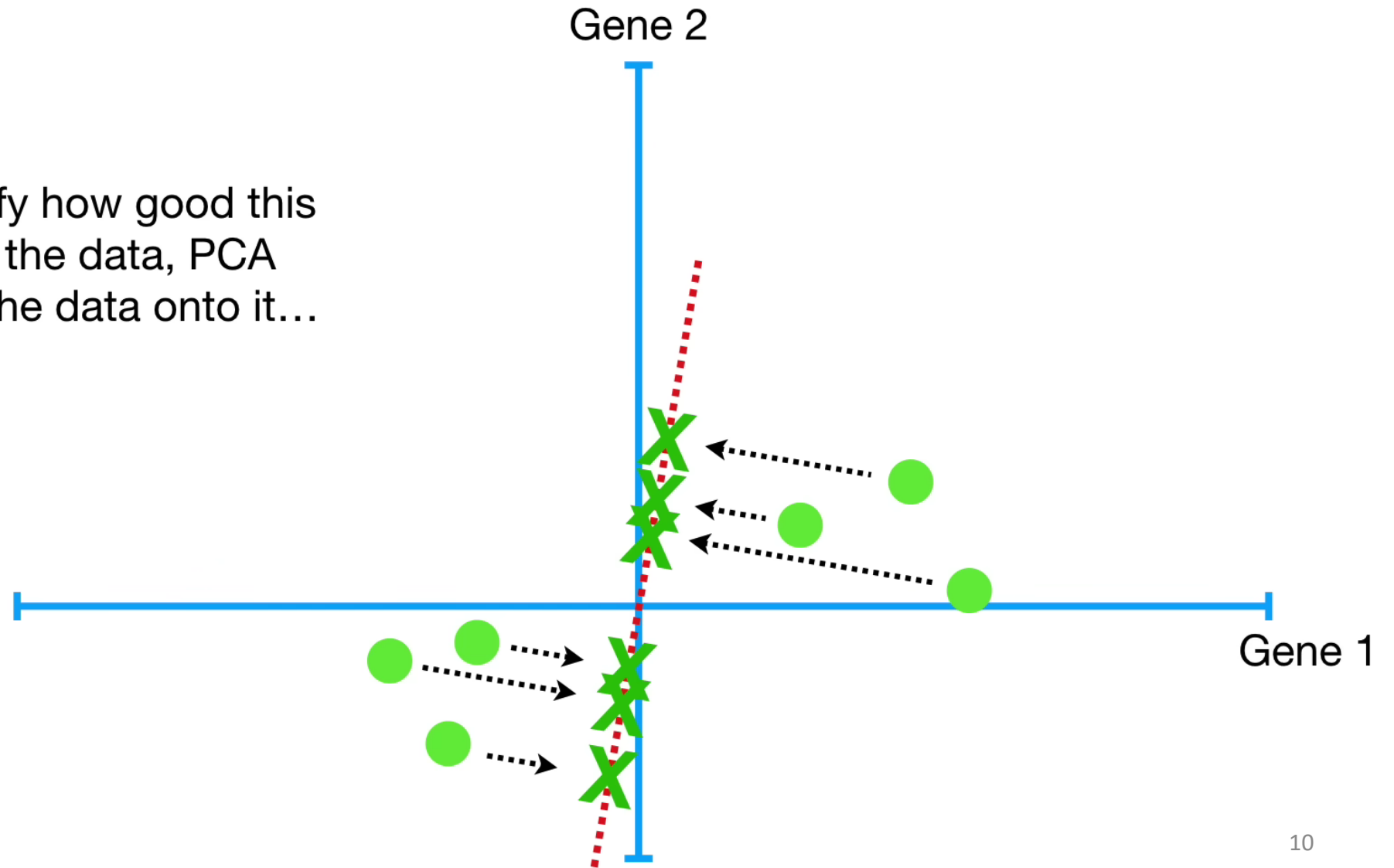
# PCA

Now we'll shift the data so that the center is on top of the origin (0,0) in the graph.

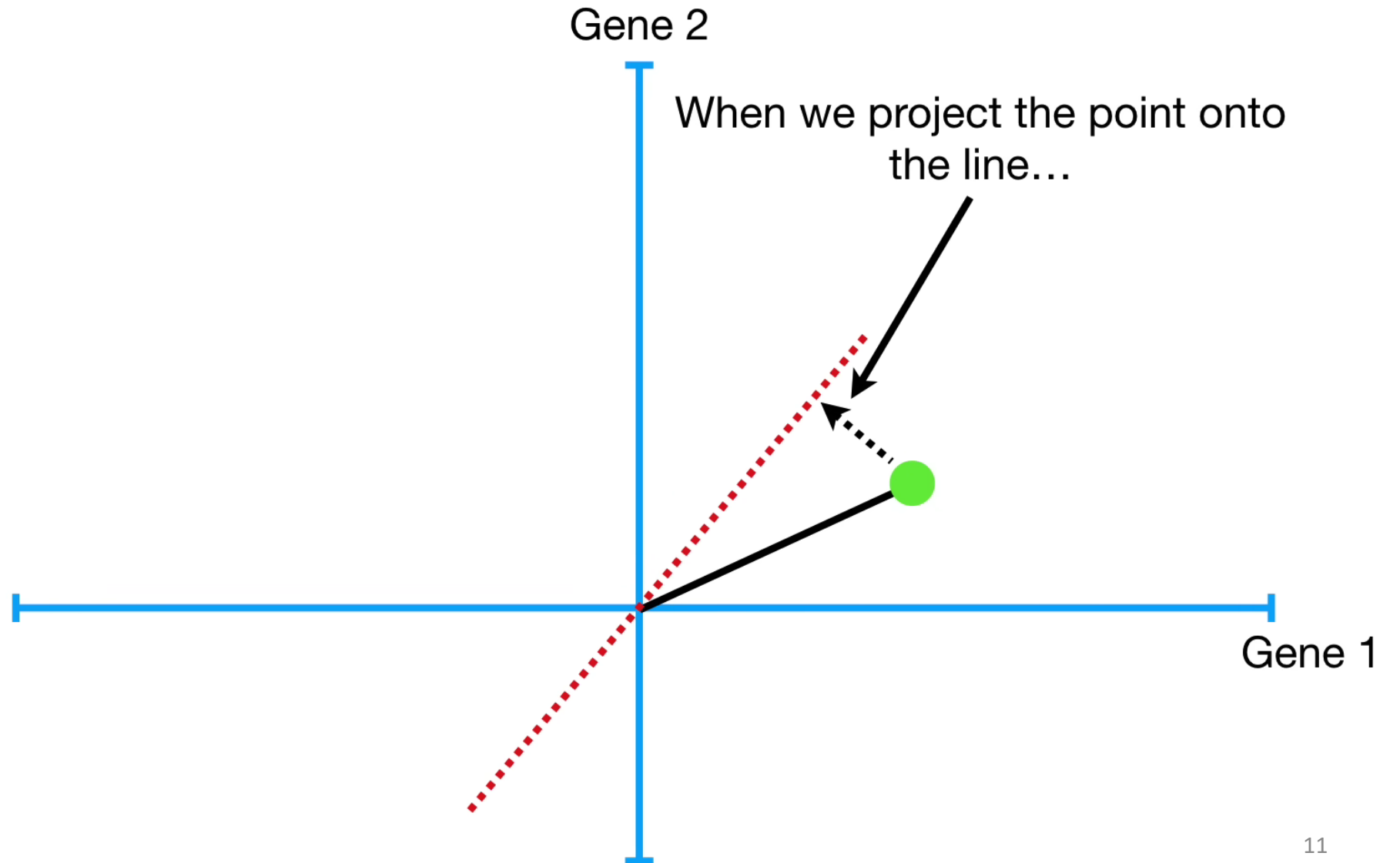


# PCA

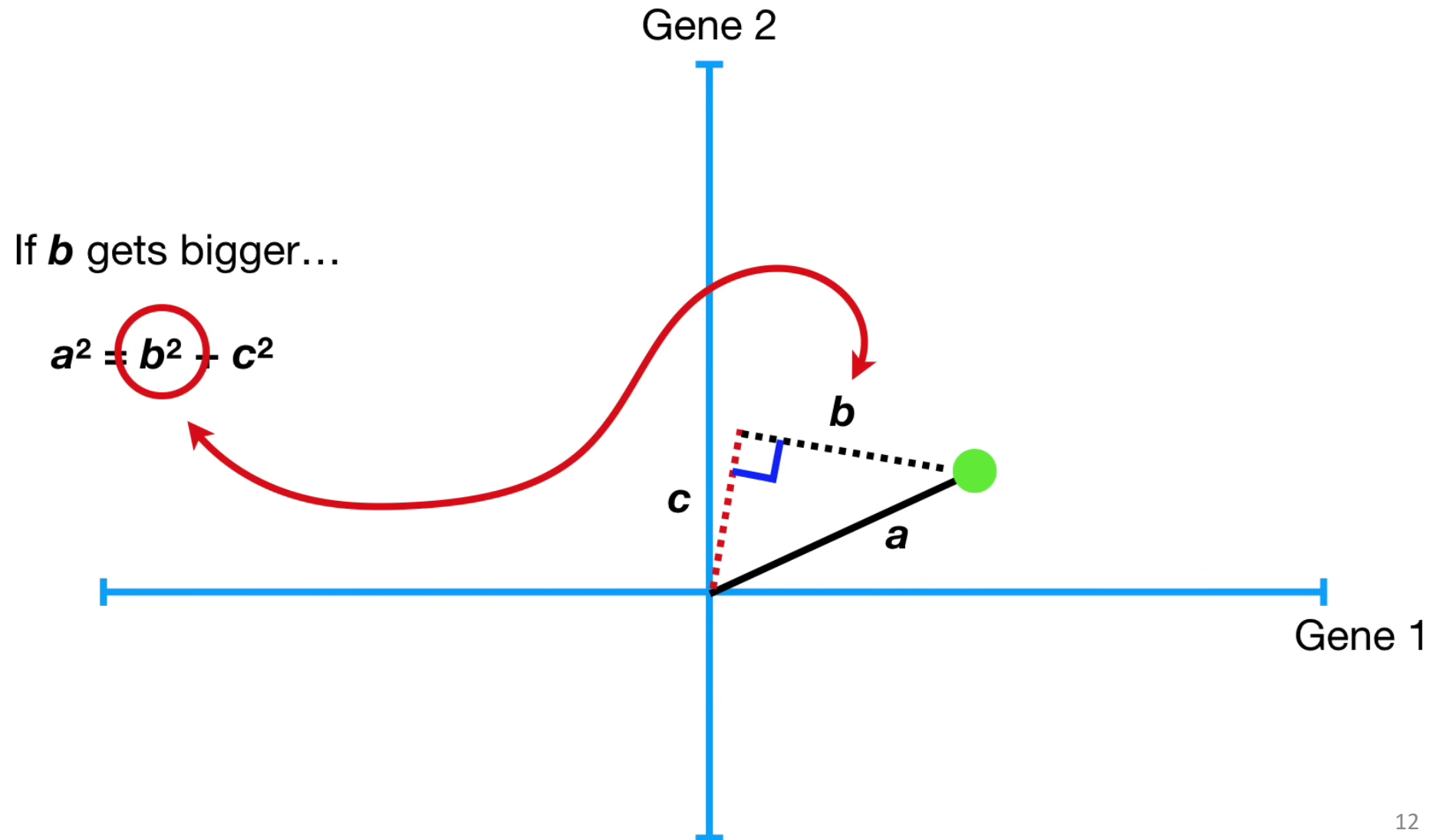
To quantify how good this line fits the data, PCA projects the data onto it...



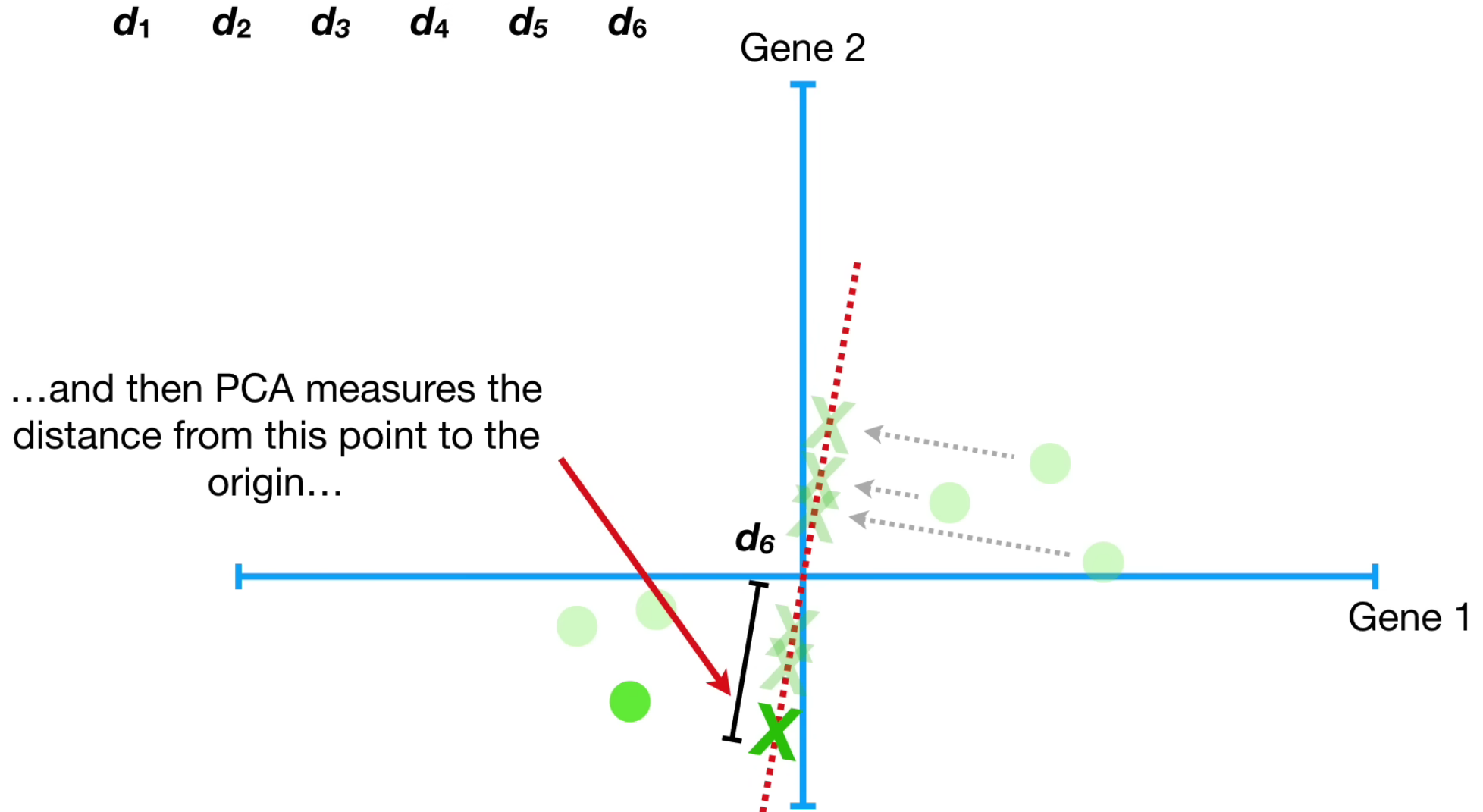
# PCA



# PCA



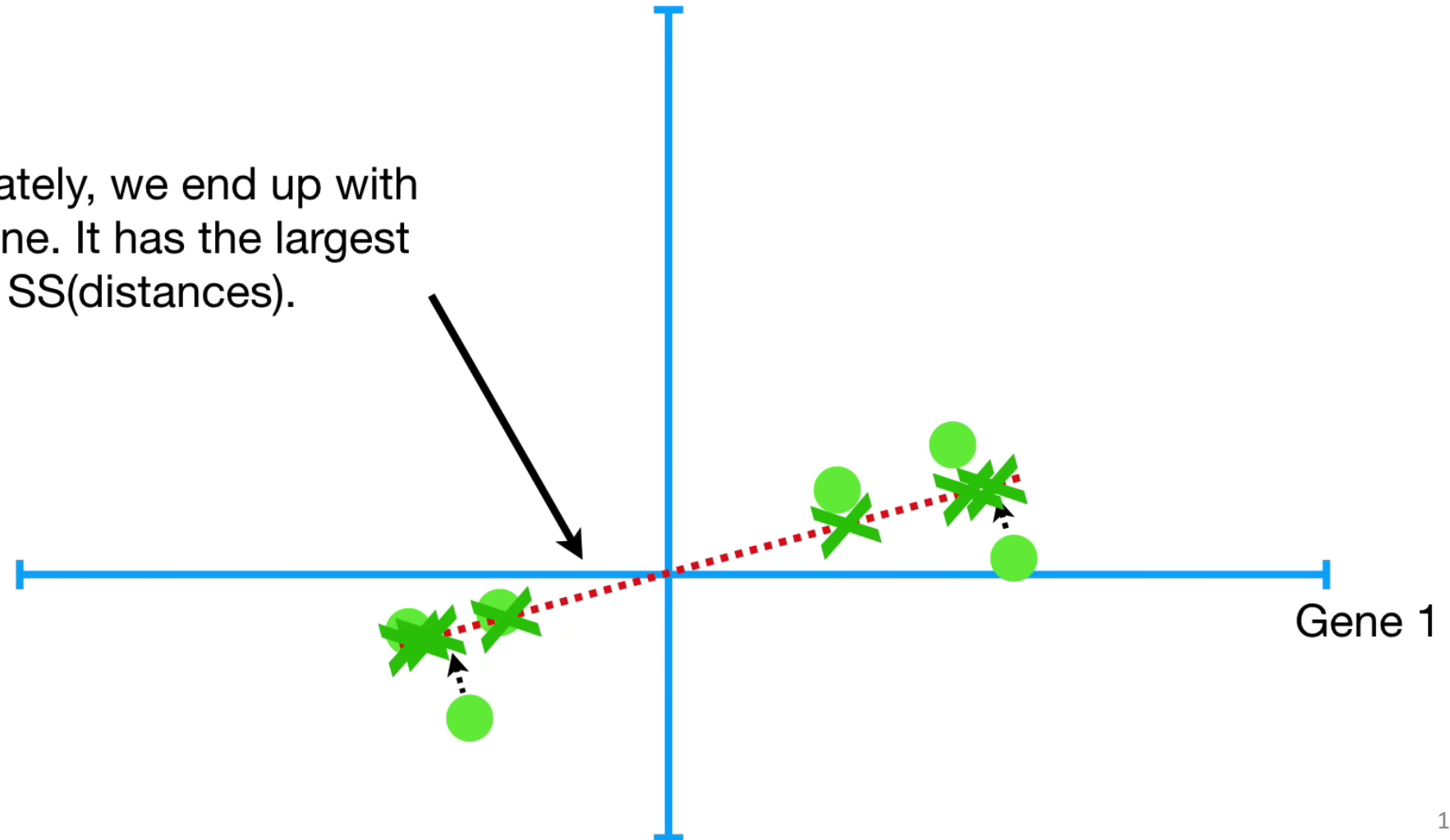
# PCA



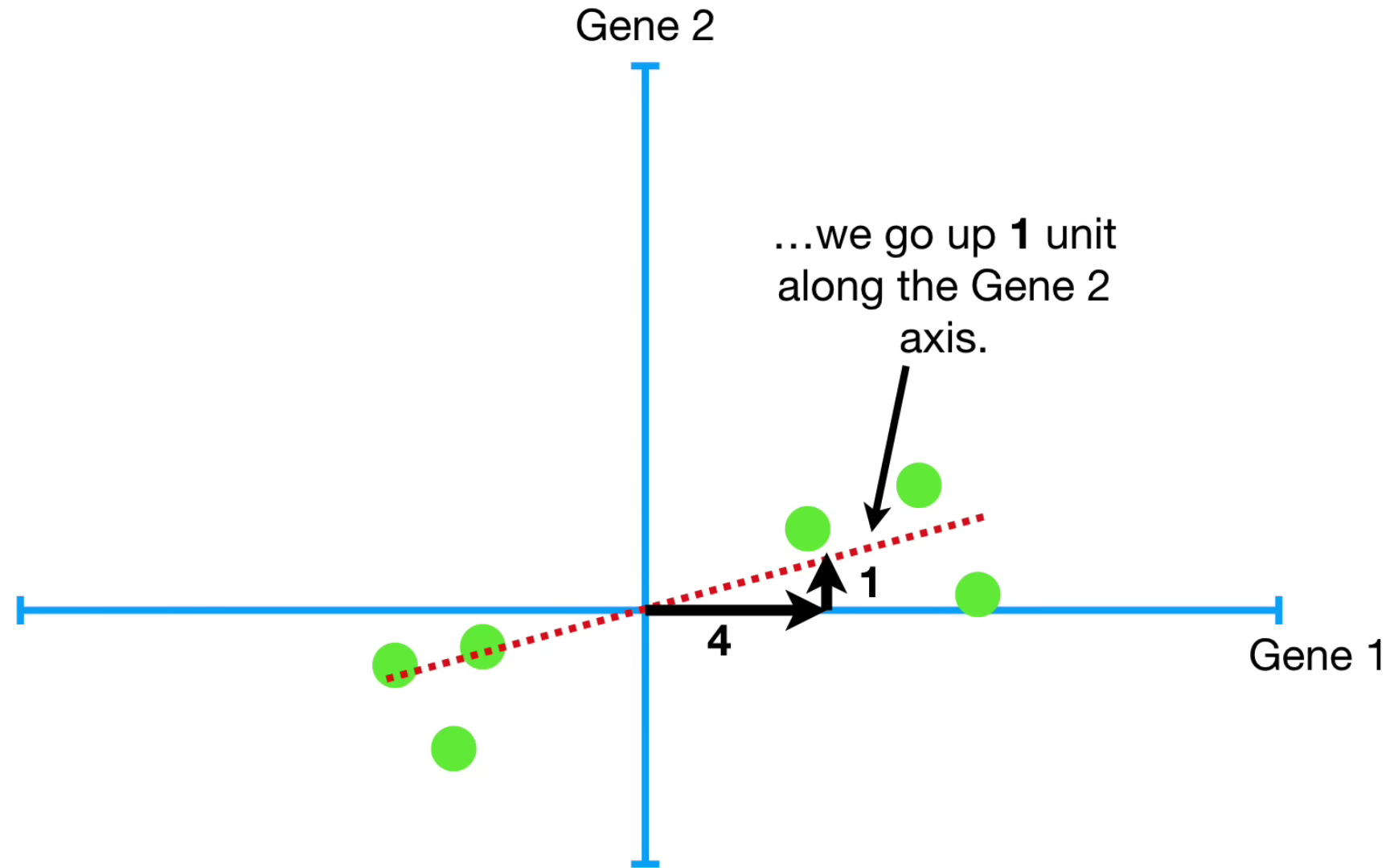
# PCA

$$d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2 + d_6^2 = \text{sum of squared distances} = \text{SS}(\text{distances})$$

Ultimately, we end up with this line. It has the largest SS(distances).



# PCA



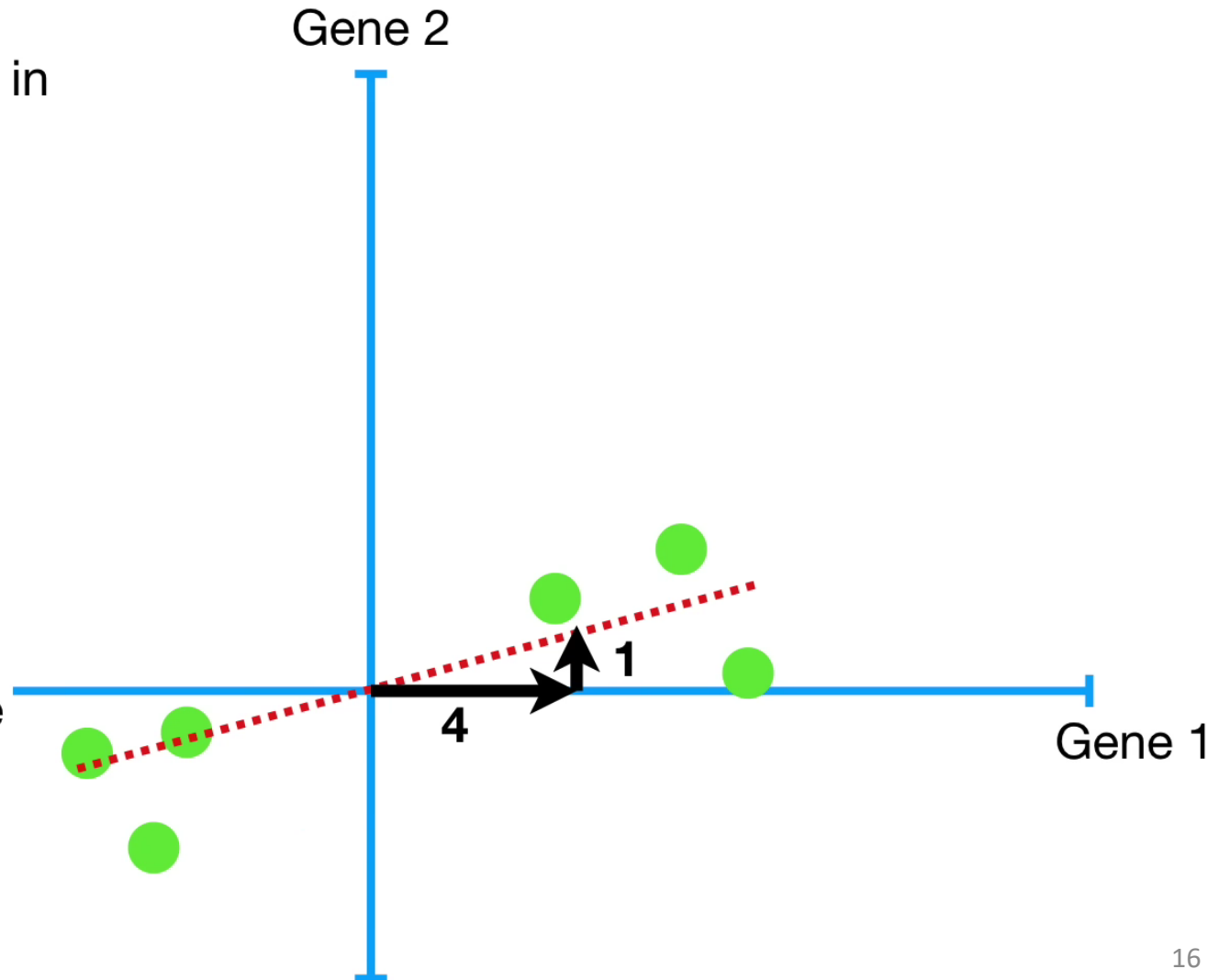
# PCA

One way to think about PC1 is in terms of a cocktail recipe...

## To make PC1

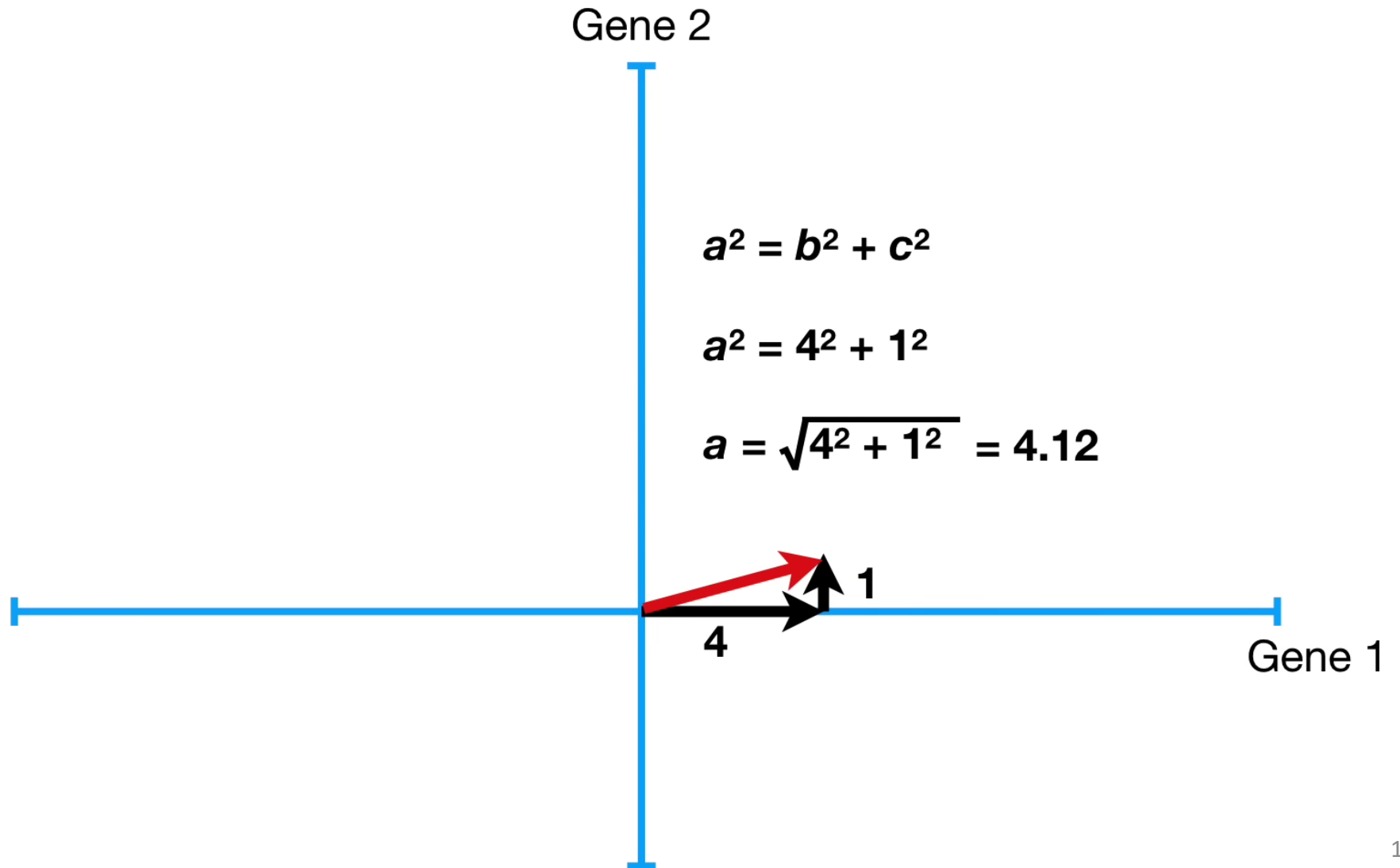
Mix **4** parts Gene 1  
with **1** part Gene 2

The ratio of Gene 1 to Gene 2 tells you that Gene 1 is more important when it comes to describing how the data are spread out..





# PCA



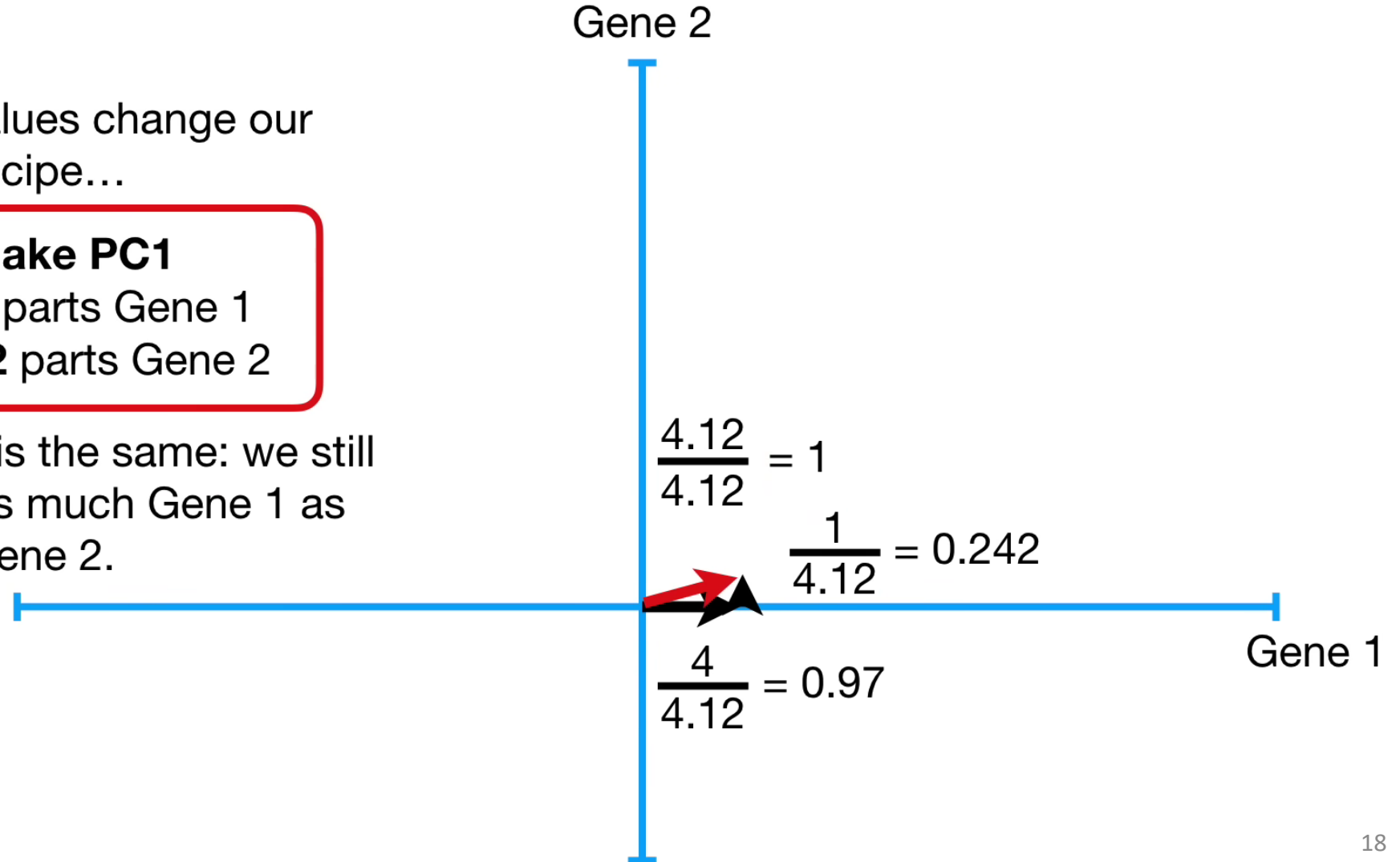
# PCA

The new values change our recipe...

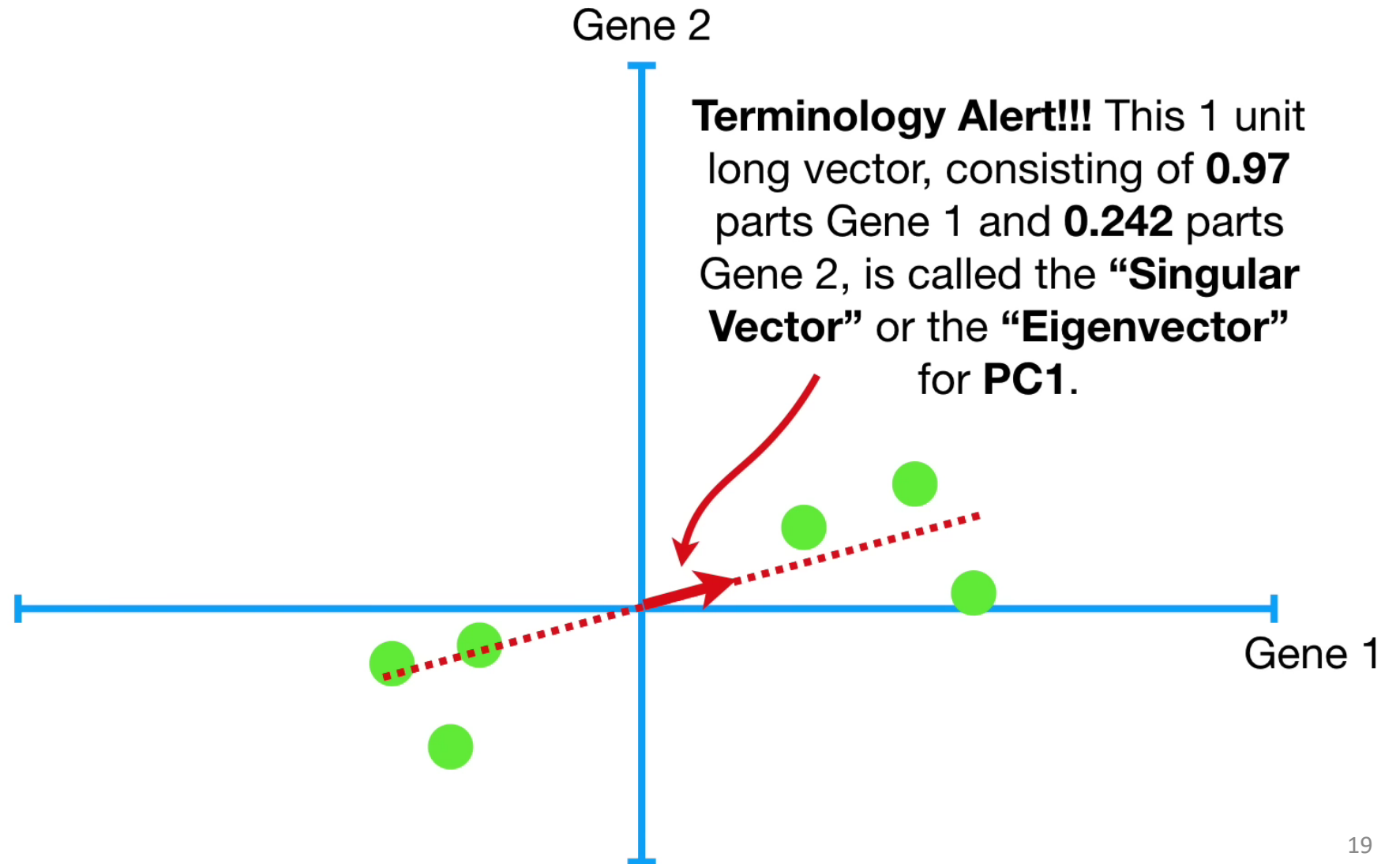
## To make PC1

Mix **0.97** parts Gene 1  
with **0.242** parts Gene 2

...but the ratio is the same: we still  
use 4 times as much Gene 1 as  
Gene 2.



# PCA

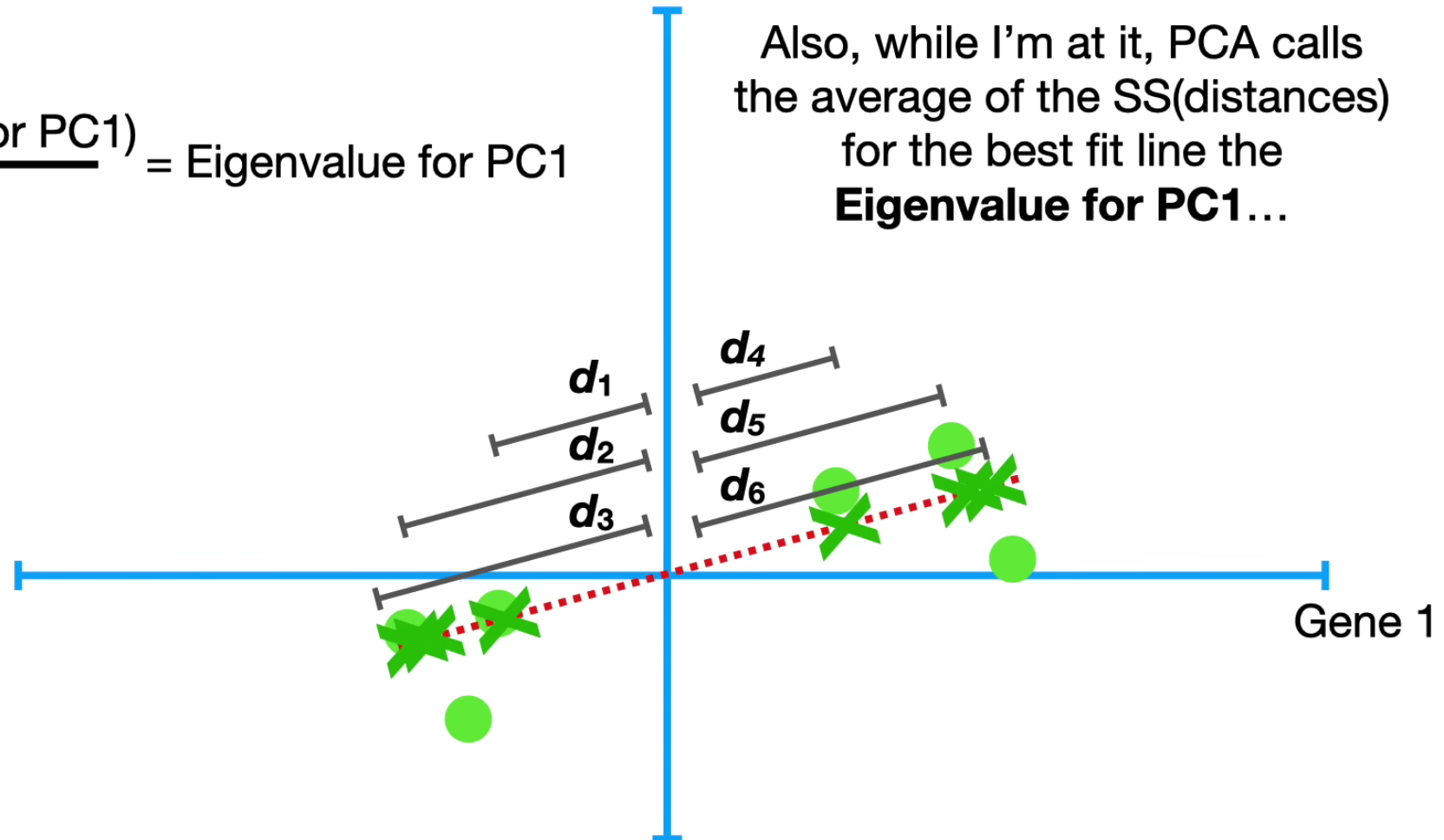


# PCA

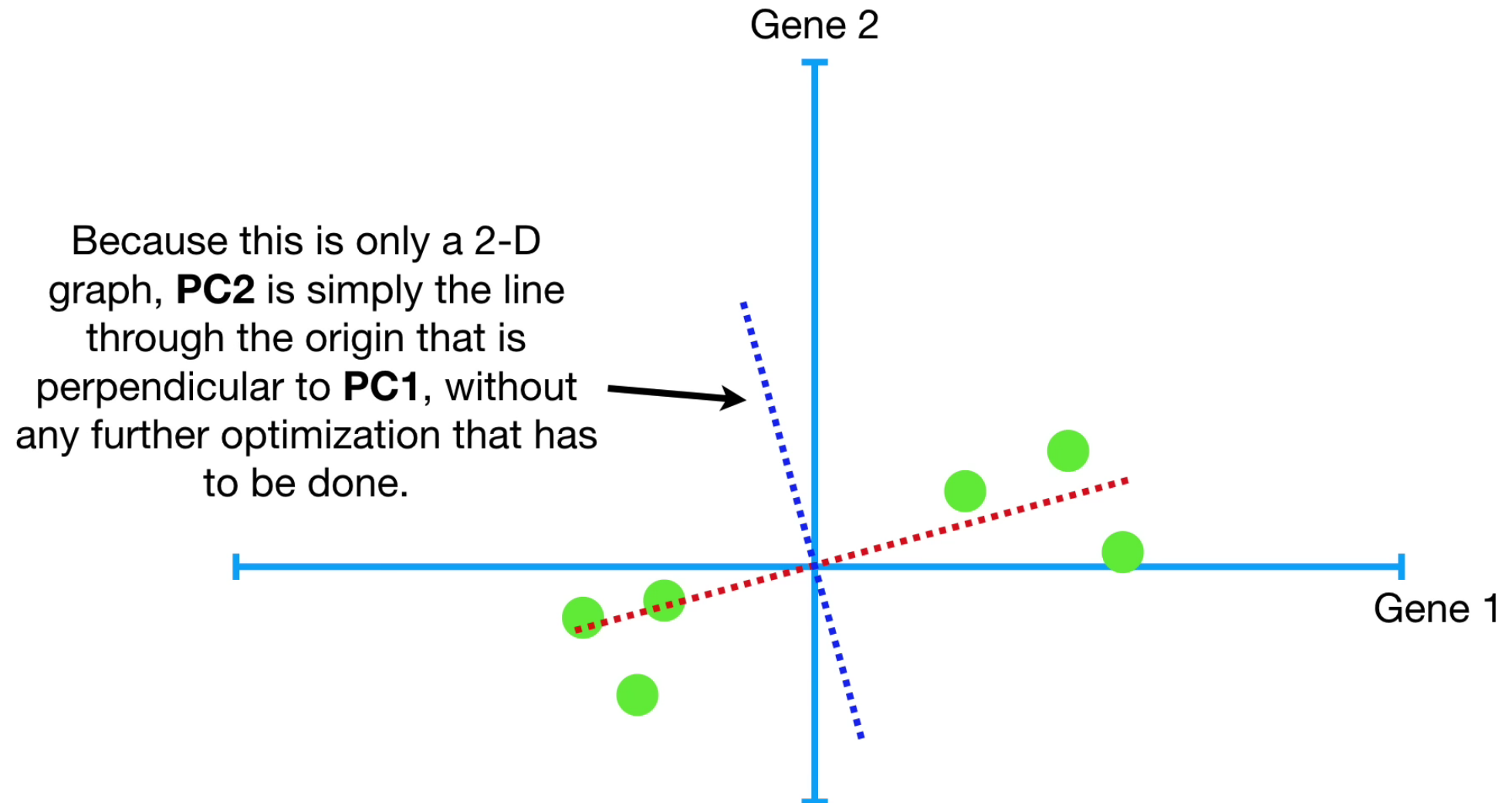
$$d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2 + d_6^2 = \text{sum of squared distances} = \text{SS}(\text{distances})$$

$$\frac{\text{SS}(\text{distances for PC1})}{n - 1} = \text{Eigenvalue for PC1}$$

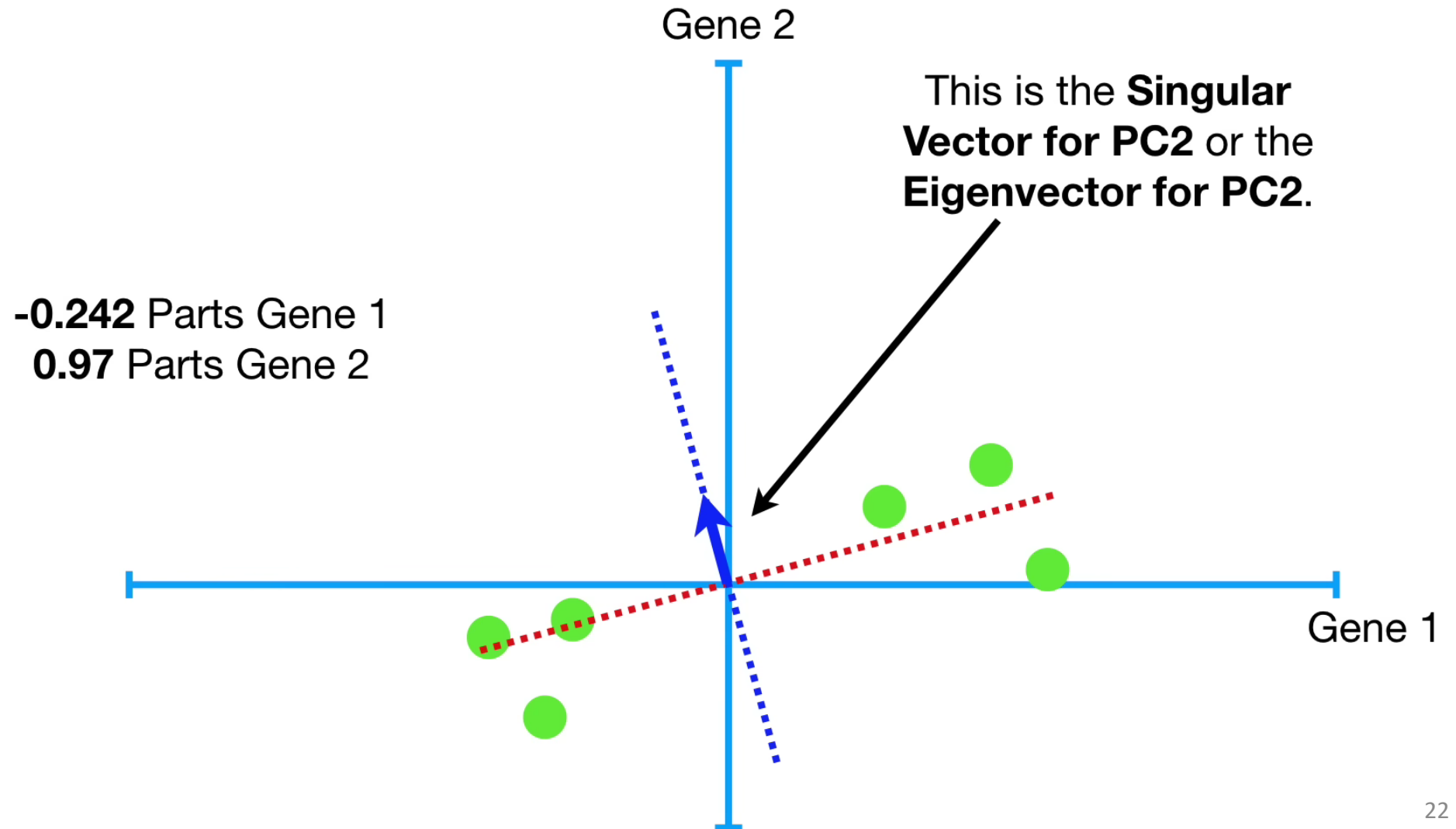
Also, while I'm at it, PCA calls  
the average of the SS(distances)  
for the best fit line the  
**Eigenvalue for PC1...**



# PCA



# PCA



# PCA

For the sake of the example, imagine that the Variation for **PC1 = 15**, and the variation for **PC2 = 3**.

That means that the total variation around both PCs is **15 + 3 = 18...**

$$\frac{SS(\text{distances for PC1})}{n - 1} = \text{Variation for PC1}$$

$$\frac{SS(\text{distances for PC2})}{n - 1} = \text{Variation for PC2}$$

PC2 ...and that means PC1 accounts for **15 / 18 = 0.83 = 83%** of the total variation around the PCs.

