

Clustering. Unsupervised Learning

硬聚类：每个点只可能是一种聚类.

Maria-Florina Balcan

04/06/2015

Reading:

- Chapter 14.3: Hastie, Tibshirani, Friedman.

Additional resources:

- Center Based Clustering: A Foundational Perspective.
Awasthi, Balcan. Handbook of Clustering Analysis. 2015.

Clustering, Informal Goals

Goal: Automatically partition **unlabeled** data into groups of similar datapoints.

Question: When and why would we want to do this?

Useful for:

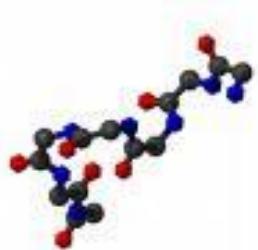
- Automatically organizing data. 数据管理.
- Understanding hidden structure in data.
- Preprocessing for further analysis.
 - Representing high-dimensional data in a low-dimensional space (e.g., for visualization purposes).

Applications (Clustering comes up everywhere...)

- Cluster news articles or web pages or search results by topic.

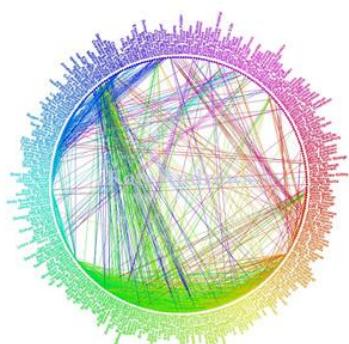


- Cluster protein sequences by function or genes according to expression profile.

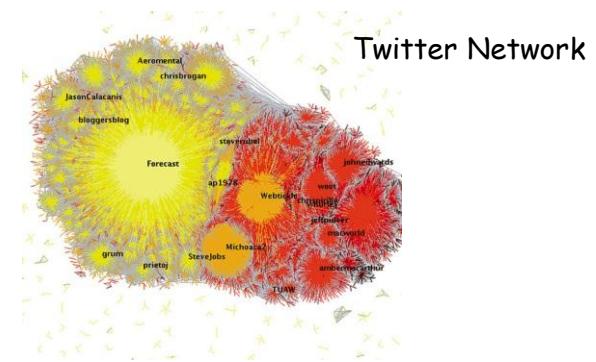


-MTEGGFDPCCECISHERTMARLINILRQSRAYCTNTECIRELPGP--SGDSG-	-ISITVILMAWMVIAVLLFLLRPPHLR-----GFSLPCKP--SSPHS--GQVPPAPPVG--	99
-MTEGGFDPCCECISHERTMARLINILRQSRAYCTNTECIRELPGP--SGDSG-	-ISITVILMAWMVIAVLLFLLRPPNLR-----GFSLPCKP--SSPHS--GQVPPAPPVG--	99
-MTEGGFDPCCECISHERTMARLINILRQSRAYCTNTECIRELPGP--SGDSG-	-ISITAILMWMMVIAVLLFLLRPPNLR-----GFSLPCKP--SSPHS--GQVPPAPPVG--	99
-MTEGGFDPCCECISHERTMARLINILRQSRAYCTNTECIRELPGP--SGDSG-	-ISITAILMWMMVIAVLLFLLRPPNLR-----GFSLPCKP--SSPHS--GQVPPAPPVG--	99
-MTEGGFDPCCECISHERAMARMLINILRQSRAYCTDTRECLRELPGP--SGDSG-	-ISITVILMAWMVIAVLLFLLRPPNLR-----GFSLPCKP--SSPHS--GQVPPAPPVG--	99
-MARGKFDPCCECISCHERAMARMLINILRQSRAYCTDTRECLRELPGP--SGDSG-	-ISITVILMAWMVIAVLLFLLRPPHLR-----GFSLPCKP--SSPHS--GQVPPAPPVG--	99
-MVEGGFDPCCECISCHERAMARMFINILRQSQSYCTNTECIRELPGP--SGDSG-	-ISITVILMAWMVIAVLLFLLRPPNLR-----GFSLPCKP--SSPHS--GQVPPAPPVG--	99
-MVEGGFDPCCECISCHERAMARMFINILRQSQSYCTNTECIRELPGP--SGDSG-	-ISITVILMAWMVIAVLLFLLRPPNLR-----GFSLPCKP--SSPHS--GQVPPAPPVG--	99
-MTEGGFDPCCECIYSHEFAMARMLINILRQSQSYCTNTECIRELPGP--SGDSG-	-ISITVILMAWMVIAVLLFLLRPPNLR-----GFSLPCKP--SSPHS--GQVPPAPPVG--	99
-MAGGGFDPCCECISHEFAMARMLINILRQSQSYCTDTRECLRELPGP--SGDSG-	-ISITVILMAWMVIAVLLFLLRPPNLR-----GSSLPCKP--SSPHS--GQDPPAPPVG--	99
-MAGGGFDPCCECVCSHEFAMARMLINILRQSQSYCTDTRECLRELPGP--SGDSG-	-ISITVILMAWMVIAVLLFLLRPPNLR-----GSSLPCKP--SSPHS--GQDPPAPPVG--	99
-MAGGGFDPCCECVCSHEFAMARMLINILRQSQSYCTDTRECLRELPGP--SGDSG-	-ISITVILMAWMVIAVLLFLLRPPNLR-----GSNLPCKP--SSPHS--GQDPPAPPVG--	99

- Cluster users of social networks by interest (community detection).



Facebook network



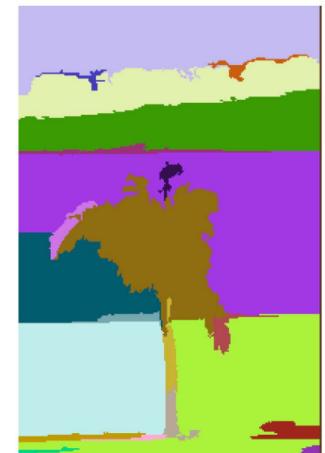
Twitter Network

Applications (Clustering comes up everywhere...)

- Cluster customers according to purchase history.



- Cluster galaxies or nearby stars (e.g. Sloan Digital Sky Survey)



- Image segmentation (clustering pixels)

Clustering

Today:

- Objective based clustering
- Hierarchical clustering
- Mention overlapping clusters

[March 4th: EM-style algorithm for clustering for mixture of Gaussians (specific probabilistic model).]

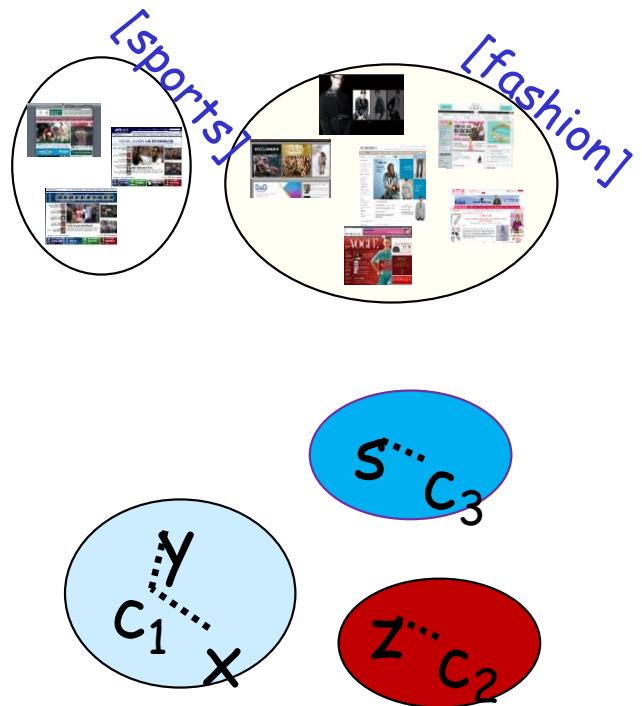
Objective Based Clustering

Input: A set S of n points, also a **distance/dissimilarity** measure specifying the distance $d(x,y)$ between pairs (x,y) .

E.g., # keywords in common, edit distance, wavelets coef., etc.

Goal: output a **partition** of the data.

- k-means: find center pts c_1, c_2, \dots, c_k to
minimize $\sum_{i=1}^n \min_{j \in \{1, \dots, k\}} d^2(x^i, c_j)$
- k-median: find center pts c_1, c_2, \dots, c_k to
minimize $\sum_{i=1}^n \min_{j \in \{1, \dots, k\}} d(x^i, c_j)$
- K-center: find partition to minimize the maxim radius



Euclidean k-means Clustering

Input: A set of n datapoints x^1, x^2, \dots, x^n in \mathbb{R}^d
target #clusters k

Output: k representatives $c_1, c_2, \dots, c_k \in \mathbb{R}^d$

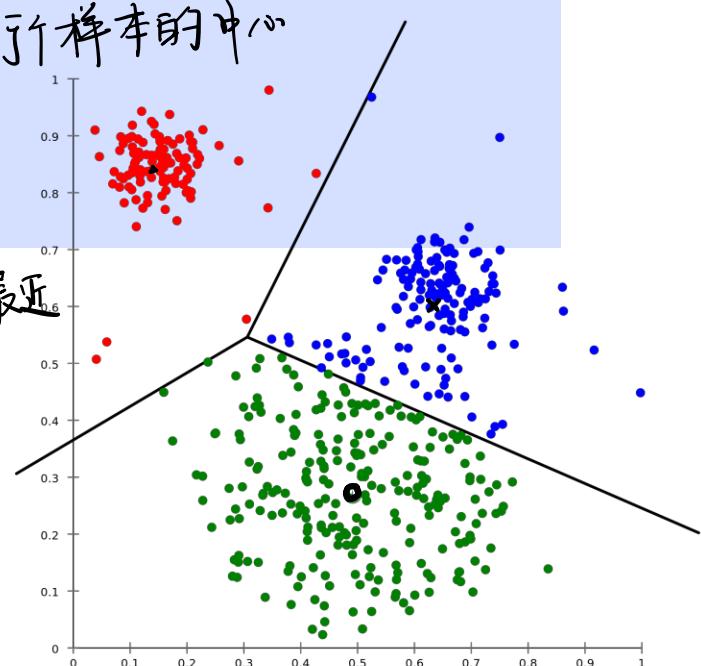
Objective: choose $c_1, c_2, \dots, c_k \in \mathbb{R}^d$ to minimize

$$\sum_{i=1}^n \min_{j \in \{1, \dots, k\}} \|x^i - \underline{c_j}\|^2$$

→ 第 i 个样本的中心

但最开始未知 c_j . 未知 c_j 和哪个点最近

\Rightarrow 非凸且 NP-hard \Rightarrow 依赖初始化.



Euclidean k-means Clustering

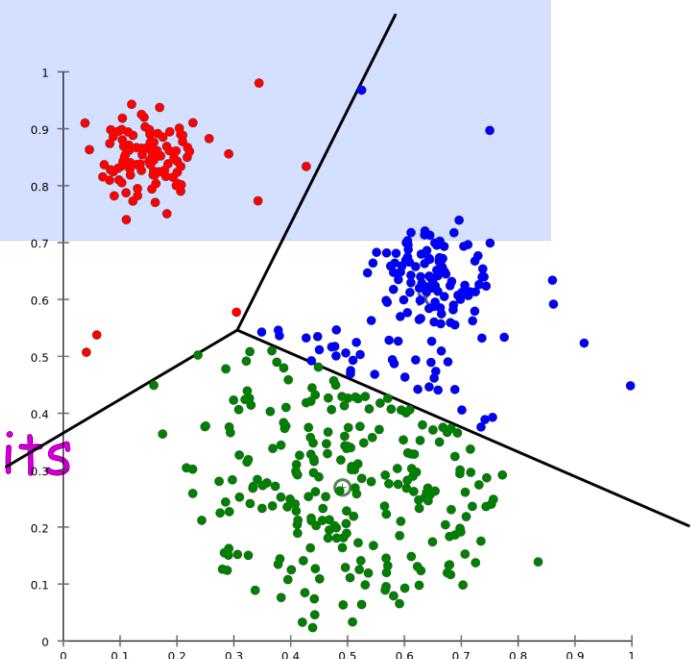
Input: A set of n datapoints x^1, x^2, \dots, x^n in \mathbb{R}^d
target #clusters k

Output: k representatives $c_1, c_2, \dots, c_k \in \mathbb{R}^d$

Objective: choose $c_1, c_2, \dots, c_k \in \mathbb{R}^d$ to minimize

$$\sum_{i=1}^n \min_{j \in \{1, \dots, k\}} \|x^i - c_j\|^2$$

Natural assignment: each point assigned to its
closest center, leads to a Voronoi partition.



简化目标函数: Euclidean k-means Clustering

定义 r_{ij} : x^i 属于 C_j 的指示变量, 只能取 0 或 1.

\rightarrow 全为 1 的 k 行列向量.

\Rightarrow 矩阵 $R: \{0, 1\}^{n \times k}$ n 行 k 列 的矩阵. 约束: $R \times \underline{1}_k = \underline{1}_n$

则目标函数可以化成 $\min_{R, C} \sum_{j=1}^k r_{ij} \|x^i - c_j\|^2$ st. $R \times \underline{1}_k = \underline{1}_n, R \in \{0, 1\}^{n \times k}$

还是一个 NP-hard 且非凸. 依赖于一个 c_j 的初始化.

维度是 d

有 k 类



矩阵化 (去掉乘和符号): $\min_{R, C} \|X - RC\|^2$ C : 所有的 c_j 组成的矩阵, $k \times d$.

$$R: n \times k \quad C: k \times d$$

i -th row: $r_i \cdot C = c_j$

只有第 j 个元素是 1 (第 i 行是 1 其它均为 0)

只有一个聚类.

只有 $k=1$ 时是非闭式解. (+均值)

Euclidean k-means Clustering

Input: A set of n datapoints x^1, x^2, \dots, x^n in \mathbb{R}^d
target #clusters k

Output: k representatives $c_1, c_2, \dots, c_k \in \mathbb{R}^d$

Objective: choose $c_1, c_2, \dots, c_k \in \mathbb{R}^d$ to minimize

$$\sum_{i=1}^n \min_{j \in \{1, \dots, k\}} \|x^i - c_j\|^2$$

Computational complexity:

NP hard: even for $k = 2$ [Dagupta'08] or
 $d = 2$ [Mahajan-Nimborkar-Varadarajan09]

There are a couple of easy cases...



An Easy Case for k-means: k=1

Input: A set of n datapoints x^1, x^2, \dots, x^n in \mathbb{R}^d

Output: $c \in \mathbb{R}^d$ to minimize $\sum_{i=1}^n \|x^i - c\|^2$

Solution: The optimal choice is $\mu = \frac{1}{n} \sum_{i=1}^n x^i$

Idea: bias/variance like decomposition

$$\frac{1}{n} \sum_{i=1}^n \|x^i - c\|^2 = \|\mu - c\|^2 + \frac{1}{n} \sum_{i=1}^n \|x^i - \mu\|^2$$

Avg k-means cost wrt c

Avg k-means cost wrt μ

Quiz

So, the optimal choice for c is μ .

Common Heuristic in Practice: The Lloyd's method

[Least squares quantization in PCM, Lloyd, IEEE Transactions on Information Theory, 1982]

Input: A set of n datapoints x^1, x^2, \dots, x^n in \mathbb{R}^d

Initialize centers $c_1, c_2, \dots, c_k \in \mathbb{R}^d$ and
clusters C_1, C_2, \dots, C_k in any way.

Repeat until there is no further change in the cost.

E-step 给定当前的 c 判断属于哪个类别。

- For each j : $C_j \leftarrow \{x \in S \text{ whose closest center is } c_j\}$
- M-step MLE 优化 θ

k-means, 假设：所有 σ 都是一样的。（三个聚类形状是一样的，且是一个凸集）。

若是非均匀分布：⑥ 则要用 Spectral Clustering Ck-means 和 GM 均无法行驶完全)

Common Heuristic in Practice: The Lloyd's method

[Least squares quantization in PCM, Lloyd, IEEE Transactions on Information Theory, 1982]

Input: A set of n datapoints x^1, x^2, \dots, x^n in \mathbb{R}^d

Initialize centers $c_1, c_2, \dots, c_k \in \mathbb{R}^d$ and
clusters C_1, C_2, \dots, C_k in any way.

Repeat until there is no further change in the cost.

- For each j : $C_j \leftarrow \{x \in S \text{ whose closest center is } c_j\}$
- For each j : $c_j \leftarrow \text{mean of } C_j$

Holding c_1, c_2, \dots, c_k fixed,
pick optimal C_1, C_2, \dots, C_k
保证了每一步都减小了目标函数

$$\text{E-step: } C^{(t)} = \arg \min_{C} \|x - DC^{(t-1)}\|^2$$

Holding C_1, C_2, \dots, C_k fixed,
pick optimal c_1, c_2, \dots, c_k

$$\text{M-step: } R(C, C) = \min_{R} \|x - RC\| \quad \text{with } R(C, C^{(t)}) C^{(t-1)} \leq Q(C, C^{(t)})$$

t-th iteration: $R \leftarrow \arg \min_R \|X - R\|_F$ if $X = R$, $R \leftarrow X$
M-step:
 $C^{(t)} \leftarrow \arg \min_C \|X - R^{(t)} C\|_F^2 \Rightarrow Q(R^{(t)}, C^{(t)}) \leq Q(R^{(t)}, C^{(t-1)})$, (交替优化)

Common Heuristic: The Lloyd's method

Input: A set of n datapoints x^1, x^2, \dots, x^n in \mathbb{R}^d

Initialize centers $c_1, c_2, \dots, c_k \in \mathbb{R}^d$ and
clusters C_1, C_2, \dots, C_k in any way.

Repeat until there is no further change in the cost.

- For each j : $C_j \leftarrow \{x \in S \text{ whose closest center is } c_j\}$
- For each j : $c_j \leftarrow \text{mean of } C_j$

Note: it always converges.

- the cost always drops and
- there is only a finite #s of Voronoi partitions
(so a finite # of values the cost could take)

Initialization for the Lloyd's method

Input: A set of n datapoints x^1, x^2, \dots, x^n in \mathbb{R}^d

Initialize centers $c_1, c_2, \dots, c_k \in \mathbb{R}^d$ and
clusters C_1, C_2, \dots, C_k in any way.

Repeat until there is no further change in the cost.

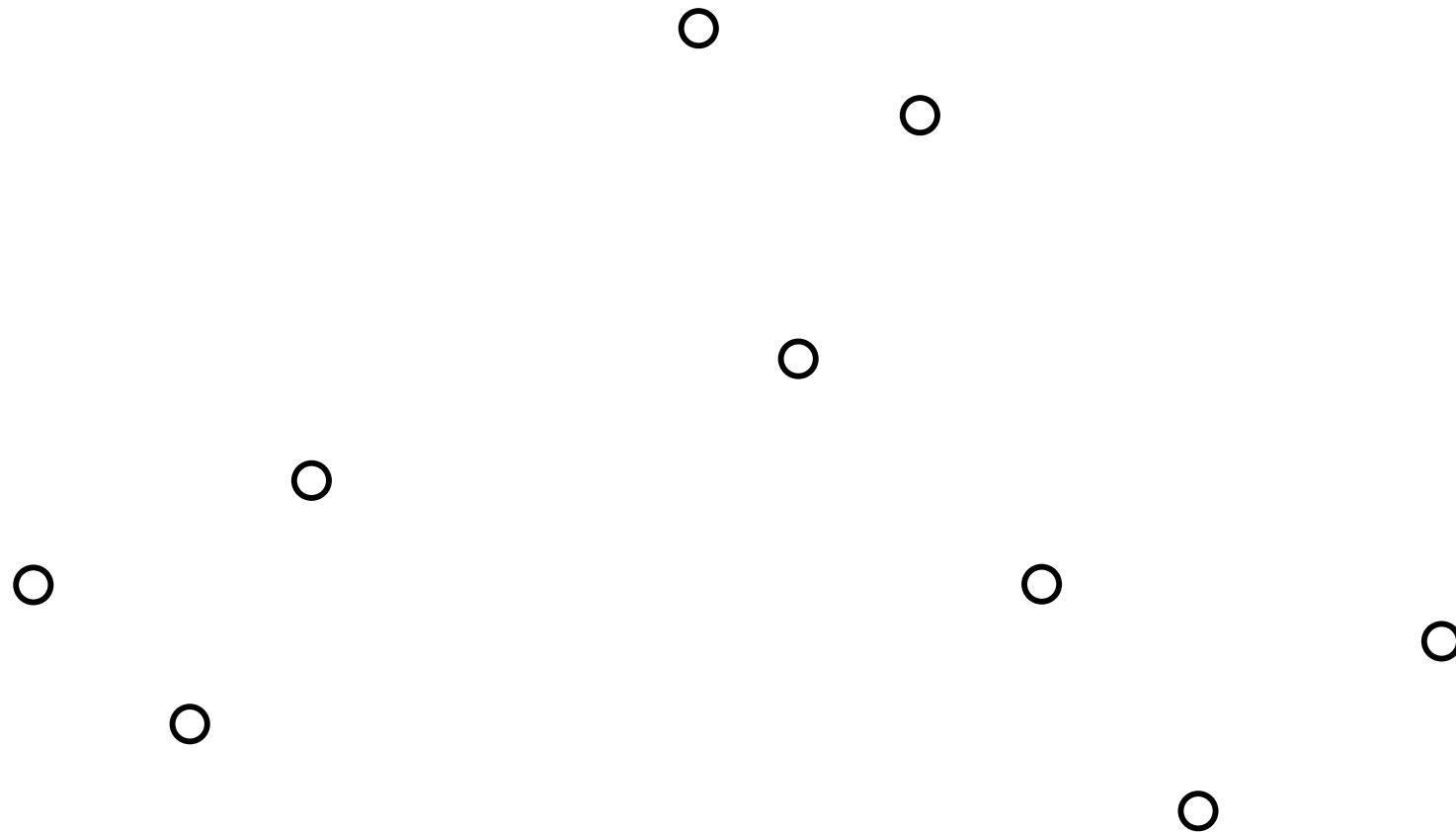
- For each j : $C_j \leftarrow \{x \in S \text{ whose closest center is } c_j\}$
- For each j : $c_j \leftarrow \text{mean of } C_j$

- Initialization is crucial (how fast it converges, quality of solution output)
- Discuss techniques commonly used in practice
 - Random centers from the datapoints (repeat a few times)
随机初始化
 - Furthest traversal 最远点遍历.
 - K-means ++ (works well and has provable guarantees)

Lloyd's method: Random Initialization

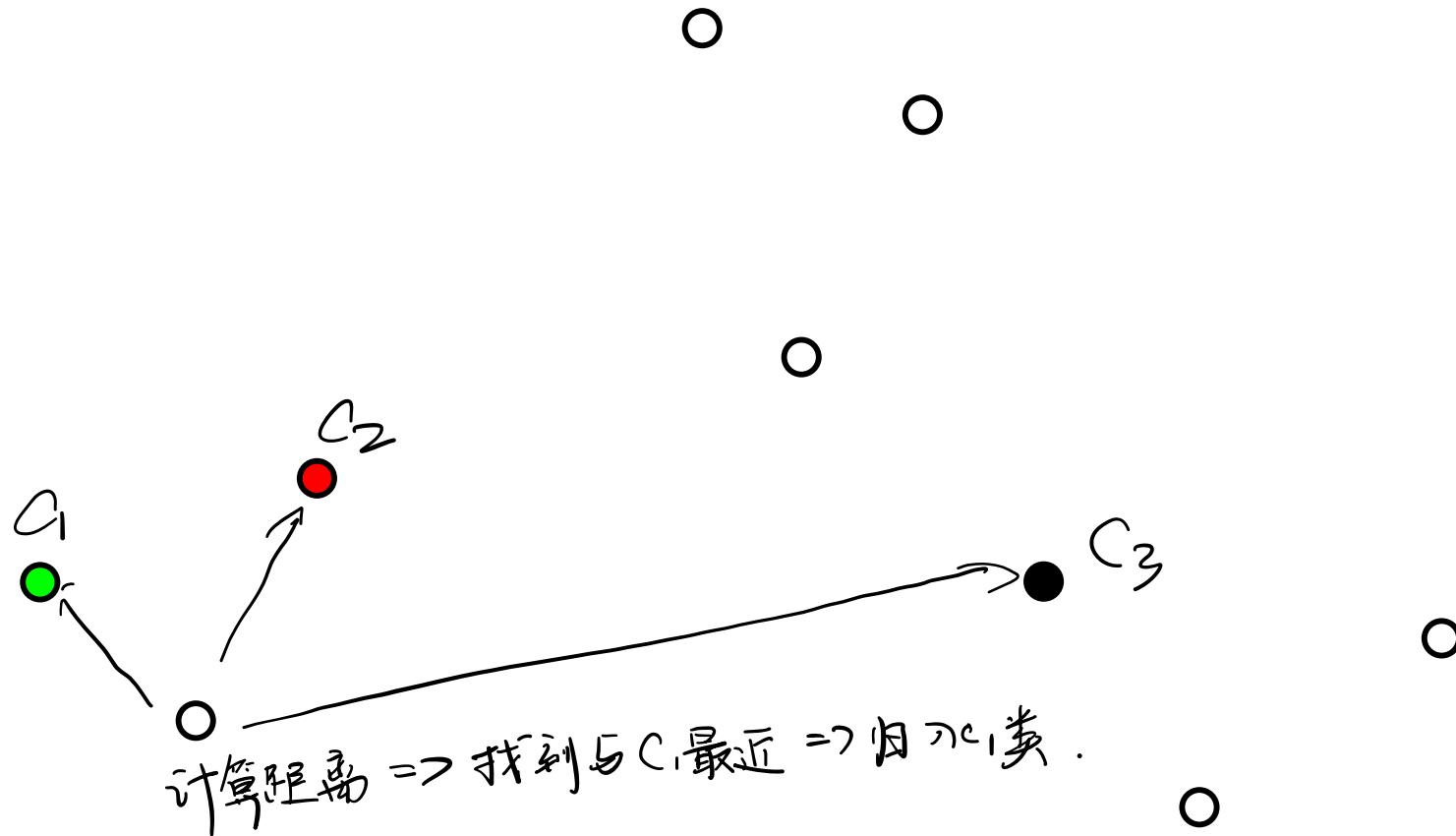
Lloyd's method: Random Initialization

Example: Given a set of datapoints



Lloyd's method: Random Initialization

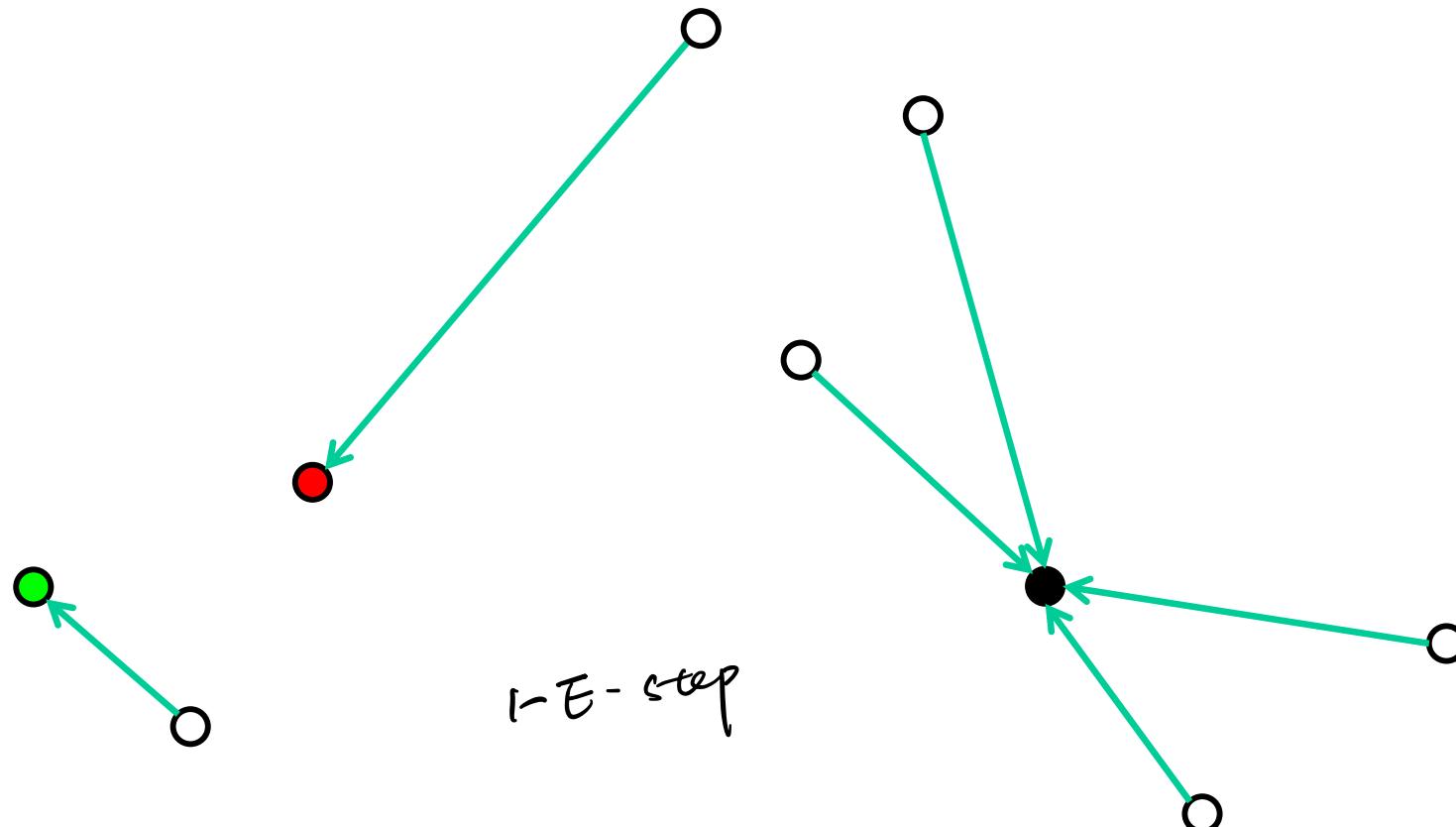
Select initial centers at random



Random Initialize

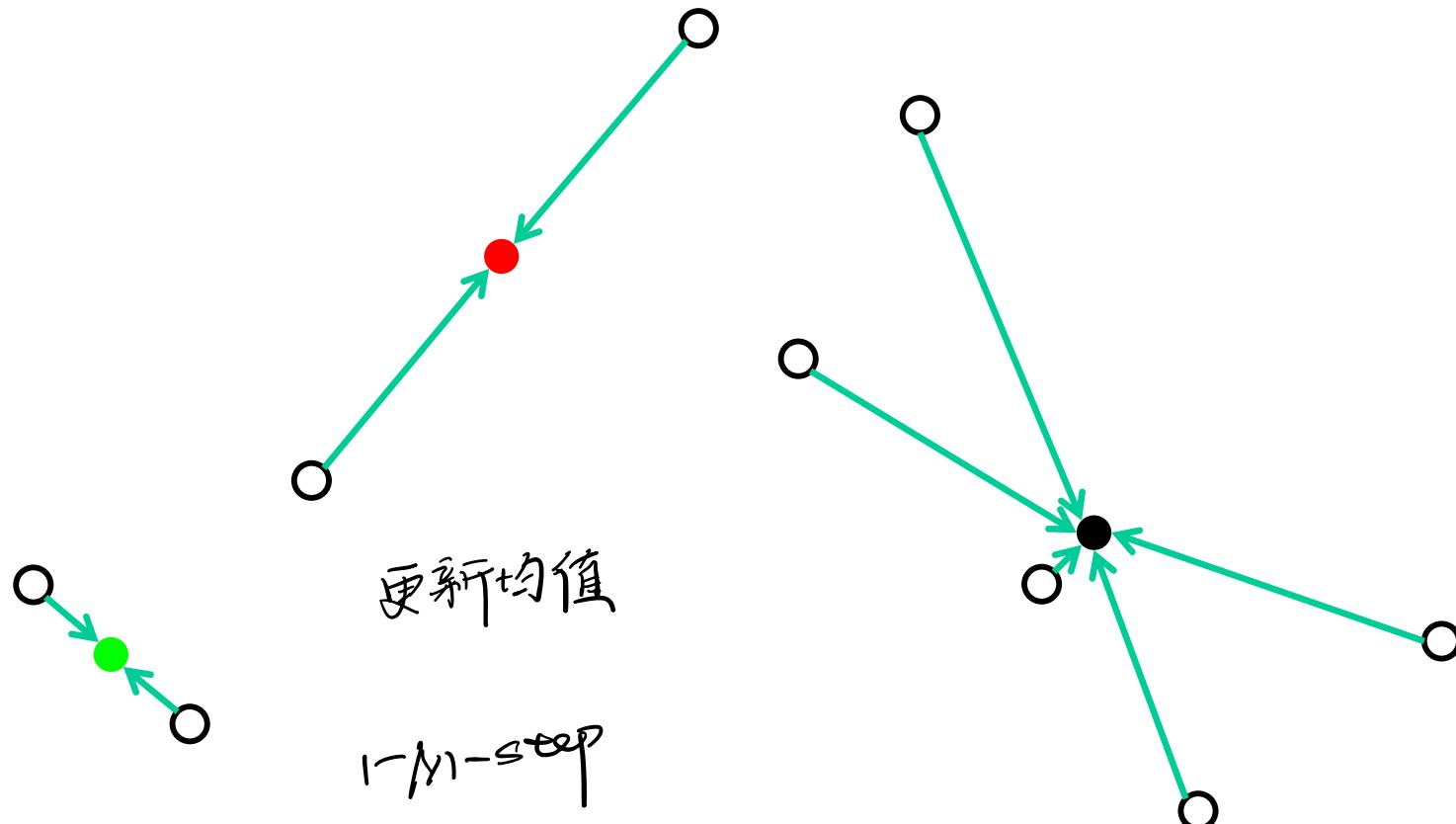
Lloyd's method: Random Initialization

Assign each point to its nearest center



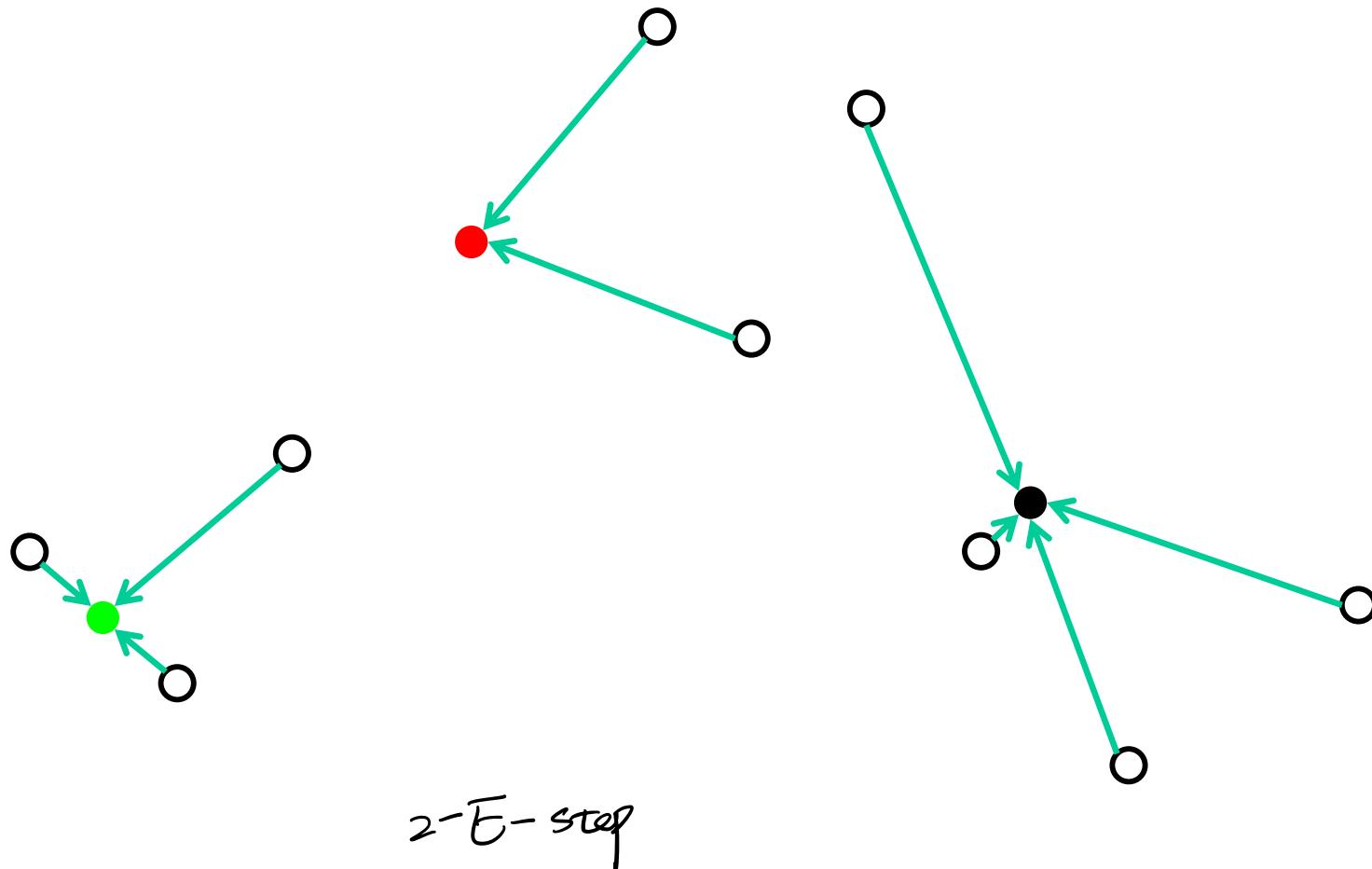
Lloyd's method: Random Initialization

Recompute optimal centers given a fixed clustering



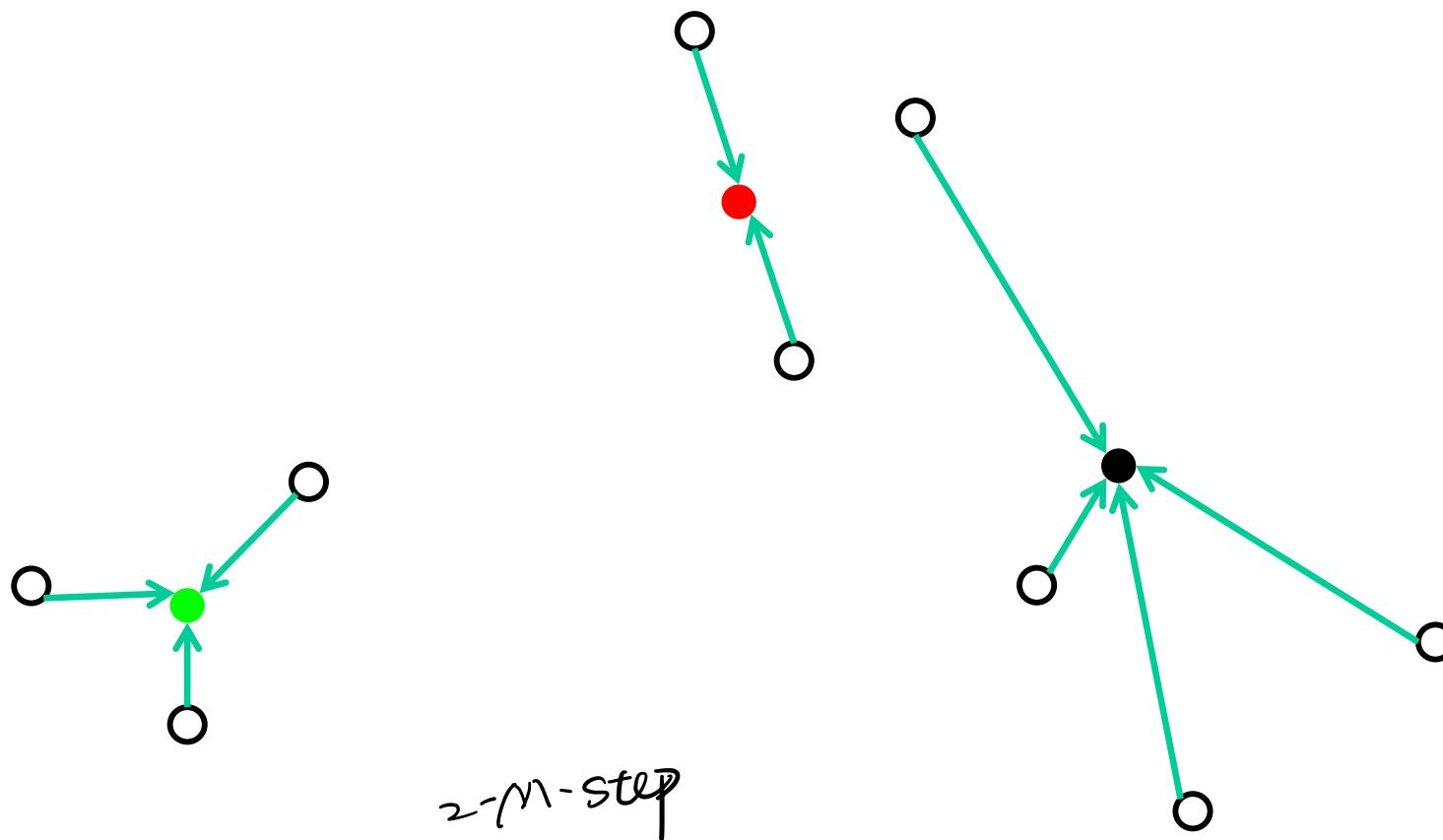
Lloyd's method: Random Initialization

Assign each point to its nearest center



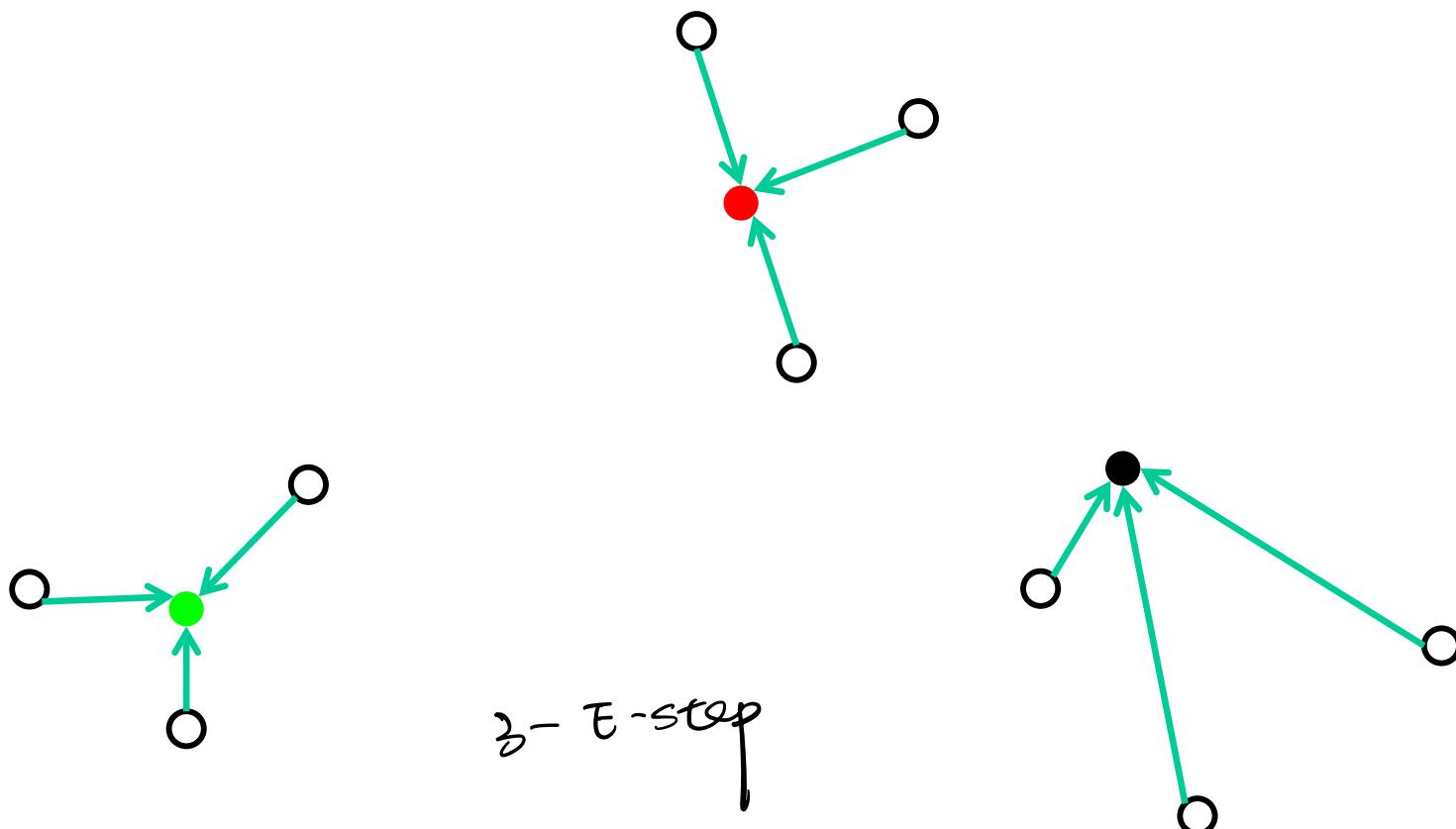
Lloyd's method: Random Initialization

Recompute optimal centers given a fixed clustering



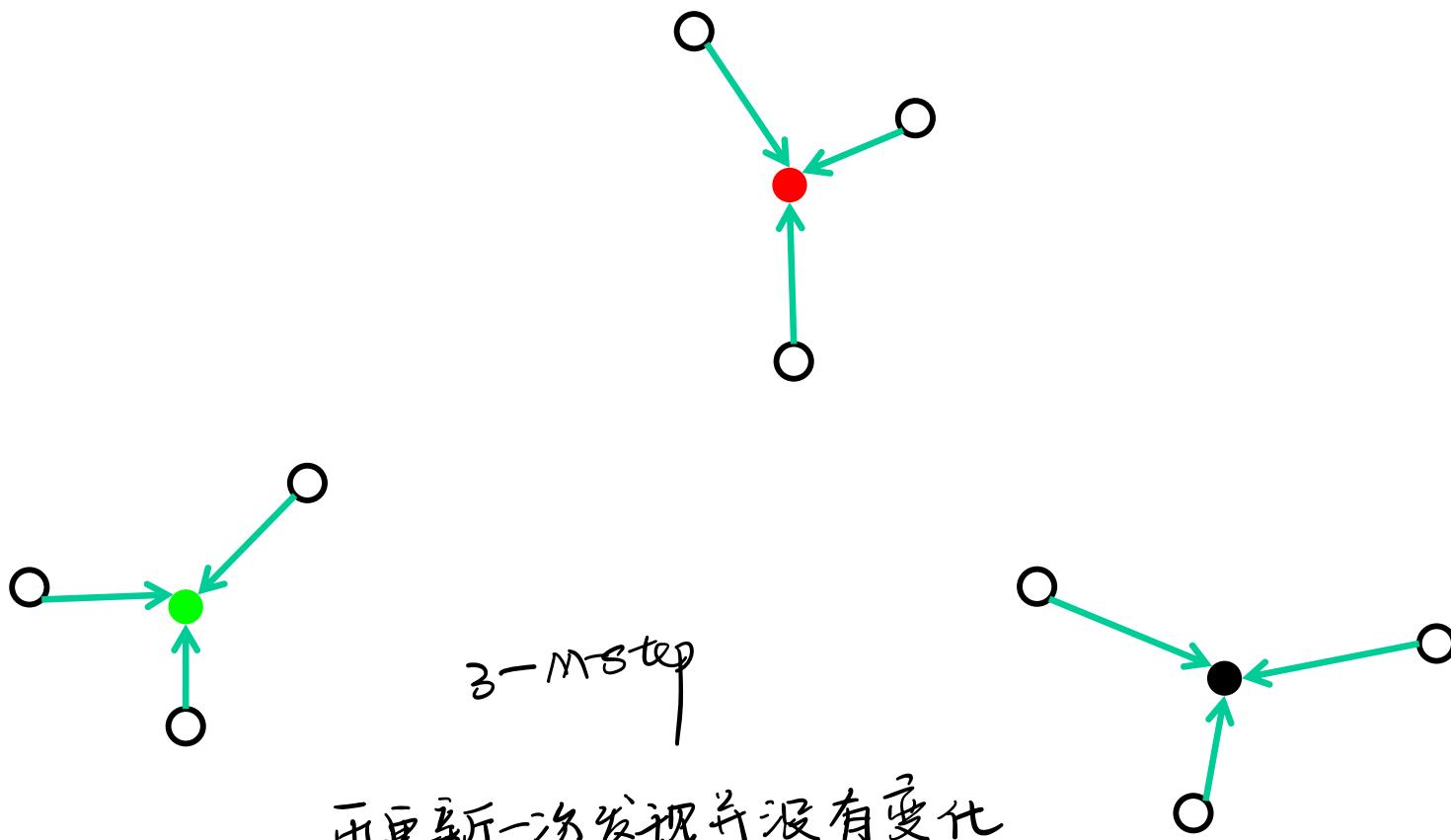
Lloyd's method: Random Initialization

Assign each point to its nearest center



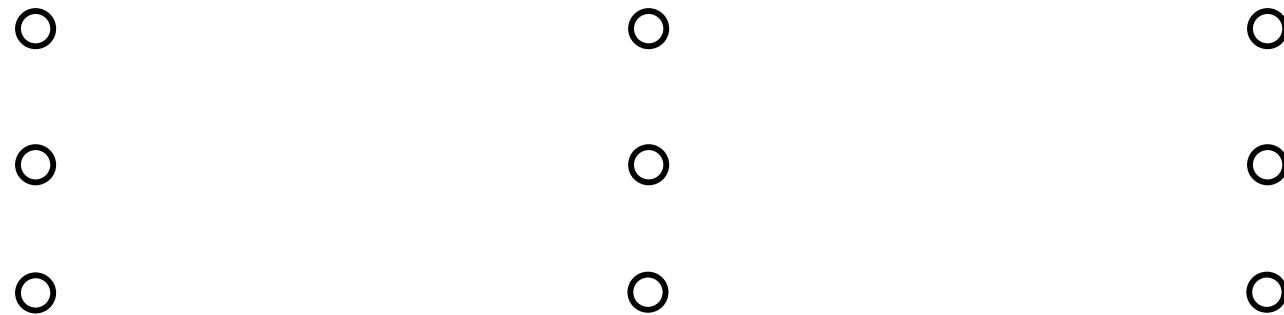
Lloyd's method: Random Initialization

Recompute optimal centers given a fixed clustering



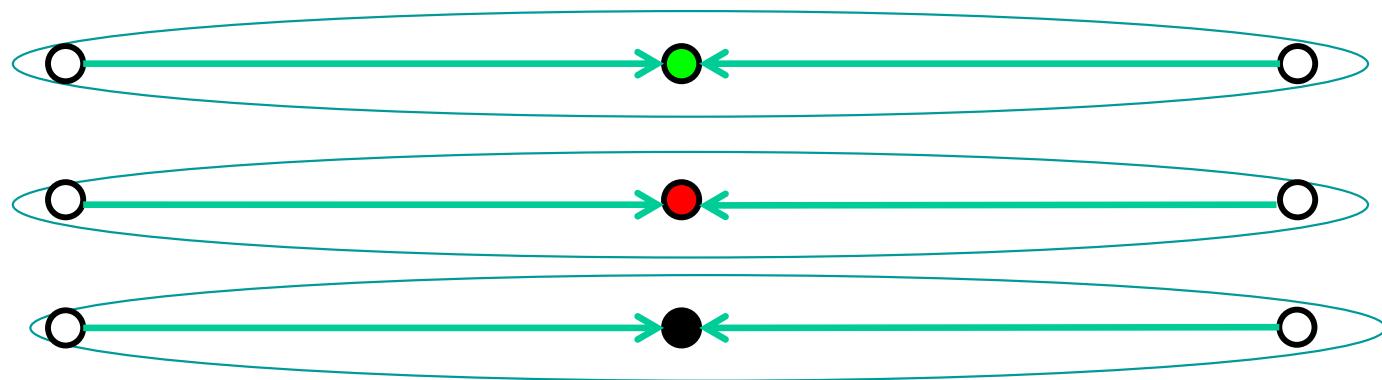
Get a good quality solution in this example.

Lloyd's method: Performance



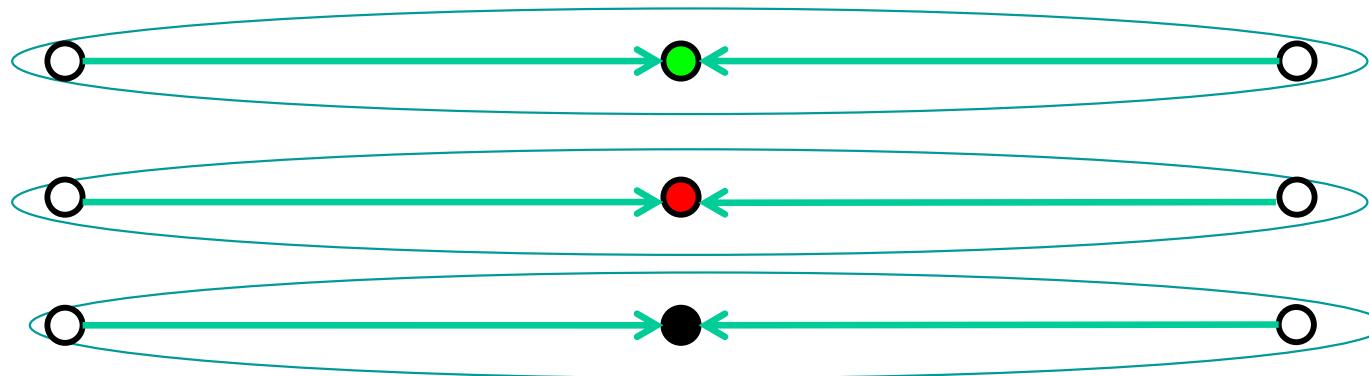
It always converges, but it may converge at a local optimum that is different from the global optimum, and in fact could be arbitrarily worse in terms of its score.

Lloyd's method: Performance

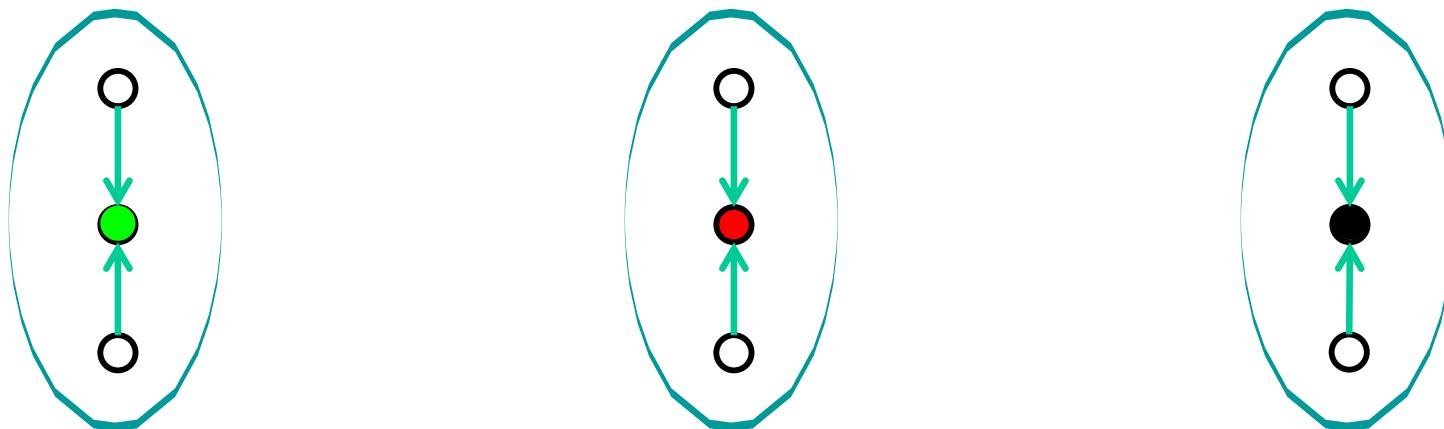


Local optimum: every point is assigned to its nearest center and every center is the mean value of its points.

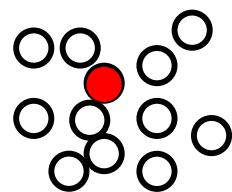
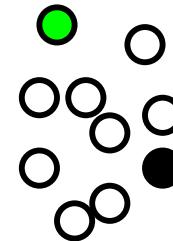
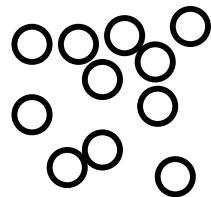
Lloyd's method: Performance



.It is arbitrarily worse than optimum solution....

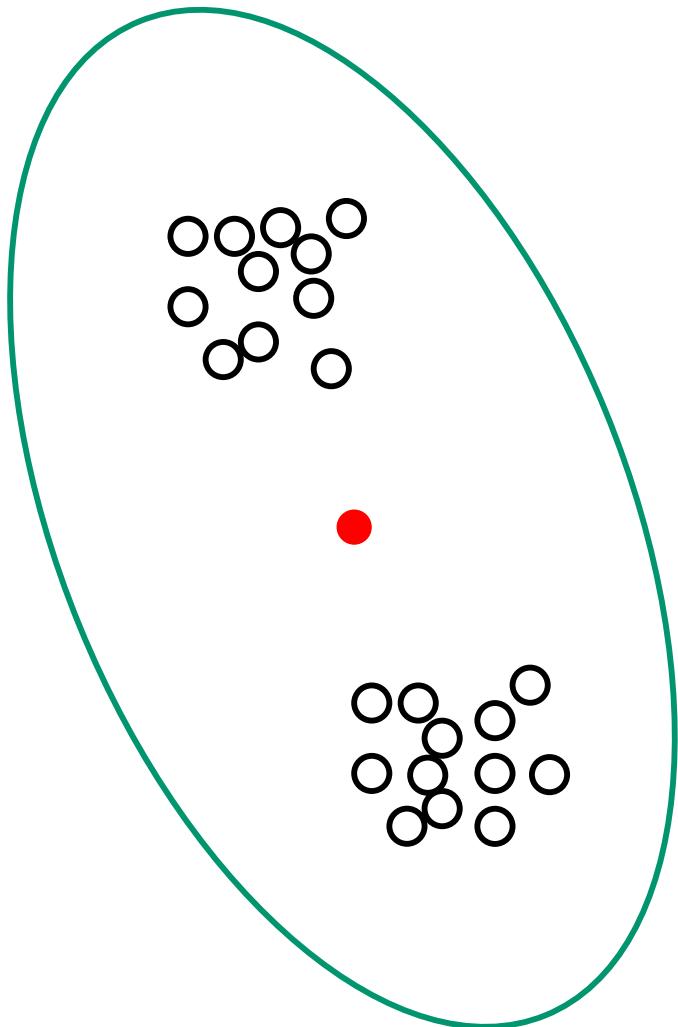


Lloyd's method: Performance



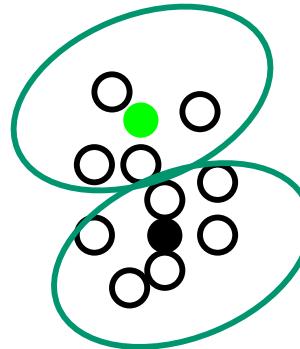
This bad performance, can happen even with well separated Gaussian clusters.

Lloyd's method: Performance



收敛于局部最优解.

但不是全局最优解.



This bad performance, can happen even with well separated Gaussian clusters.

Some Gaussian are combined.....



Lloyd's method: Performance

- If we do random initialization, as k increases, it becomes more likely we won't have perfectly picked one center per Gaussian in our initialization (so Lloyd's method will output a bad solution).

认为 k 个 Gaussian Distribution 相同

- For k equal-sized Gaussians, $\Pr[\text{each initial center is in a different Gaussian}] \approx \frac{k!}{k^k} \approx \frac{1}{e^k}$ → 好的初始化的概率.
• Becomes unlikely as k gets large. 随机初始化概率太小，故不用.

Another Initialization Idea: Furthest Point Heuristic

最远点遍历启发式.

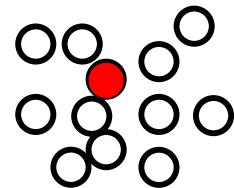
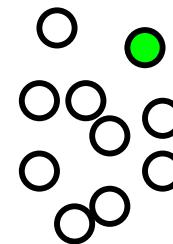
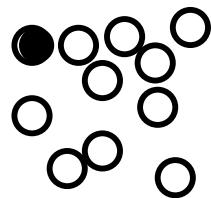
第一次是 Random. 然后每次, 确定 c_1 后

Choose c_1 arbitrarily (or at random). c_2 是距 c_1 最远的点.
 c_3 是距 $c_1 c_2$ 最远的点.

- For $j = 2, \dots, k$ 缺点: 有异常值受影响大
优: 理想情况下优秀.
 - Pick c_j among datapoints x^1, x^2, \dots, x^d that is farthest from previously chosen c_1, c_2, \dots, c_{j-1}

Fixes the Gaussian problem. But it can be thrown off by outliers....

Furthest point heuristic does well on previous example



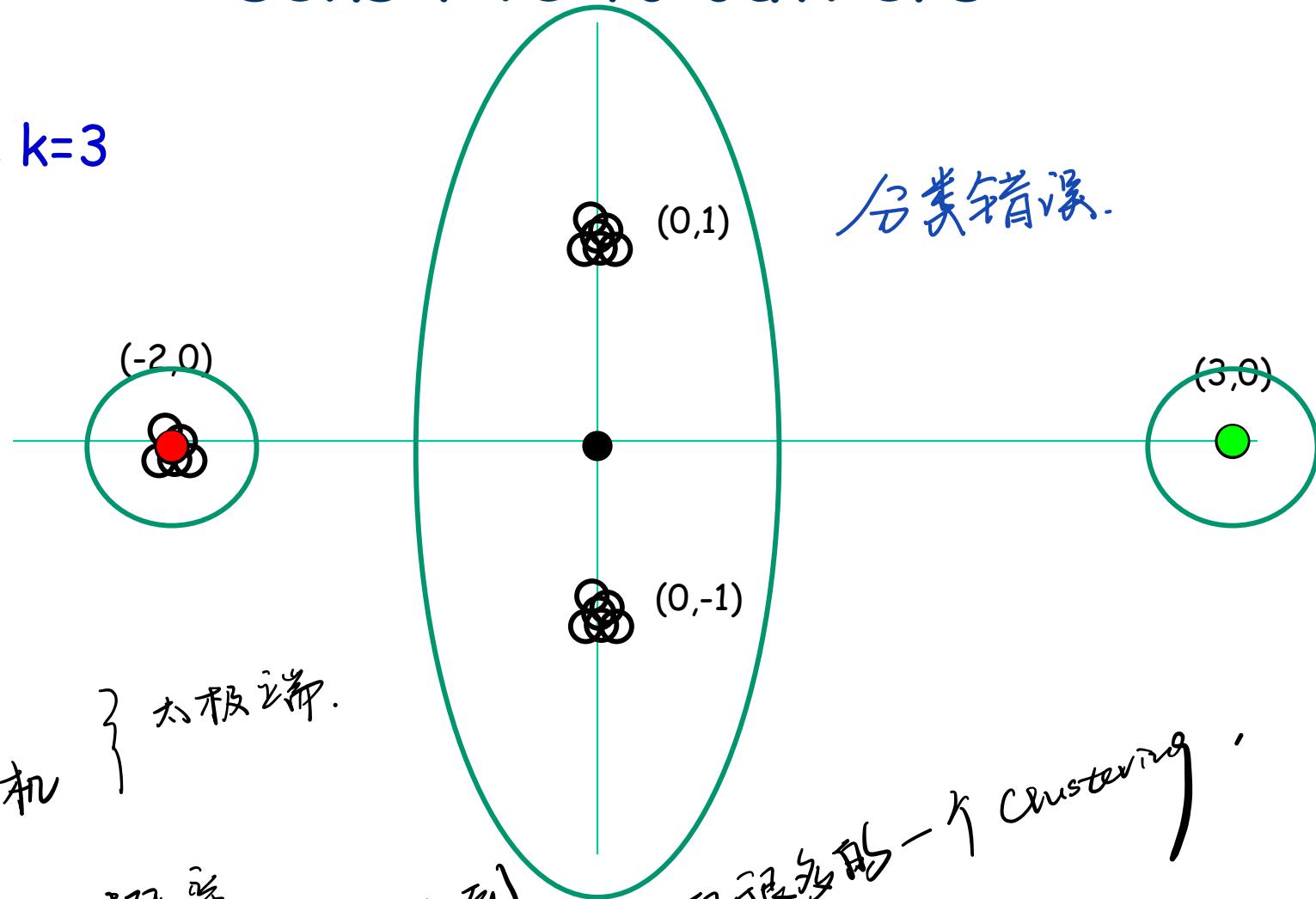
Furthest point initialization heuristic sensitive to outliers

Assume $k=3$



Furthest point initialization heuristic sensitive to outliers

Assume $k=3$



优化：每距离矩阵计算。
可能造出很多的 clustering。
⇒ 最远点，依然很可能选到
但是也可能选中其它次远且很多的

K-means++ Initialization: D² sampling [AV07]

- Interpolate between random and furthest point initialization
- Let $D(x)$ be the distance between a point x and its nearest center. Choose the next center proportional to $D^2(x)$.

- Choose c_1 at random.
- For $j = 2, \dots, k$
 - Pick c_j among x^1, x^2, \dots, x^n according to the distribution

$$\Pr(c_j = x^i) \propto \min_{j' < j} \|x^i - c_{j'}\|^2 / D^2(x^i)$$

在该分布下抽取 c_j

Theorem: K-means++ always attains an $O(\log k)$ approximation to optimal k-means solution in expectation.

以 $O(\log k)$ 的速率逼近

Running Lloyd's can only further improve the cost.

K-means++ Idea: D^2 sampling

- Interpolate between random and furthest point initialization
- Let $D(x)$ be the distance between a point x and its nearest center. Choose the next center proportional to $D^\alpha(x)$.

$$\|x^i - C_j\|_2^\alpha = C \cdot |x^i - C_j|^\alpha$$

- $\alpha = 0$, random sampling

无穷范数：选择最大值。

- $\alpha = \infty$, furthest point

(Side note: it actually works well for k-center)

- $\alpha = 2$, k-means++ 欧氏距离。→是随机与最远点的权衡。

Side note: $\alpha = 1$, works well for k-median

K-means ++ Fix



K-means++/ Lloyd's Running Time

- K-means ++ initialization: $O(nd)$ and one pass over data to select next center. So $O(nkd)$ time in total.
- Lloyd's method 迭代 T 次 : $O(nkd + Tnkdt)$

Repeat until there is no change in the cost.

- For each j : $c_j \leftarrow \{x \in S \text{ whose closest center is } c_j\}$
- For each j : $c_j \leftarrow \text{mean of } C_j$

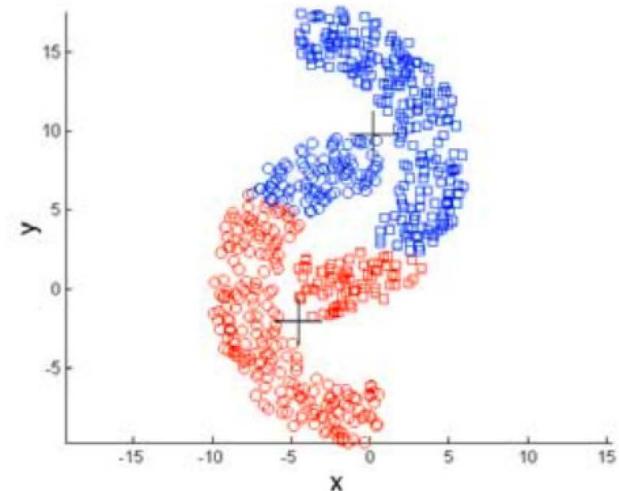
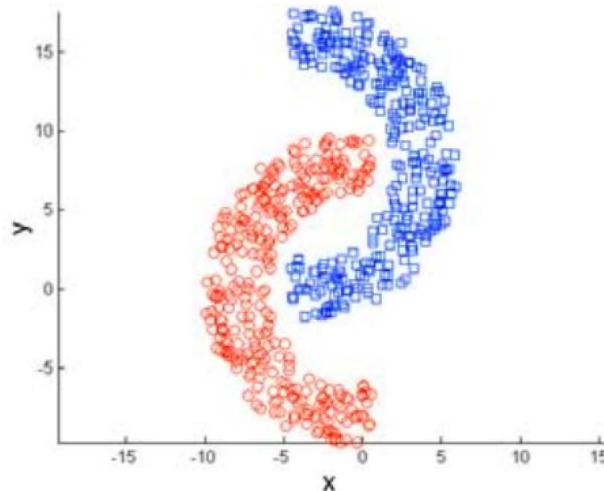
Each round takes time $O(nkd)$.

最差时 T 为指教次数 (迭代次数)

- Exponential # of rounds in the worst case [AV07].
- Expected polynomial time in the smoothed analysis model!

K-means++/ Lloyd's Summary

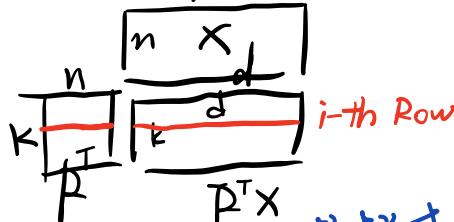
- K-means++ always attains an $O(\log k)$ approximation to optimal k-means solution in expectation.
- Running Lloyd's can only further improve the cost.
- Exponential # of rounds in the worst case [AV07].
- Expected polynomial time in the smoothed analysis model!
- Does well in practice.



Kernel K-means

$$\text{矩阵化目标函数} = \min_{RC} \|X - RC\|_F^2 \quad \text{s.t. } R \in \{0, 1\}^{n \times k}, \quad R \times \mathbb{1}_k = \mathbb{1}_n$$

$$P^T X :$$



假设 P_i^T 中有 m_i 1 \Rightarrow 从 X 中找到 m 个样本点.

→ 对角矩阵, 描述样本个数.

$$\begin{array}{c} (P^T R)^{-1} R^T X \\ \frac{R^T n \times d}{k \times k} \quad \frac{k \times n \times d}{k \times k} \\ \frac{k \times d}{k \times d} \end{array} \Rightarrow \begin{array}{c} d \\ \frac{k \times d}{k \times d} \end{array} \quad \text{i-th Row}$$

就是 C 矩阵.

$$\begin{aligned} Q(R) &= \|X - R(P^T R)^{-1} R^T X\|_F^2 \\ &= \underbrace{\langle X, X \rangle}_{\text{常数}} - 2 \underbrace{\langle X, R(P^T R)^{-1} R^T X \rangle}_{+ \underbrace{\langle R(P^T R)^{-1} R^T X, R(P^T R)^{-1} R^T X \rangle}_{= - \operatorname{Tr}(X^T R(P^T R)^{-1} R^T X)}} \\ &= - \operatorname{Tr}(X^T R(P^T R)^{-1} R^T X) \end{aligned}$$

$$\min_R Q(R) = \max_R \operatorname{Tr}(X^T R(P^T R)^{-1} R^T X) = \max_R \operatorname{Tr}(R(P^T R)^{-1} R^T \underline{X X^T})$$

内积
可在这里套 kernel

$$R(P^T R)^{-1} R^T = \underbrace{R(P^T R)^{-\frac{1}{2}}}_{\text{定义为 } F} \underbrace{(P^T R)^{-\frac{1}{2}} R^T}_{F^T} = F F^T$$

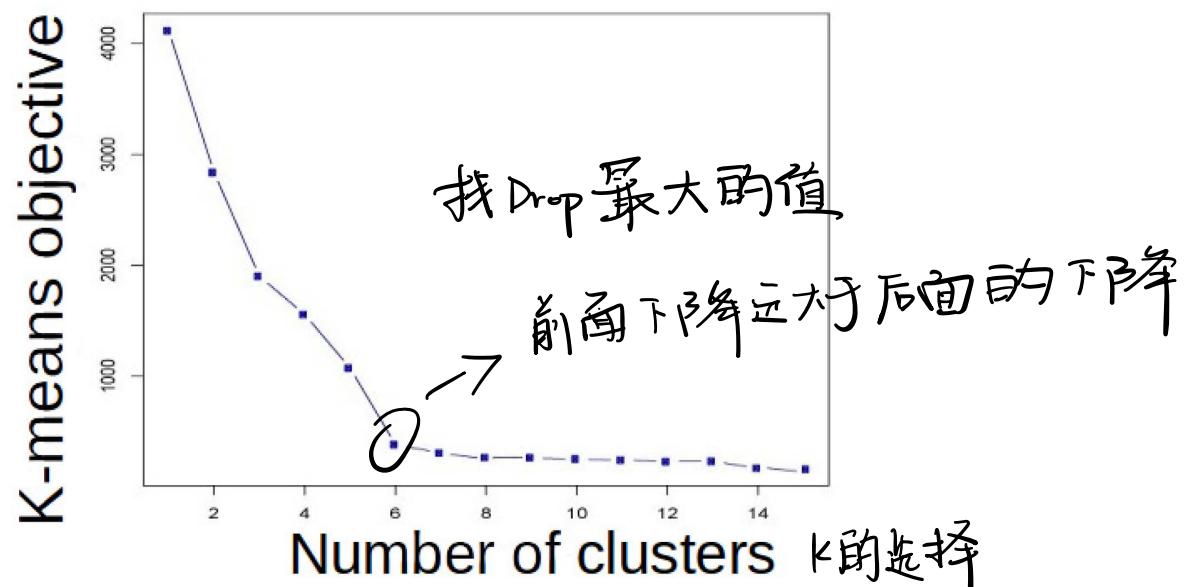
$$\Rightarrow \max_F \operatorname{Tr}(F X \underline{X^T} F^T) \quad \text{而 } \underline{F^T F} = I \quad (\text{特征值分解, 可给出闭式解})$$

去掉 $\{0, 1\}$ 的约束

$$F^T F = (P^T R)^{-\frac{1}{2}} R^T R (P^T R)^{-\frac{1}{2}} = (P^T R)^{\frac{1}{2}} (P^T R)^{-\frac{1}{2}} = I.$$

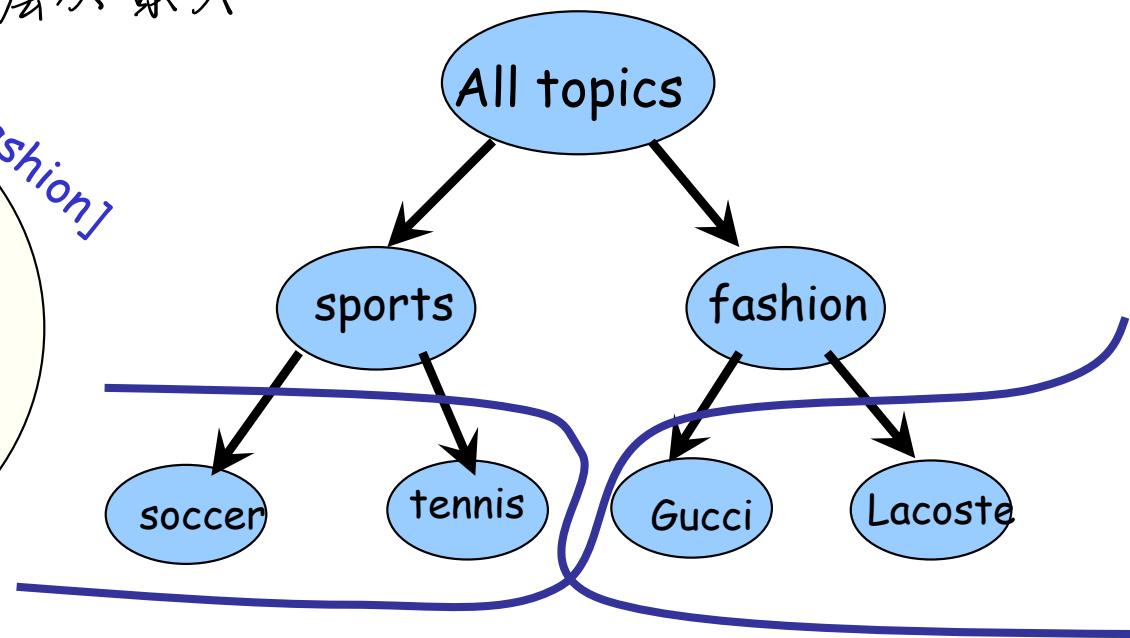
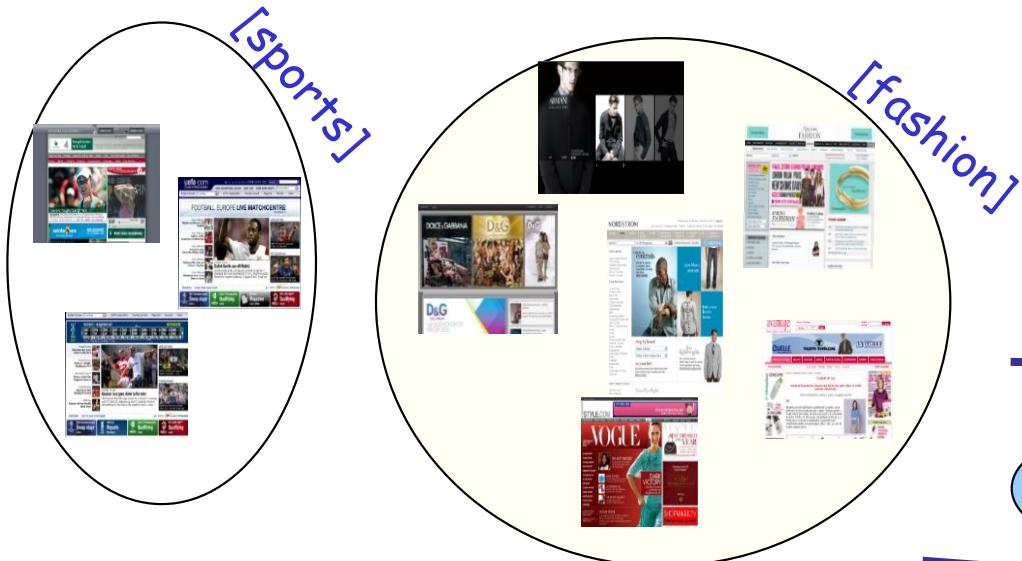
What value of k???

- Heuristic: Find large gap between $k - 1$ -means cost and k -means cost.
- According to information criteria (AIC, BIC, etc.)
- Try hierarchical clustering.



Hierarchical Clustering

层次聚类

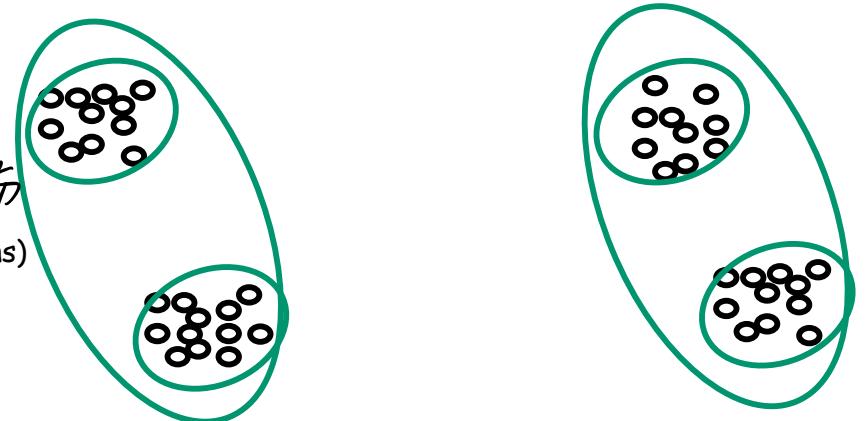


- A hierarchy might be more natural.
- Different users might care about different levels of granularity or even prunings.

Hierarchical Clustering

Top-down (divisive) 自上而下

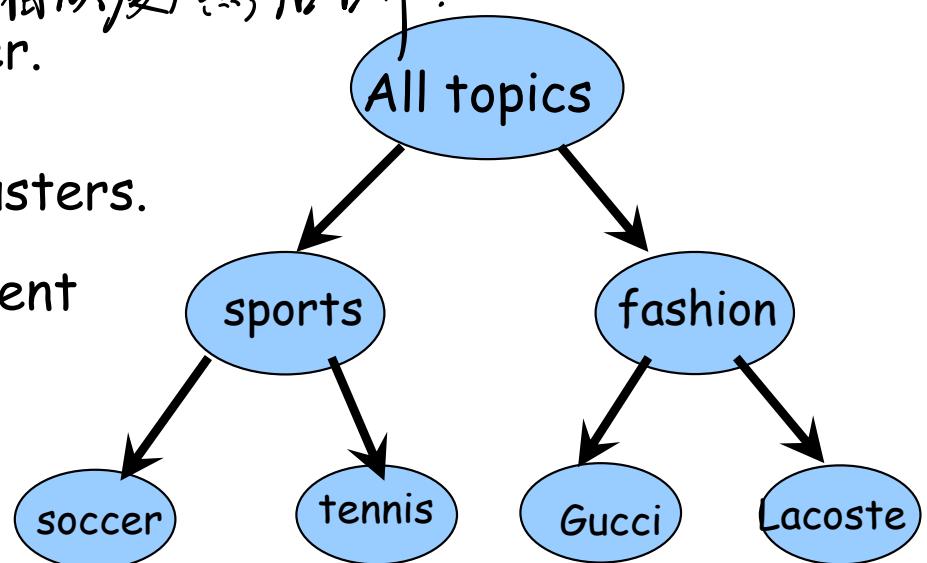
- Partition data into 2-groups (e.g., 2-means)
- Recursively cluster each group.



Bottom-Up (agglomerative)

最开始分成几个子集(相似度),然,后合并.

- Start with every point in its own cluster.
- Repeatedly merge the "closest" two clusters.
- Different defs of "closest" give different algorithms.

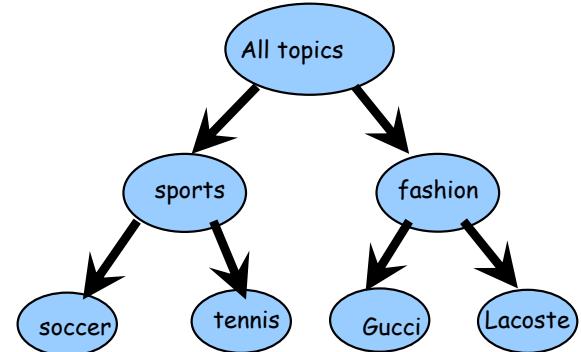


Bottom-Up (agglomerative)

Have a **distance** measure on pairs of objects.

$d(x,y)$ - distance between x and y

E.g., # keywords in common, edit distance, etc



- Single linkage:

作为两集合距离。
$$\text{dist}(A, B) = \min_{x \in A, x' \in B'} \text{dist}(x, x')$$
 两集合间最小值

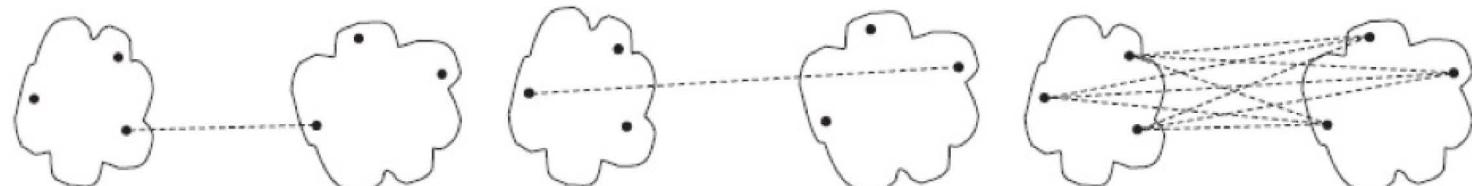
- Complete linkage:

$$\text{dist}(A, B) = \max_{x \in A, x' \in B'} \text{dist}(x, x') \quad \text{最大值}$$

- Average linkage:

$$\text{dist}(A, B) = \text{avg}_{x \in A, x' \in B'} \text{dist}(x, x') \quad \text{均值}$$

- Ward's method



(a) MIN (single link.)

(b) MAX (complete link.)

(c) Group average.

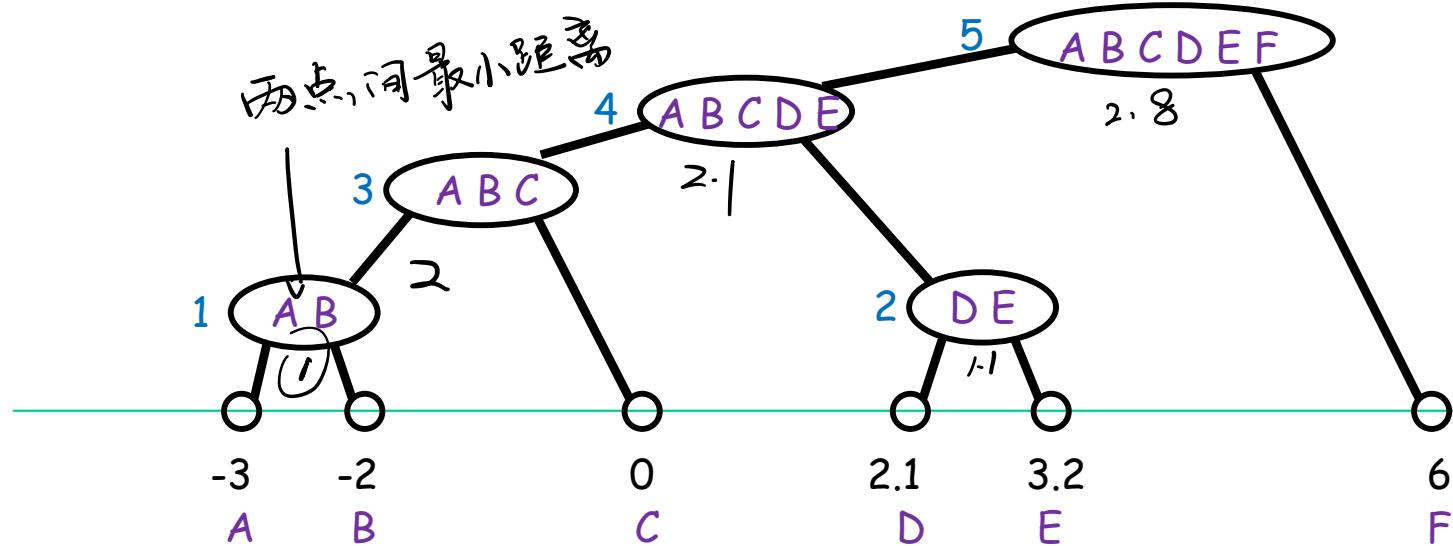
Single Linkage

Bottom-up (agglomerative)

- Start with every point in its own cluster.
- Repeatedly merge the “closest” two clusters.

Single linkage: $\text{dist}(A, B) = \min_{x \in A, x' \in B} \text{dist}(x, x')$

Dendrogram



Single Linkage

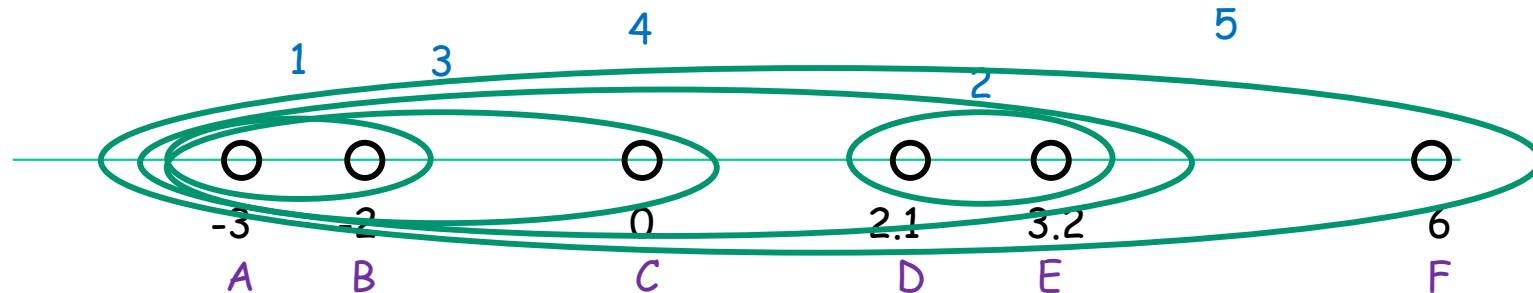
Bottom-up (agglomerative)

- Start with every point in its own cluster.
- Repeatedly merge the “closest” two clusters.

Single linkage: $\text{dist}(A, B) = \min_{x \in A, x' \in B} \text{dist}(x, x')$

One way to think of it: at any moment, we see connected components of the graph where connect any two pts of distance $< r$.

Watch as r grows (only $n-1$ relevant values because we only merge at value of r corresponding to values of r in different clusters).



Complete Linkage

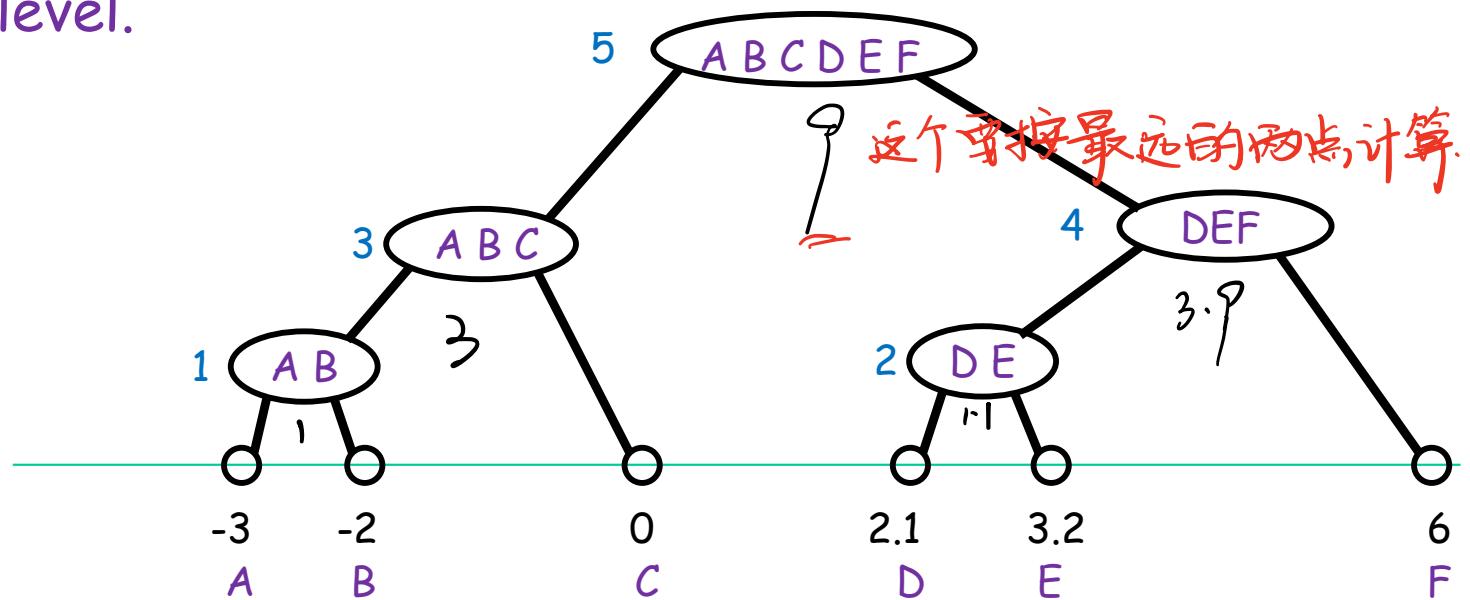
计算最大距离.

Bottom-up (agglomerative)

- Start with every point in its own cluster.
- Repeatedly merge the "closest" two clusters.

Complete linkage: $\text{dist}(A, B) = \max_{x \in A, x' \in B} \text{dist}(x, x')$

One way to think of it: keep max diameter as small as possible at any level.



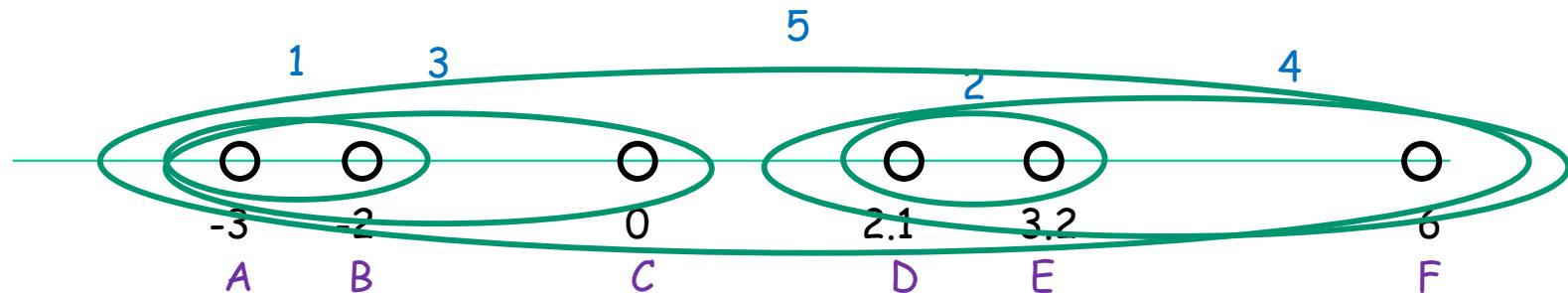
Complete Linkage

Bottom-up (agglomerative)

- Start with every point in its own cluster.
- Repeatedly merge the “closest” two clusters.

Complete linkage: $\text{dist}(A, B) = \max_{x \in A, x' \in B} \text{dist}(x, x')$

One way to think of it: keep max diameter as small as possible.



Ward's Method

Bottom-up (agglomerative)

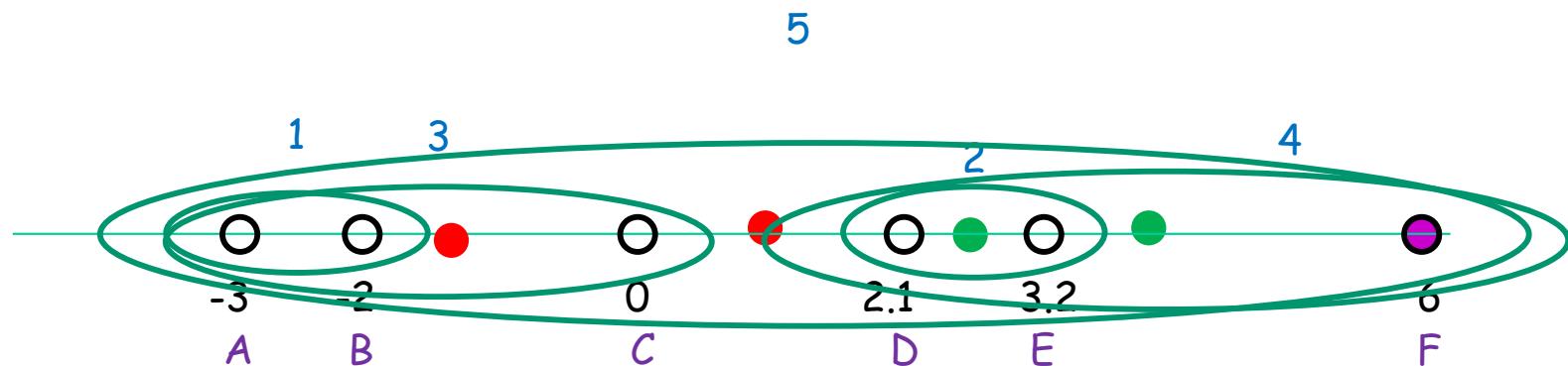
- Start with every point in its own cluster.
- Repeatedly merge the “closest” two clusters.

Ward's method: $\text{dist}(C, C') = \frac{|C| \cdot |C'|}{|C| + |C'|} \|\text{mean}(C) - \text{mean}(C')\|^2$

→ 样本个数作为权重
均值距离.

Merge the two clusters such that the increase in k-means cost is as small as possible.

Works well in practice.



Running time

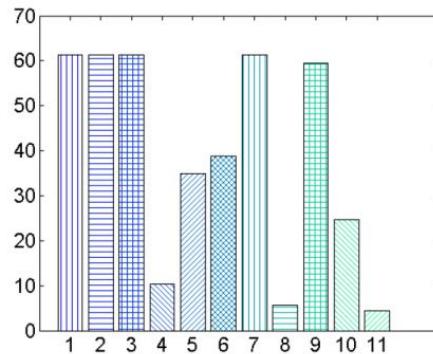
- Each algorithm starts with N clusters, and performs $N-1$ merges.
- For each algorithm, computing $dist(C, C')$ can be done in time $O(|C| \cdot |C'|)$. (e.g., examining $dist(x, x')$ for all $x \in C, x' \in C'$)
- Time to compute all pairwise distances and take smallest is $\underline{\underline{O(N^2)}}$.
- Overall time is $O(N^3)$.

In fact, can run all these algorithms in time $O(N^2 \log N)$.

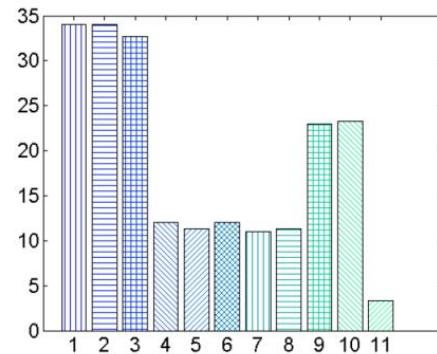
See: Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze, Introduction to Information Retrieval, Cambridge University Press. 2008. <http://www-nlp.stanford.edu/IR-book/>

Hierarchical Clustering Experiments

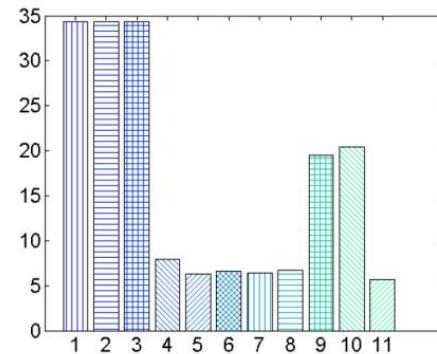
[BLG, JMLR'15]



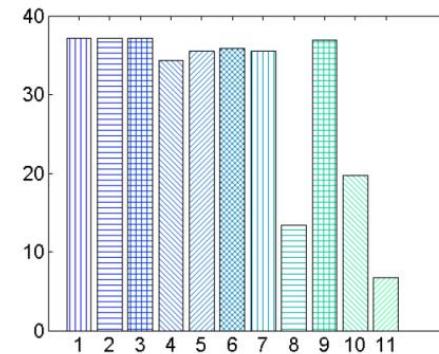
(a) Wine



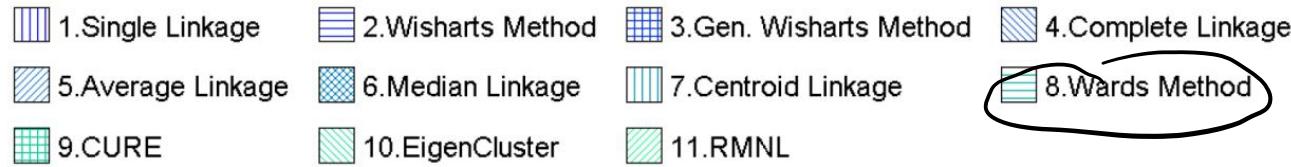
(b) Iris



(c) BCW



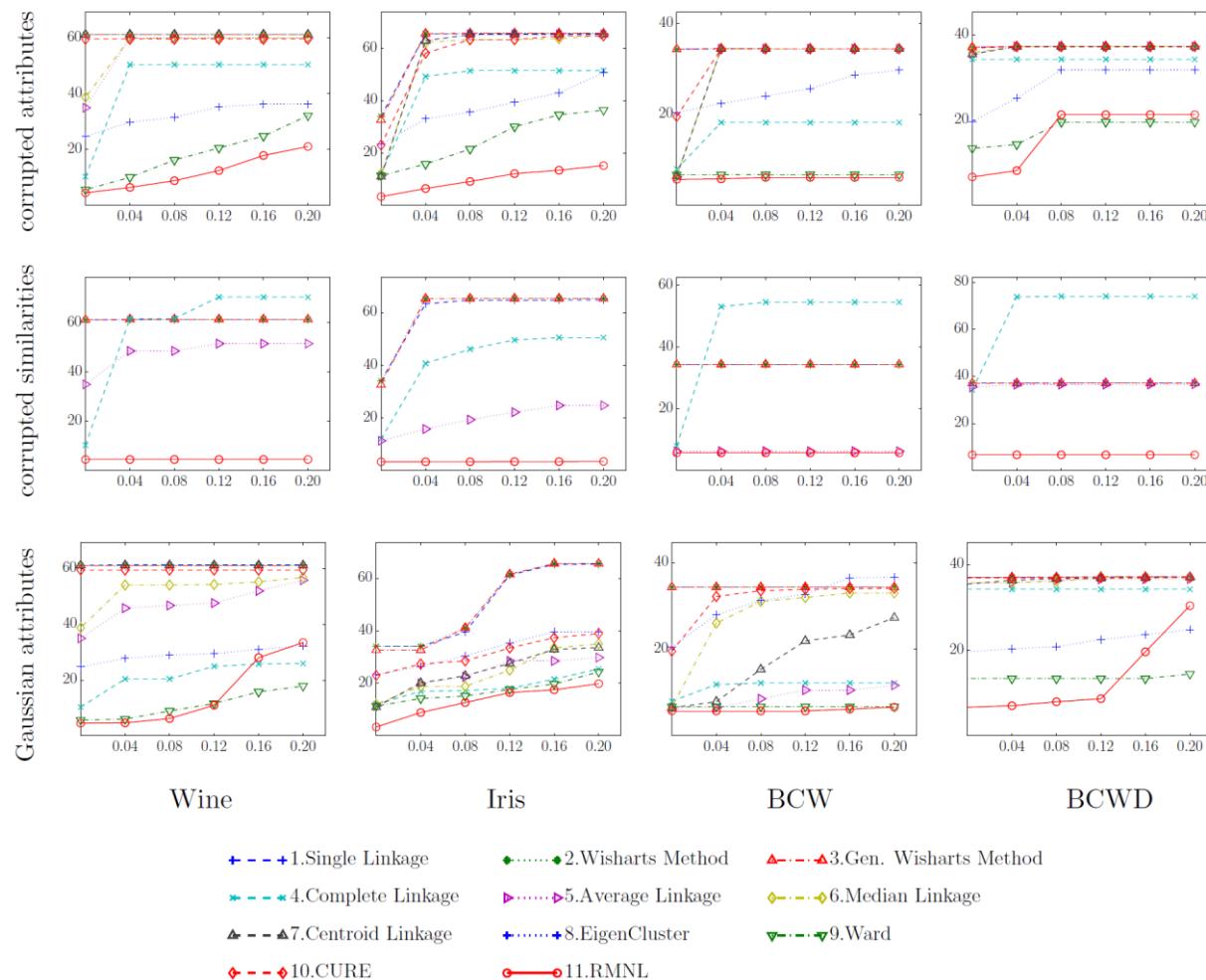
(d) BCWD



Ward's method does the best among classic techniques.

Hierarchical Clustering Experiments

[BLG, JMLR'15]



Ward's method does the best among classic techniques.

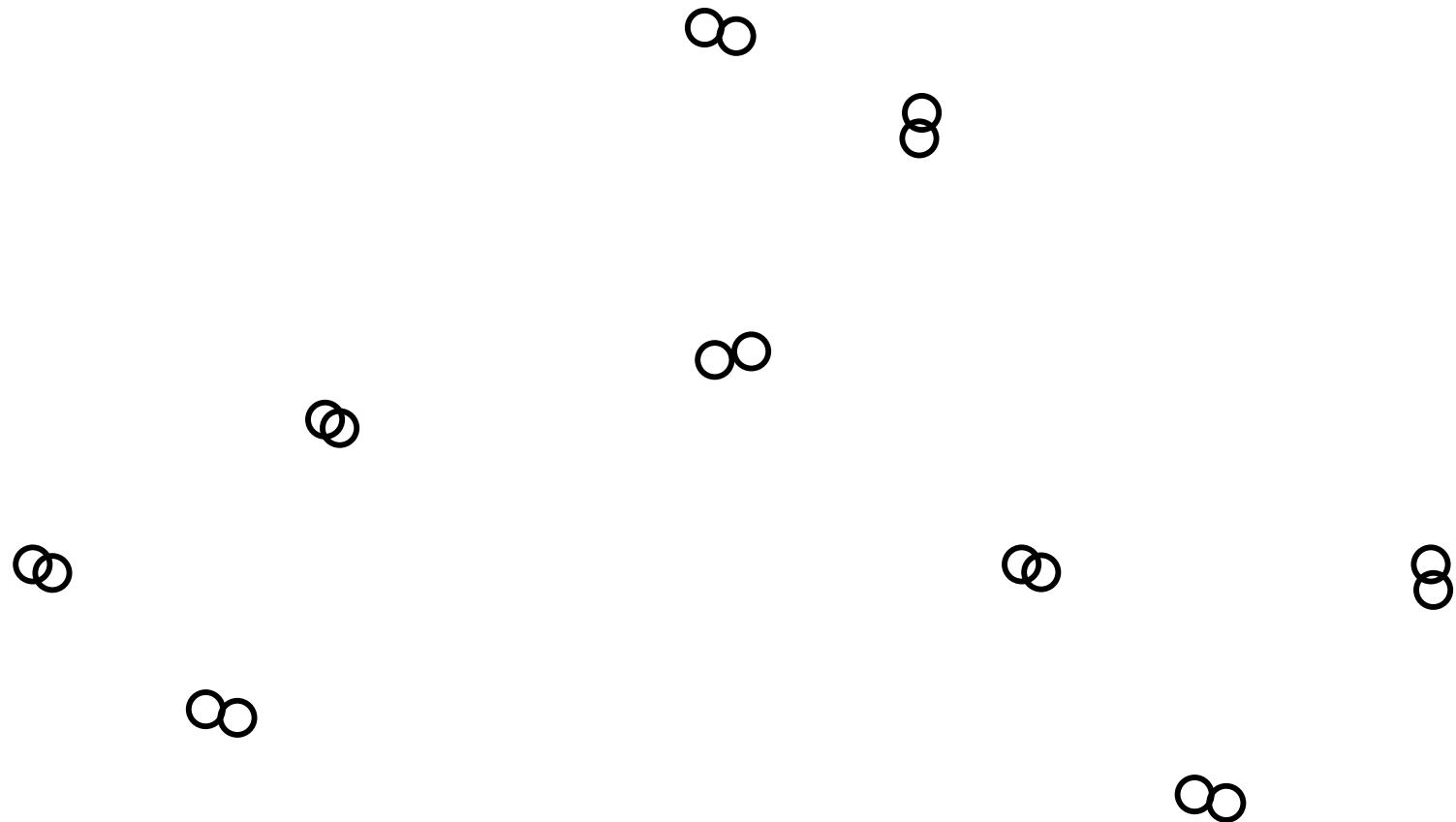
What You Should Know

- Partitional Clustering. k-means and k-means ++
 - Lloyd's method
 - Initialization techniques (random, furthest traversal, k-means++)
- Hierarchical Clustering.
 - Single linkage, Complete Linkage, Ward's method

Additional Slides

Smoothed analysis model

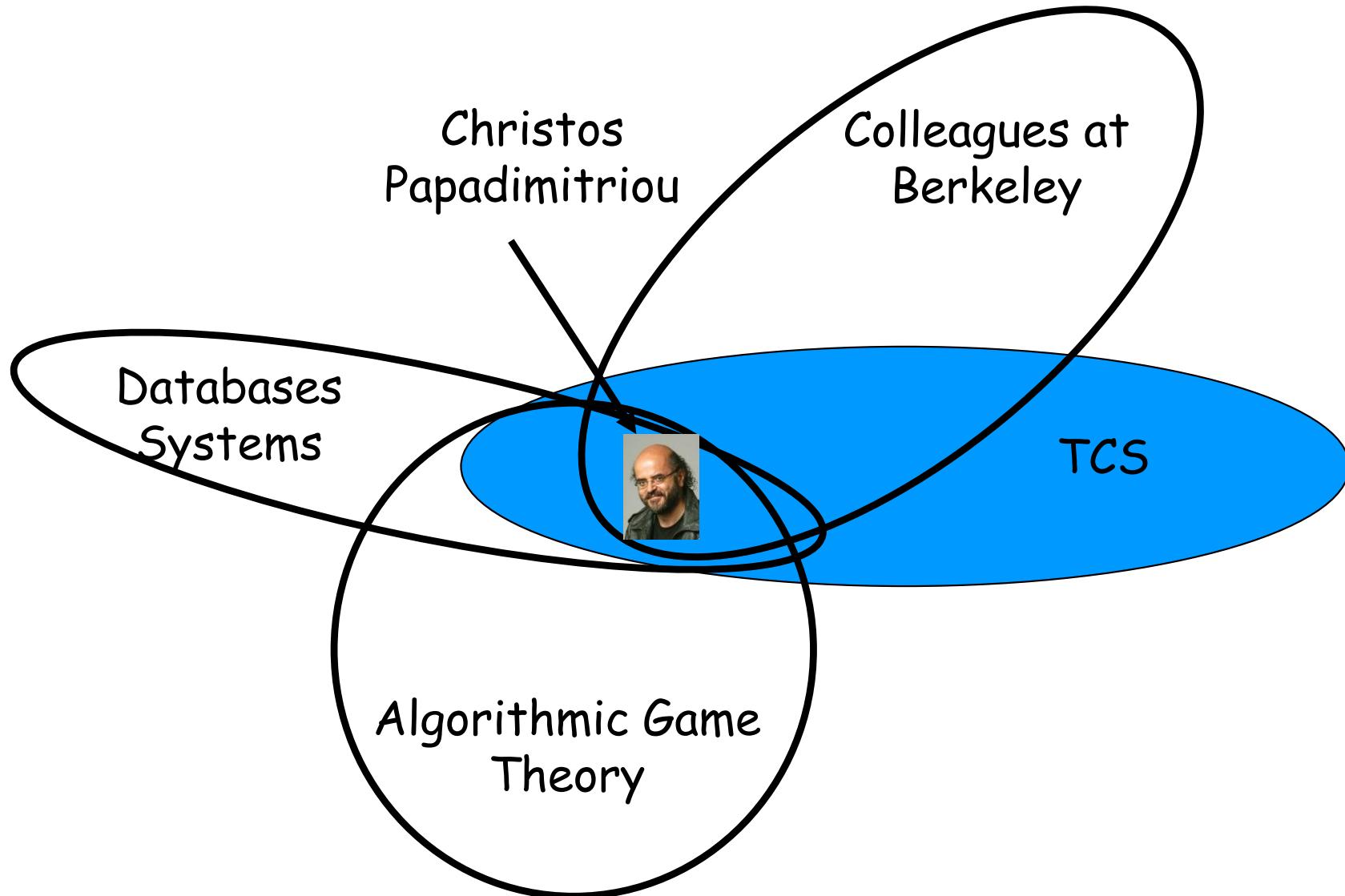
- Imagine a worst-case input.
- But then add small Gaussian perturbation to each data point.



Smoothed analysis model

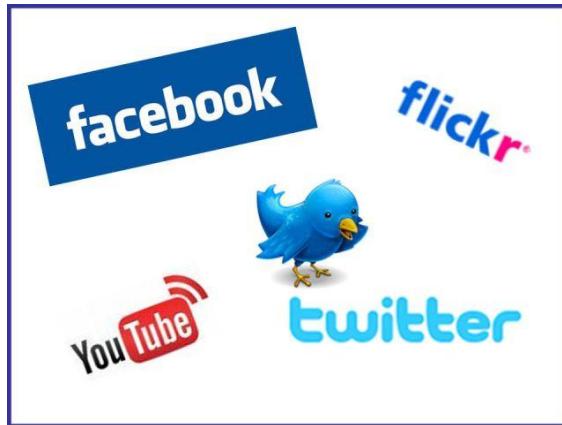
- Imagine a worst-case input.
- But then add small Gaussian perturbation to each data point.
- Theorem [Arthur-Manthey-Roglin 2009]:
 - $E[\text{number of rounds until Lloyd's converges}]$ if add Gaussian perturbation with variance σ^2 is polynomial in $n, 1/\sigma$.
 - The actual bound is : $O\left(\frac{n^{34}k^{34}d^8}{\sigma^6}\right)$
- Might still find local opt that is far from global opt.

Overlapping Clusters: Communities



Overlapping Clusters: Communities

- Social networks
- Professional networks



- Product Purchasing Networks, Citation Networks, Biological Networks, etc

Overlapping Clusters: Communities

