Nama    : Arfara Yema Samgusdian
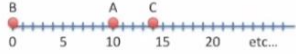
NIM     : 1103202004

Kelas   : TK-44-G4

## 1. Principal Component Analysis (PCA) Clearly Explained (2015)

2. ROC and AUC explained

$$\text{True Positive Rate} = \text{Sensitivity} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

——Actual——

|  |  | Is Obese | Is Not Obese |
|---|---|---|---|
| Predicted | Is Obese | 4 | 4 |
|  | Is Not Obese | 0 | 0 |

**ROC and AUC, Clearly Explained!**

And depending on how many **False Positives** I'm willing to accept, the optimal threshold is either this one…

True Positive Rate (Sensitivity)

0, 0    False Positive Rate (1 - Specificity)   1

**ROC and AUC, Clearly Explained!**

StatQuest wit… ✓

The **red** method… …is better than the **blue** method.

True Positive Rate (Sensitivity)

0, 0    False Positive Rate (1 - Specificity)   1

▶ ▶| 🔊   15:49 / 16:16 · Summary of concepts ›        ⬤ CC ⚙ 🖼 ▢ ⛶

**ROC and AUC, Clearly Explained!**

StatQuest wit… ✓   Join   🔔 Subscribed ⌄   👍 30K   👎   Share   ⬇ Download …

## 3. StatQuest: K-nearest neighbors, Clearly explained

Step 2: Add a new cell, with unknown category, to the PCA plot. We don't know this cell's category because it was taken from another tumor where the cells were not properly sorted.

**StatQuest: K-nearest neighbors, Clearly Explained**

Step 3: We classify the new cell by looking at the nearest annotated cells. (i.e. the "nearest neighbors").

If the "K" in "K-nearest neighbors" is equal to 1, then we only use the nearest neighbor to define the category.

In this case, the category is **GREEN**.

If K=11, we would use the 11 nearest neighbors.

In this case, the category is still **GREEN**.

2:05 / 5:30 • K-NN applied to scatterplot data

**StatQuest: K-nearest neighbors, Clearly Explained**

If K=11 and the new cell is between two (or more) categories, we simply pick the category that "gets the most votes".

In this case....

7 nearest neighbors are RED.
3 nearest neighbors are ORANGE.
1 nearest neighbor is GREEN.



**StatQuest: K-nearest neighbors, Clearly Explained**

StatQuest wit... ✓    Join    🔔 Subscribed ∨    👍 9.5K 👎    ➤ Share    ⬇ Download    ...

---

# A few thoughts on picking a value for "K"

- There is no physical or biological way to determine the best value for "K", so you may have to try out a few values before settling on one. Do this by pretending part of the training data is "unknown".

- Low values for K (like K=1 or K=2) can be noisy and subject to the effects of outliers.

- Large values for K smooth over things, but you don't want K to be so large that a category with only a few samples in it will always be out voted by other categories.

4:53 / 5:30 • Thoughts on how to pick 'K' >

**StatQuest: K-nearest neighbors, Clearly Explained**