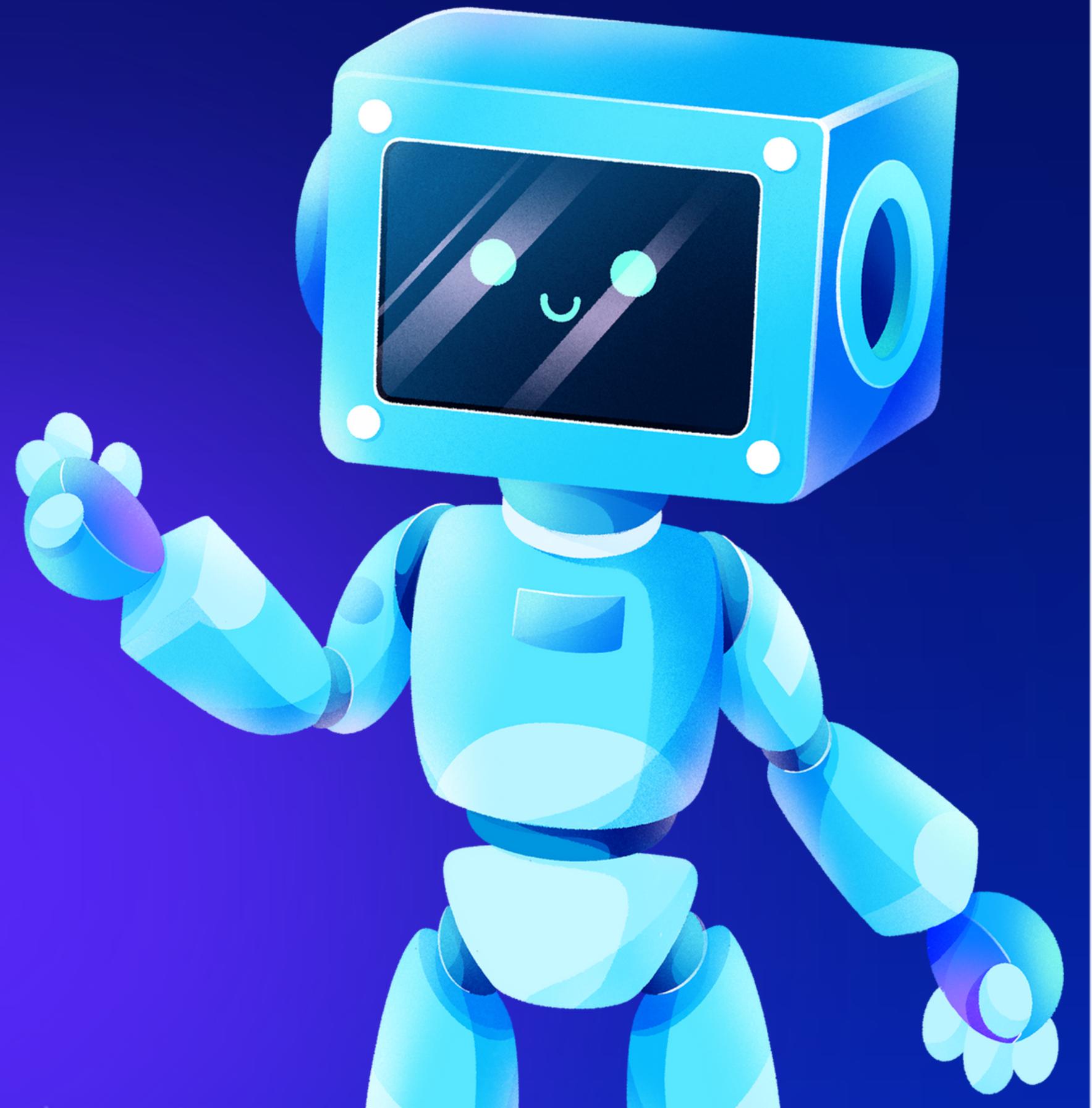




ARTIFICIAL INTELLIGENCE

PCA (PHÂN TÍCH THÀNH PHẦN CHÍNH)

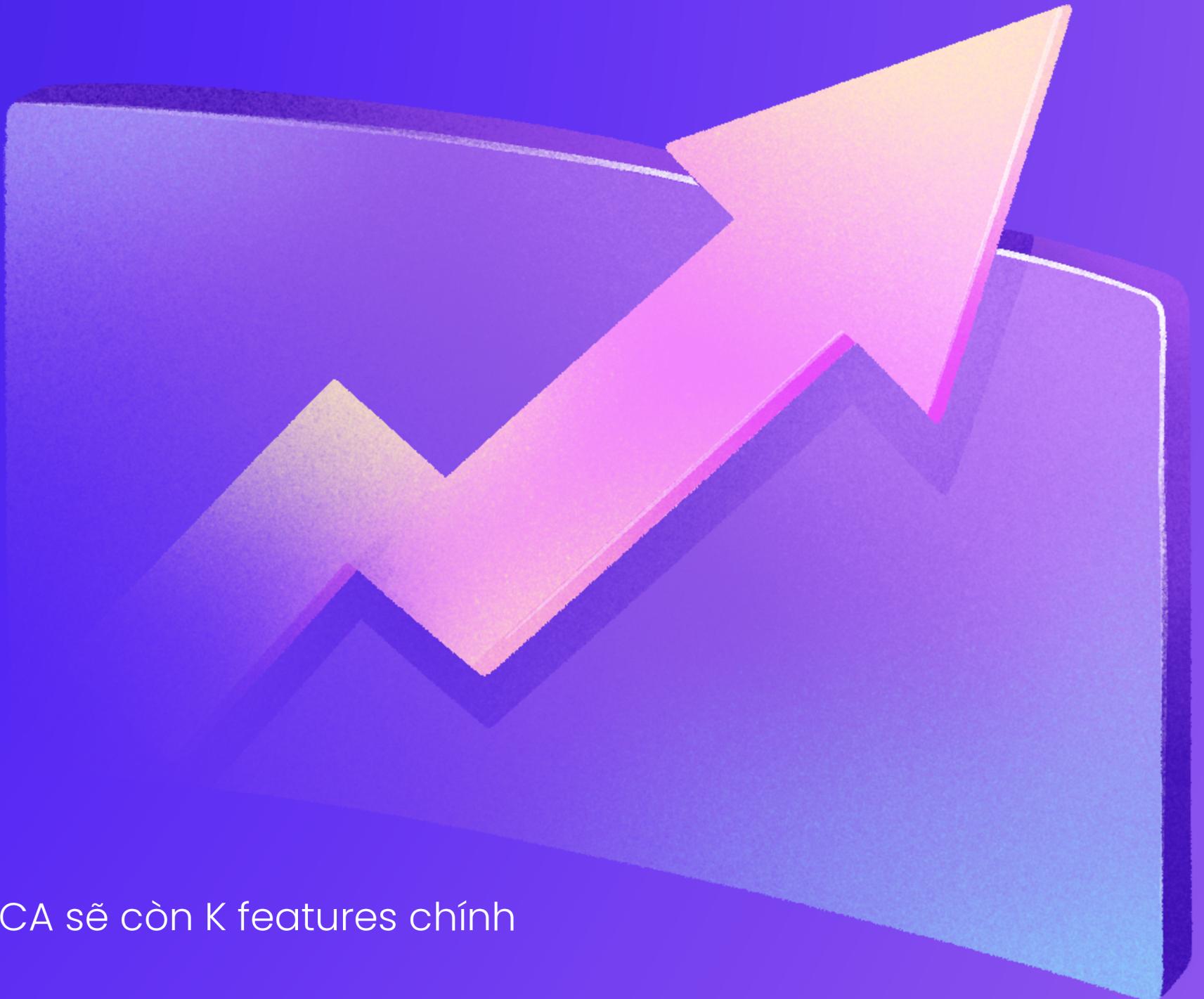


PCA LÀ GÌ?

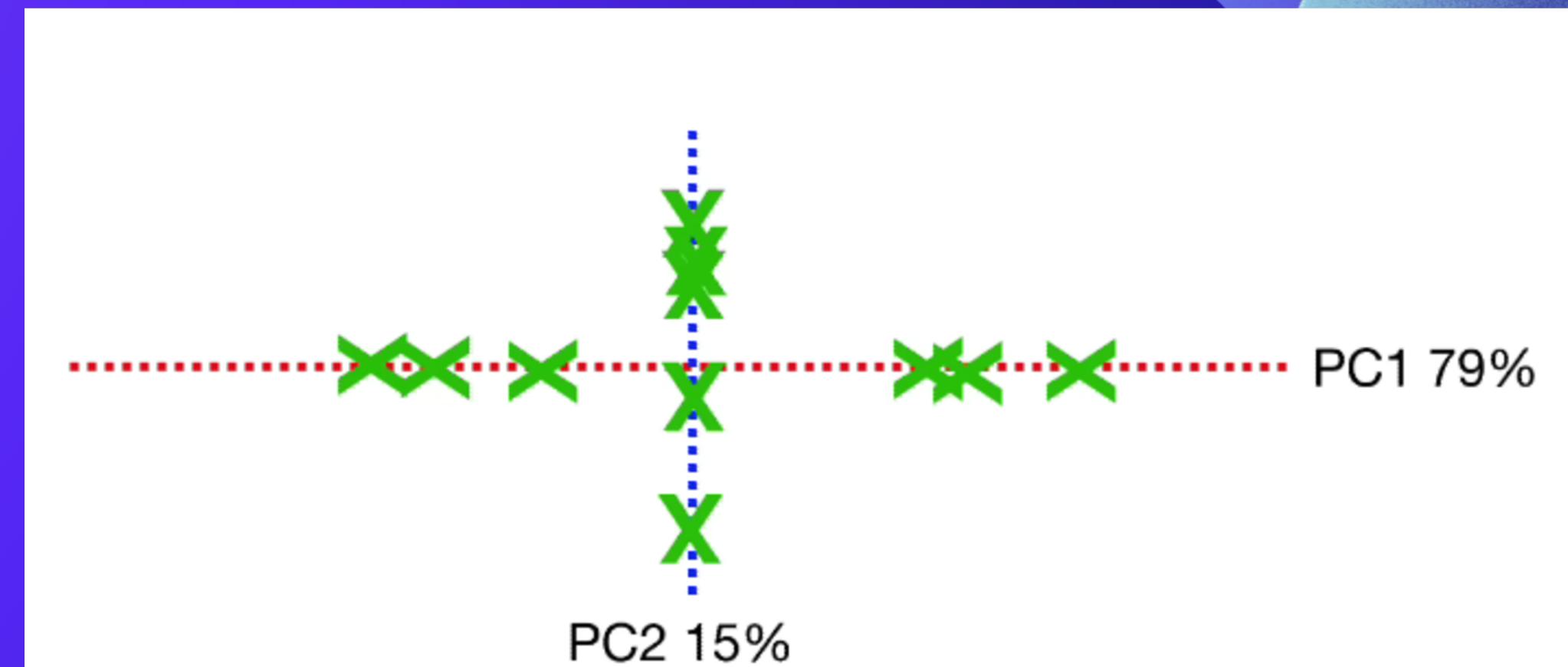
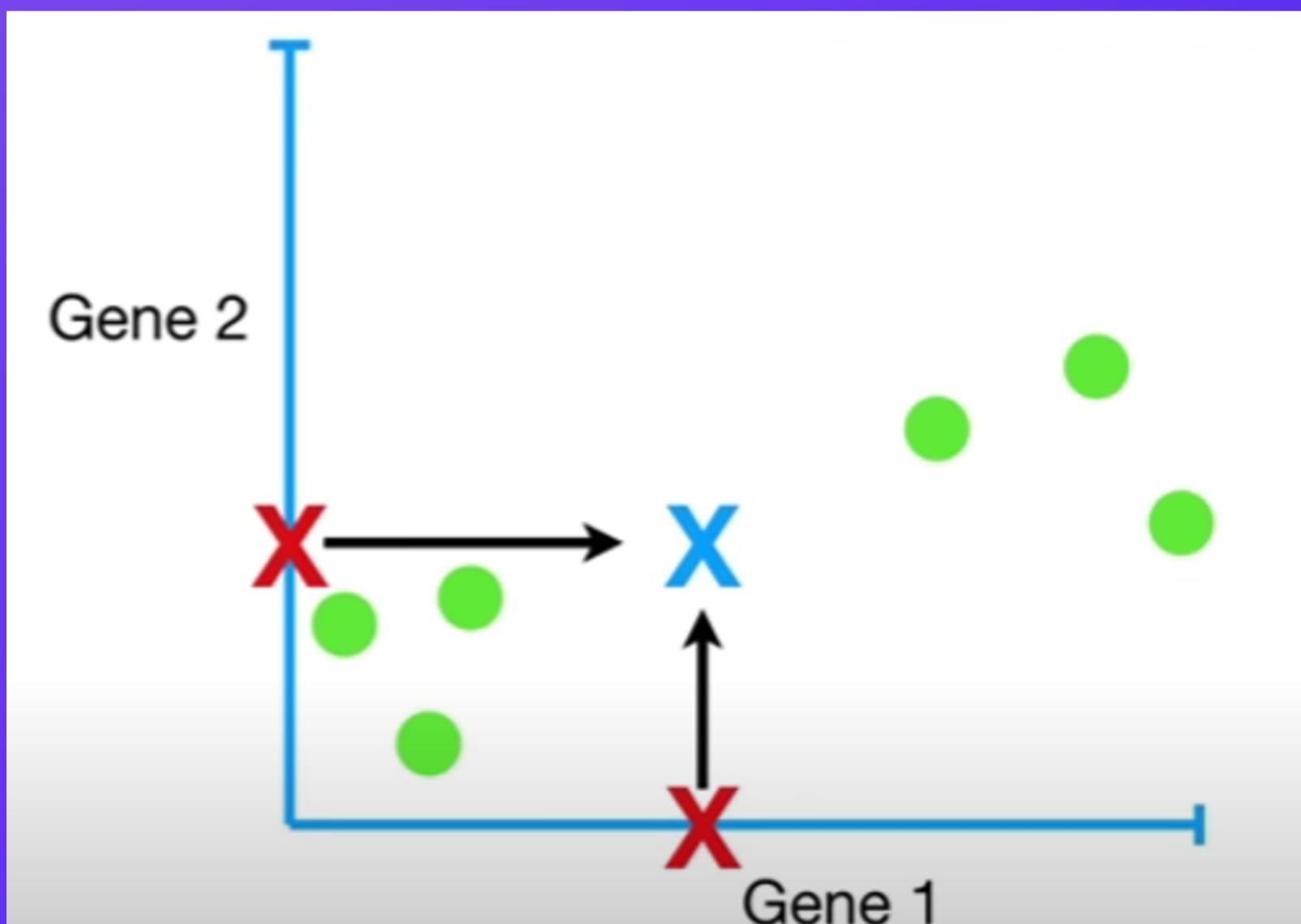
- Phân tích thành phần chính hoặc PCA là một kỹ thuật để giảm kích thước của tập dữ liệu lớn.
- Việc làm như trên sẽ giúp cho chúng ta: giảm chiều dữ liệu mà vẫn giữ được đặc trưng chính, chỉ mất đi “chút ít” đặc trưng, tiết kiệm thời gian, chi phí tính toán, dễ dàng visualize dữ liệu hơn

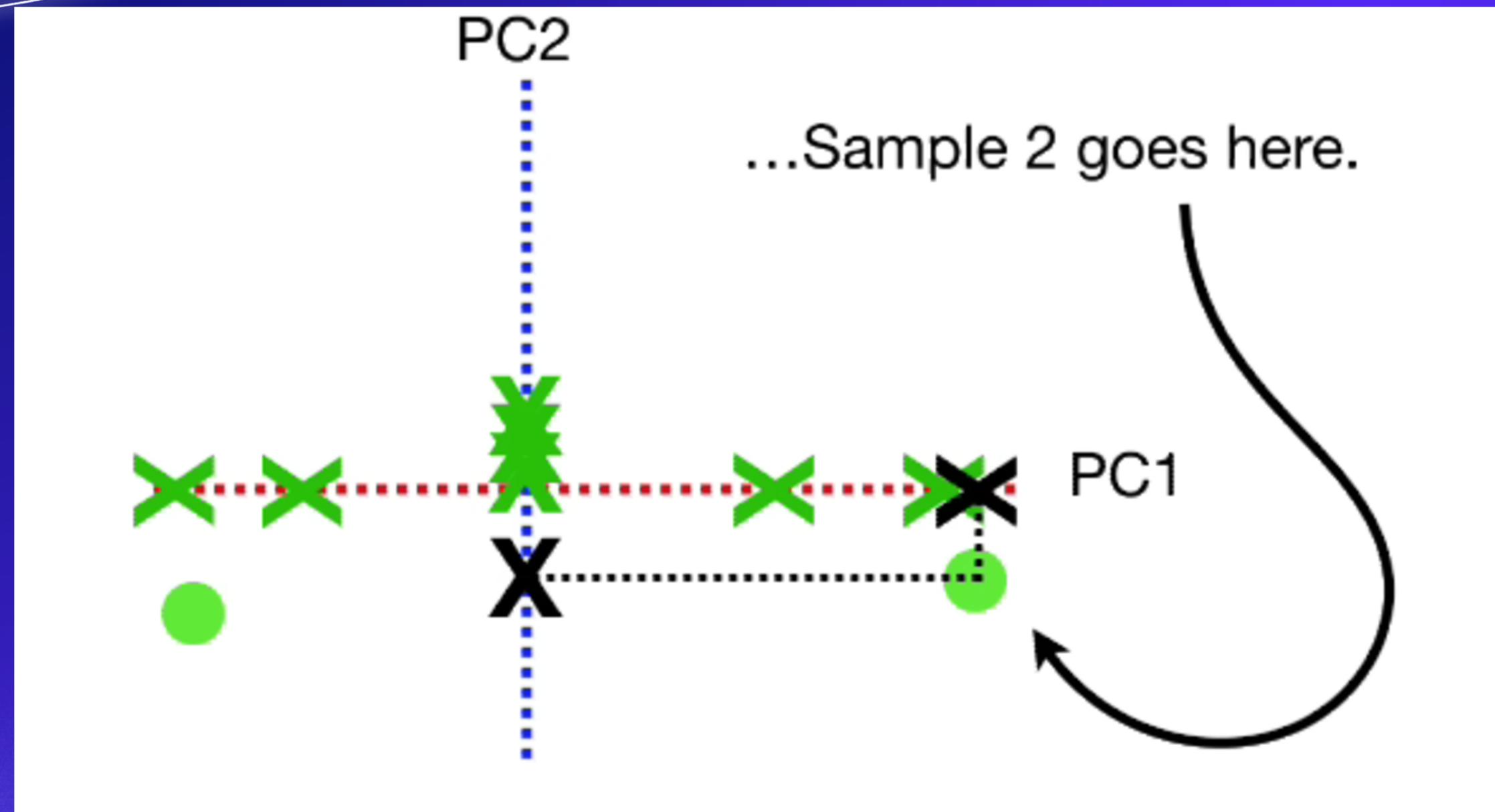
Ví dụ

Dữ liệu của chúng ta có N features thì sau khi áp dụng PCA sẽ còn K features chính mà thôi ($K < N$)



MỤC TIÊU PCA

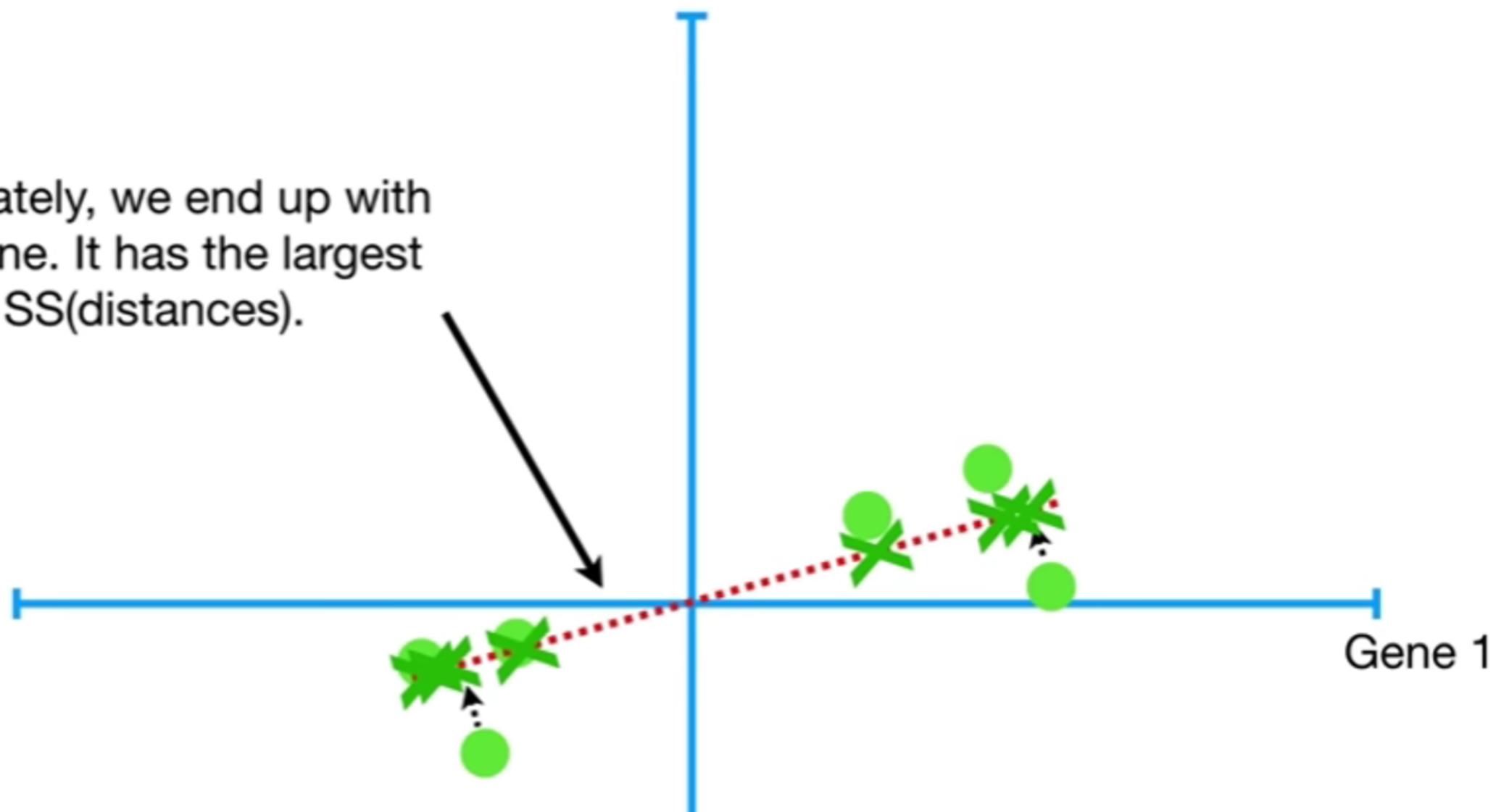




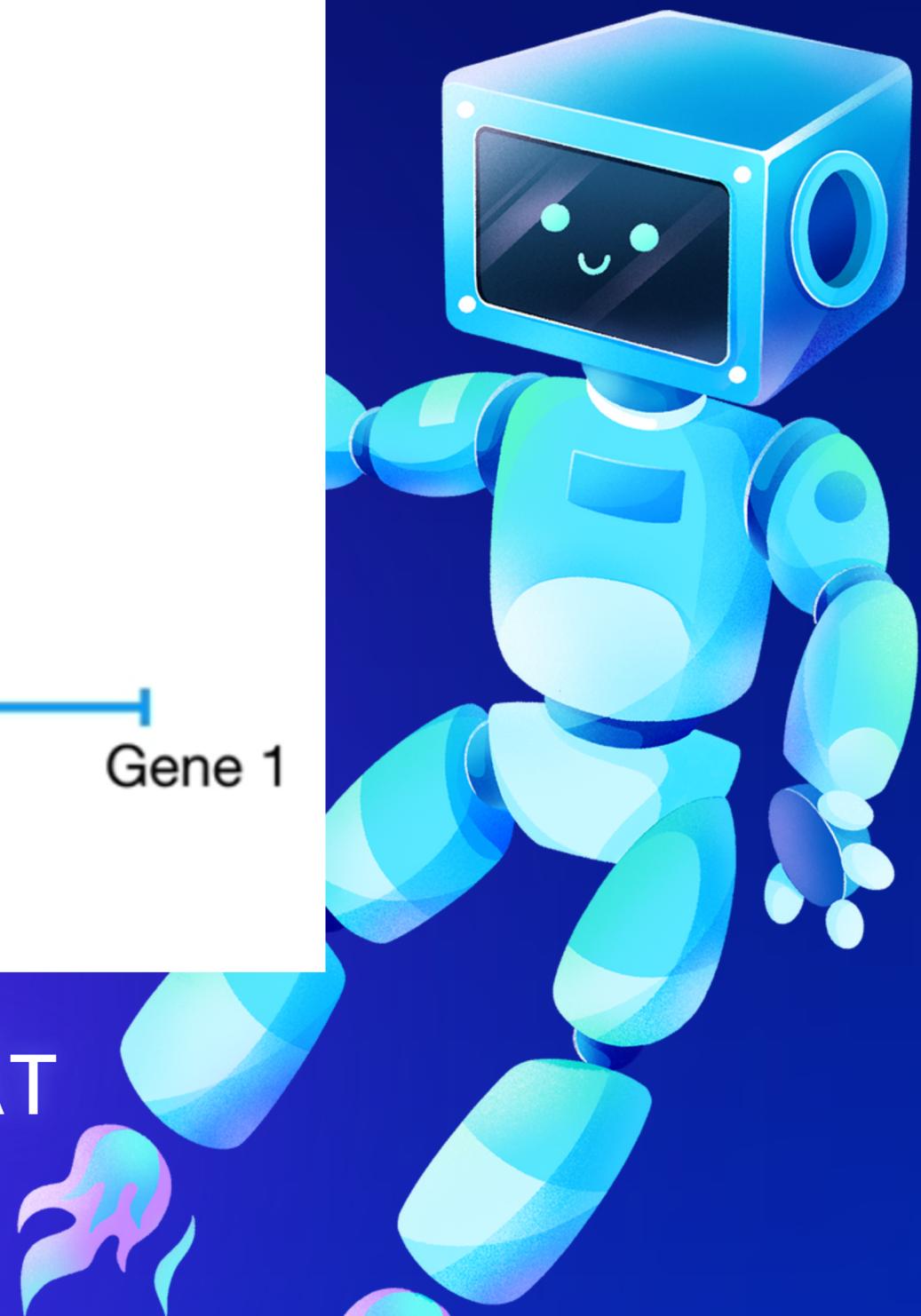
CÁCH XÁC ĐỊNH PCA

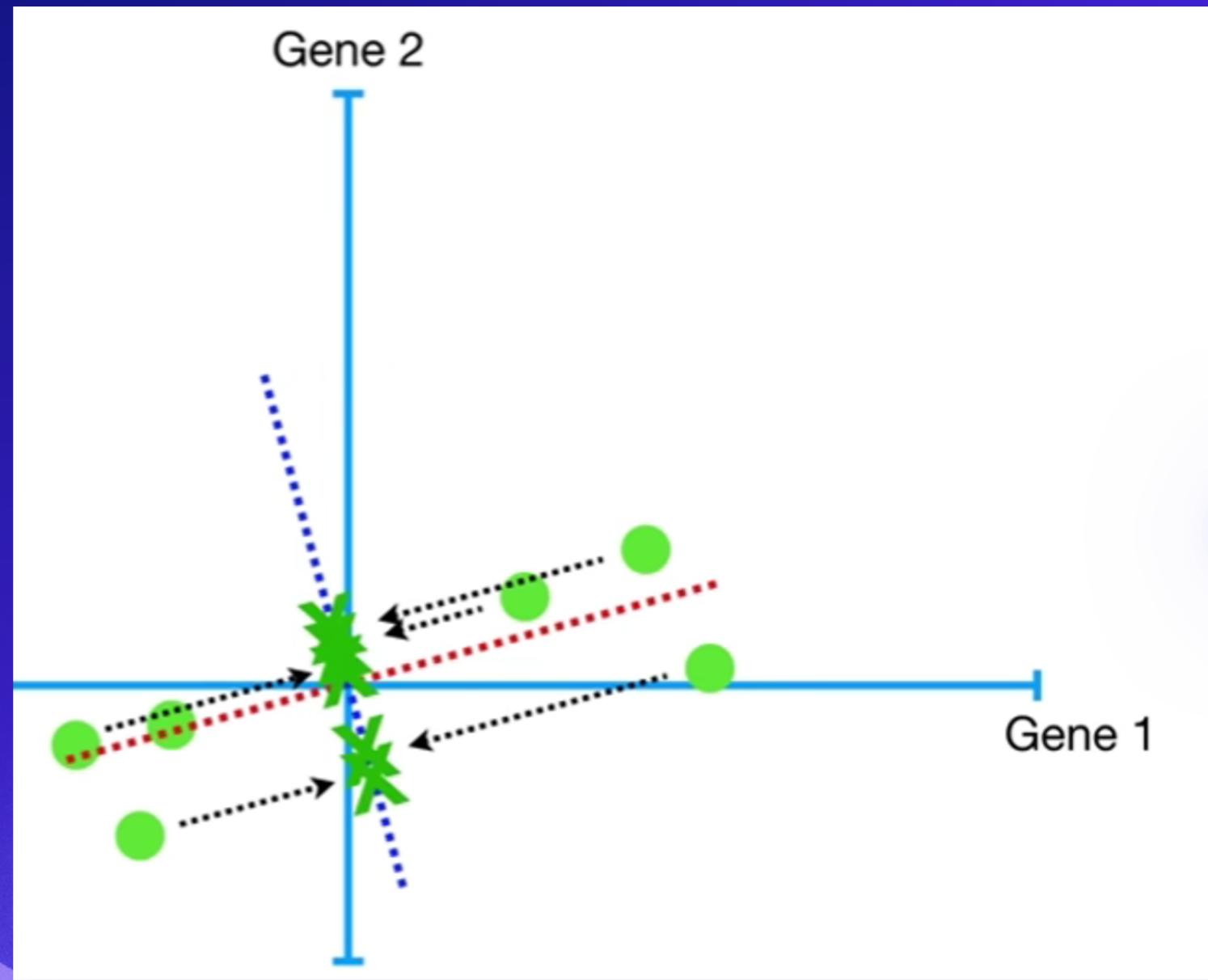
$$d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2 + d_6^2 = \text{sum of squared distances} = \text{SS(distances)}$$

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



SAO CHO SS(DISTANCE) LỚN NHẤT



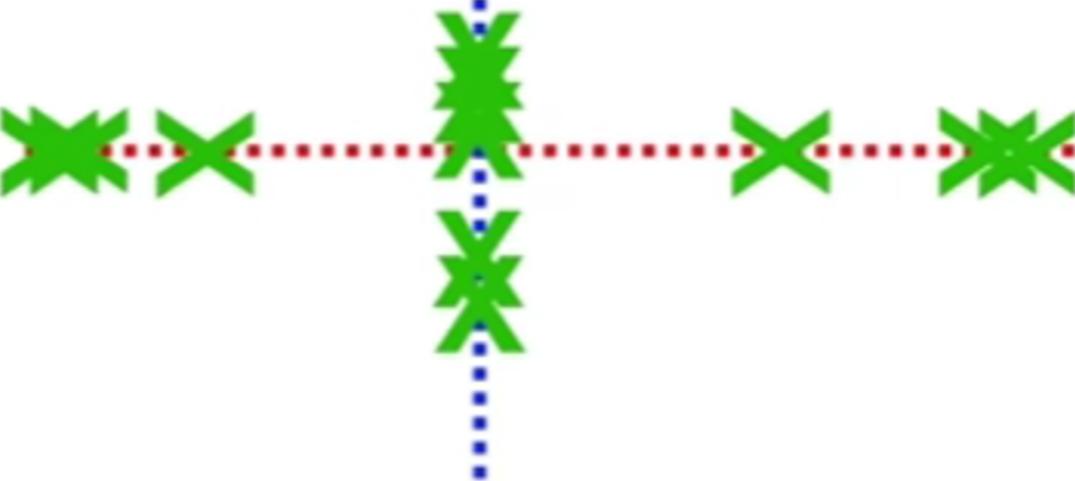


====>>

projected points
samples go in
^ plot.

PC2

PC1



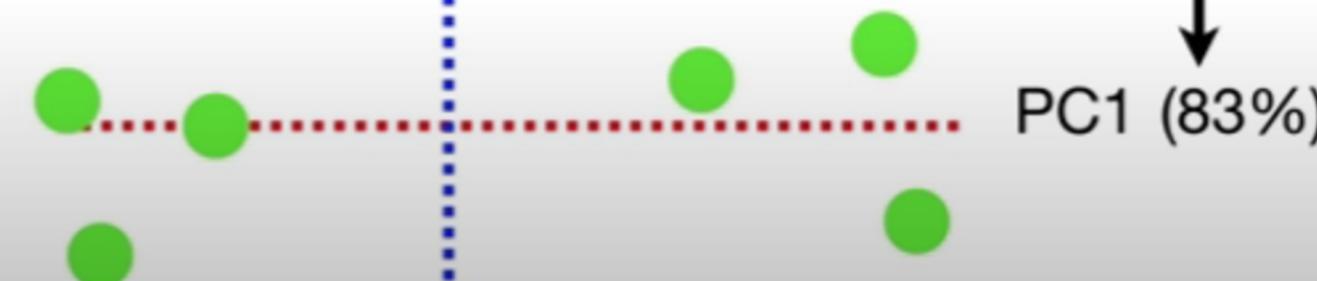
CÁCH TÍNH TỈ LỆ (%) CỦA CÁC PCA

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

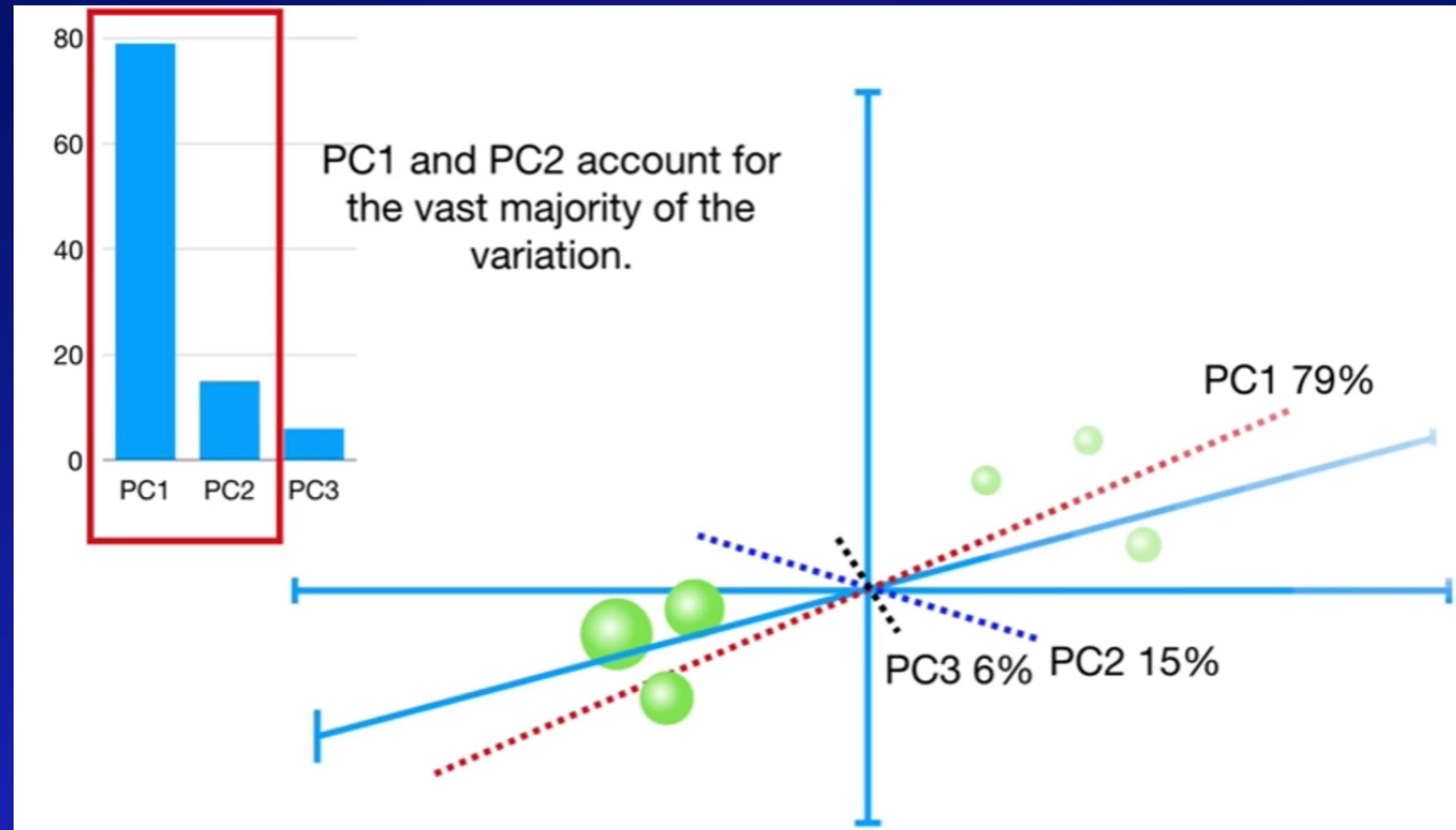
$$\frac{SS(\text{distances for PC1})}{n - 1} = \text{Variation for PC1}$$
$$\frac{SS(\text{distances for PC2})}{n - 1} = \text{Variation for PC2}$$

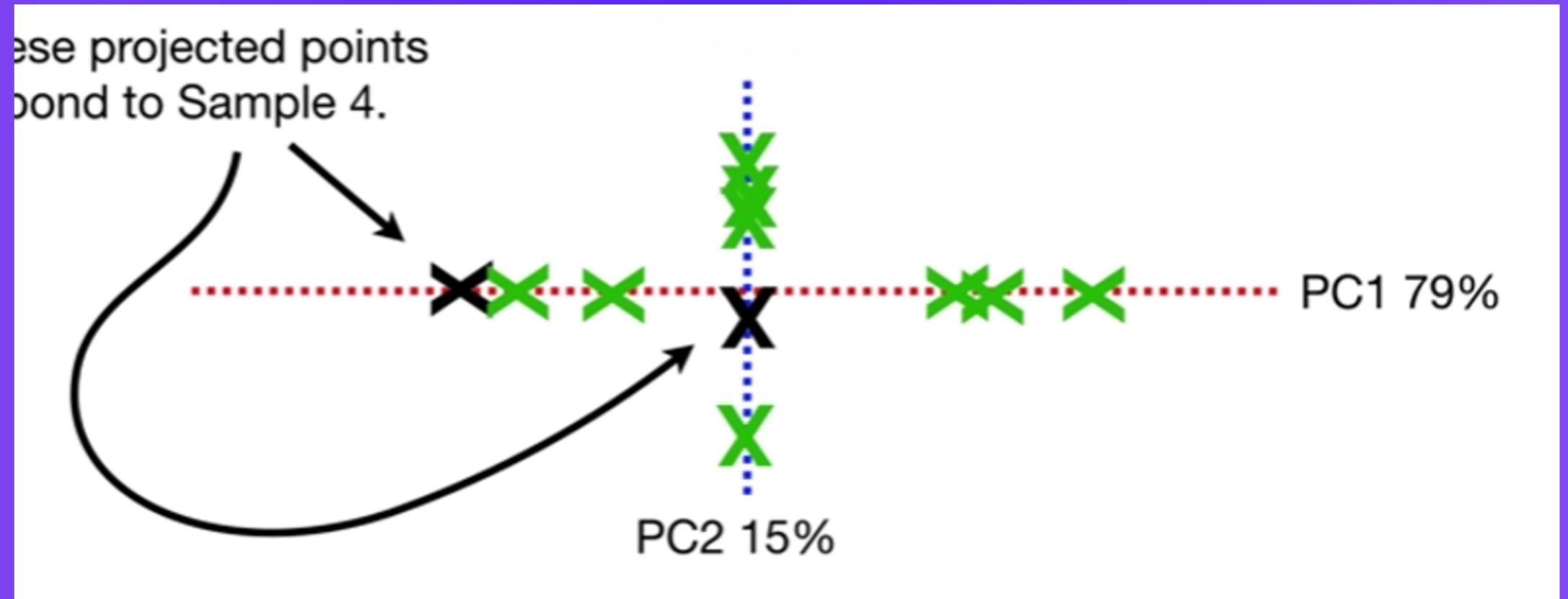
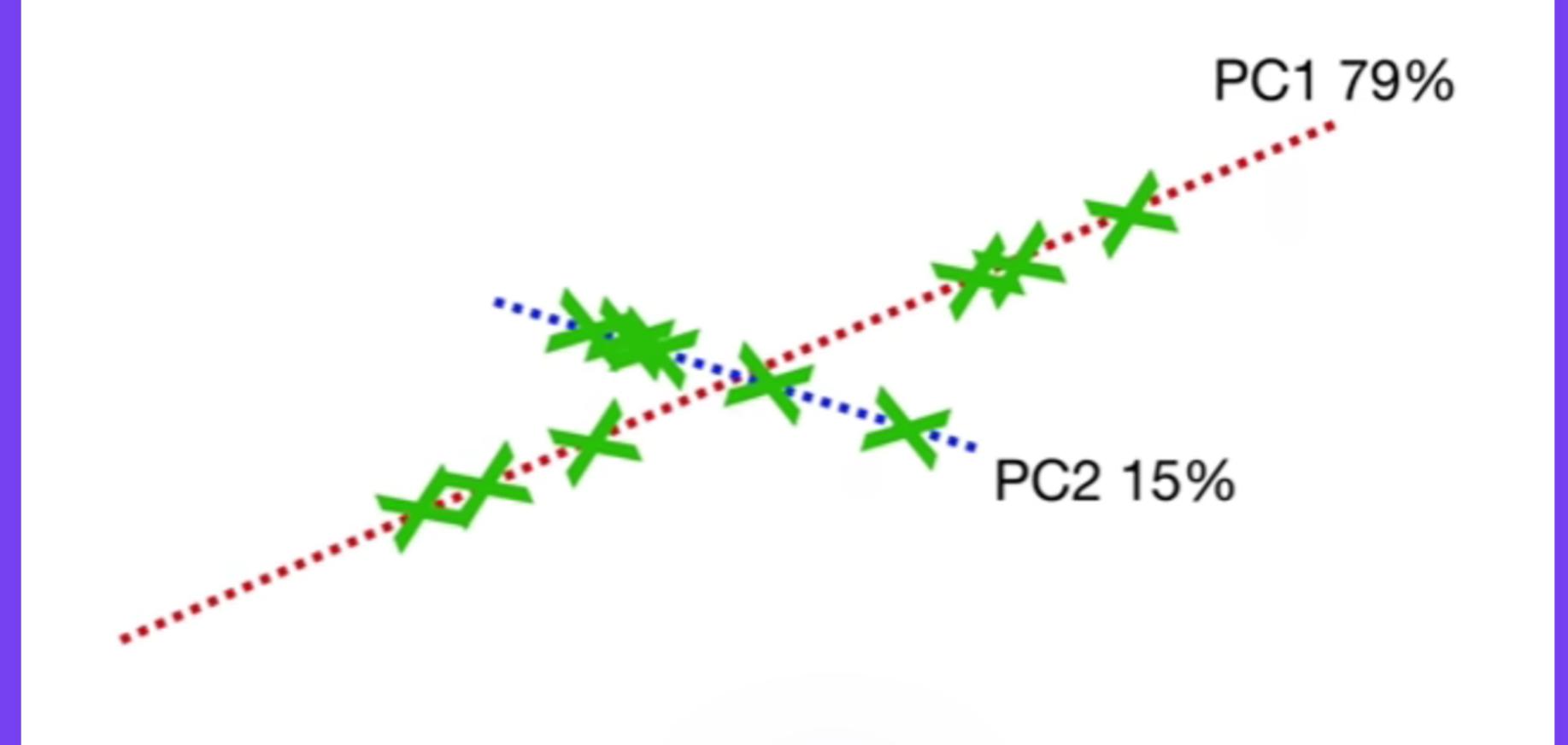
That means that the total variation around both PCs is **15 + 3 = 18...**

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



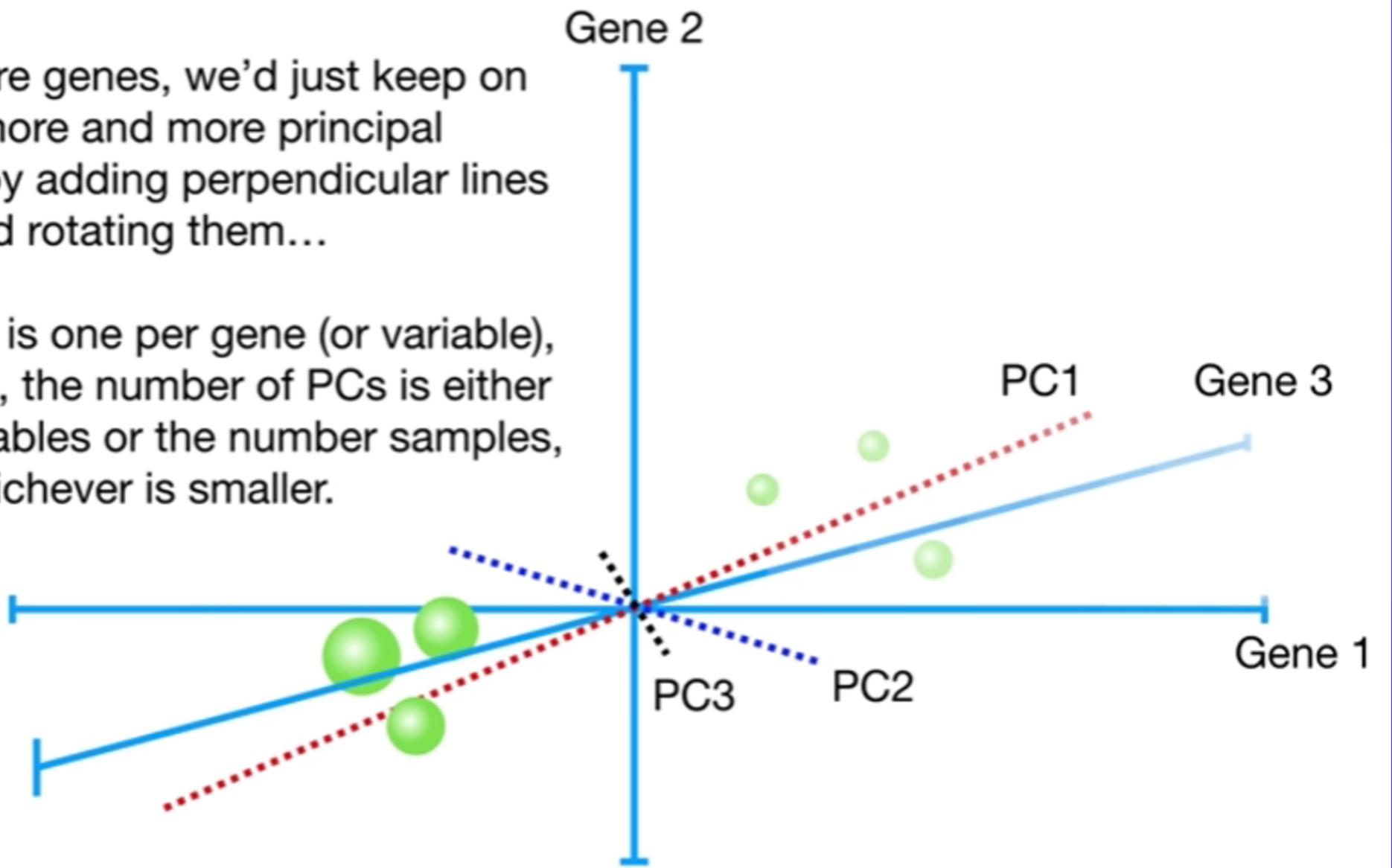
ĐỐI VỚI DỮ LIỆU NHIỀU CHIỀU HƠN





If we had more genes, we'd just keep on finding more and more principal components by adding perpendicular lines and rotating them...

In theory there is one per gene (or variable), in practice, the number of PCs is either the number of variables or the number samples, whichever is smaller.



THANK YOU!

