

# Lecture 4: Nearest Neighbours II

15 January 2020

Lecturer: Dr. Kanat Tangwongsan

Scribe: Kanokpon &amp; Kanokpon

Input: a collection  $D = \{P_1, P_2, P_3, \dots, P_n\} \subset \chi$  of  $n$  points  
 a dist function  $d : \chi * \chi \rightarrow \mathbb{R}_+ \cup \{0\}$

Query: given a query  $q \in x$  find (all) points near  $q$

Two ingredients in the above formulation:

1. a (dis)similarity measure.
2. an efficient d/s & algo baseline:  $O(n)$  probes

Similarity = -dist

## 1 Euclidian distance ( $l_2$ distance)

$$\begin{aligned}\vec{x}, \vec{y} &\in \mathbb{R}^d \\ d(\vec{x}, \vec{y}) &= ||x - y||_{l_2} \\ &= \sqrt{(\vec{x} - \vec{y})^T (\vec{x} - \vec{y})} \\ &= \sqrt{\sum_{i=1}^d (x_i - y_i)^2}\end{aligned}$$

$$\text{For } p > 0, ||x - y||_{l_p} = \sqrt[p]{\sum (x_i - y_i)^p}$$

$p = 0$  : How many non-zero coordinate  $||\vec{x} - \vec{y}||_{l_0}$

$p = \infty$  :  $\max x_i - y_i$

$p = 1$  : Manhattan Dist

## 2 Jaccard

$$\begin{aligned}S, T &\subseteq u \\ J(S, T) &= \frac{|S \cap T|}{|S \cup T|}\end{aligned}$$

### 3 Cosine Similarity

$$\frac{\langle x, y \rangle}{||x|| ||y||}$$

$$\text{Angular distant} = \cos^{-1}\left(\frac{\langle x, y \rangle}{||x|| ||y||}\right)$$

## 4 Low-dimensional Space

### 4.1 d = 1

- balance search tree
- sorted array
- $O(\log n)$  time/query
- $O(n)$  space
- $O(n \log n)$  preprocessing

### 4.2 d = 2

- Voronoi
- KD tree
- OCT tree

#### **Voronoi**

For each point in  $D$ , construct (and store) the region that contains all of its closest neighbors

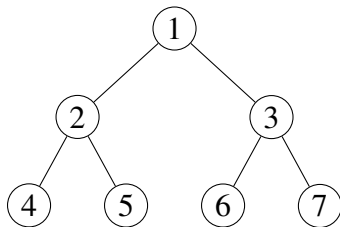
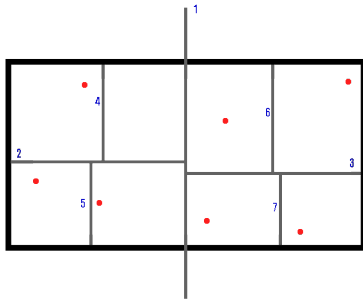
$$VR_i = \{x \in \chi | d(x_i, p_i) \leq d(x_j, p_j) \forall j \neq i\}$$

- $n \log n$  time to build
- $O(n)$  space
- $O(\log n)$  query(via point location)

#### **KD trees**

- space partitioning technique
- alternate b/w vertical and horizontal cuts
  - find a median split point

- recurse until 1 point
- yield a BST



Rectangle range query  
find all points of  $D$  in  $D \cap R$

Start at root & recurse

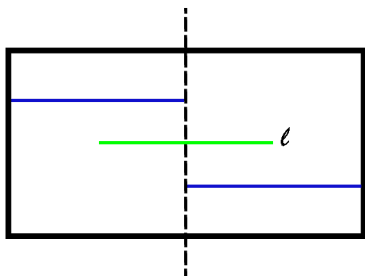
case i: This node region  $\subseteq R$   
return everything in this subtree

case ii: this node's region  $\cap R = \emptyset$   
return  $\emptyset$

case iii: this node's region  $\cap R \neq \emptyset$   
recurse & return the union of the answers from 2 subtrees

Total query cost =  $\overbrace{\# \text{ nodes visited}}^{O(\sqrt{n})} + \overbrace{\# \text{ pts reported}}^s$

Observation:  $\# \text{ nodes visited} \leq \# \text{ nodes where regions are intersected by an edge of } R$



$$f(n) \leq 3 + 2f\left(\frac{n}{4}\right)$$

$$f(n) \leq O(\sqrt{n})$$

- space:  $1 + 2 + 2^2 + \dots + 2^{\log n} = O(n)$
- preprocessing:  $O(n \log n) \Leftarrow n \log n(\text{sorting}) + [T(n) = 2T(n/2) + O(n)]$
- query:  $O(\sqrt{n} + s)$

### 4.3 $d > 2$

- Voronoi
  - has  $\Omega(n^{\lceil \frac{d}{2} \rceil})$  sites
  - $O(n^d)$  space,  $(d + \log n)^{O(1)}$  query
- KD-tree
  - query:  $O(dn^{1-\frac{1}{d}} + s)$  [still prefer if  $d \leq 50$ ]

## 5 Aim

- #1 Dimensionality reduction JL
- #2 Hash & Approximate LSH
- #3 Data-dependent schemes PCA

## 5.1 Johnson-Lindenstrauss

JL-Lemma states that any  $n$  points in high dimensional euclidian space can be mapped onto  $k$  dimensions where  $k \geq O(\frac{1}{\varepsilon^2} \log n)$  without distorting the euclidian distance between any two more than a factor

Find a mapping  $f : u^D \rightarrow V^d$  s.t. (embedded space)  $d_V(f(x), f(y)) \in (1 \pm \varepsilon)d_u(x, y)$

**Lemma 1 (Johnson-Lindenstrauss)**

Let  $0 < \varepsilon < \frac{1}{2}$  Given any set of points

$\{p_1, p_2, \dots, p_n\} \subseteq \mathbb{R}^D$  there exists a mapping  $f : \mathbb{R}^D \mapsto \mathbb{R}^k$

with  $k = O(\frac{1}{\varepsilon^2} \log n)$  such that  $\|f(p_i) - f(p_j)\|_2^2 \in (1 \pm \varepsilon)\|p_i - p_j\|_2^2$

$$M_X = \begin{bmatrix} M_1 \\ \vdots \\ M_k \end{bmatrix}_{K \times D} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad ; M_{n,i} \sim N(0, 1) \quad (1)$$

$$= \begin{bmatrix} \langle M_1, x \rangle \\ \vdots \\ \langle M_k, x \rangle \end{bmatrix} \quad (2)$$

$$f(x) = \frac{1}{\sqrt{k}} M_x \quad (3)$$

**Lemma 2 (Distributional JL)**

Let  $0 < \varepsilon < \frac{1}{2}$  If  $f$  is constructed

per above with  $k = 8\varepsilon^{-2} \ln \frac{2}{\delta}$  with  $x \in \mathbb{R}^D$ , with  $\|x\|_2 = 1$  then

$\Pr[\|f(x)\|_2^2 \in 1 \pm \varepsilon] \geq 1 - \delta$

$$\Delta = f\left(\frac{\overbrace{P_i - P_j}^x}{\|p_i - p_j\|}\right) \quad (4)$$

$$= \frac{1}{\sqrt{k}} M(\dots) \quad (5)$$

$$= \frac{1}{\|p_i - p_j\|} (f(p_i) - f(p_j)) \quad (6)$$

$$\|\Delta\|_2^2 = \frac{1}{\|p_i - p_j\|_{l_2}^2} \|f(p_i) - f(p_j)\|_2^2 \quad (7)$$

$$(8)$$

$$\text{Let } x = \frac{P_i - P_j}{\|p_i - p_j\|} \quad (9)$$

$$\|x\|_2 = 1 \quad (10)$$

$$\frac{\|f(p_i) - f(p_j)\|_2^2}{\|p_i - p_j\|_2^2} = \|\Delta\|_2^2 \in 1 \pm \varepsilon \quad \text{Proved Lemma 1} \quad (11)$$

$$(12)$$

$$\text{Let } \delta = \frac{2}{n^2} \quad (13)$$

$$k = \frac{8}{\varepsilon^2} \ln \frac{2}{\frac{2}{n^2}} \quad (14)$$

$$= \frac{8^2}{\varepsilon} \ln n^2 \quad (15)$$

for pair  $(i, j)$ ,  $\Pr[\Delta \notin 1 \pm \varepsilon] < \frac{1}{n^2}$

$\binom{n}{2}$  pairs total  $\Rightarrow \Pr[\text{at least one has } \Delta \notin 1 \pm \varepsilon]$

$\stackrel{\text{union bound}}{<} \binom{n}{2} \frac{1}{n^2} \leq \frac{1}{2}$

**w.p.**  $\geq \frac{1}{2}$ , all pairs are good

$$A(x) = \frac{1}{\sqrt{k}} M_x$$

matrix  $M_{K \times D}, N_{i,j} (0, 1)$

Let  $0 < \varepsilon < \frac{1}{2}$  If  $A$  is as defined above with  $k = 8\varepsilon \ln(\frac{2}{\delta})$  and  $x \in \mathbb{R}^D$  with  $\|x\|_2 = 1$

then,  $\Pr[\|A(x)\|_2^2 \in 1 \pm \varepsilon] \geq 1 - \delta$

$N(\mu, \sigma^2)$  mean  $\mu$  and variance  $\sigma^2$

$$\text{pdf: } f(x) = \frac{1}{2\pi\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

then:

$$\begin{aligned} CG_1 &\sim N(c\mu_1, c^2\sigma_1^2) \\ G_1 + G_2 &\sim (\mu_1 + \mu_2, \sigma_1^2, \sigma_2^2) \end{aligned}$$

$$A(x) = \begin{bmatrix} M_1 \\ \vdots \\ M_k \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix} = \frac{1}{\sqrt{k}} \begin{bmatrix} \langle M_1 x \rangle & -y_1 \\ \langle M_2 x \rangle & -y_2 \\ \vdots & \vdots \\ \langle M_k x \rangle & -y_K \end{bmatrix}$$

$$y_i = [G_1, G_2, \dots, G_2] \begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix} = \sum_{j=1}^D G_j x_j, \text{ where } G_j \sim N(0, 1)$$

$$y_i \sim N(0, \underbrace{\|x\|_2^2}_{=1}) = N(0, 1)$$

$$Z = \|A(x)\|_2^2 = A(x)^T A(x) = \frac{1}{k} \sum y_i^2 \rightarrow \text{chi}^2 \text{ distribution with } k \text{ doF}$$

$$\mathbb{E}[Z] = \frac{1}{k} \sum \mathbb{E}[y_i^2] = 1 = \|x\|_2^2$$

$$\mathbb{E}[\|A(p) - A(q)\|_2^2] = \mathbb{E}[\|A(p - q)\|_2^2] = \|p - q\|_2^2$$

$$\begin{aligned} \Pr[Z > 1 + \varepsilon] \forall t &> 0 \\ &= \Pr[e^{tkz} \geq e^{tk(1+\varepsilon)}] \\ &\leq \frac{\mathbb{E}[e^{tkz}]}{e^{tk(1+\varepsilon)}} ; e^{tkz} = e^{tk\frac{1}{k}\sum y_j^2} = \Pi e^{ty_j^2} \\ &= \Pi \frac{\mathbb{E}[e^{ty_j^2}]}{e^{tk(1+\varepsilon)}} \end{aligned}$$

$$\textbf{claim1: } \mathbb{E}\left[e^{tG^2}\right] = \frac{1}{\sqrt{1-2t}} \text{ for } t < \frac{1}{2}$$

**claim2:**  $\frac{1}{e^t \sqrt{1-2t}} \leq e^{\frac{t^2}{1-2t}}$

$$= \left( \frac{1}{e^{t(1+\varepsilon)} \sqrt{1-2t}} \right)^k$$

$$\leq \exp \left\{ \frac{kt^2}{1-2t} - k + \varepsilon \right\}$$

use  $t = \frac{\varepsilon}{4} \leq e^{\frac{-k\varepsilon^2}{8}}$

Recall  $\varepsilon < \frac{1}{2}$

choose  $k = \frac{8}{\varepsilon^2} \left( \frac{2}{\delta} \right) = \frac{\delta}{2}$

$\Pr [Z < 1 - \varepsilon] < \frac{\delta}{2}$