

# Nearest neighbor searching

Preprocess points  $P = \{p_1, \dots, p_n\}$

Later, given point  $q$

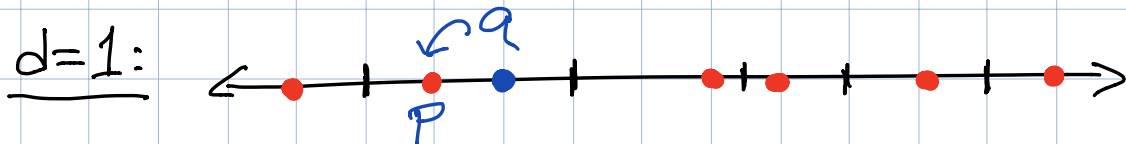
quickly find  $\min \{d(p, q) \mid p \in P\}$

$\uparrow$  distance

Euclidean distance in  $\mathbb{R}^d$

Fast alg / small space when  $d \leq 2$

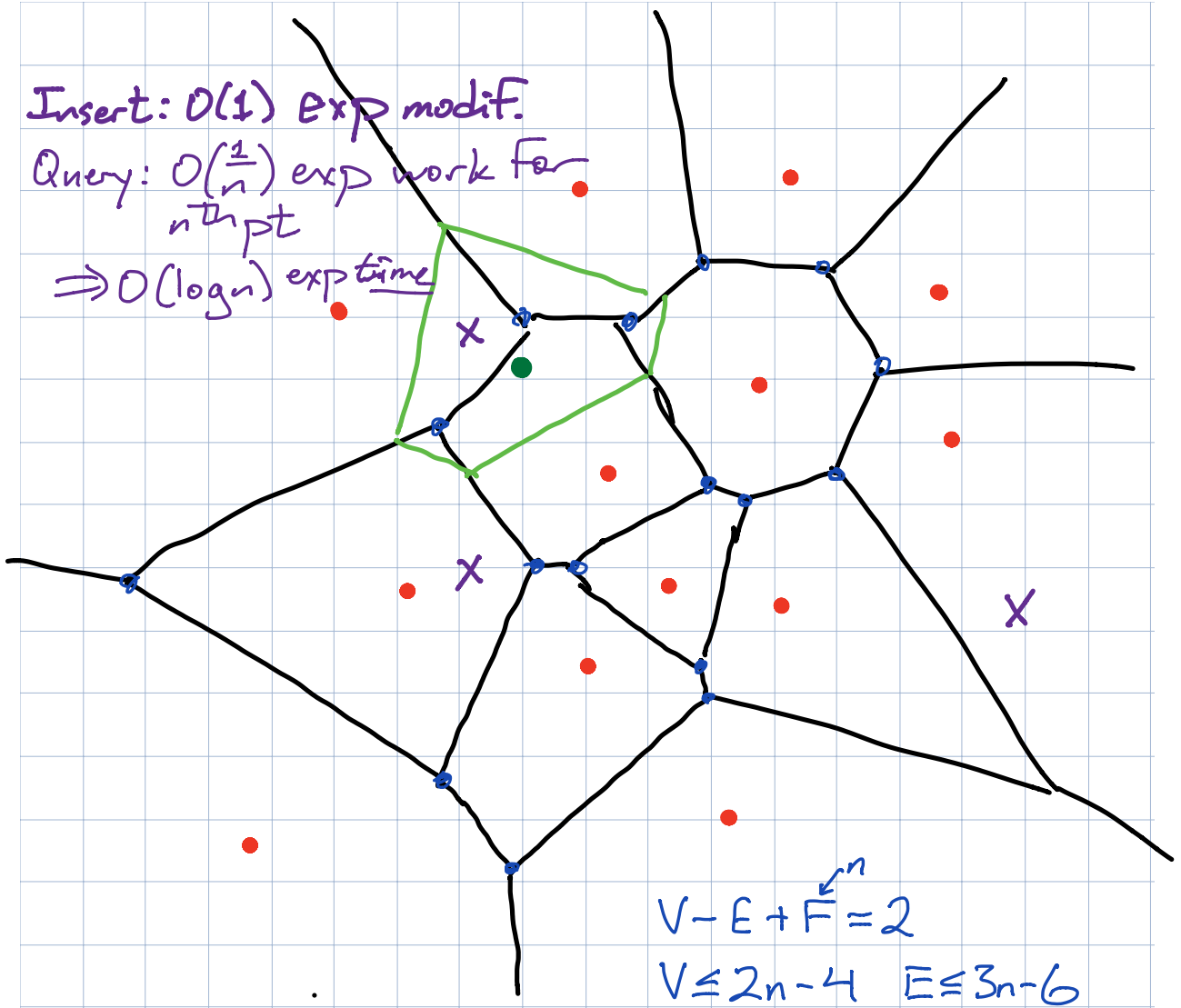
Brute force:  $O(dn)$   
time



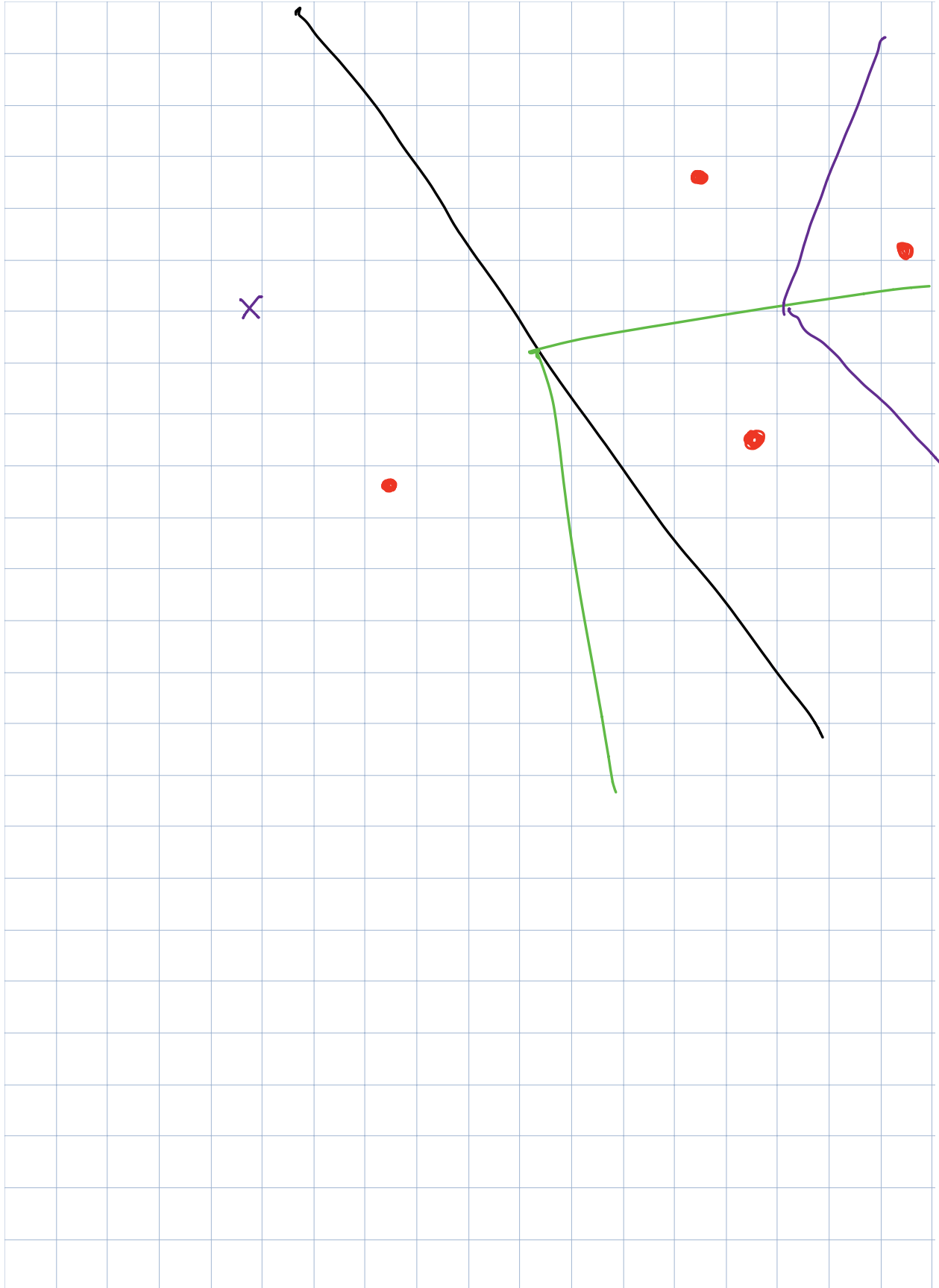
Binary search,  $O(\log n)$  time  $O(n)$  space

$d=2$ : Voronoi diagram [1906]

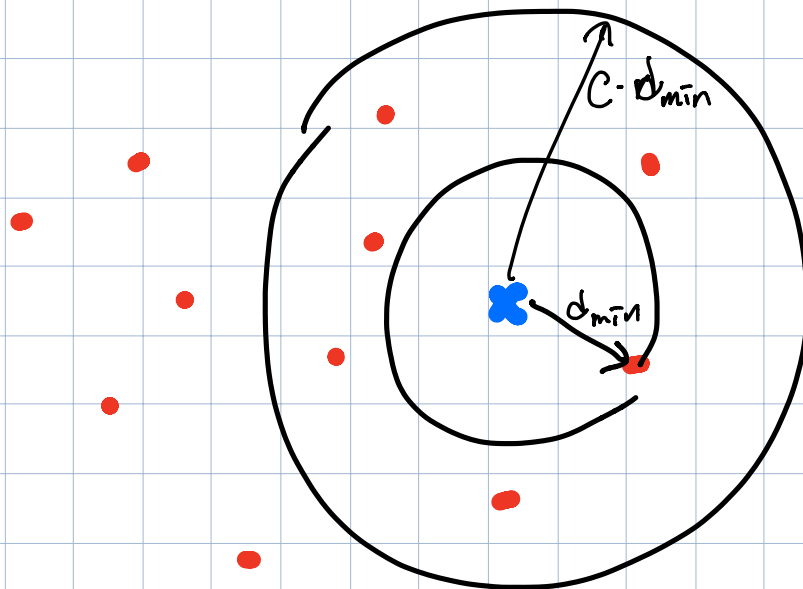
Insert:  $O(1)$  Exp modif.  
Query:  $O(\frac{1}{n})$  exp work for  
nth pt  
 $\Rightarrow O(\log n)$  exp time



$$V - E + F = 2$$
$$V \leq 2n - 4 \quad E \leq 3n - 6$$



Approx. nearest neighbor query:



Return any point  $p_i \in P$  s.t.  
 $d(p_i, q) \leq c \cdot \min\{d(p, q) \mid p \in P\}$   
with prob  $1 - \delta$ .

Locality-sensitive hashing

Approx. near-neighbor search

Given  $P$  to preprocess

Given point  $q$ , distances  $r, R$

If there is a point  $p \in P$  s.t.  $d(p, q) \leq r$   
return a point  $p' \in P$  s.t.  $d(p', q) \leq R$   
with prob  $1 - \delta$ .

[Har-Peled  
2001]

A family  $\mathcal{H}$  of hash functions is *locality sensitive* if

$$\text{IF } d(p, q) < r \Rightarrow \Pr[h(p) = h(q)] \geq \Pi$$

$$\text{IF } d(p, q) > R \Rightarrow \Pr[h(p) = h(q)] \leq \pi$$

---

"Points" are  $d$ -bit vectors

$$\text{dist}(p, q) = \#1\text{'s in } p \oplus q \quad \text{Hamming dist}$$

$$p = 000011101 \quad \text{dist} = 3$$

$$q = 001001100$$

$$h_i(p) = p_i = \text{ith bit in } p \quad \mathcal{H} = \{h_i \mid 1 \leq i \leq d\}$$

$$\Pr[h(p) = h(q)] = \frac{d - \text{dist}(p, q)}{d} = 1 - \frac{d(p, q)}{d}$$

$$\Pi = 1 - \frac{r}{d} \quad \pi = 1 - \frac{R}{d}$$

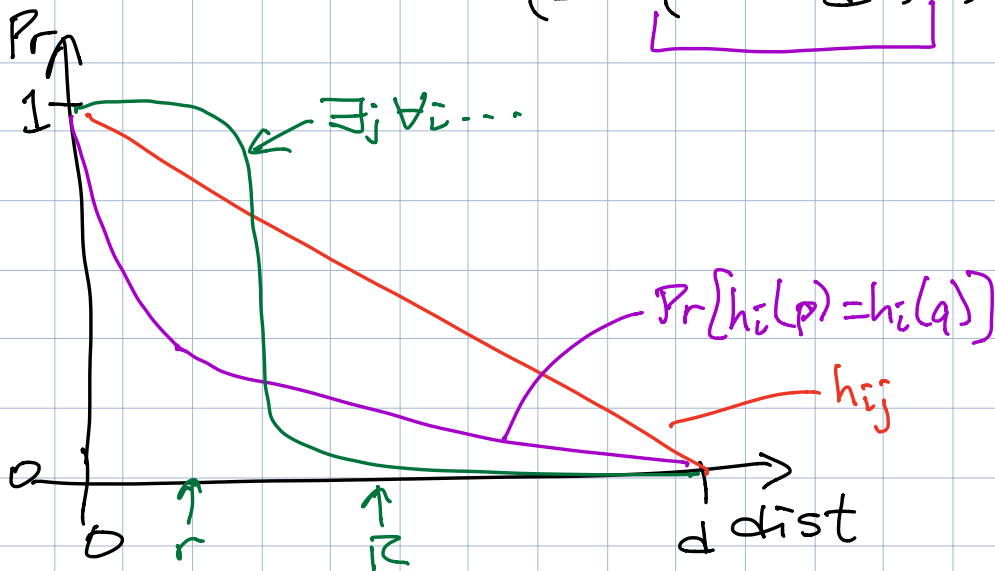
Pick  $k$  functions  $h_{ij} \in \mathcal{H}$  independently  
 $1 \leq i \leq k \quad 1 \leq j \leq l$

$$\text{For all } j: h_j(p) = (h_{1j}(p), \dots, h_{kj}(p)) \in \{0, 1\}^k$$

Define "collide" = For some  $j$ ,  $h_j(p) = h_j(q)$

$$\exists j \forall i \ h_{ij}(p) = h_{ij}(q)$$

$$\begin{aligned} \Pr[\text{collision}] &\leq 1 - \left(1 - \Pr[h_{ij}(p) = h_{ij}(q)]\right)^k \\ &= 1 - \left(1 - \left(1 - \frac{\text{dist}}{d}\right)^k\right)^l \end{aligned}$$



for each  $j$ :  
 For each  $p \in P$   
 store  $p$  at  $T_j[h_j(p)]$

query( $q$ )

for each  $j$ :  
 $S \leftarrow T_j[h_j(q)]$   
 check every pt in  $S$   
 if close pt found return it  
 stop after checking  $2l$  points

