# PCA

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

Gene 2

Gene 1

Gene 2

...or it can try to find the line that **maximizes** the distances from the projected points to the origin.

Gene 1

Gene 2

...while **these distances get larger when the line fits better.**

Gene 1

$d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2 + d_6^2$ = sum of squared distances = SS(distances)
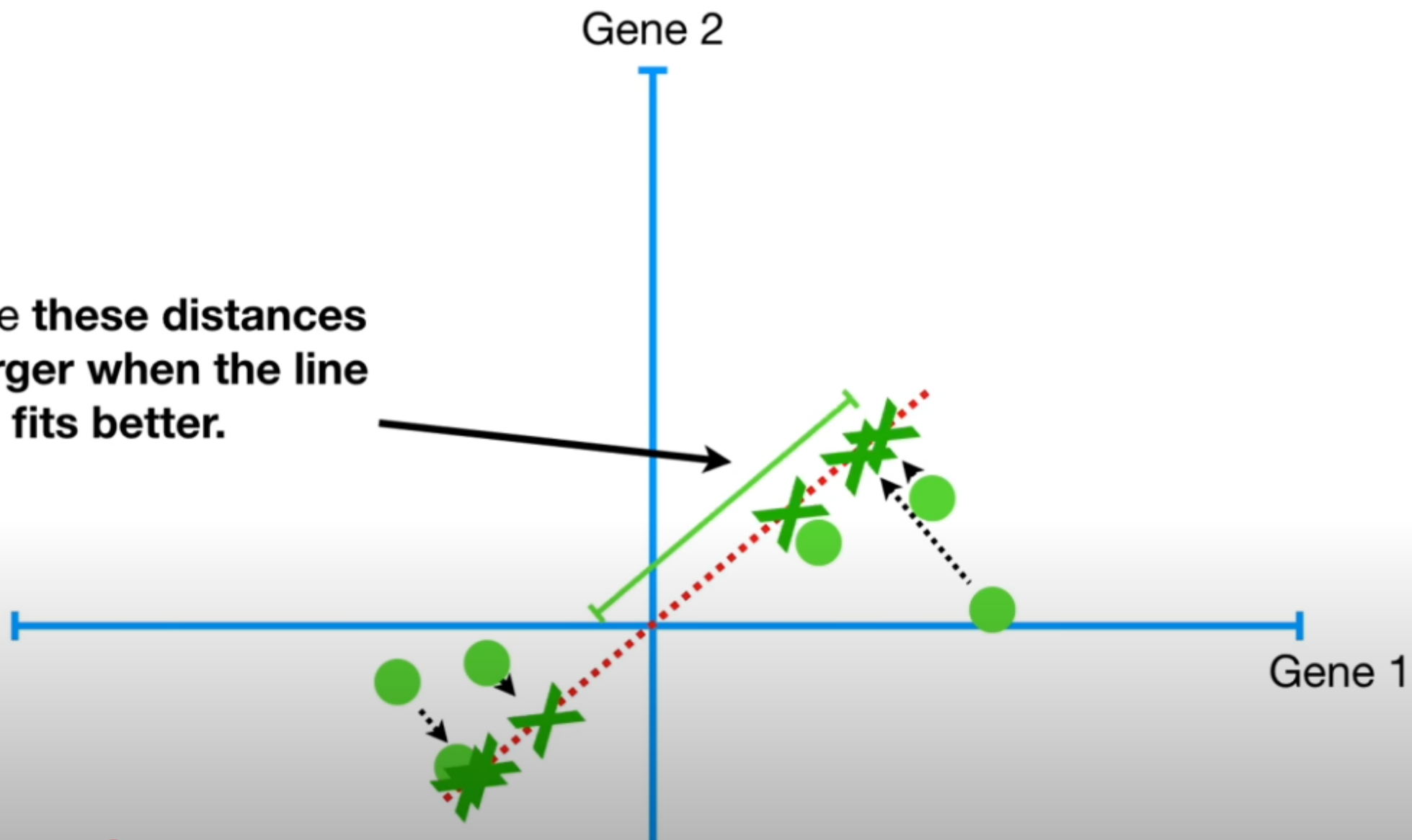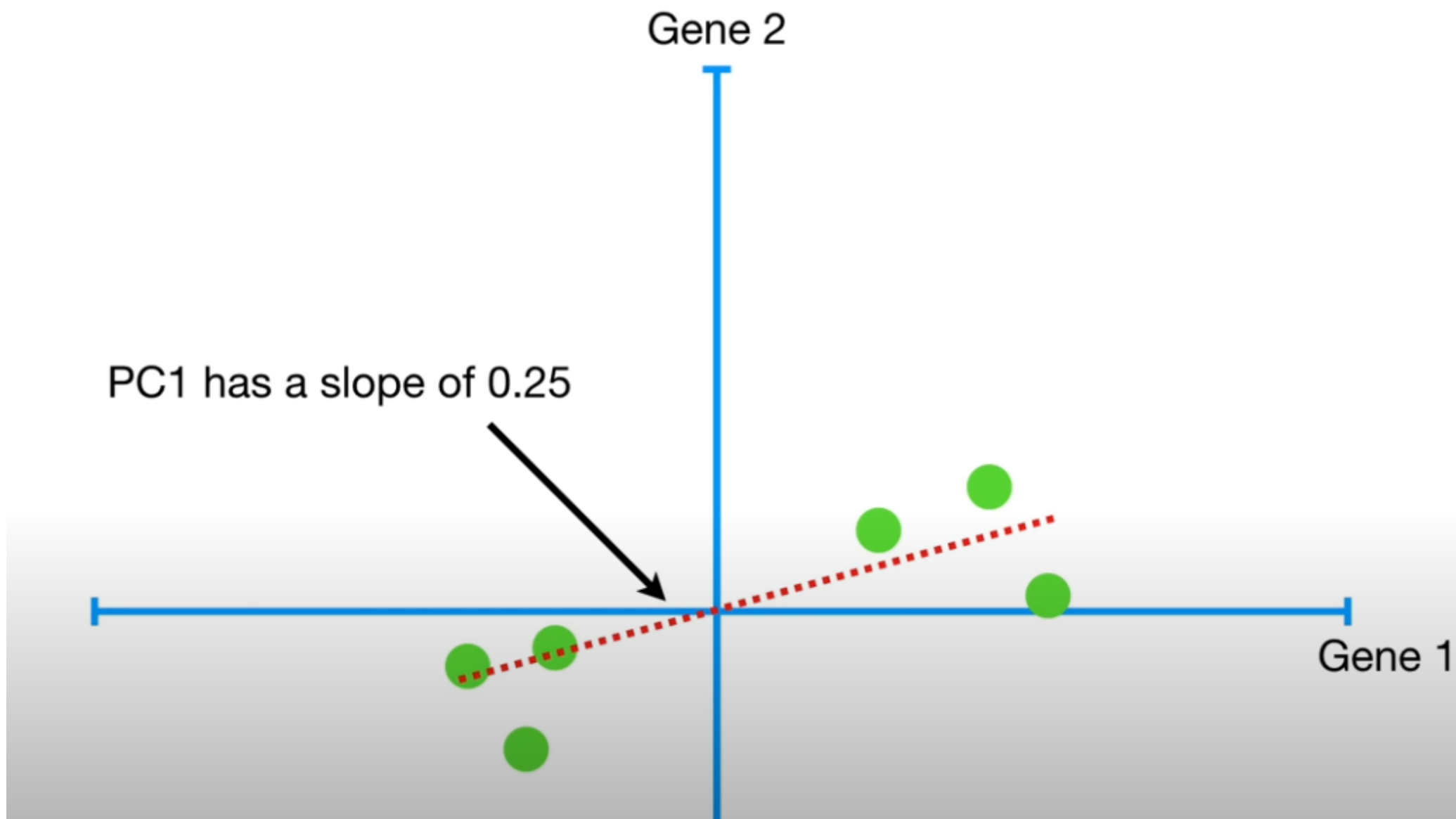
$$\frac{SS(\text{distances for PC1})}{n-1} = \text{Eigenvalue for PC1}$$

$$\sqrt{SS(\text{distances for PC1})} = \text{Singular Value for PC1}$$

...and the square root of the SS(distances) is called the **Singular Value for PC1.**

Gene 1

...so Sample 6 goes here.

PC2

PC1

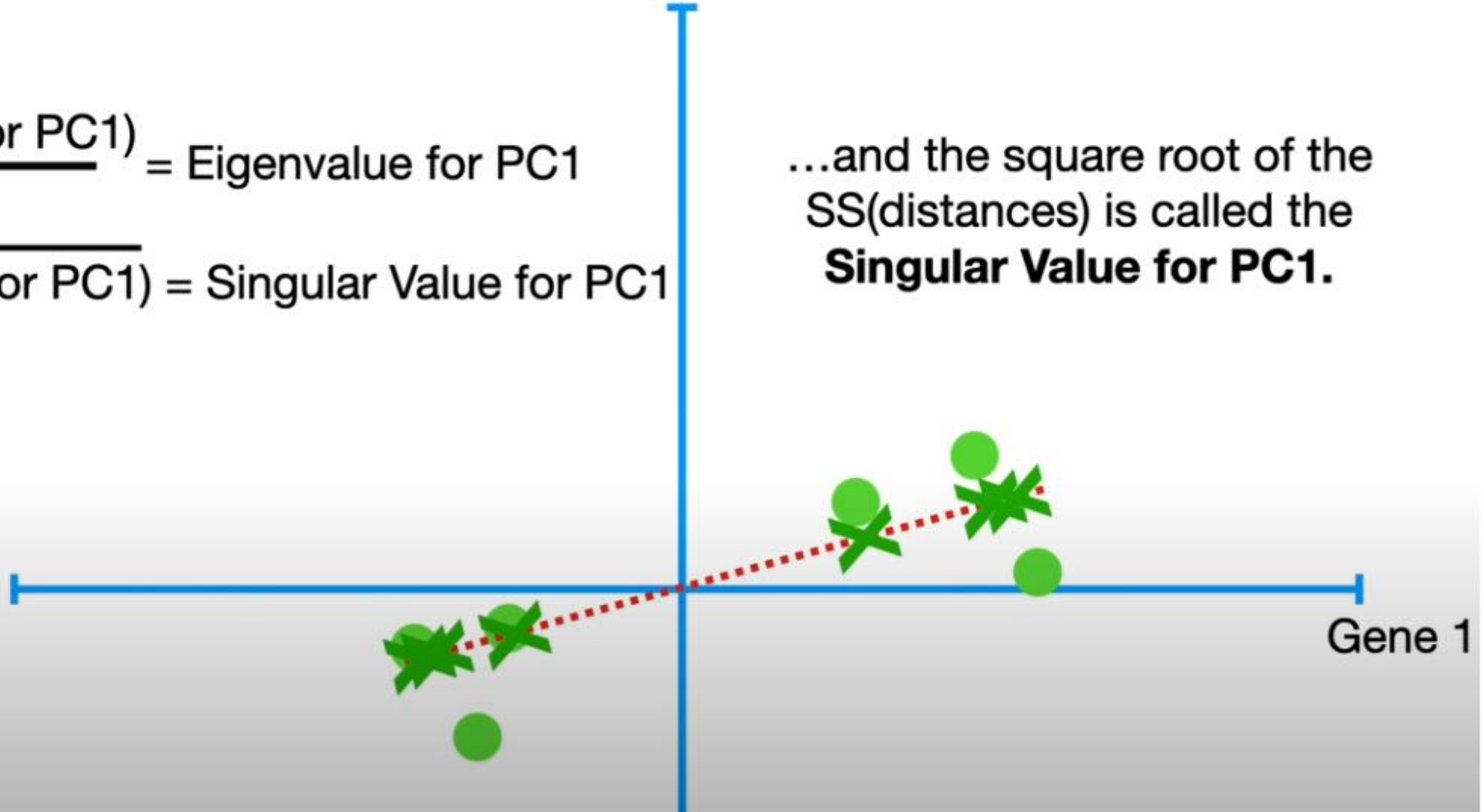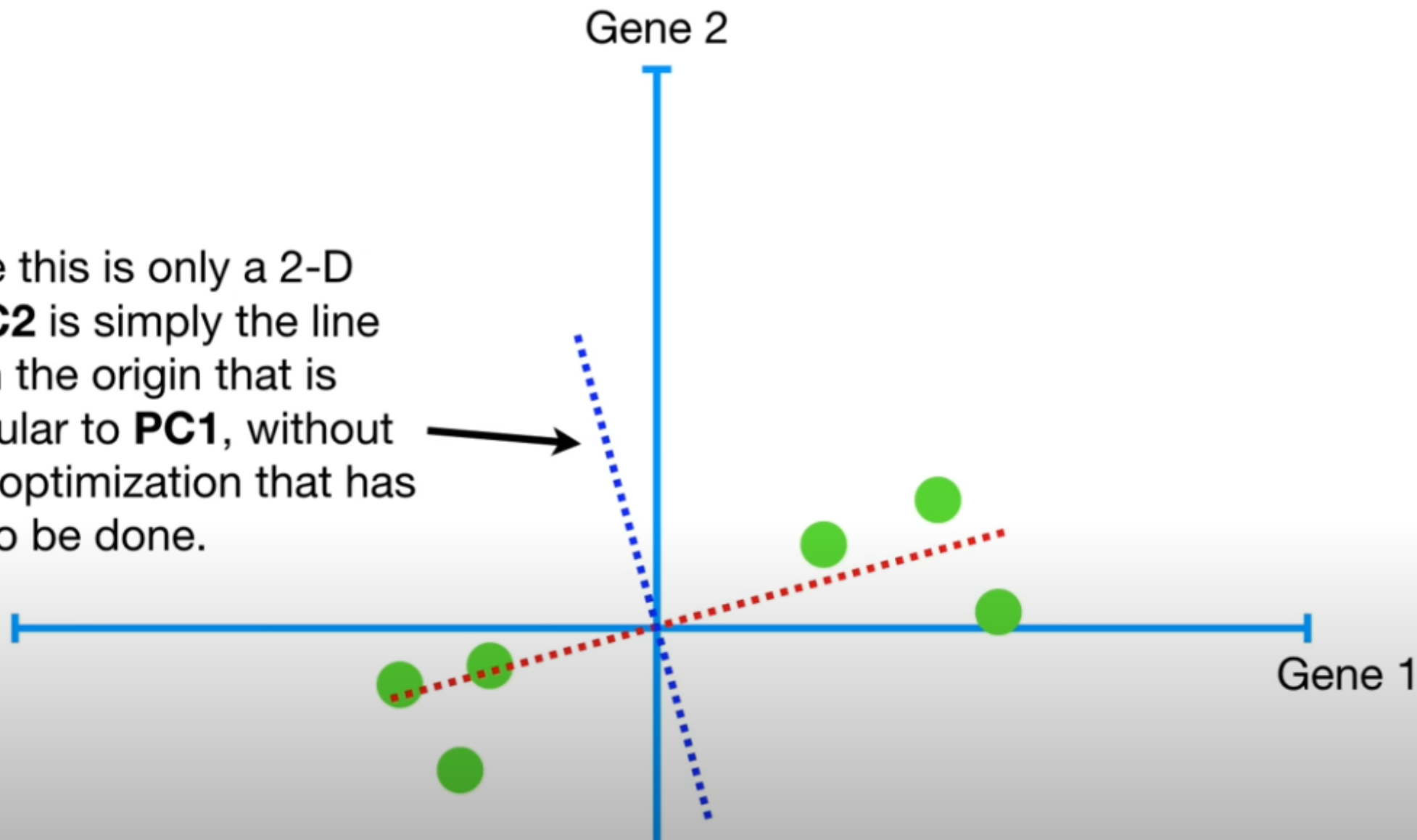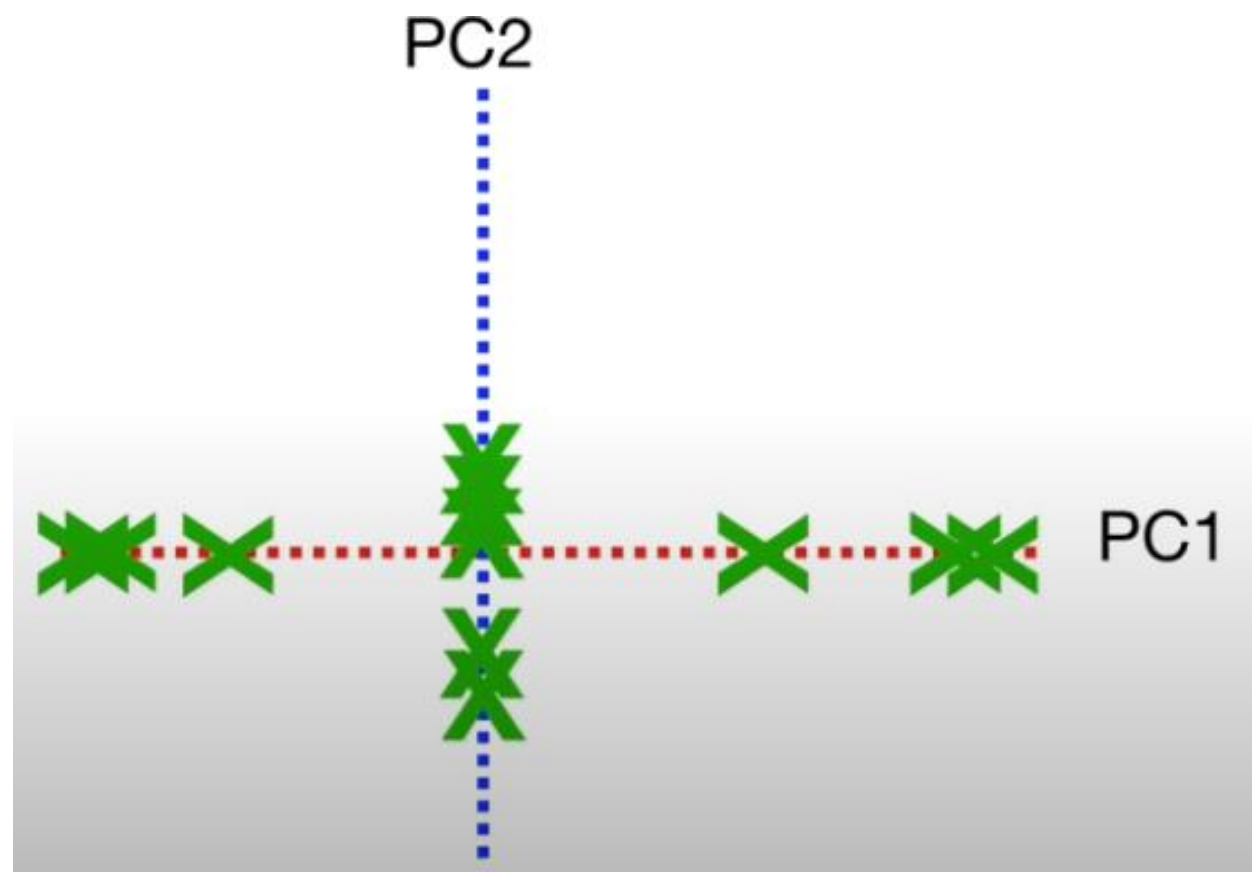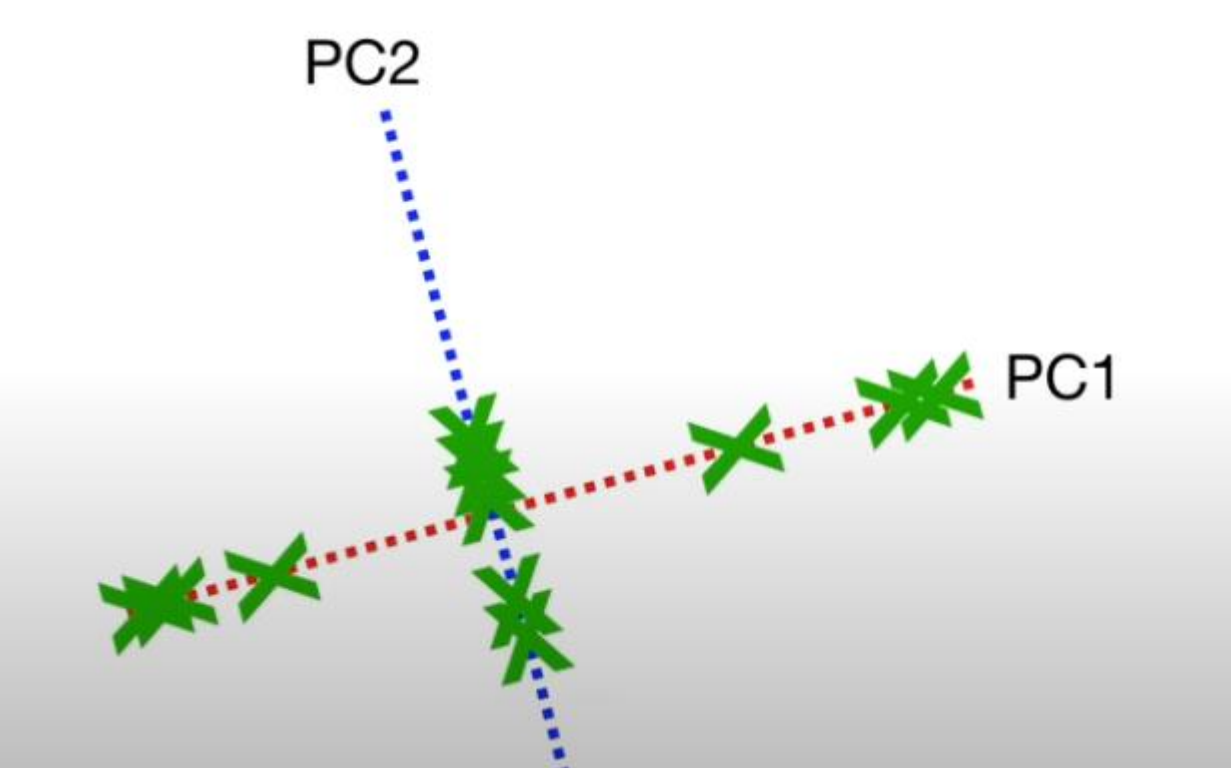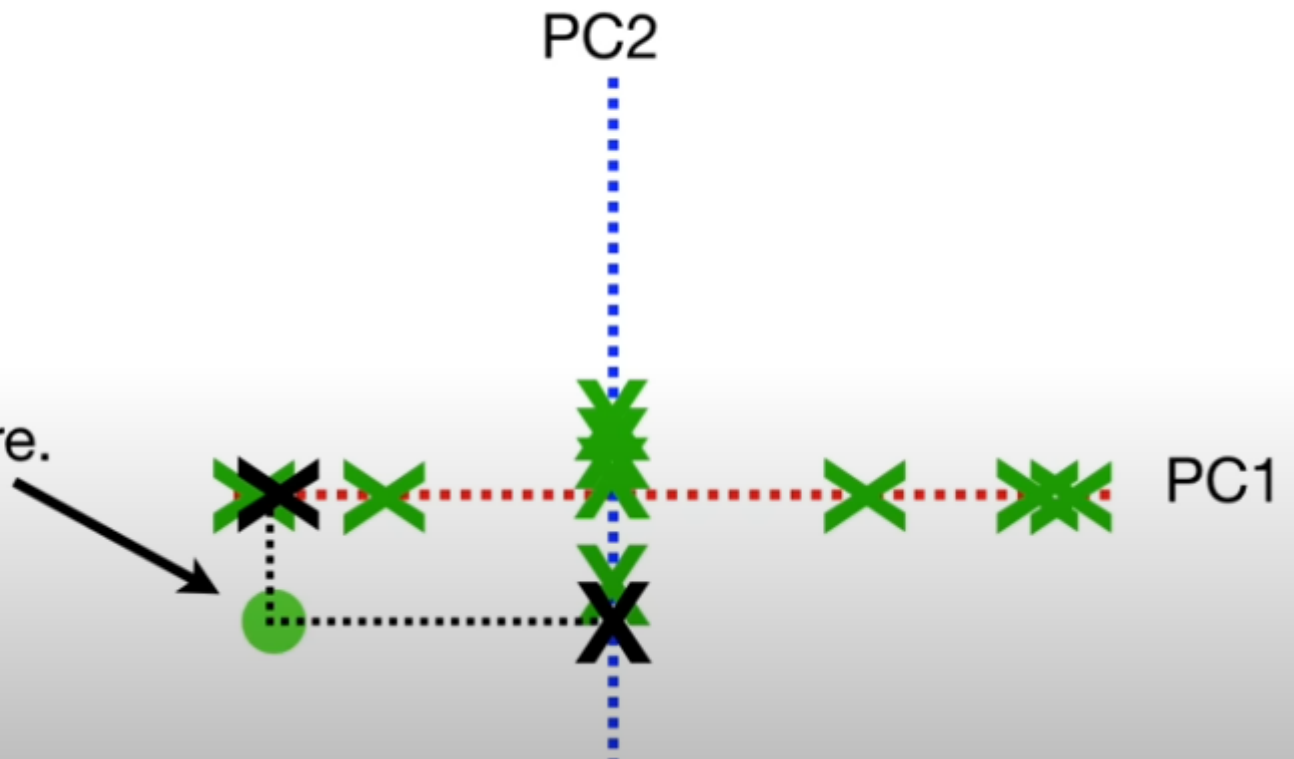For the sake of the example, imagine that the Variation for **PC1 = 15**, and the variation for **PC2 = 3**.

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

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

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

PC2

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



PC1 (83%)

**TERMINOLOGY ALERT!!!!** A **Scree Plot** is a graphical representation of the percentages of variation that each PC accounts for.
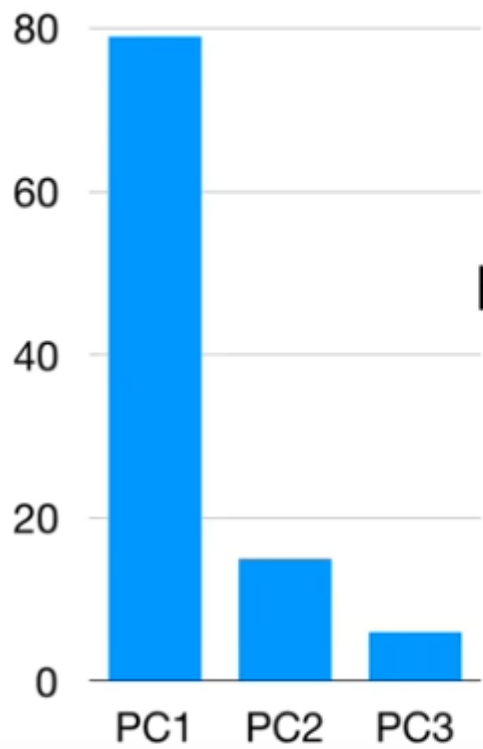
3D

Here's the scree plot...