# 1. K-means

Groups similar profiles into clusters

→ requires training data
→ can adapt to user patterns
Pros: Adapts   Cons: LOTS data
needed

## HYBRID Approach:

- SBERT initial embeddings
- K-means to cluster
  similar profiles
- Cosine-similarity → final
  match
  (inside clusters)

$O(n^2) \rightarrow O(n \times k)$
         K-means

→ efficient

→ SBERT embeds allows good
compat.

MongoDB  - usr id
         - name          } everything

→ embed - bio
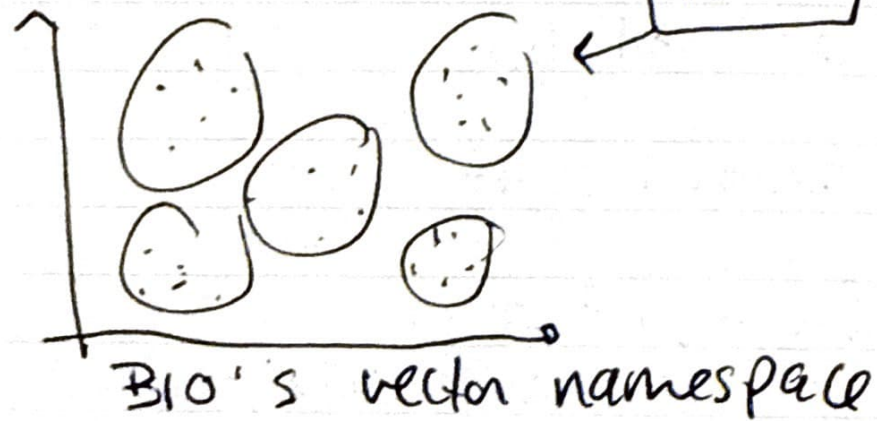
→ embed - pref

Pinecone ⎡ - id
         ⎢ - ~~vector~~ values = vector
         ⎢ - metadata = { vector type
         1          faculty
        second          edh
         ⎢          ...
         ⎣            .. }

⎡ *vector type changes
⎢ 1 index
⎣ 2 namespaces

# ① k-means : periodically

- pull embedding from ALL

- assign user bio → Clusters id



BIO's vector namespace

- store cluster-id in Mongo

# ② filter cluster id only



find simil. score

$$\pi_1 = user1.bio \quad v.s \quad user2.pref$$
$$\pi_2 = user2.bio \quad v.s \quad user1.pref$$
$$\pi = (\pi_1 \times \pi_2)/2$$

→ Meta_data filtering w/ weights

→ final score $= \alpha \cdot \pi + \beta \cdot Meta$

→ Sort by final scores

→ Provide top 3 best matches
↑
∞ N