

② → k-means clustering •
→ k-nearest neighbors -

• Exploratory Data Analysis (EDA)

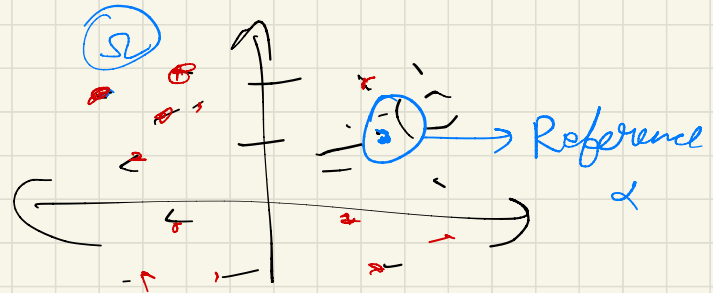
Virus $\xrightarrow{\text{RNA}}$ RT-PCR \rightarrow Data

5000 samples

NYC, Pregnant women \rightarrow covid testing

" Clustering "

Cluster \rightarrow Group



$(x_1, x_2, \dots, x_n) \in \mathbb{R}^d$ Partition $\rightarrow P_1, P_2, P_3, \dots, P_k$
 $\hookrightarrow \textcircled{S}$

$$P_1 \cup P_2 \cup \dots \cup P_k = S, \quad P_i \subseteq S \quad \forall i \in \{1, 2, \dots, k\}$$

$$\text{For any } i \neq j, \quad P_i \cap P_j = \emptyset$$

$\alpha, \beta, \gamma, \Omega, \textcircled{N/A}$
 $\uparrow \quad \uparrow \quad \uparrow \quad \uparrow \quad \uparrow$

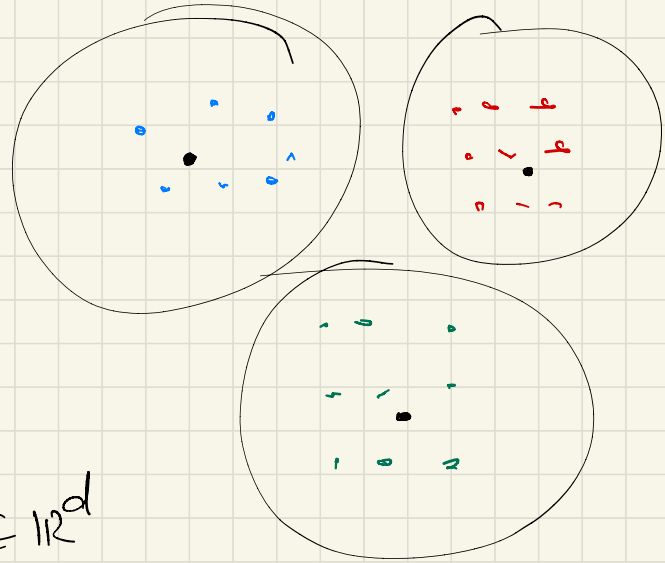
Unsupervised.

• Pandas.

"k"

"centroid"

k-means, k-medians, k-centers,
DBSCAN, k-medoids



Objective fn: $x_1, x_2, \dots, x_n \in \mathbb{R}^d$

Goal: $c_1, c_2, \dots, c_k \in \mathbb{R}^d$

C_i = set of all points in dataset that closest
to c_i Mean-squared error

$$* \text{cost}(x, C) = \sum_{i=1}^n (x_i - \underline{c(i)})^2, \quad c(i) = \min_{j=1}^k \|x_i - c_j\|^2$$

- NP-hard to exactly minimize this cost.

$$\min_{c \in \mathbb{R}^d} \sum_{i=1}^n (x_i - c)^2$$

$$\hookrightarrow \sum_{i=1}^n (\overbrace{x_i - \mu} + \overbrace{\mu - c})^2$$

mean
centroid

$$= \sum_{i=1}^n (x_i - \mu)^2 + (\mu - c)^2 + 2(x_i - \mu)(\mu - c)$$

$$= \underbrace{\sum_{i=1}^n (x_i - \mu)^2} + \underbrace{\sum (\mu - c)^2} + 2(\mu - c) \sum_{i=1}^n (x_i - \mu)$$

$$\text{cost}(c) = \underbrace{\sum_{i=1}^n (x_i - \mu)^2} + \underbrace{(\mu - c)^2 n}$$

$$(n)$$

$$(\mu)$$

$$\mu = \frac{\sum_{i=1}^n x_i}{n}$$

$$\text{cost}(C_1, C_2) = \sum_{i=1}^n \underbrace{(x_i - c_i)^2}$$

• Lloyd's method / "k-means" algorithm

- ① Assign : For every ^{data} point, assign it to its closest center
- ② Re-center : 'New centres' are the means of the clusters chosen in step ①.