

KDD 2021 Tutorial

High-Dimensional Similarity Query Processing for Data Science

Jianbin Qin

Shenzhen Institute of Computing
Sciences
Shenzhen University

Wei Wang

Hong Kong University of
Science and Technology
(Guangzhou)

Chuan Xiao

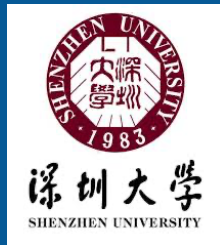
Osaka University and
Nagoya University

Ying Zhang

University of
Technology Sydney

Yaoshu Wang

Shenzhen Institute of Computing
Sciences
Shenzhen University



Outline

2

- Introduction
- Exact Query Processing
- Approximate Query Processing
- Selectivity Estimation
- Open Problems

Introduction

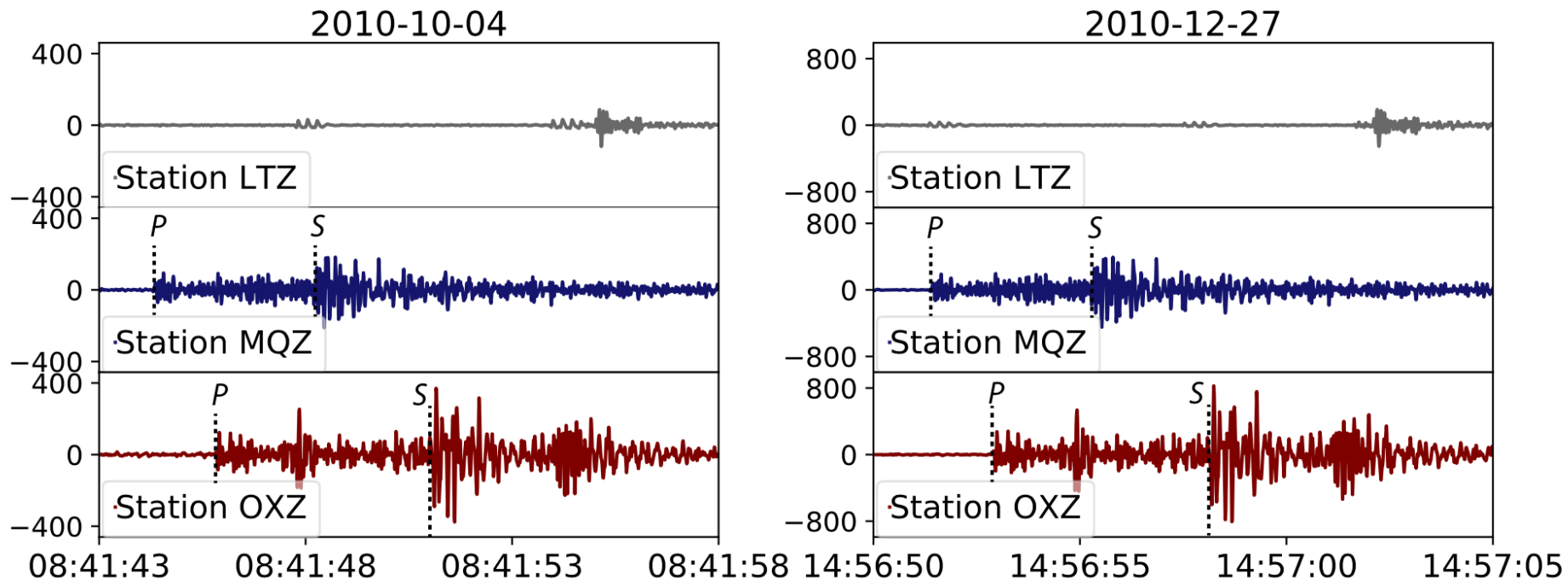
3

- High-dimensional data is abundant
 - ▣ Traditional sources:
 - **Time-series** [EZPB19], scientific applications
 - Document, multimedia, strings, feature vectors
 - ▣ New data sources:
 - Embedding from deep learning models
- Growing size and complexity
 - ▣ Web, social network, IoT
 - ▣ NOAA (USA) collects 100TB sensing data / day for weather forecasting
 - ▣ A variety of similarity/distance functions concerned

Example: Scientific Applications

4

- High-dimensional data in huge volumes in scientific domains [YHEB17, RYBE+18]

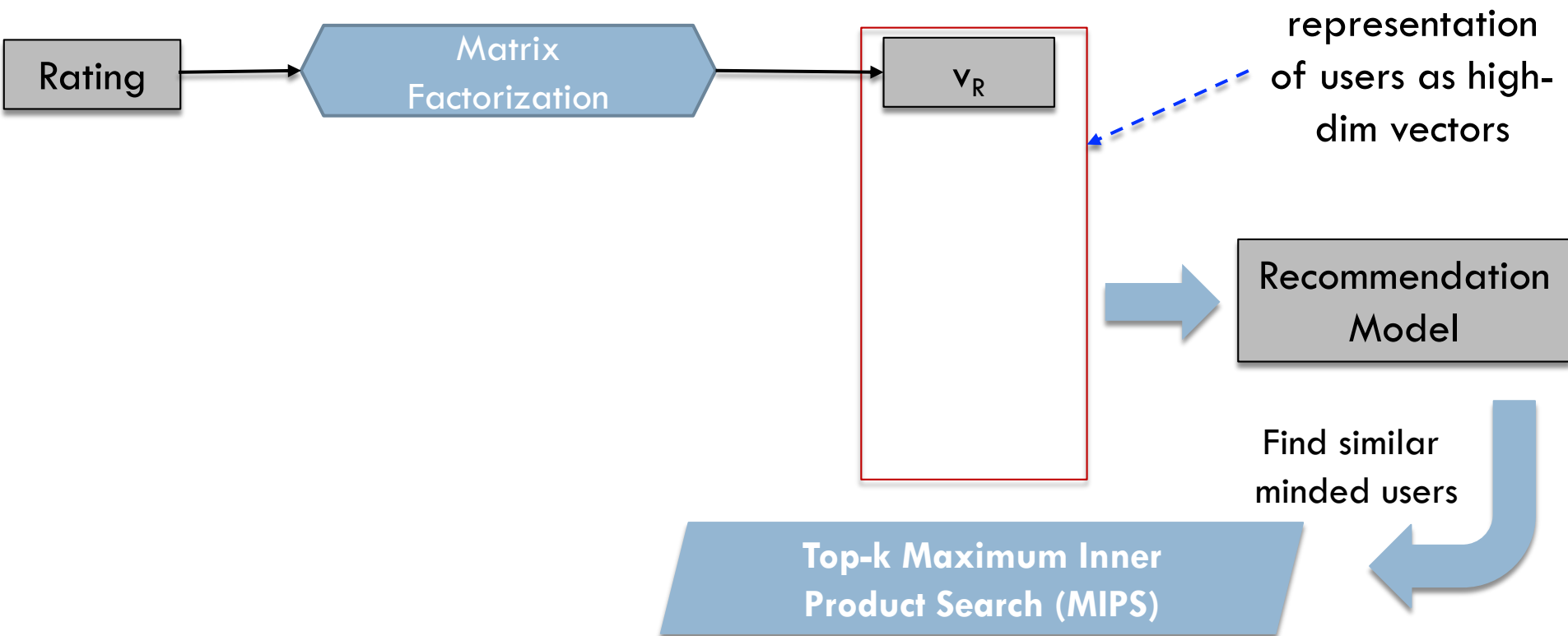


Research Question: Whether the magnitude 4.7 earthquake in Arkansas 2011 was caused by wastewater injection

Example: Embedding Vectors

5

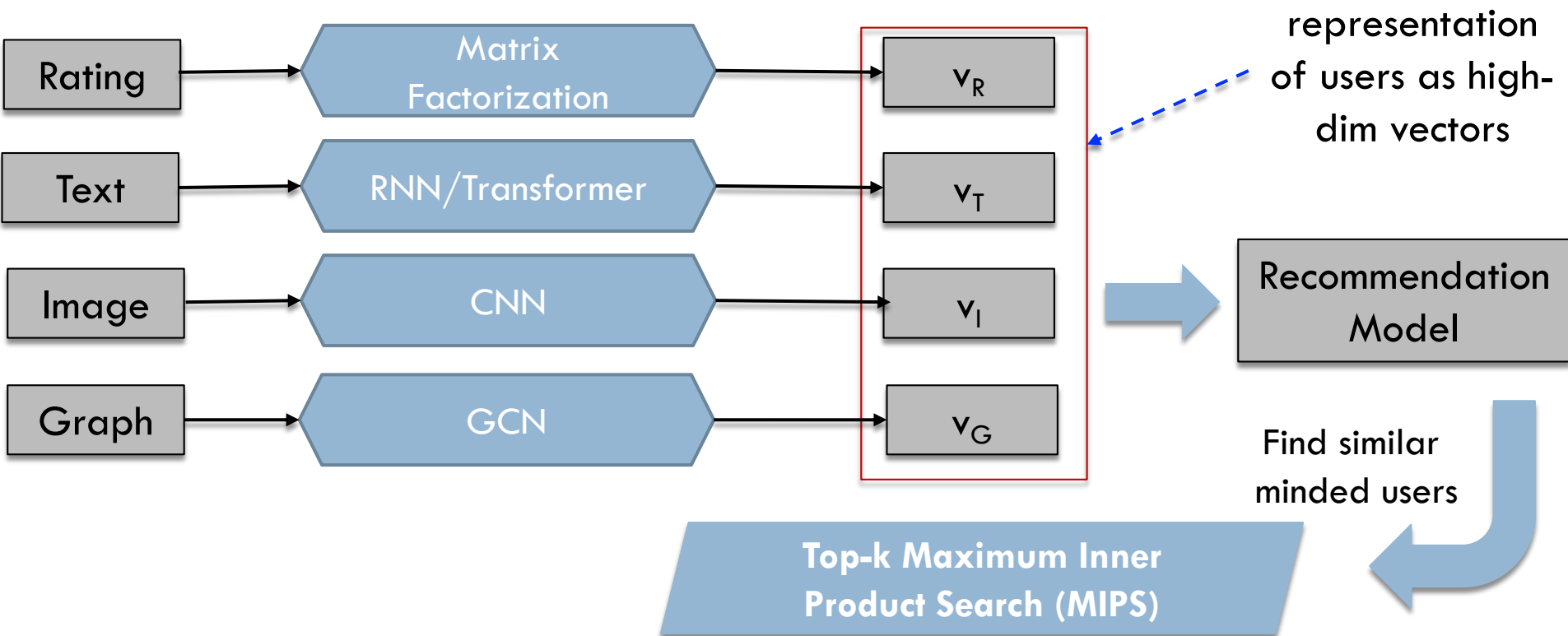
- DL presents a **unified** and **engineering-friendly** way to handle various information sources
 - ▣ Representation learning: e.g., **embedding**



Example: Embedding Vectors

6

- DL presents a **unified** and **engineering-friendly** way to handle various information sources
 - ▣ Representation learning: e.g., **embedding**



Example: Usage in Machine/Deep Learning

7

□ Kernel trick

- φ : mapping **low-dim** feature vectors to **high-dim** vectors

$$\langle \varphi(x), \varphi(x') \rangle = \mathcal{K}(x, x')$$

□ Feature hashing trick

- φ : **random** mapping **high-dim** feature vectors to **low-dim** vectors

$$\mathbb{E} [\langle \varphi(x), \varphi(x') \rangle] = \langle x, x' \rangle$$

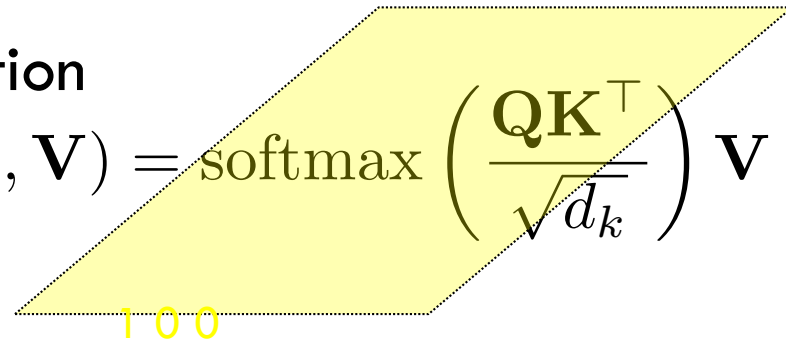
Improves efficiency, scalability and sometimes effectiveness

Example: Usage in Machine/Deep Learning

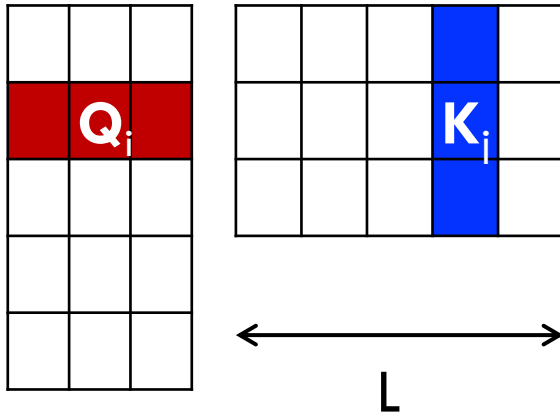
8

- Reformer [KKL20]
 - Speed up self-attention

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \right) \mathbf{V}$$



Time
flies
like
an
arrow



Find batch top-k \mathbf{K}_i 's for each \mathbf{Q}_i

Scale to long sequences, $O(L \log L)$
instead of $O(L^2)$

Usage in Machine/Deep Learning

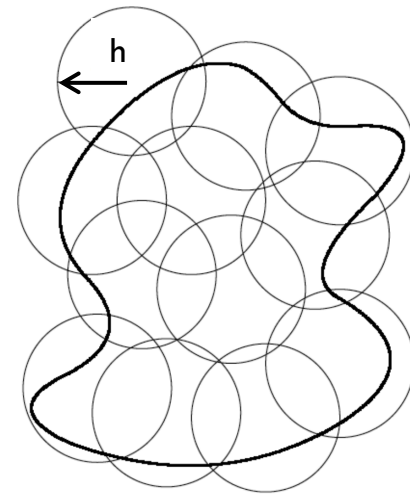
9

□ Q-learning with nearest neighbor [SX18]

▣ Idea:

- quantization of the **state space** X into $\{c_i\}_{i=1}^N$
- (non-parametric) kernel ridge regression for new (x, a) values

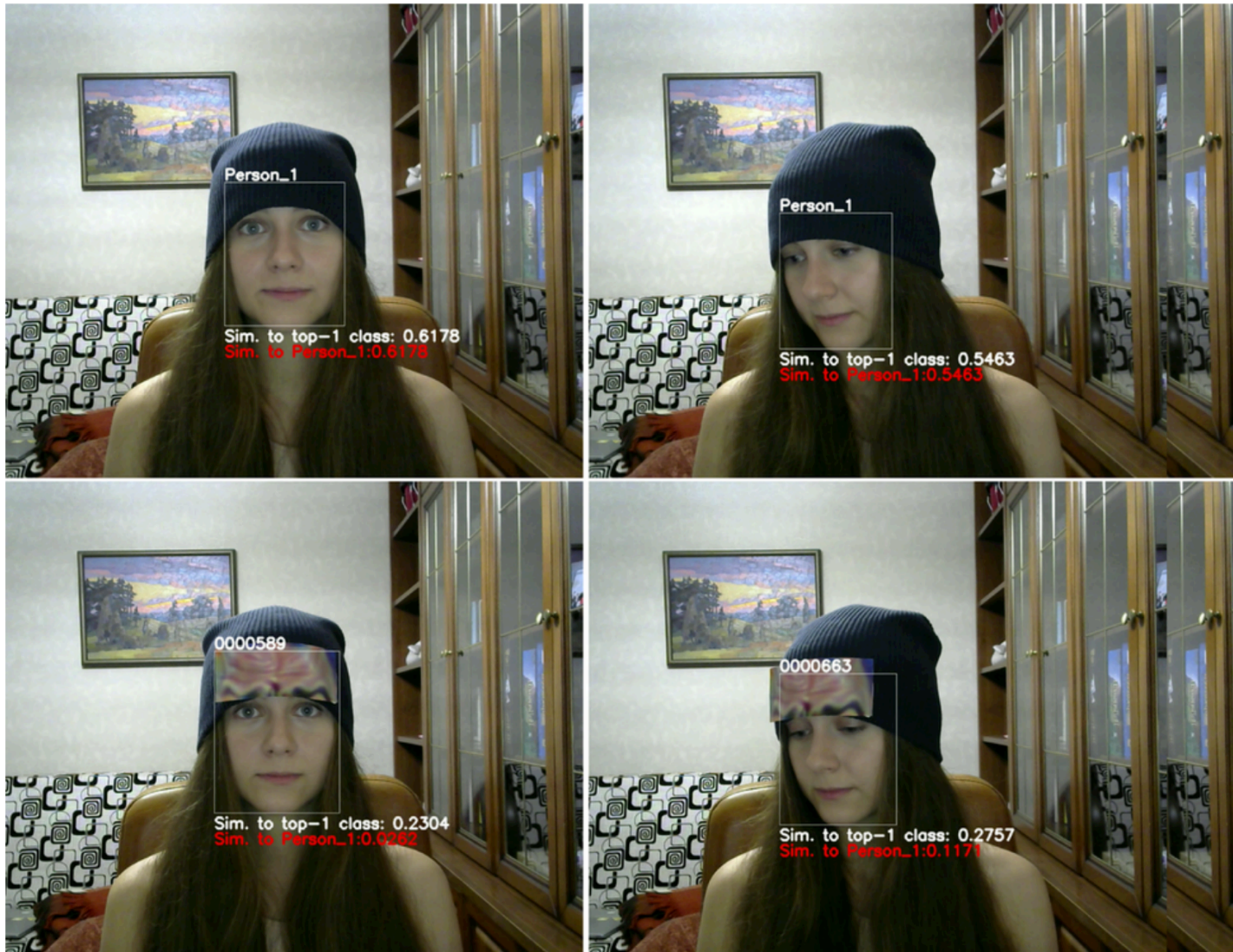
$$\begin{aligned}\hat{q}(x, a) &= \sum_{i=1}^n K(x, c_i) q(c_i, a) \\ &= \sum_{K(c, x) \leq h} K(x, c) q(c, a)\end{aligned}$$



Reinforcement Learning has been used in several DB problems, including Neo query optimizer [MNMZ+19]

Example: Adversarial Machine Learning

10



$\text{sim} \geq 0.54$

Adversarial sticker
on the forehead

$\text{sim} \leq 0.28$

[KP19]

Example: Adversarial Machine Learning

11

- Local intrinsic dimensionality (LID) is an important feature to detect adversarial examples [MLWE+18]

$$\widehat{\text{LID}}(x) = - \left(\frac{1}{k} \sum_{i=1}^k \log \frac{r_i(x)}{r_k(x)} \right)^{-1}$$

- High-dimensional geometry explains the existence of adversarial examples [GMFS+18]

require kNN queries

kNN queries are also useful in

- outlier/novelty detection
- kNN classification
- zero/few-shot learning

Example: Few-shot Learning

12

- Classify test images where there are few learning examples in the training data for each class
- Nearest neighbor classifier with learned embedding outperforms sophisticated methods [WCWM19]

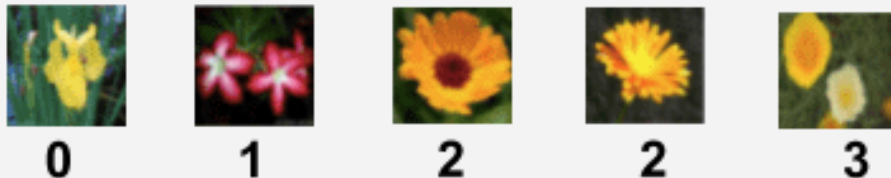
Support



Query



Support



Query



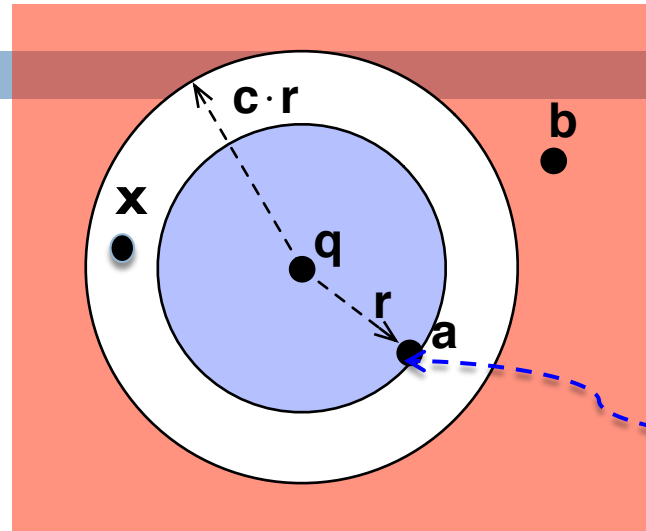
Problem Definitions

14

- Database object & the query
 - d-dimensional point/vectors $\in \mathbb{R}^d$
- Distance or similarity functions
 - $\text{dist}(u, v)$ L_p distance ($0 < p \leq 2, \infty$), Hamming dist, edit dist ...
 - $\text{sim}(u, v)$ cosine similarity/inner product, Jaccard
- Query types
 - k-nearest neighbor queries (kNN)
 - range queries
 - conjunctive queries
 - similarity/distance join queries (top-k, range, closest pair, containment, ...)

NN and kNN

15



d-dimensional space

$$D = \{a, b, x\}$$

$$NN = a$$

$$\text{top2 NN} = \{a, x\}$$

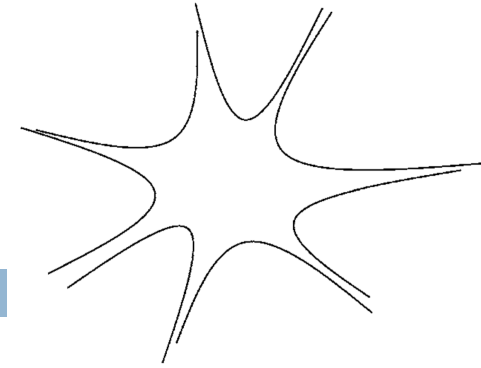
$$1.3\text{-ANN} = \{a, x\}$$

dist() is typically L_2 distance

- Nearest Neighbor (of q): o^*
 - ▣ $\text{dist}(o^*, q) = \min \{ \text{dist}(o, q), o \in D \}$
 - ▣ Generalizes to k-NN
- c -Approximate NN: o
 - ▣ $\text{dist}(o, q) \leq c * \text{dist}(o^*, q)$

Challenges / 1

high-dimensional
convex body



16

□ Non-intuitive high-dimensional Geometry

▣ Sampling uniformly within a unit hypercube → samples are within a thin ε 'shell'

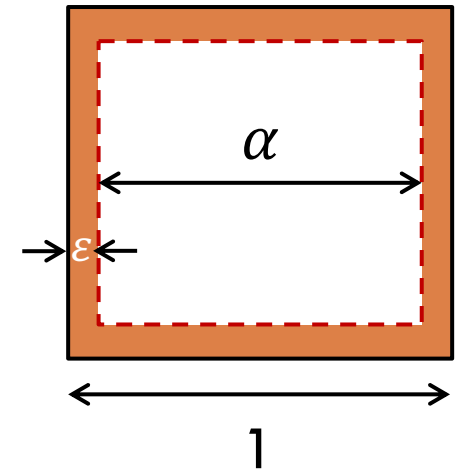
▣ $\text{Vol}(r) = \alpha^d \approx e^{-2\varepsilon d} \rightarrow 0$ ($\alpha < 1$)

▣ Angle between two vectors

▣ random Radamacher vectors →

$$\Pr \left[|\cos(\theta_{x,y})| > \sqrt{\frac{\log c}{d}} \right] < \frac{1}{c}$$

orthogonal w.h.p



Challenges / 1

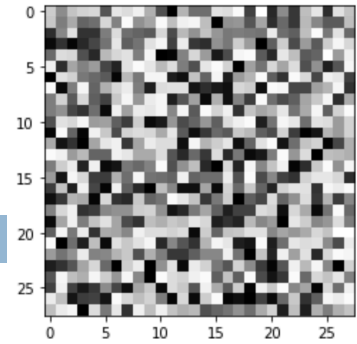
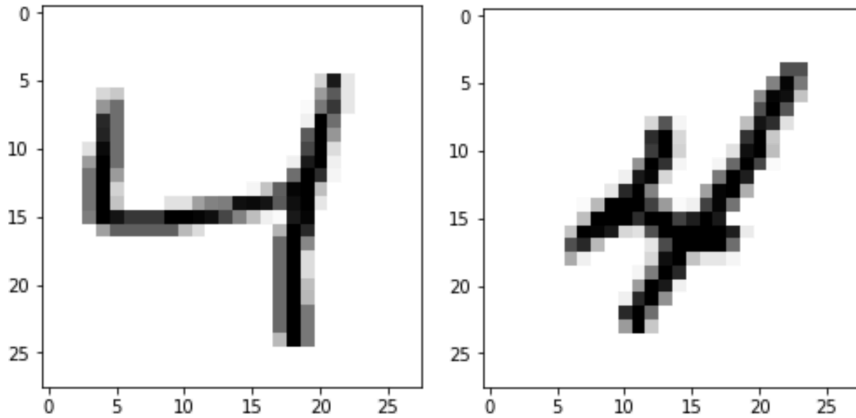
17

- Curse of Dimensionality / Concentration of Measure
 - ▣ Under some assumptions, $\text{maxdist}(q, D)/\text{mindist}(q, D)$ converges to 1
 - Key assumption: independent distribution in each dimension
 - k-NN is still meaningful for real datasets
 - ▣ Hard to find algorithms sub-linear in n (# of points) and polynomial in d (# of dimensions)
 - ▣ Approximate version (c-ANN) is not much easier

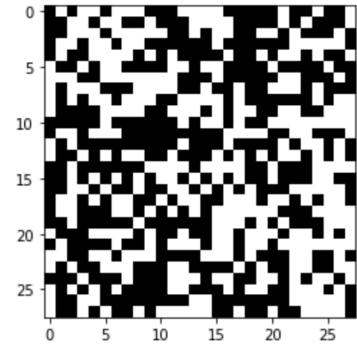
Challenges /2

18

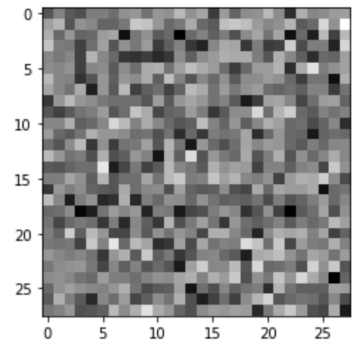
- No idea of the distribution of real data
- Manifold hypothesis



uniform



Radamacher



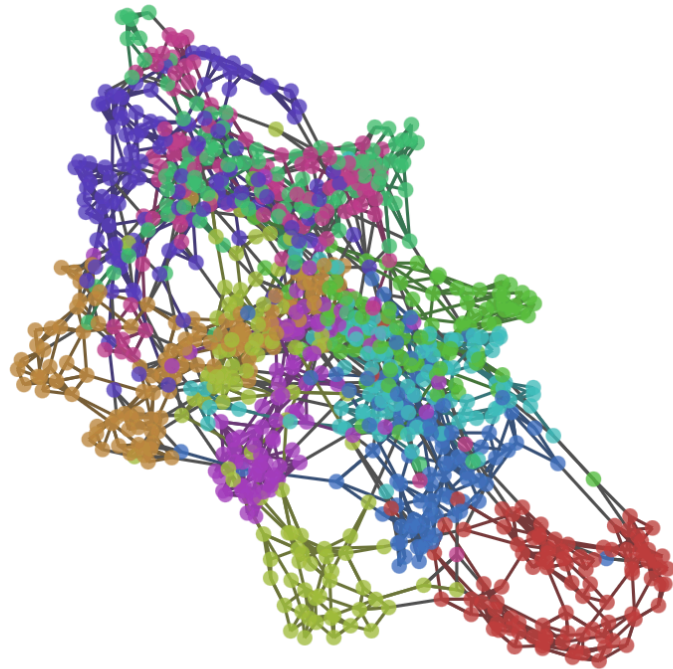
Gaussian

Challenges /2

19

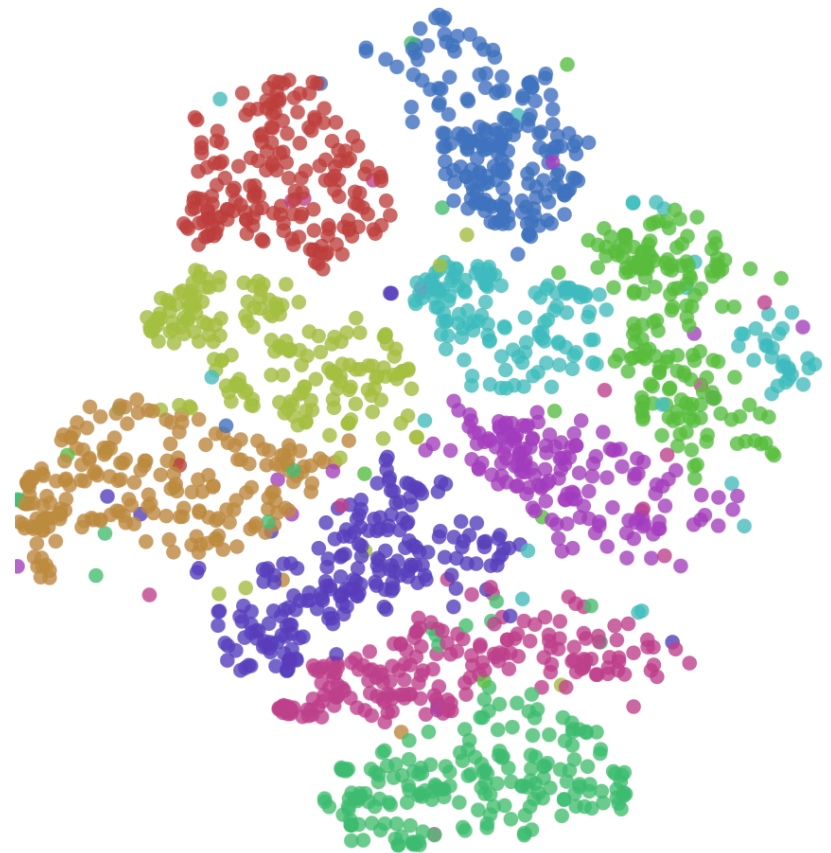
□ Distribution of Real Data

▣ Manifold hypothesis



play

Visualizing MNIST as a Graph



A t-SNE plot of MNIST

Challenges /3

20

- Large data size
 - ▣ 1KB for a single point with 256 dims → 1B pts = 1TB
 - ~100 SIFT vectors per image
 - ▣ High-dimensionality (e.g., documents → millions of dimensions)
- Variety of distance/similarity functions
 - ▣ Less of an issue in the DL era