

Lecture 7

K-means

Find k points
 m_1, m_2, \dots, m_k

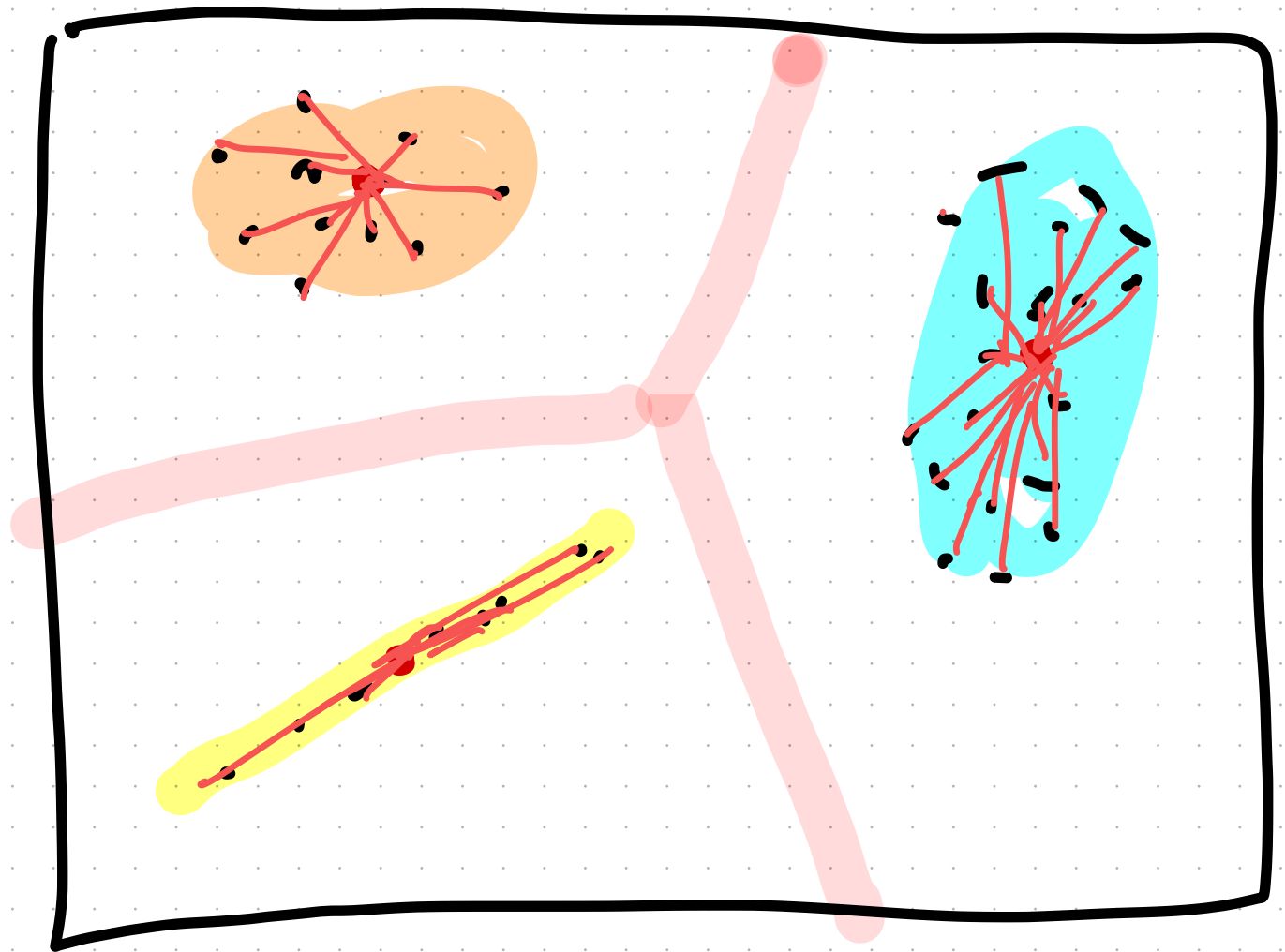
That minimize

$$\sum d(x_i, \underset{m}{\text{closest}}(x_i))$$

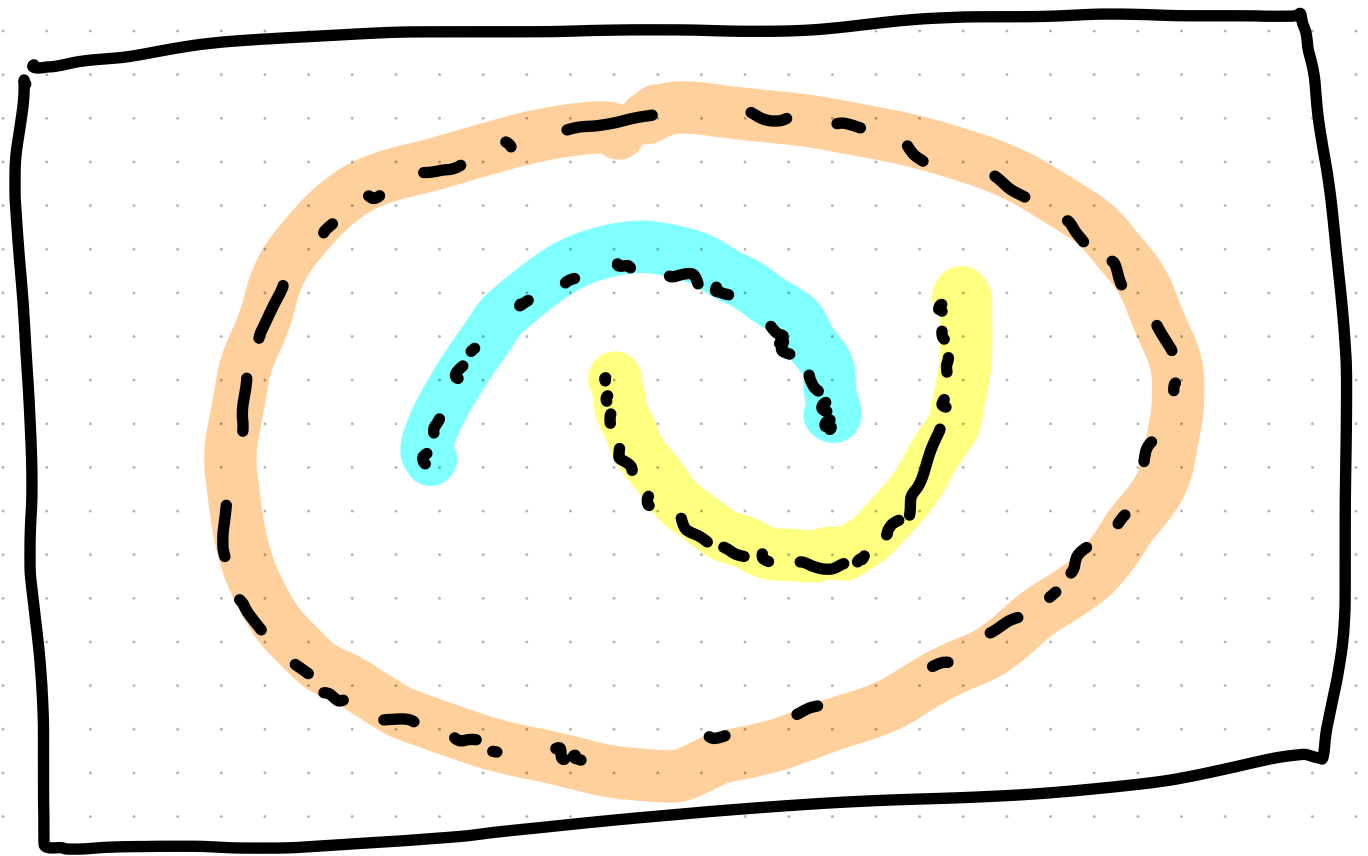
NP-Complete:

we think there
is no $O(n^c)$

algorithm for any c

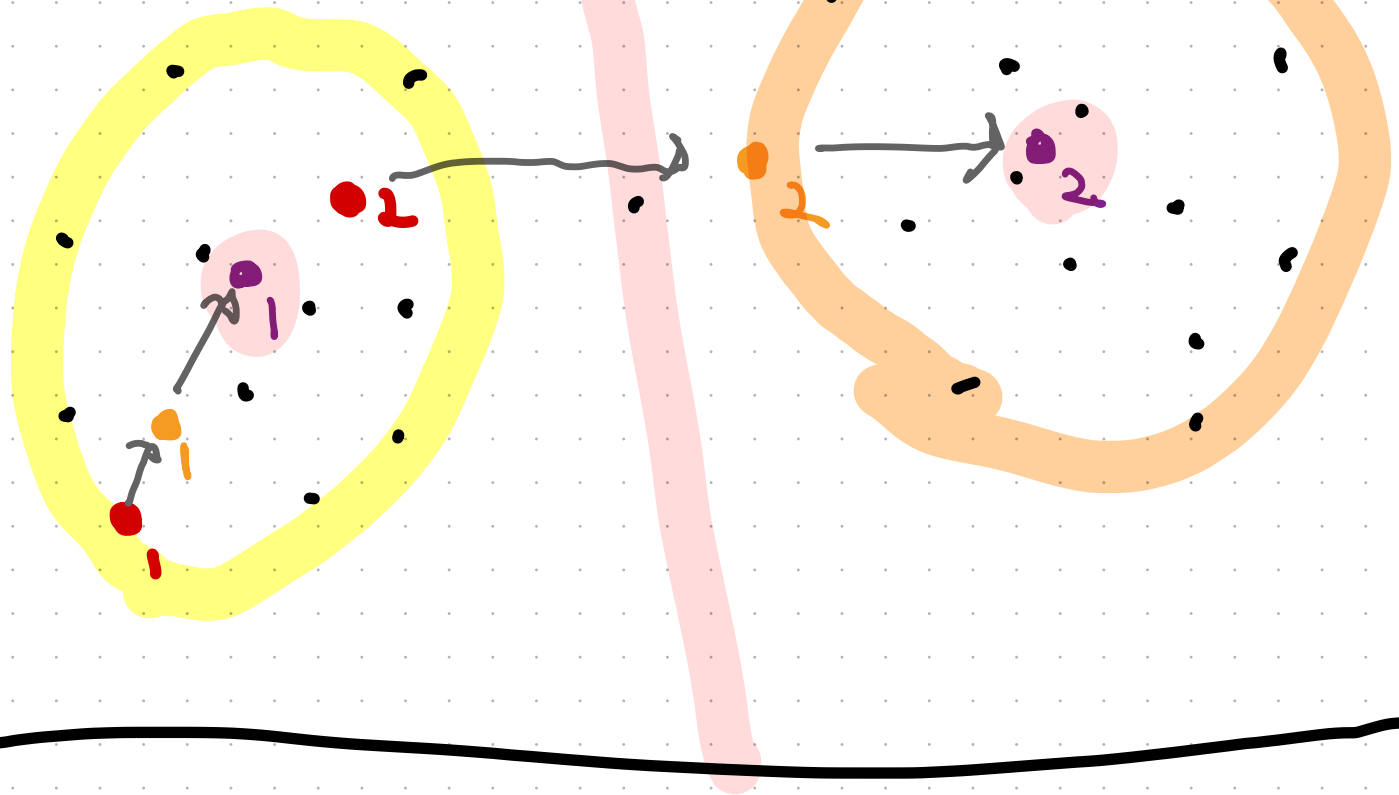


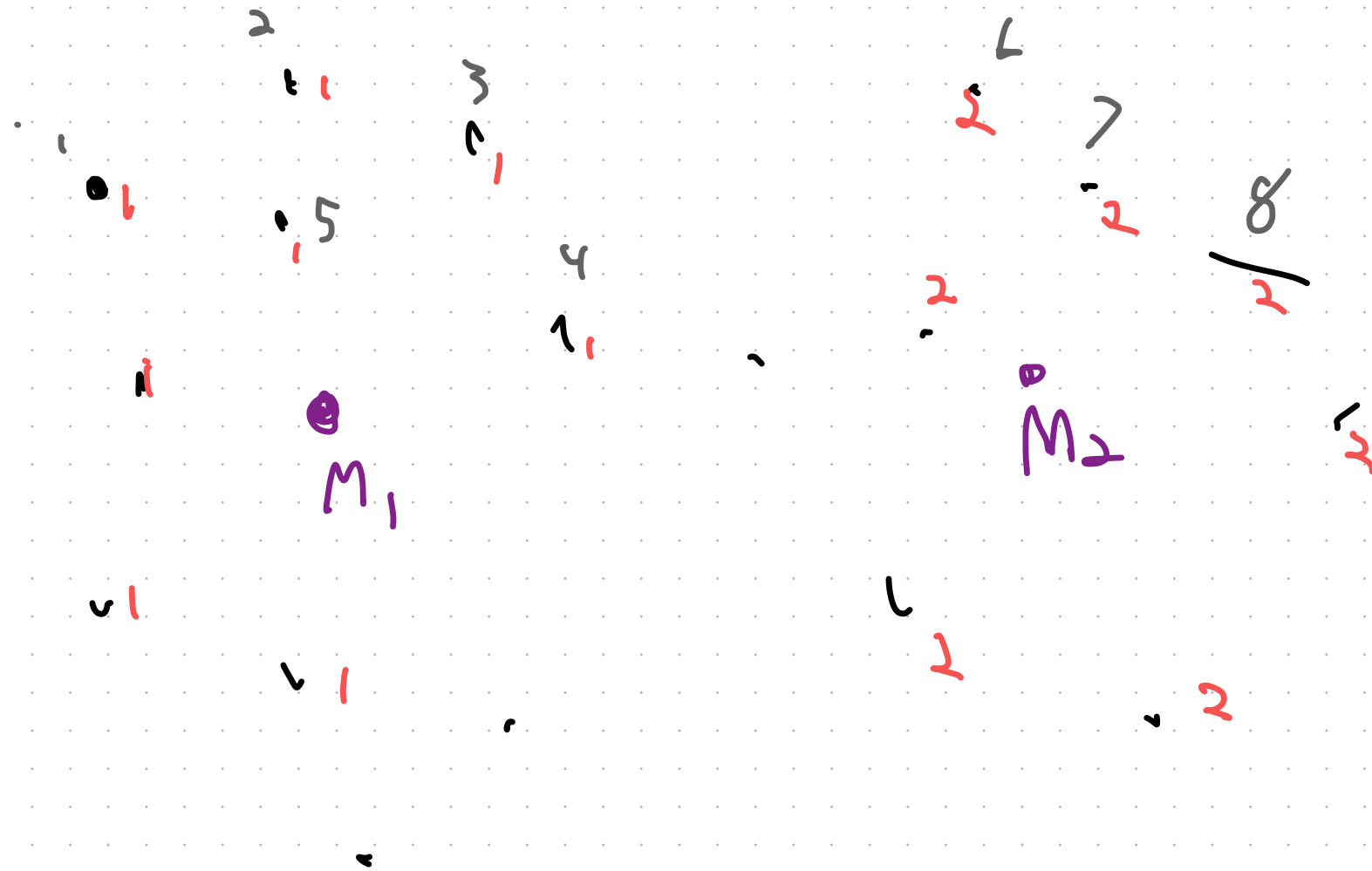
Clustering



$k=2$

For each point
group according
to closest mean
Compute the
centroid of
each group





P:

Runtime

n : Number of data points

d : Dimension

i : # of iterations to converge

t : # of different trials

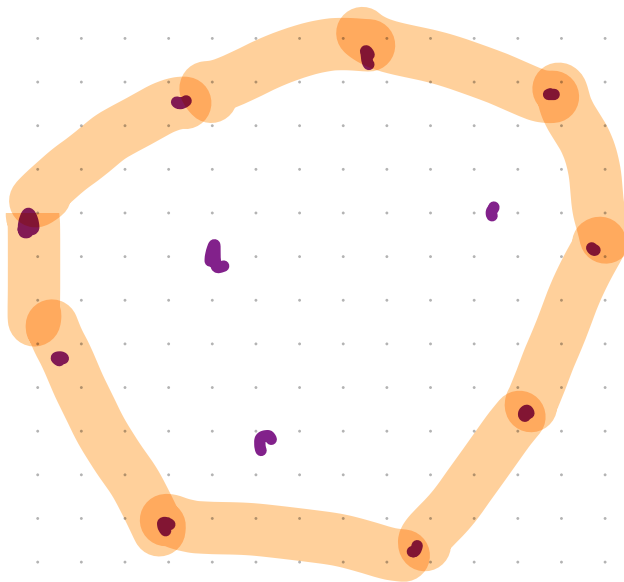
k : # of means

$$O(nkd): \begin{cases} n \cdot k \cdot d \\ n \cdot k \\ n \cdot d \end{cases}$$

for each point find closest mean
for each mean group closest point
Compute centroid

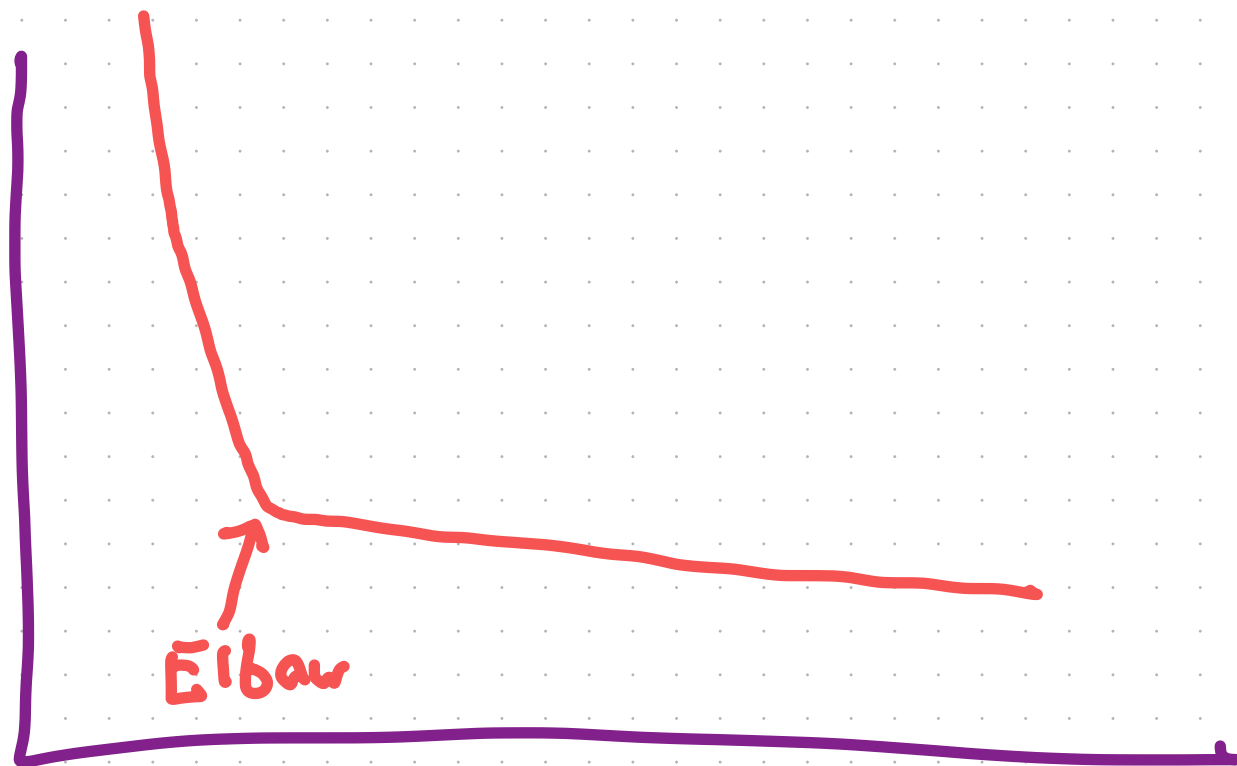
Total: $O(Nkd \underbrace{ti}_{\text{Not too big}})$

"Curse of dimensionality"



$$n^{L\frac{d}{2}+1}$$

quality



K

Dimensionality Reduction

Points in high dimension



Points in a lower dimension

Preserve distance

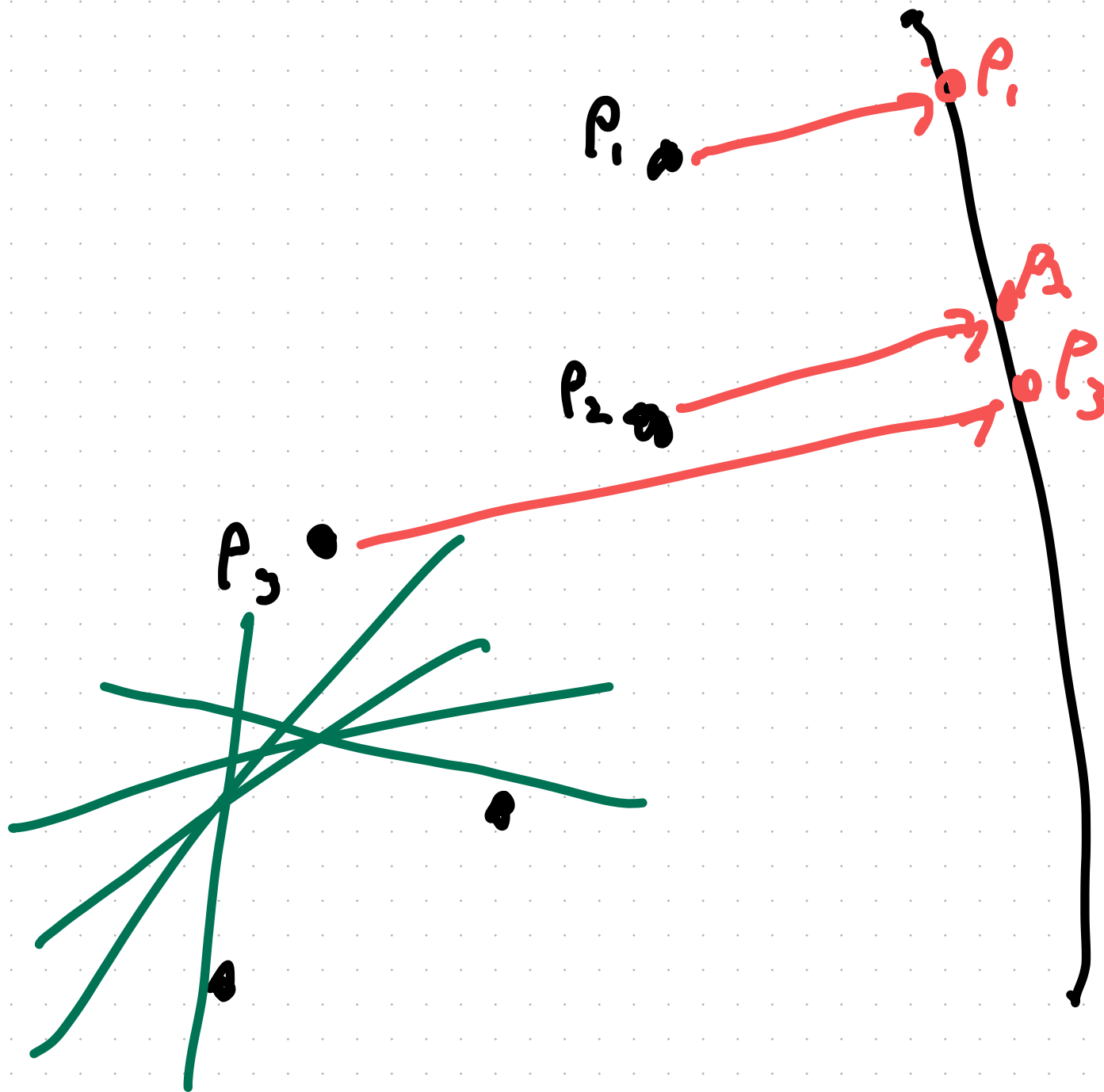
Johsan - Lindenstrauss (JL)

Have a set of s points in \mathbb{R}^d
 d dimension

Then there is a function $JL(x)$ such that
 $JL(x)$ is K dimensional
and $K \geq \frac{20 \ln n}{\epsilon^2}$

for any x, y in S :

$$(1-\epsilon)d(JL(x), JL(y)) \leq d(x, y) \leq (1+\epsilon)d(JL(x), JL(y))$$



$$3 \left\{ \begin{array}{c} \overbrace{\left[\begin{array}{ccccc} - & - & - & - & - \\ - & - & - & - & - \\ - & - & - & - & - \end{array} \right]}^5 \left[\begin{array}{c} 3 \\ 2 \\ 1 \\ 4 \\ 9 \end{array} \right] \end{array} \right\} 5 = \begin{bmatrix} \cdot \\ \cdot \\ 0 \\ \cdot \\ \cdot \end{bmatrix}$$

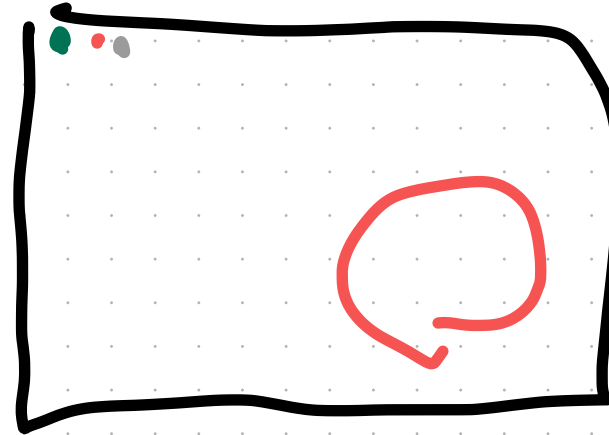
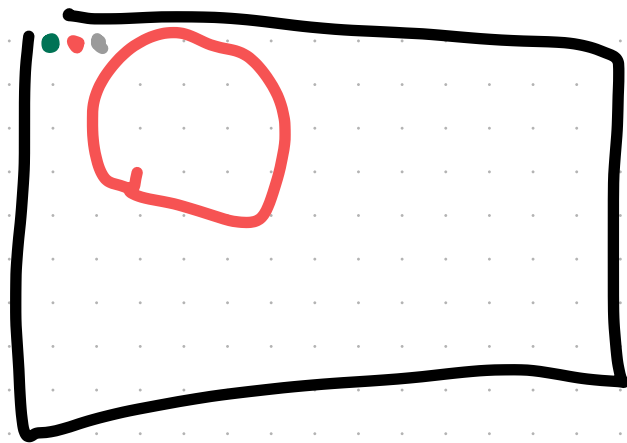
~~① Use a normal distribution~~

Pick -1 1

with 50% chance

if $k \geq \frac{4 + 28}{\frac{\epsilon^2}{2} + \frac{\epsilon}{3}} \ln n$

Then with $1/\epsilon$ works
with prob $\geq 1 - \delta$



Image



Feature selection



few thousand dimensional
vector

Distance Estimation

Approx Nearest Neighbor

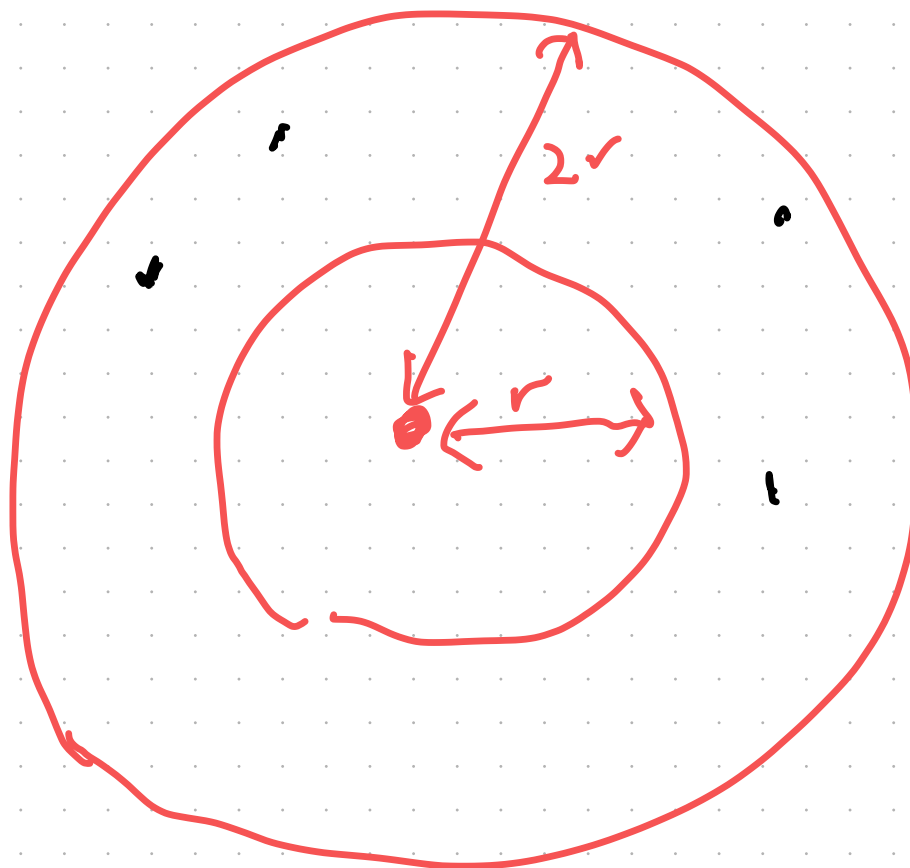


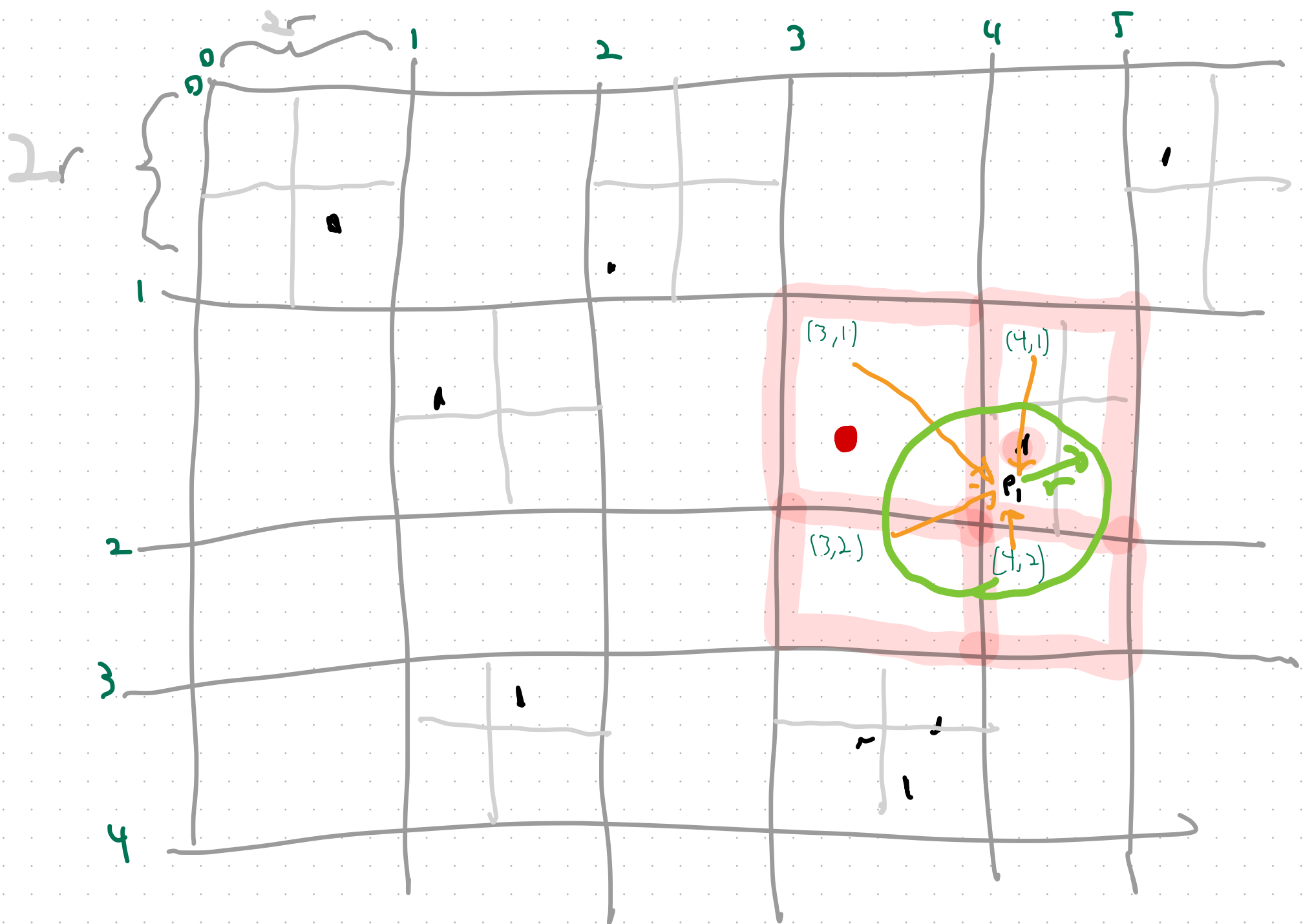
if there is a point
within distance r find it

if there is a point
 $r \dots 2r$
perhaps find it

if there are no points
 $\leq 2r$ report that







$$2^d$$

$$JL \quad d = \log N$$

$$2^{\log N} = N$$

