

Approximate Nearest Neighbors (+SIFT)

Large Dataset P $|P|=n=1$ million

(Q1) Which pairs $p, p' \in P$ are close

(Q2) Given query g find closest point $p \in P$ to g .

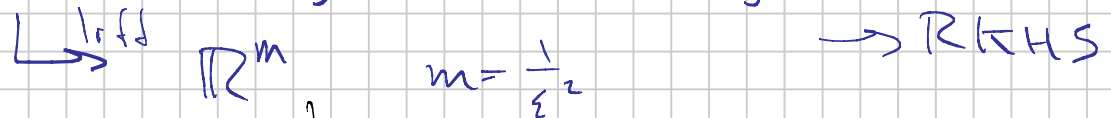
$$p^* = \underset{p \in P}{\operatorname{argmin}} d(p, g) = \phi_P(g)$$

High Dimensional Euclidean Data

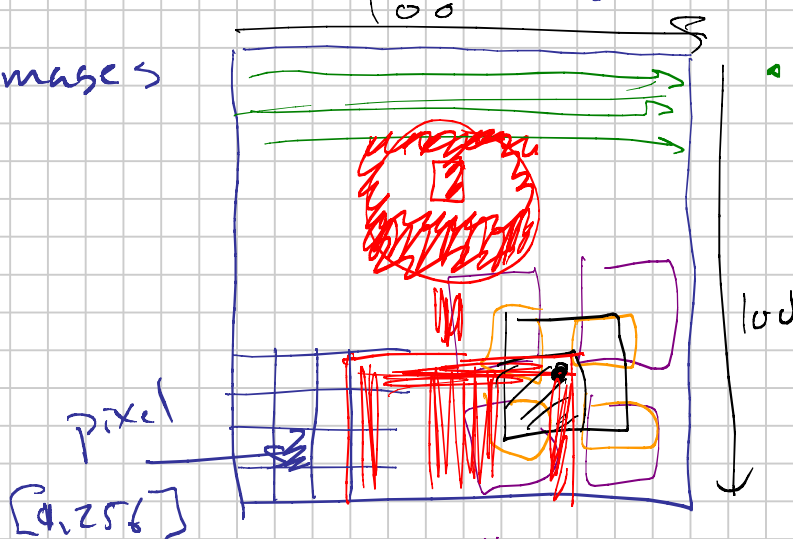
- time series: price, stocks, temperature

1 2 3
5 15 16 18 ... 35
10

- machine learning kernel $K(x, y)$



- Images



- vector $100 \times 100 \rightarrow d = 10,000$ (if pre-aligned)

- SIFT (David Lowe)



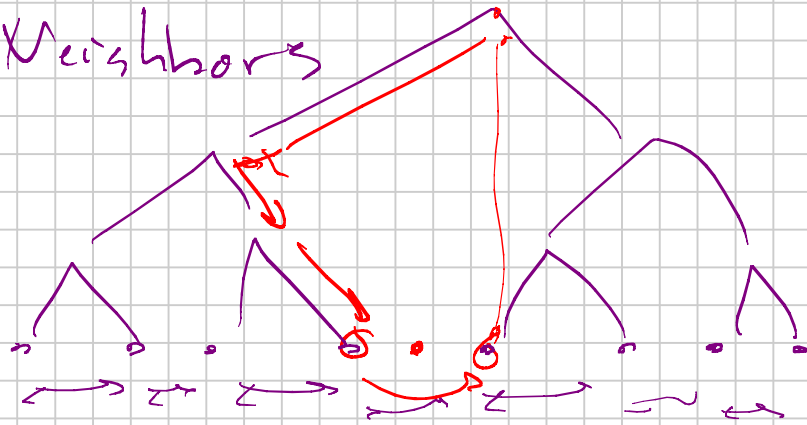
8 coord

$4 \times 4 = 16$ scale

invariant to scale, rotation, $8 \cdot 16 = 128$ dimensional shift

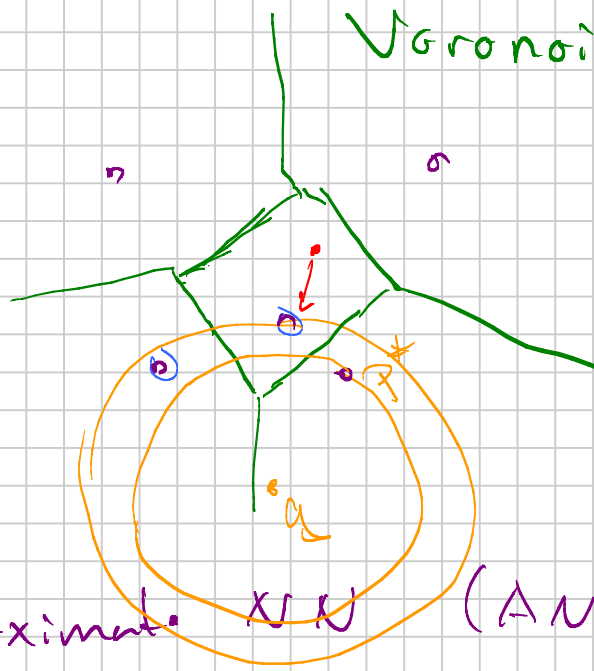
Nearest Neighbors

$P \subset \mathbb{R}^1$



$\log n$

$P \subset \mathbb{R}^2$



Voronoi Diagram

size = $O(n)$
 query = $O(\log n)$
 construction = $O(n \log n)$

$P \subset \mathbb{R}^d$

Approximate NN (ANN)

d-dim
 size $O(n^{\lceil d/2 \rceil})$

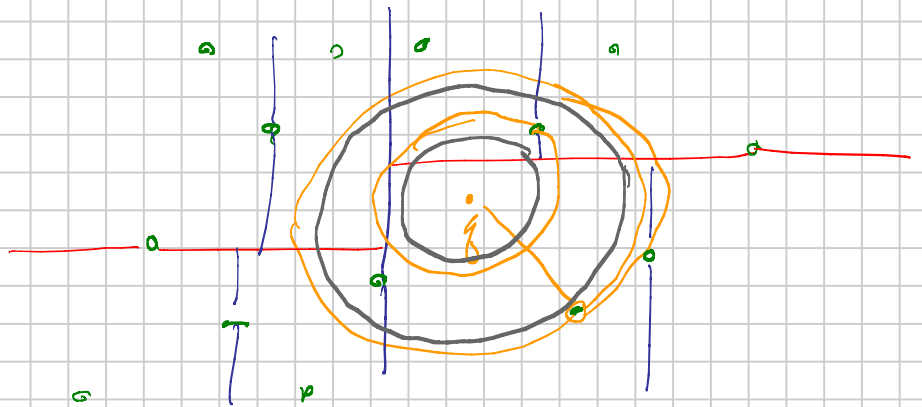
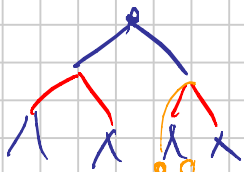
$$p^* = \underset{P \in P}{\operatorname{argmin}} d(p, q) = \phi_P(q)$$

$\epsilon \in (0, 1)$

Goal: Find \hat{p} s.t. $d(\hat{p}, q) \leq (1+\epsilon) d(p^*, q)$

Medium Dimensions $d \in [3-12]$

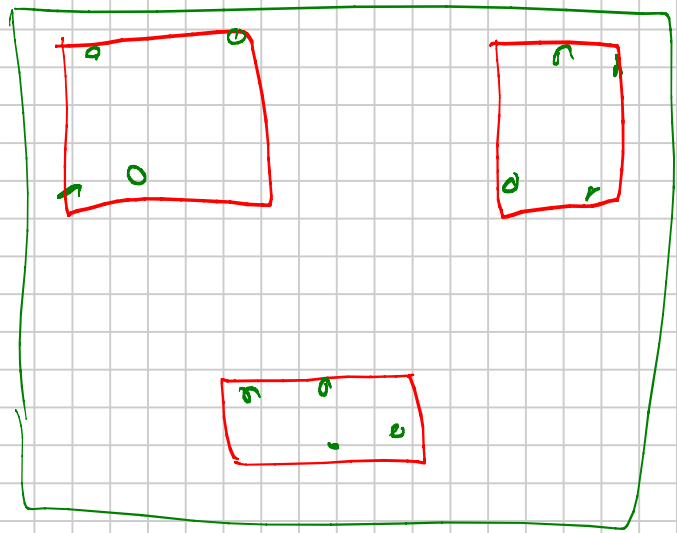
sk-d tree



ANN

• R-tree

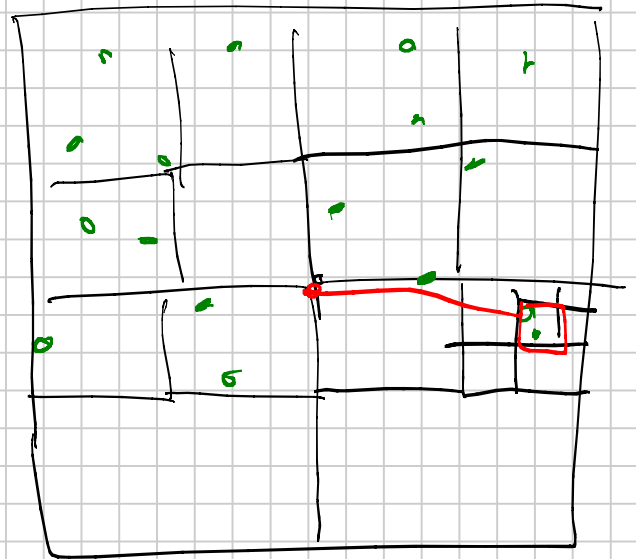
R = rectangle



(compressed)
• Quad-tree

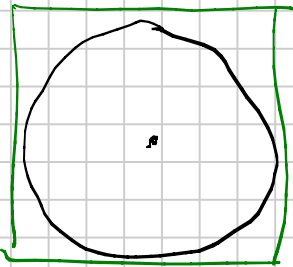
$$O(n \log n + \frac{1}{\epsilon} d)$$

$$O(\log n + \frac{1}{\epsilon} d)$$



High-Dimensional $d > 20$

"curse of dimensionality"



$B(d, \text{rad}=1)$ ball in \mathbb{R}^d

$$\text{Vol}(B(d, 1)) = \frac{\pi^{d/2} \text{radius}^d}{\Gamma(d/2 + 1)} \approx \frac{\pi^{d/2}}{(d/2)!} \leq 1$$

$$[-1, 1]^d$$

$$\text{Vol}(\text{cube}(d, \text{rad}=1)) = 2^d$$

• ANN $d \in [2, 20?]$

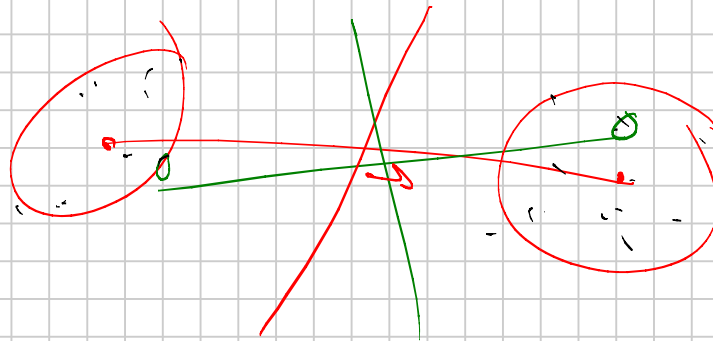
• LSH $n \rightarrow$ very big

• Mods to kd-tree $d \approx 160-200$

inherently low dimensional
bounded doubling dimension



• Build kd tree but smarter splits



- cluster 2 clusters
- Random Rotations do few picks best