

## Nearest Neighbor Search

### Retrieval as k-nearest neighbor search

#### Retrieving documents of interest

##### Document retrieval

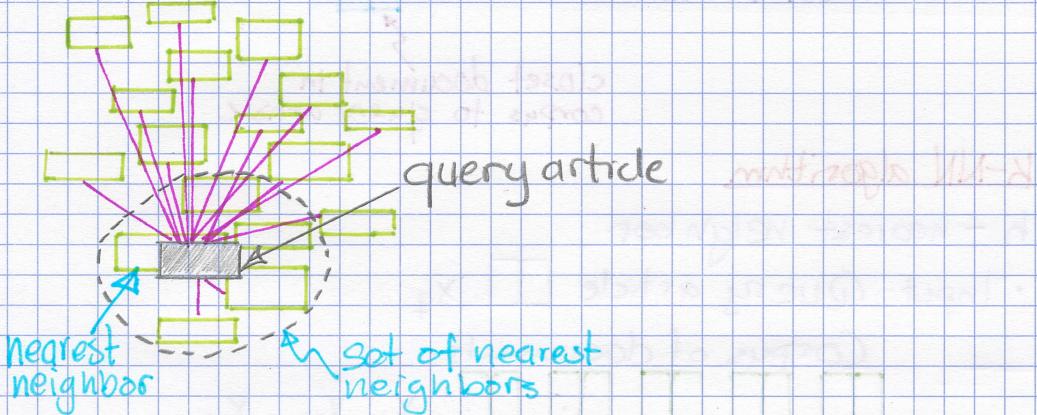
- Currently reading article you like
- Goal: Want to find similar article

##### Challenges

- How do we measure similarity?
- How do we search over articles?

#### 1-NN search for retrieval

Space of all articles, organized by similarity of text



#### 1-NN algorithm

##### 1 - Nearest neighbor

- Input: Query article  $\boxed{\quad}$ :  $x_q$

Corpus of documents ( $N$  docs)

$\boxed{\quad} \boxed{\quad} \boxed{\quad} \boxed{\quad} \boxed{\quad} \boxed{\quad} : x_1, x_2, \dots, x_N$

- Output: Most similar article  $\boxed{\quad} \leftarrow x^{NN}$

Formally:  $x^{NN} = \min_{x_i} \text{distance}(x_q, x_i)$

### 1-NN algorithm

Initialize  $\text{Dist2NN} = \infty$ ,  $\boxed{\quad} = \emptyset$

For  $i = 1, 2, \dots, N$

Compute:  $\delta = \text{distance}(\boxed{\quad}, \boxed{\quad})$

If  $\delta < \text{Dist2NN}$

set  $\boxed{\quad} = \boxed{x_i}$

set  $\text{Dist2NN} = \delta$

Return most similar document

closest document

$\boxed{\quad} = \emptyset$

$x_i$

query document  
 $x_q$

$\boxed{\quad}$

document  $i$   
from corpus

$\boxed{\quad}$

$x^{NN}$

closest document in  
corpus to query article

### K-NN algorithm

$k$  - nearest neighbor

• Input: Query article  $\boxed{\quad} : x_q$

Corpus of documents

$\boxed{\quad} \quad \boxed{\quad} \quad \boxed{\quad} \quad \boxed{\quad} \quad \boxed{\quad} : x_1, x_2, \dots, x_N$

• Output: List of  $k$  similar articles

$\boxed{\quad} \quad \boxed{\quad} \quad \boxed{x_i} \quad \boxed{\quad} \quad \boxed{\quad} \quad \boxed{\quad} \quad \boxed{\quad} \quad \boxed{\quad}$

Formally

$$x^{NN} = \{x^{NN_1}, \dots, x^{NN_k}\}$$

for all  $x_i$  not in  $x^{NN}$ ,  $\text{distance}(x_i, x_q) \geq$

$$\geq \max_{x^{NN_j}, j=1, \dots, k} \text{distance}(x^{NN_j}, x_q)$$

K-NN algorithm

Initialize  $\text{Dist2kNN} = \text{sort}(\delta_1, \dots, \delta_k)$  ← list of sorted distances



sort first k documents by distance to query doc

= sort( $\boxed{\quad}_1, \boxed{\quad}_2, \dots, \boxed{\quad}_k$ ) ← list of sorted docs

For  $i = h+1, \dots, N$

Compute:  $\delta = \text{distance}(\boxed{\quad}_i, \boxed{\quad}_q)$

If  $\delta < \text{Dist2kNN}[k]$  ← distance to  $k^{\text{th}}$  NN (furthest NN in set)

find  $j$  such that  $\delta > \text{Dist2kNN}[j-1]$  but  
but  $\delta < \text{Dist2kNN}[j]$   
remove furthest article and shift queue.

inserting new article

$$\boxed{\quad}_{[j+1:k]} = \boxed{\quad}_{[j:k-1]}$$

$$\text{Dist2kNN}[j+1:k] = \text{Dist2kNN}[j:k-1]$$

Set  $\text{Dist2kNN} = \delta$  and  $\boxed{\quad}_{[j]} = \boxed{\quad}_i$

Return  $k$  most similar articles



closest  $k$  docs to query doc



Critical elements of NN search

Item (e.g., doc) representation

$$x_q \leftarrow \boxed{\quad}$$

Measure of distance between items.

$$\delta = \text{distance}(x_i, x_q)$$