

## Similarity Search Queries

### 1. Range queries

Given a query point (vector)

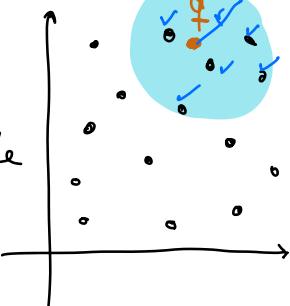
$q$

Find all the points in the vicinity of it

Table T:

Each black dot represents a tuple (vector) in Table T

Table T



Formally

Find

$$\{t \in T \mid d(q, t) \leq r\}$$

## SQL Extension

```
SELECT *
FROM T
WHERE d(q) <= r
```

e.g., Pgvector standard for  $\ell_2$ -norm

```
Select *
from T
where
```

Column Names

Embedding  $\leftrightarrow [3, 5, 8, -1] < 5$

$d$        $q$        $r$

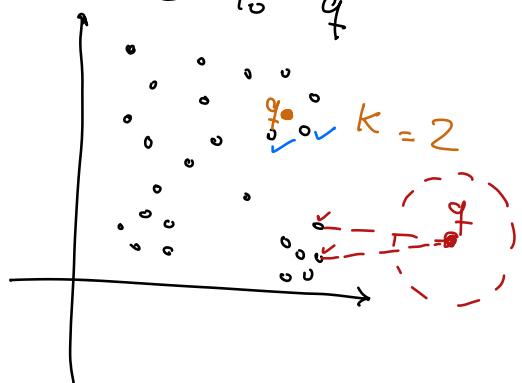
### 2 - nearest neighbor queries

Given a query point

$q$

, and a value  $k$

Find the  $k$  nearest points in the Table to  $q$



Formally

Find

$$k \operatorname{argmin}_{t \in T} d(t, q)$$

## SQL Extension

```
SELECT *
From T
order by d(q)
Limit k
```

In, Pgvector, Euclidean ( $\ell_2$ -norm)  
Select \* from T  
Order by

Embedding  $\leftrightarrow [-1, 3, -1, 2]$   
Limit 10  
 $\hookrightarrow k$

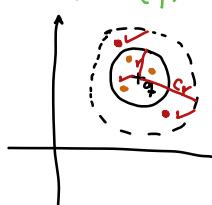
## Efficient Processing of K-NN queries

Exact: Given a value  $k$ , Find the  $k$ -nearest neighbor of  $Q$   
 $\rightarrow S = \underset{x \in T}{\text{argmin}} d(q, x)$

Approximate: Returns a set  $S'$  which may not be the exact KNN, but  $S'$  is "close" enough

measuring Approximation

(1) by distance  
 $\forall x' \in S' \quad \text{dist}(q, x') \leq C \cdot \text{dist}(x, q)$



(2) by recall  
 Recall@ $K$ : what Percentage of the returned Points belong to the KNN Set

$$\text{Recall@}K = \frac{|S \cap S'|}{K}$$

Question: How to find KNN exactly?

Step 1:

Compute  $d(q, x)$  for all points in the Table  $\Theta(nd)$

Number of Rows  
Number of Columns

Step 2:  
 apply the extended version of the median Finding Alg. on the scores to find the top- $K$  (K-NN)

$\Theta(n)$

$\Rightarrow$  Total time Complexity:  $\Theta(nd)$

X

- This is very slow, particularly when  $n$  is large

$\Rightarrow$  Question: Can we preprocess the database & Create Indices that enable Sublinear query answering at least approximately?

Indexing the Database for fast ANN query answering

- To Partition The Data  
↳ How?

- Build Some Data Structures on top of the Data

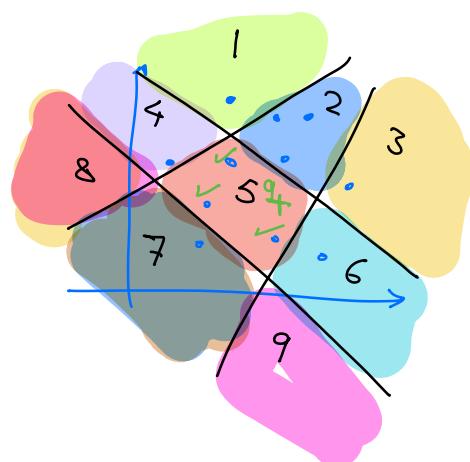
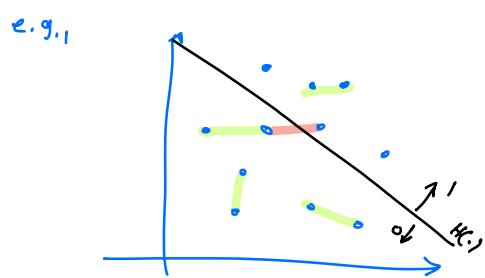
Data Structures

- Tables (Partition the Data)
- Trees
- Graphs
- Hybrid Data Structures

## Local Sensitive Hashing (LSH)

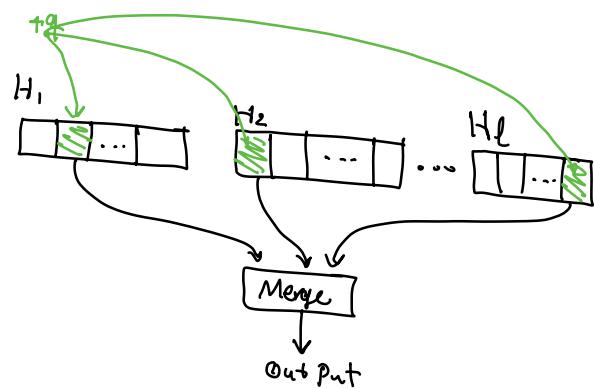
A hash function  $H$  is in class of LSH, if

- $\forall a, b, \text{dist}(a, b) \leq r$   
 $P(H(a) = H(b)) \geq P_1$
  - $\forall a, b, \text{dist}(a, b) \geq Cr$   
 $P(H(a) = H(b)) < P_2$
- $$P_1 > P_2$$



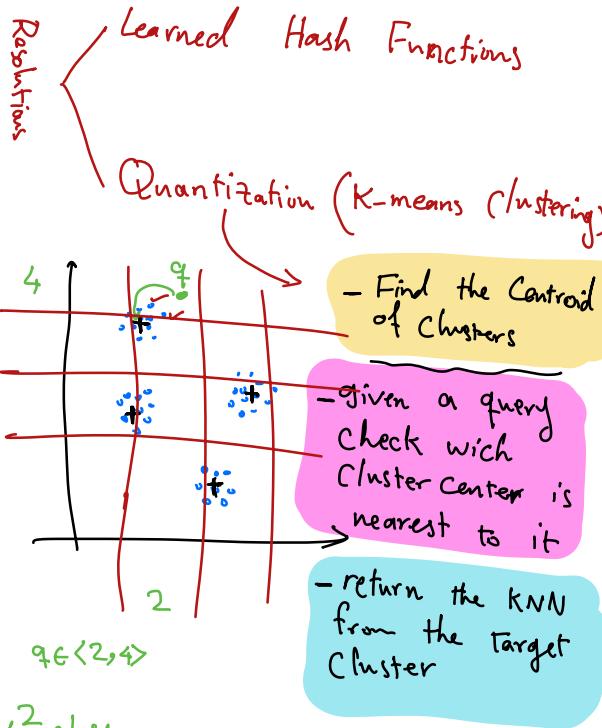
Using One Hash function,  
it is likely that we miss  
Some of the NNs

\*Resolution: Use multiple hash functions



This is called LSH

- \* When Data Points are **not evenly distributed** in the Space (when many points are in close vicinities), LSH loses its power.
- ↳ LSH does not Prune the Search Space in these cases.

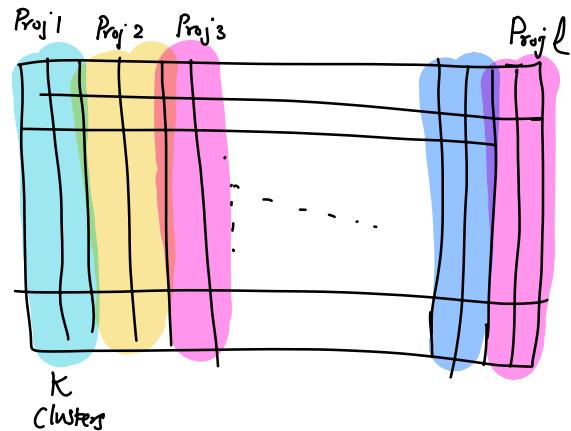


\* Issue 1:  $K$  Should be  $\mathcal{O}(\frac{n}{\log n})$  for this approach to have  $\mathcal{O}(\log n)$  query time

\* Issue 2: Curse of Dimensionality

Resolution:

Projection on lower Dimensions



$$\Rightarrow \text{No. Clusters} = k^l$$

$$k^l = \frac{n}{\log n}$$

$$\Rightarrow l = \log_k n - \log \log n$$

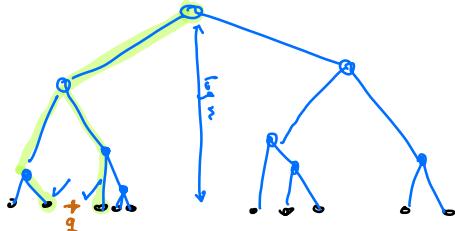
→ Given a query Point:

- Step 1: to find the Cluster Centers that  $q$  falls into

$$q \in \langle c_1, c_2, \dots, c_l \rangle$$

$c_i$ : is the closest Cluster Center to  $q$ , On projection  $i$

- Step 2: Check the Points within the same Cluster & return the top-k



To Construct the k-NN Tree

Construct Tree ( $S, i$ )

if  $|S| \leq T$   $\leftarrow$  Size Threshold  
return  $S$

Rule  $\leftarrow$  Select Rule ( $S, i$ )

The rule for Partitioning  $S$

Left  $\leftarrow \{x \in S \mid \text{Rule}(x) = \text{False}\}$

Right  $\leftarrow \{x \in S \mid \text{Rule}(x) = \text{True}\}$

return (rule, Construct Tree (Left,  $i-1$ ),  
Construct Tree (Right,  $i-1$ ))

(A) Round Robin on Columns

Select Rule<sub>RR</sub> ( $S, i$ )

$j \leftarrow$  median of  $S$  based on  
Column( $i \bmod d$ )

Rule  $\leftarrow (x[i] \leq j)$

e.g.

$S$	$i=1$
$t_1$	3
$t_2$	-1
$t_3$	5
$t_4$	8
$t_5$	0

$\{ -1, 0, 3 \}$   
 $\{ t_2, t_5, t_1 \}$   
 $\{ 5, 8 \}$   
 $\{ t_3, t_4 \}$

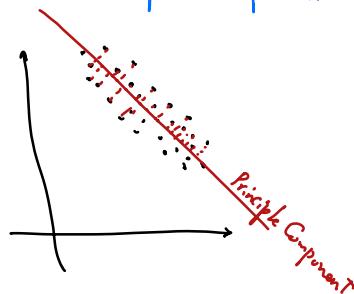
Observation: This is not a  
Binary Search Tree!

$\Rightarrow$  We initially need to check  
both Left & Right  
Branches to an answer of  
This can become  
Linear  $\Theta(nD)$

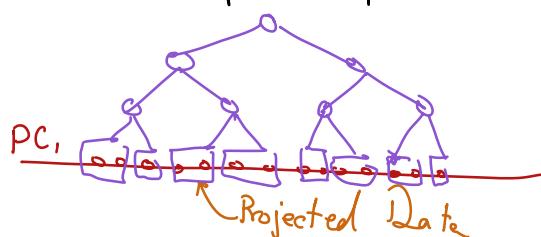
A Practical Resolution

- Use PCA (Principal Component Analysis)

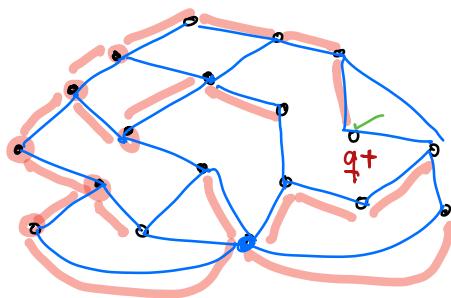
e.g.



- Instead of creating the tree based on individual columns,  
constructs the tree on  
projection(s) of data  
on principle components



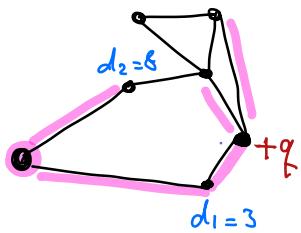
## Graph Indices



### Approach 1:

- Start from an arbitrary point in the Graph
- follow a BFS Traversal & find the NN of  $q$ .

$\Theta(nd) \times$



### Approach 2: Greedy Search (DFS)

- at each iteration, Next  $\leftarrow$  the neighbor with Smallest distance to  $q$
- Continue the search on Next
- Stop when all neighbors have larger distances

## Time Complexity of Greedy Search

$$\Theta(L \deg(u))$$

↓ length of the DFS path  
 ↓ degree of the nodes of the

## Properties of the Graph

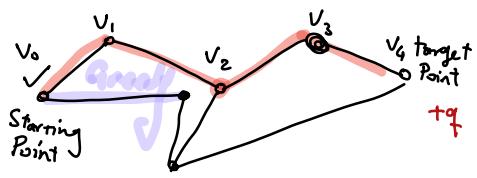
1- G Should be Sparse  
 $\deg(u) \sim \sim$  Small

2-  $l$  Should be Small.  
 after a Small number of hops, the alg. Stops

3- when the Greedy Stops a Close-to-NN Point is found

How about connecting each Point to its  $k$  nearest neighbors?

Monotonicity of the Search Path.



There should exist a monotonically decreasing Path on  $d(x, q)$  from the Starting Point to  $NN(q)$

$$\forall i \leq l$$

$$d(v_i, q) > d(v_{i+1}, q)$$

### \* Monotonic Search Network

A Graph, where the Greedy Path Satisfies the monotonicity Property.

→ if  $d$  is small ( $d \leq 3$ )

Such a graph with low degrees exist

↳ Use Voronoi Diagram

→ what if  $d$  is not small?

### Navigable Small-world Graph

↳ Should be small  $\Theta(\log n)$

NSW

→ In Practice

KNNG:

- Iteratively add the nodes to the Graph.
- for each node Connect it to its  $k$ -nearest neighbors

### HNSW (Hierarchical NSW)

