# Lecture 3: metric data

- So far, we haven't really thought about _structure_ of data.
- But usually this is what actually matters!

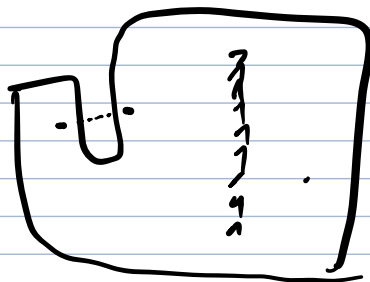

vs

vs

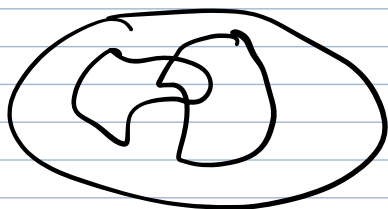Q: What is the most similar data point to mine?

Q: How similar are two data points?

High level question: **distance between data**

- points on a map
  - Euclidean ?



- Sets of things



distance = ?

__Metric spaces__: A metric space is specified by

$$(X, d)$$

a collection of points

a distance function
$$d: X \times X \longrightarrow \mathbb{R}_{\geq 0}.$$
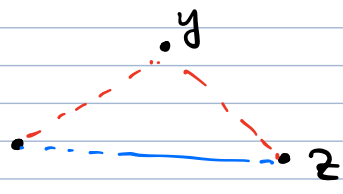
$d$ must satisfy 3 properties:

(i)  $d(x,y) = 0 \Rightarrow x = y.$   $\forall x, y$

(ii)  (Symmetry)  $d(x,y) = d(y,x)$  $\forall x, y.$

(iii) (Triangle inequality)

$$d(x,z) \leq d(x,y) + d(y,z). \quad \forall x, y, z.$$

examples of metric spaces:
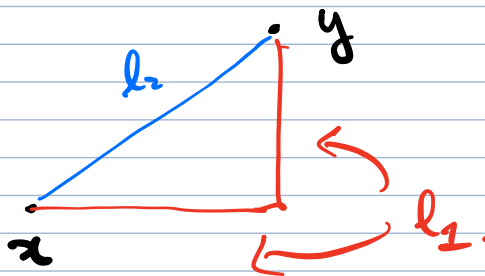
(0) $l_p$- distances. $\quad X = \mathbb{R}^d$

$$d(x,y) = \| x - y \|_p = \left( \sum_{i=1}^{d} |x_i - y_i|^p \right)^{1/p}$$

$$(p \in [1, \infty]).$$

e.g. $p = 2 \rightarrow$ Euclidean distance.

$p = 1 \rightarrow$ "Manhattan distance"

$$\| x \cdot y \|_1 = \sum_{i=1}^{d} |x_i - y_i|$$

$$p = \infty$$

$$\| x - y \|_\infty = \max \left( |x_1 - y_1|, \cdots, |x_d - y_d| \right).$$

$$= \lim_{p \to \infty} \| x - y \|_p.$$

Note: $l_2$ is very special!

 - basis independent.

 - self- dual

1). String distances. given $s_1, s_2$ strings,

edit distance:

$$d_{edit}(s_1, s_2) = \min \# \text{ insertions + deletions to go from}$$
$$\qquad s_1 \quad \text{to} \quad s_2$$

$s_1 = $ A  G  T  A  C  A

$s_2 = $ G  T  A  A  T

$d_{edit}(s_1, s_2) \leq 3$ ⟵

~ "basic" DP : $O(mn)$

$O(n^{1+\delta})$ time to get $(1+\delta)$-approx

↳ best paper FOCS 2018

[Chakraborty, Das, Goldenberg, Koucky, Saks].

~~X~~  G  T  A  C  A

↓

G  T  A  ~~X~~  A

↓

G  T  A  A

↓

G  T  A  A  **T**

<u>Hamming distance:</u>  # of replacements

$\boxed{\text{A}}$ $\boxed{\text{G}}$ $\boxed{\text{T}}$ $\boxed{\text{A}}$ $\boxed{\text{C}}$ $\boxed{\text{A}}$
$\boxed{\text{G}}$ $\boxed{\text{T}}$ $\boxed{\text{A}}$ $\boxed{\text{A}}$ $\boxed{\text{T}}$ $\boxed{\text{∅}}$

$d_{Hamming}(s_1, s_2) = 5.$

**3).** Jaccard (dis)- similarity   given S, T sets

$$J(S,T) := \frac{|S \cap T|}{|S \cup T|}$$

not a metric!

larger $J(S,T)$ ⟶ more similar S & T

$J(S,T) = 1. \Leftrightarrow S = T.$

$d_J(S,T) = 1 - J(S,T)$ ⟵ is a metric.

triangle inequality rather nonobvious but true!

Q: How similar are different metrics?

<u>metric embeddings:</u>

given two metric spaces $(X_1, d_1)$, $(X_2, d_2)$, is there a bijection $f: X_1 \to X_2$ s.t.

$$d_1(a, b) \approx d_2(f(a), f(b))$$
$$\forall a, b \in X_1 \ ?$$

or can I design $X_2$ so that $f$ exists?

Active area of study!

---

# Nearest - neighbor problem $\quad (X, d)$
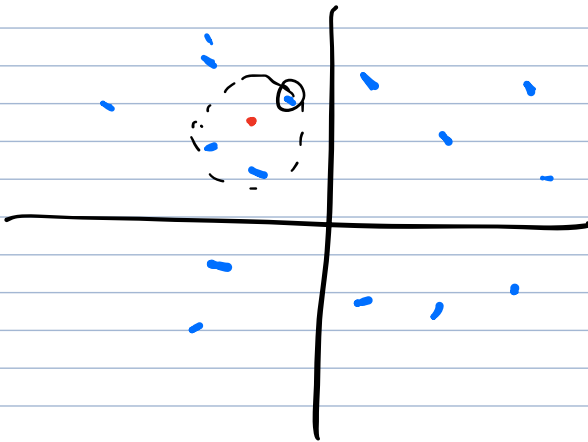
data set $\quad D \subseteq X$

preprocess $\quad D$ s.t. given a query $q \in X$,
quickly find $\quad x \in D$ s.t. $d(x, q)$ is minimized.

mostly focus on $\quad X = \mathbb{R}^k$,
$$d = \| \cdot \|_2.$$

Approx - NN: given parameter $\epsilon > 0$, return $x \in D$ s.t.

$$d(x, q) \leq (1 + \epsilon) \min_{x^* \in D} d(x, q).$$



Goal: If $|D| = n$, use space $O(n)$ and
answer queries in time $O(\log n)$.

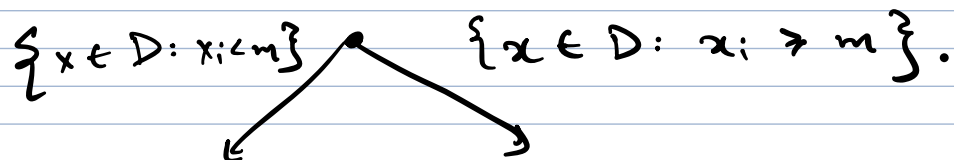| "Space partitioning" | Random projection |
|---|---|
| - kd - trees | - locality sensitive hashing |
|  | (LSH) |

# k-d tree [Bentley '75].

build a <u>search tree</u> to partition space.

many variants of this!

<u>leaf nodes:</u> Subset that contains 1 data element.

<u>non-leaf node:</u> A subset of points w/ size > 1.
Some dimension $i$, and value $m$.

$\{x \in D: x_i < m\}$      $\{x \in D: x_i \geqslant m\}$.

How to build a k-d tree: Initially, root node is all points in D.

at node v: Let $D_v$ be associated set of points to v.

    — it $|D_v| = 1$, v is leaf.

    — if $|D_v| > 1$:
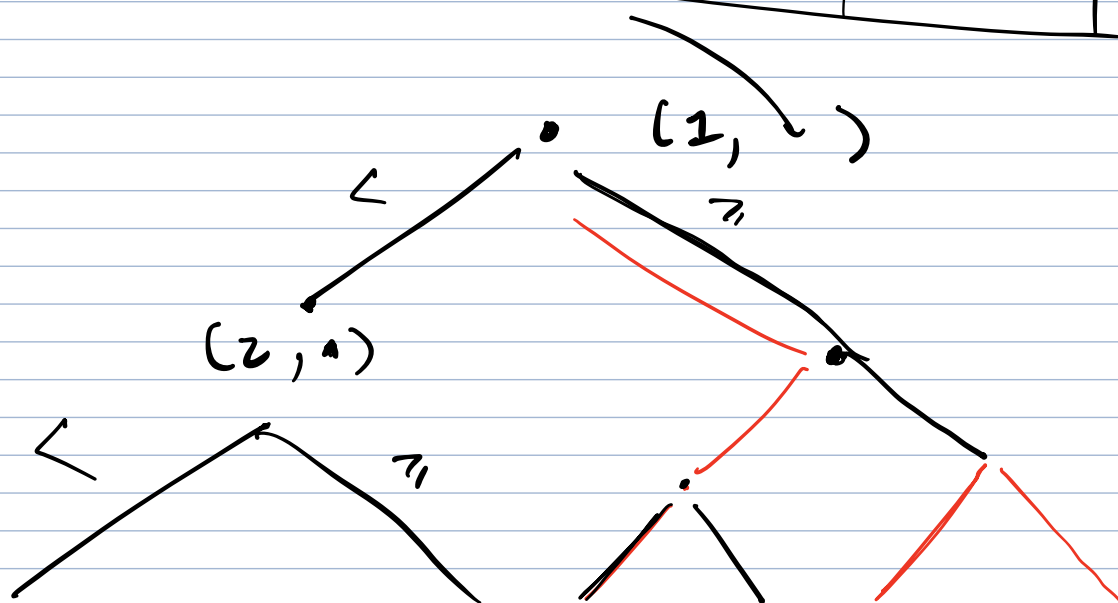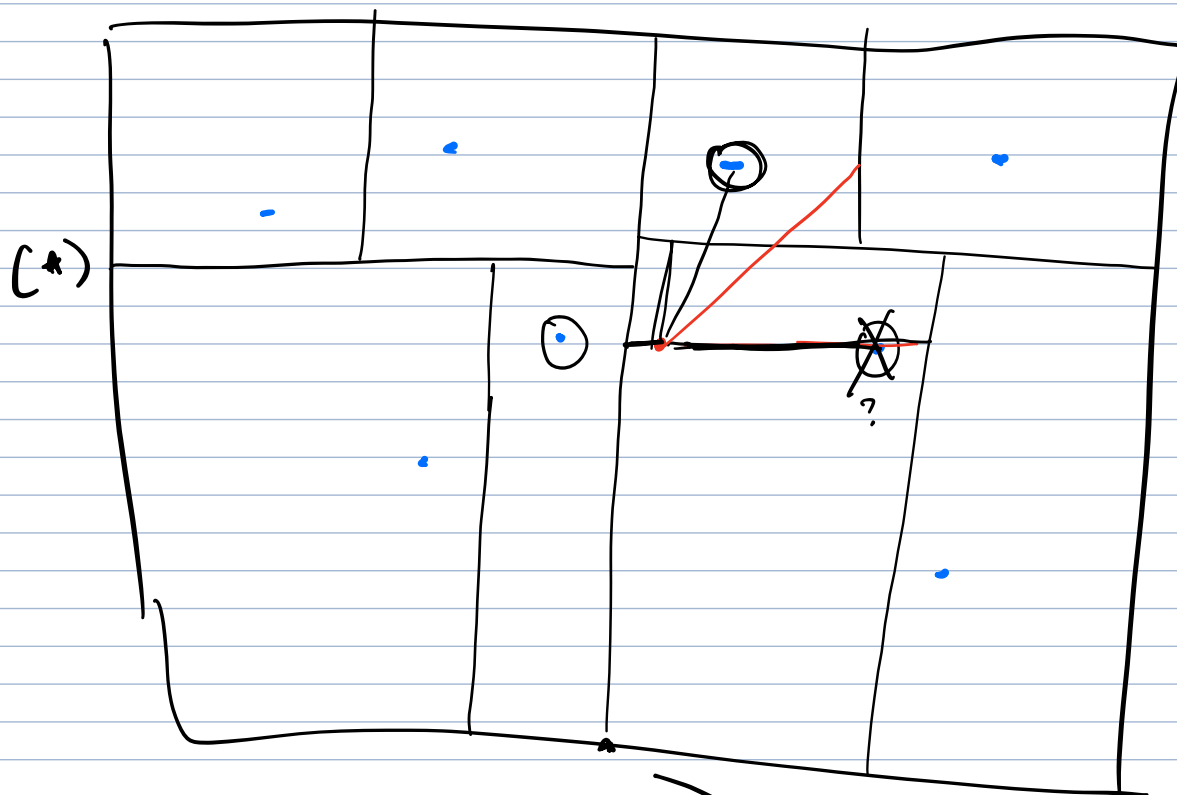
        pick dimension $i_v \in \{1, \cdots, k\}$

        Let $m_v = \underset{x \in D_v}{\text{median}}(x_{i_v})$

        store $i_v, m_v$ at v.

        left child gets $\{x \in D_v: x_{i_v} < m_v\}$

        right child gets $\{x \in D_v: x_{i_v} > m_v\}$.

        recurse

(*)

(1, )

(2, )

size of tree?  $O(n)$

depth of tree?  $O(\log n)$

How to perform lookup?

need to recurse up tree!

worst case: might have to look at full tree!

in <u>low dimensions</u>: often much faster

"on average" ∴ $O(2^k \log n)$

rule of thumb:  k-d tree useful when $\boxed{k \leq 20}$
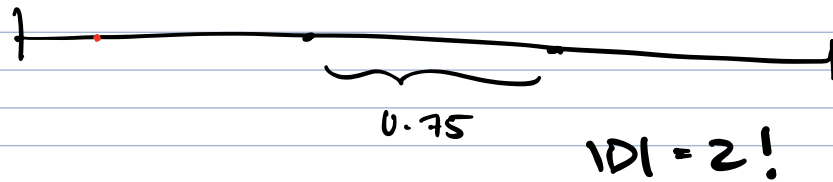
## Curse of dimensionality:

Often times,  running times for these sorts of
space-partitioning  methods  (and others)  scale exponentially
w/  dimension.

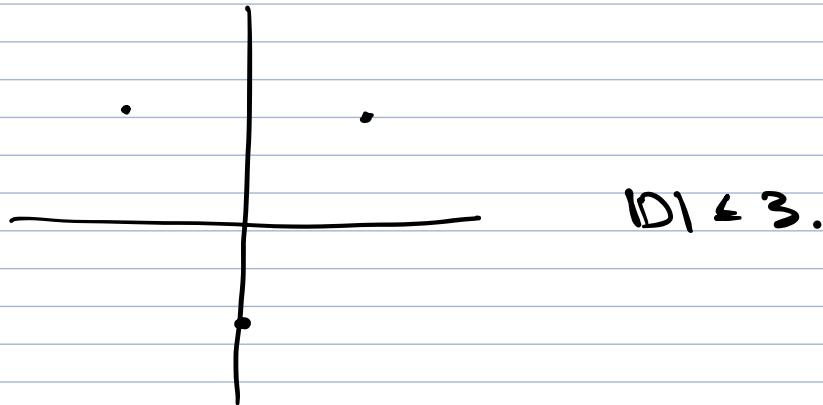high-dimensional points have lots of points w/
similar  distances!

e.g.    How  large  can  D  be s.t.
$\|x-y\|_2 \in [0.75, 1] \quad \forall x, y \in D$?

k = 1 :



0.75
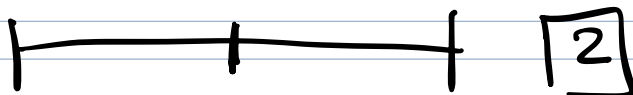
$|D| = 2!$

k = 2 :



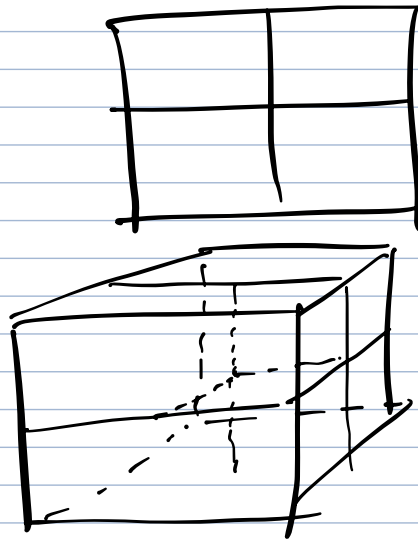$|D| \leq 3.$

k = 100 ?    A  lot!

in general:    $\exists D \subseteq \mathbb{R}^k$ s.t.  $|D| \geq 2^{ck}$

e.g.    cell  packing:  How  many  boxes of  edge length
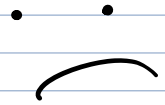1  can  fit  into  a  box  w/ edge  length 2 ?

k = 1    $\vdash\!\!-\!\!-\!\!+\!\!-\!\!-\!\!\dashv$    $\boxed{2}$

$k:$    2



4

$2^k$    :(

need   new   approach   for   high dimensions!