

* FAISS (Facebook AI Similarity Search)

- It is a library for efficient similarity search and clustering of dense vectors.
- What is similarity search?
 - Given a set of vectors x_i in dimension d, FAISS builds a data structure in ^{RAM} from it. After the structure is constructed, when given a new vector x in dimension d it performs efficiently the operation:
$$j = \operatorname{argmin}_i \|x - x_i\|$$

where $\|\cdot\|$ is the Euclidean distance (L^2)

 - In FAISS terms, the data structure is an index, an object that has an add method to add x_i vectors.
NOTE: x_i are assumed to be fixed

→ Popular Indexes:

- | | |
|----------|----------|
| (1) Flat | (3) HNSW |
| (2) LSH | (4) IVF |

x_{db} = for the database, that contains all the vectors that must be indexed

x_{qry} = ~~query~~ query vector(s) for which we need to find the nearest neighbors

K = find K-nearest neighbors

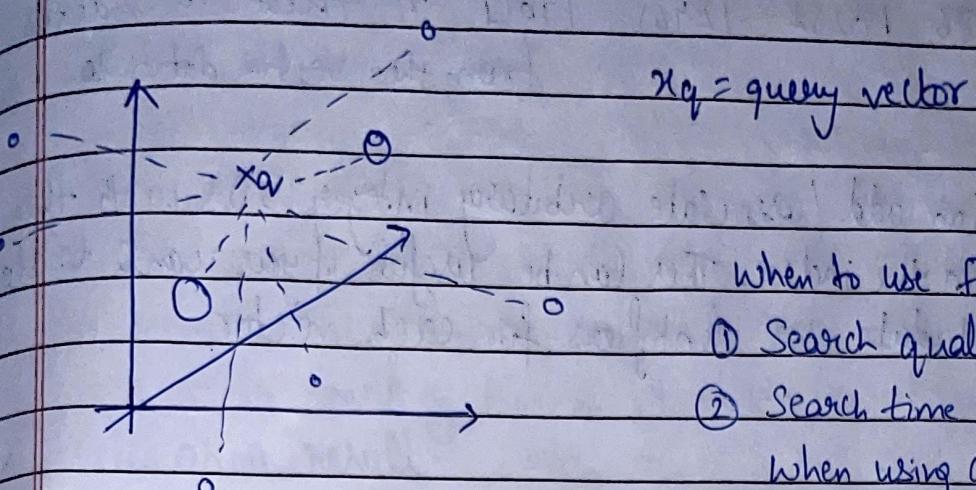
(KNN)

classmate

Date _____

Page _____

- ① Flat - Simplest version that just performs brute-force L2 distance search



x_q = query vector

When to use Flat index?

- ① Search quality is a high priority
- ② Search time doesn't matter OR
when using a small index ($< 10^4$)

~~Efficient database~~

~~stop words~~

Types:

- IndexFlatL2 - using Euclidean distance
- IndexFlatIP - using Inner Product distance

Code:

```
from sentence_transformers import SentenceTransformer
```

```
model = SentenceTransformer('model-name')
```

```
import faiss
```

```
sentence_embeddings = model.encode(data)
```

```
d = sentence_embeddings.shape[1]
```

```
index = faiss.IndexFlatL2(d)
```

```
index.is_trained # to check if index requires training. Since it's  
# flat index and is already trained, this is True
```

```
index.add(sentence_embeddings)
```

```
index.ntotal # no. of embeddings
```

```
 $x_q$  = model.encode('query sentence')
```

$K = 4$ # similar vectors need for input x_q

$D, I = \text{Index}. \text{Search}(x_q, K)$

$\text{print}(I)$

$\gg [[4586 \ 10252 \ 12465 \ 190]]$ # these are indexes
from the vector database

Note: you can add/associate arbitrary integer IDs with the
vectors in the index. This can be useful if you want to store
some metadata or identifiers for each vector.

How?

$\text{index} = \text{faiss}. \text{IndexFlatL2}(\underline{\text{df}})$

$\text{index} = \text{faiss}. \text{IndexIDMap}(\text{index})$

$\text{index}. \text{addWithIDs}(\text{embeddings}, \text{df}. \text{id}. \text{values})$

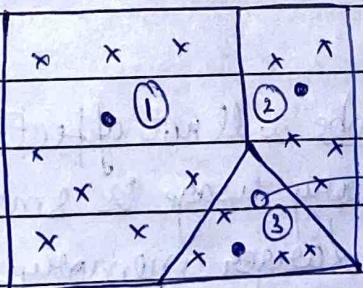
document
embeddings

dataframe whose IDs
we want to associate

② IVF \Rightarrow The Inverted File Index consists of search scope reduction through clustering.

- SVF works on Voronoi cell concept.
- Voronoi cells are regions of space that are defined by the distance to a given set of points such that each point in a cell is closer to its corresponding point than to any other point.
- One way to understand the difference b/w K-Means and voronoi cells is to think of Kmeans clustering as a process and voronoi cells as a result.
- Now, every datapoint will be contained within a cell and be assigned to that respective centroid.
- One Problem \Rightarrow if our query vector (x_q) lands near the edge of a cell, there's a good chance that its closest other datapoint is contained within a neighboring cell. we call this the [edge problem]

Eg.



query with $n\text{-probe}=1$, i.e. scope is restricted to (3) cell only

Soln: we increase $n\text{-probe}$ to > 1 , ex. to 5 now, the query will search for results in 5 nearest cells

> Python Code $d = 786$

```
n-list = 128 # no. of cells / clusters to partition data into
quantizer = faiss.IndexFlatIP(d) # dimension of vectors = d
index = faiss.IndexIVFFlat(quantizer, d, n-list)
index.train(sentence_embeddings) # we must train the index to cluster into cells
index.add(sentence_embeddings)
index.nprobe = 8 # set how many of nearest cells to search
D, I = index.search(xq, K)
```

- n-list means that we must compare our vector to more centroid vectors - but after selecting the nearest centroid's cells to search, there will be fewer vectors within each cell. So, increase n-list to prioritize search speed.
- n-probe → we find the opposite. Increase n-probe increases the search scope - thus prioritizing search quality.
- In terms of memory, nprobe will not affect this. The effect of nlist on memory usage is small too. Higher nlist means a marginally larger memory requirement.
- So, we must decide b/w greater search-quality with nprobe and faster search speed with nlist.

What is the role quantization?

① Indexing Process

- During indexing (adding vectors to IndexIVF), each input vector is assigned a quantisation index by the quantiser (which is an instance of another FAISS Index, such as IndexFlatIP, IndexFlatL2 or IndexIVPQ1, IndexPer)
- The quantisation index determines which inverted list (cell) the vector belongs to
- The group Id (quantisation Index) maps to a specific inverted list (cell), where the Id (vector Index) of the vector is stored
- The inverted list (cell) object is required only after training

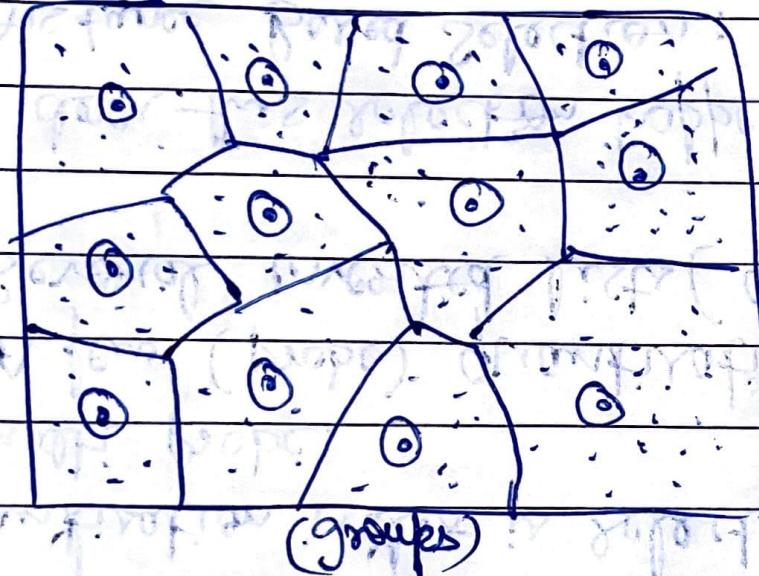
② Search Process

- During search, when a query vector is quantised using the same quantiser, it obtains a quantisation Index (group Id)
- The search process involves looking up the inverted list (cell) corresponding to this quantisation index
- The vectors within inverted list^(cell) are then considered for further distance computation

③ Multi-probe

- By default, only the inverted list corresponding to the quantisation index is selected
 - In multi-probe:
 - 3.1 A few (probe) quantisation indices are selected
 - 3.2 Several inverted lists (cells) are visited during search
- How does this selection happen?
↳ Distance Based Selection: ~~choose~~

(Indexing process)



(Typically
an integer
value)

→ Quantized Index / group Id
Vector Id
of centroid

S.No	Q.I.	S.I.
1	0	-
2	0	-
3	0	-
4	0	-
5	1	-
6	1	-
7	1	-
8	1	-
9	1	-
10	1	-
11	1	-
12	1	1

- Vector Id
- Id of centroid mapped to quantized index

Each Id is mapped to the original vector Ids within its cluster

$x_q \rightarrow$ Input query $\xrightarrow{\text{Converted to}}$ Quantized Index \rightarrow Gets the group Id
(multiple group Id in case of m-probe > 1)

(Search process)

Search in those clusters to get answer