

CSE8803/CX4803 Machine Learning in Computational Biology

Lecture 9: Learning from high-dimensional data:
UMAP, clustering

Xiuwei Zhang

School of Computational Science and Engineering

Paper presentation logistics

Presentation dates are fixed for all groups

Now start paper bidding for Phase 1 presentations!
Deadline Thursday 11:59pm.

Bidding link:

https://gatech.co1.qualtrics.com/jfe/form/SV_0Vbl3SXv1pCTlk

List of papers:

https://docs.google.com/document/d/1RJDWddTV3hqnc6YGzSXGpKCIYgjL5k_GtLarOU54s/edit

Your name

Your group ID

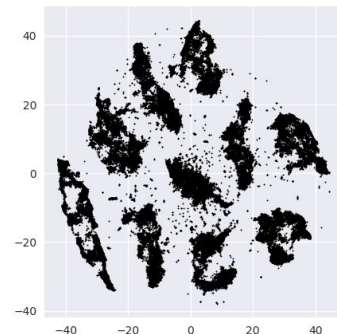
Your 1st choice paper

Number of points for your 1st choice paper

t-SNE

Why its output (the low-dimensional representations of data) is more often used for visualization than for further downstream analysis (clustering, etc)?

- It doesn't learn a function from high-dim data to low-dim, so if new data points come we can't directly convert it to low-dim.
- It doesn't directly preserve distance but rather preserves the neighborhood of every datapoint. So the distance between points in low-dim space can't be interpreted as close representation of the original distance. Distance-based clustering methods should be used with caution to the output of tSNE. In particular, larger distances are not preserved.
- It's output is often good for human eyes (also considering the effect of varying parameters), but not good for automatic clustering methods like k-means.
- Still controversial.



UMAP Uniform Manifold Approximation and Projection

MDS: preserving distance

tSNE: preserving neighborhood

UMAP: preserving graph topology

Compared to t-SNE, UMAP seems to be

- faster
- deterministic

Very similar intuition but
different mathematical
framework

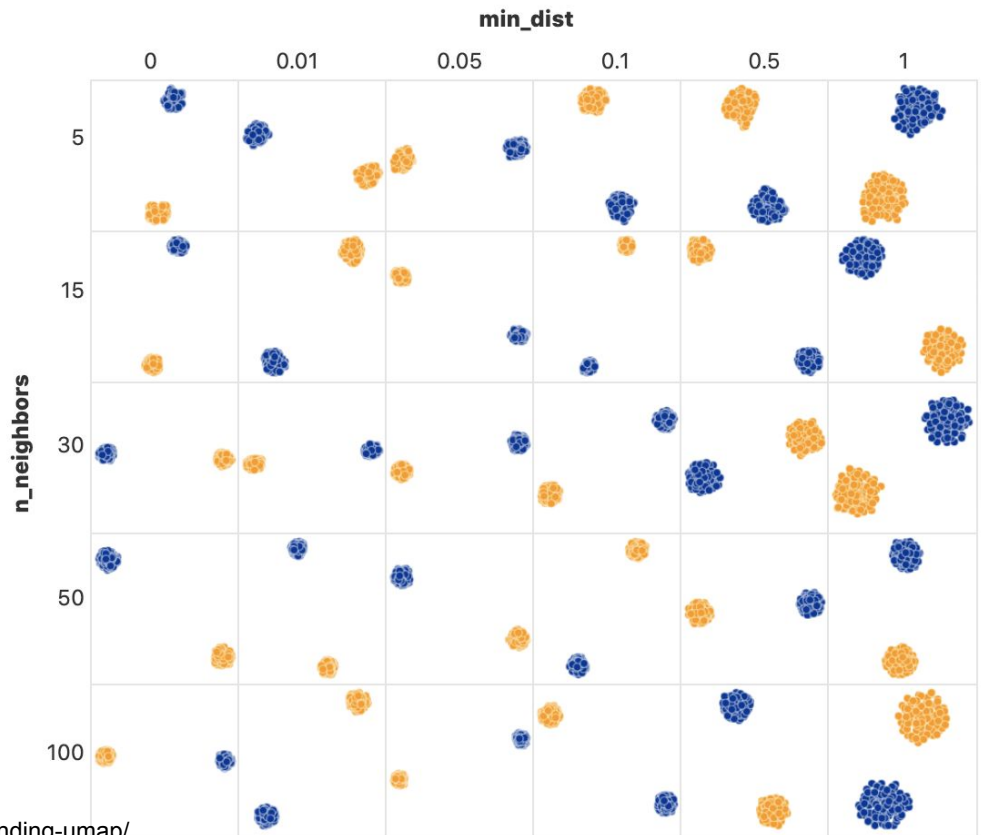
UMAP Uniform Manifold Approximation and Projection

Hyper-parameters

N_neighbors: the number of nearest neighbors to consider

Min_dist: minimum distance apart that points are allowed to be in the low dimensional representation

Two clusters with equal numbers of points, but different variances within the clusters.

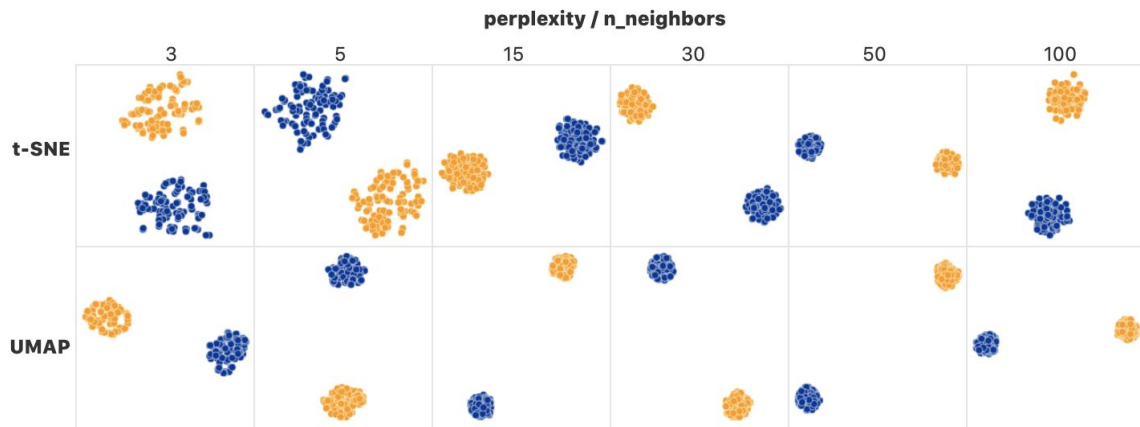


UMAP Uniform Manifold Approximation and Projection

Hyper-parameters

N_neighbors: the number of nearest neighbors to consider

Min_dist: minimum distance apart that points are allowed to be in the low dimensional representation



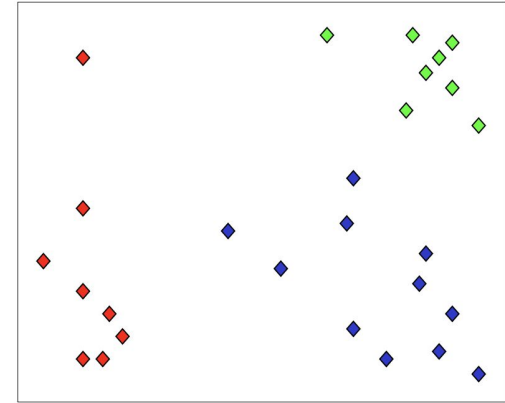
Two clusters with equal numbers of points.

<https://pair-code.github.io/understanding-umap/>

Clustering algorithms

Clustering - unsupervised learning methods

- Organizing data into clusters such that there is
 - high intra-cluster similarity
 - low inter-cluster similarity
- Informally, finding natural groupings among objects.
- Applications
 - Organizing data into clusters provides information about the internal structure of the data
 - Sometimes the partitioning is the goal
 - Knowledge discovery in data



K=2



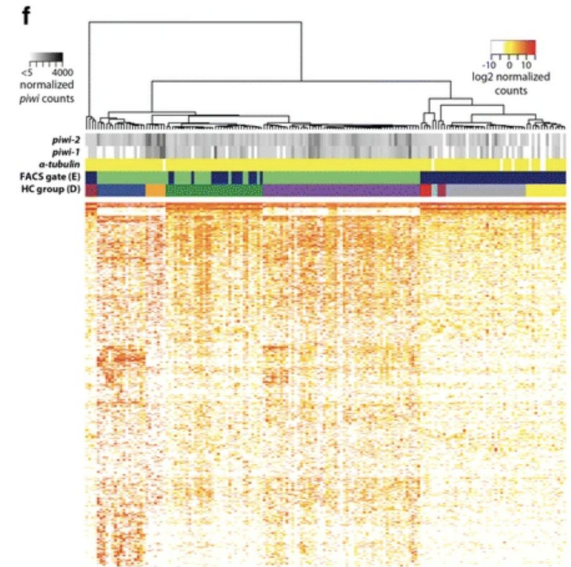
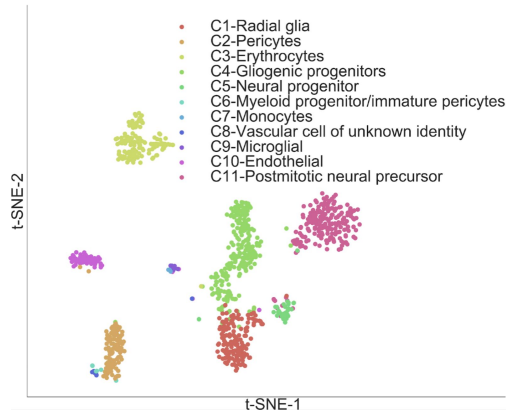
Original



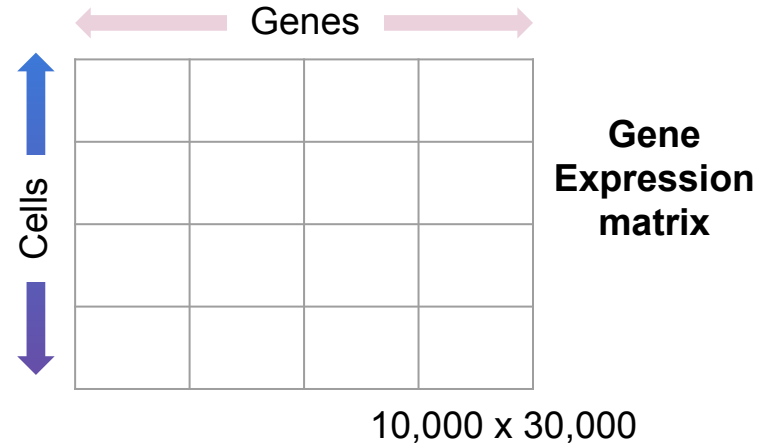
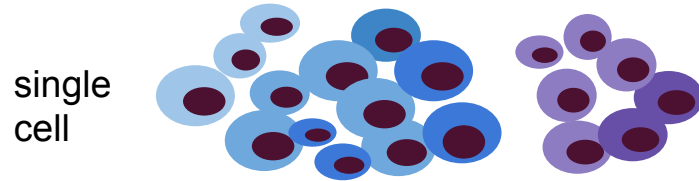
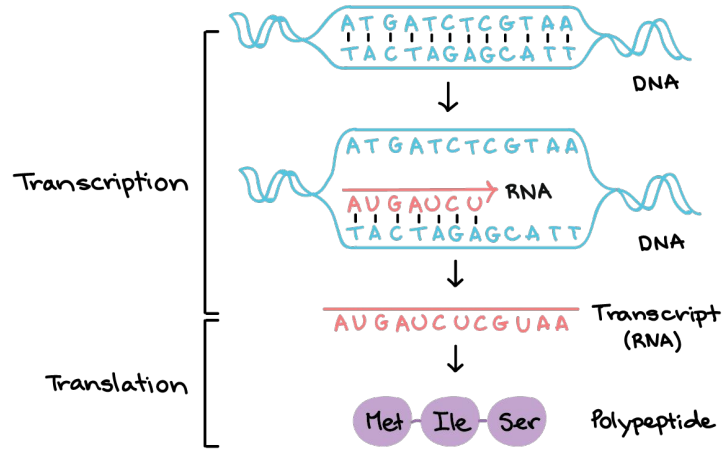
Clustering algorithms

Clustering - unsupervised learning methods

- K-means (and fuzzy k-means)
- Hierarchical clustering
- Spectral clustering
- Graph-based : Louvain, Leiden clustering



Some biological background: Gene expression analysis



k-means

Given the number of clusters k
An iterative clustering algorithm.

Initialize:

Pick k random points as cluster centers

Alternate:

1. Assign data points to closest cluster center
2. Change the cluster centers to the average of its assigned points

Stop:

When no points' assignments change

k-means

Given the number of clusters k
An iterative clustering algorithm.

Initialize:

Pick k random points as cluster centers

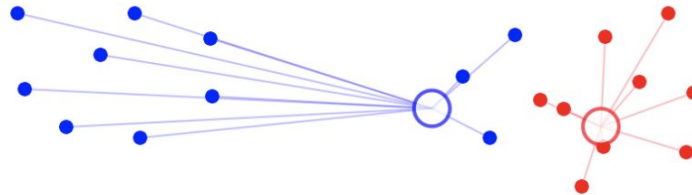
Alternate:

1. Assign data points to closest cluster center

2. Change the cluster centers to the average of its assigned points

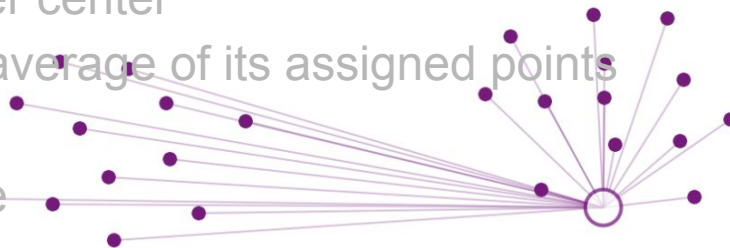
Stop:

When no points' assignments change



Initial locations of centers

Points are assigned to
closest centers



k-means

Given the number of clusters k
An iterative clustering algorithm.

Initialize:

Pick k random points as cluster centers

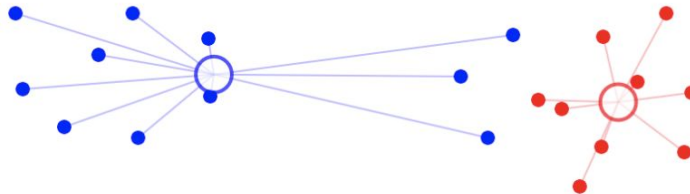
Alternate:

1. Assign data points to closest cluster center

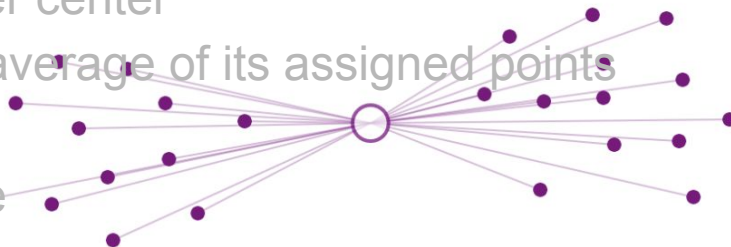
2. Change the cluster centers to the average of its assigned points

Stop:

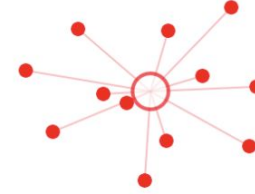
When no points' assignments change



Centers updated as the
average of all points of that
label



k-means



Labels for each point
updated

Given the number of clusters k
An iterative clustering algorithm.

Initialize:

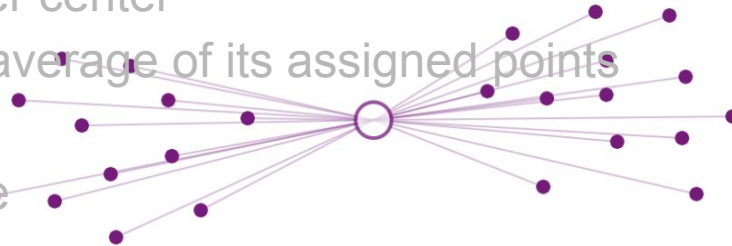
Pick k random points as cluster centers

Alternate:

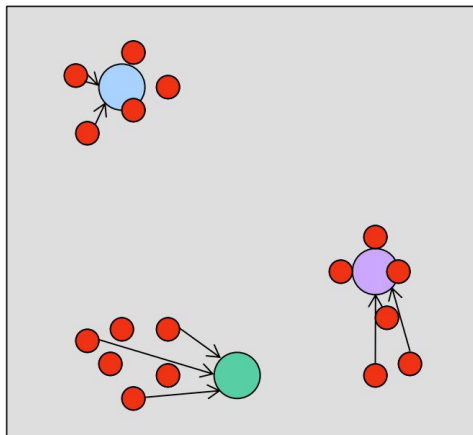
1. Assign data points to closest cluster center
2. Change the cluster centers to the average of its assigned points

Stop:

When no points' assignments change



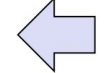
K-means as an “expectation maximization” process



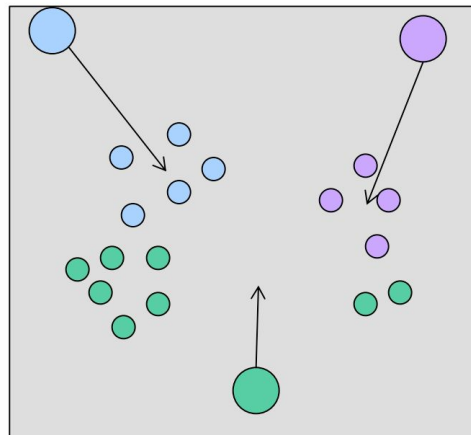
Re-assign each point \mathbf{x}_i
to **nearest center** k
→ Minimize distance from \mathbf{x}_i to μ_k :

$$d_{i,k} = (\mathbf{x}_i - \mu_k)^2$$

(“M”)



(“E”)



Update center μ_k to the
mean of the points assigned to it:

$$\mu_k = \sum_{\mathbf{x}_i \text{ with label } k} \frac{\mathbf{x}_i}{|\mathbf{x}^k|}$$

where: $|\mathbf{x}^k| = \# \mathbf{x}_i \text{ with label } k \rightarrow n_k$

K-means as an “expectation maximization” process

$$\text{COST}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_n) = \sum_{\mu_k} \sum_{\mathbf{x}_i \text{ with label } k} (\mathbf{x}_i - \mu_k)^2$$

The M step:

With the current the cluster assignment (labels) of points

$$\begin{aligned} \sum_{\mathbf{x}_i \text{ with label } k} (\mathbf{x}_i - \mu_k)^2 &= \sum_{\mathbf{x}_i \text{ with label } k} (\mathbf{x}_i^2 - 2\mathbf{x}_i \mu_k + \mu_k^2) \\ &= \sum_{\mathbf{x}_i \text{ with label } k} \mathbf{x}_i^2 - 2 \mu_k \sum_{\mathbf{x}_i \text{ with label } k} \mathbf{x}_i + n_k \mu_k^2 \end{aligned}$$

Take derivative on μ :

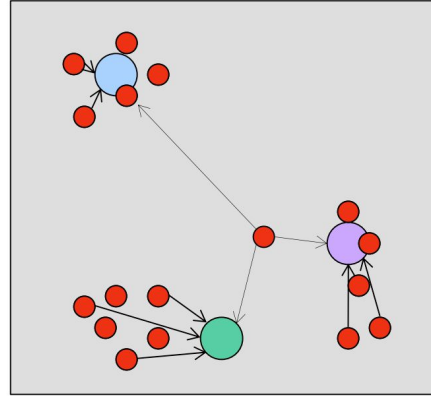
$$2 \sum_{\mathbf{x}_i \text{ with label } k} \mathbf{x}_i = 2 n_k \mu_k$$

$$\mu_k = 1/n_k \sum_{\mathbf{x}_i \text{ with label } k} \mathbf{x}_i$$

Fuzzy k-means

What if some points are half-way between two cluster centers?

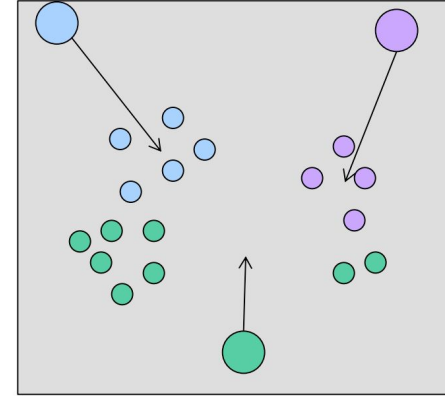
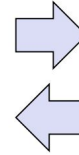
Assign partial weights.



Re-assign each point \mathbf{x}_i to all centers, weighted by distance

→ For each point calculate the probability of membership for each category K:

$$P(\text{label } K \mid \mathbf{x}_i, \boldsymbol{\mu}_k)$$



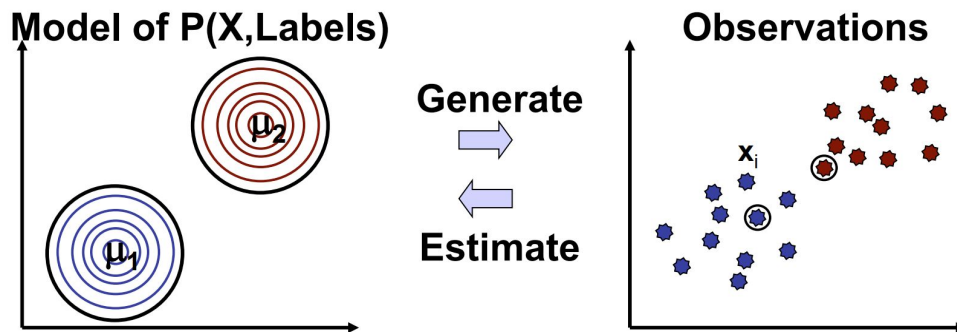
Update center $\boldsymbol{\mu}_k$ to the **weighted mean** of the points assigned to it:

$$\boldsymbol{\mu}_k = \frac{\sum_{\mathbf{x}_i \text{ with label } j} \mathbf{x}_i P(\boldsymbol{\mu}_k \mid \mathbf{x}_i)^b}{\sum_{\mathbf{x}_i \text{ with label } j} P(\boldsymbol{\mu}_k \mid \mathbf{x}_i)^b}$$

Regular K-Means is a special case of fuzzy k-means where:

$$P(\text{label } K \mid \mathbf{x}_i, \boldsymbol{\mu}_k) = \begin{cases} 1 & \text{if } \mathbf{x}_i \text{ is closest to } \boldsymbol{\mu}_k \\ 0 & \text{otherwise} \end{cases}$$

K-means as a generative model



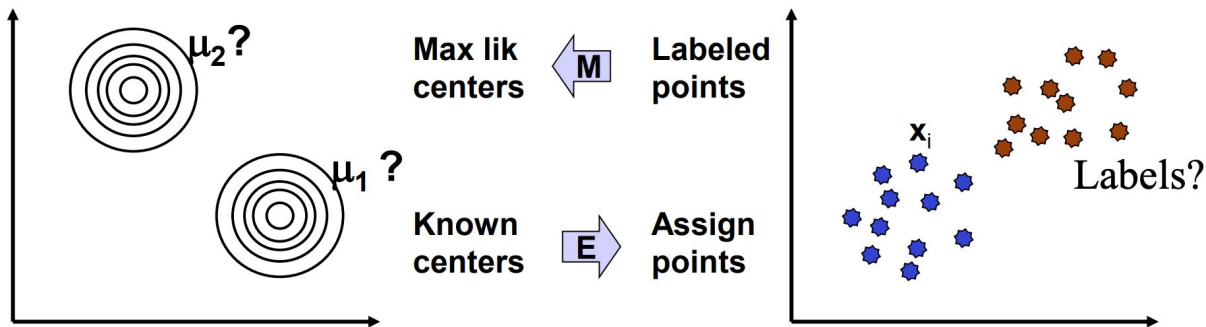
Consider that each point is generated from a Gaussian distribution with unit variance.

$$P(x_i | \mu_k) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{(x_i - \mu_k)^2}{2} \right\}$$

Problem: find cluster centers and labels for data points such that the total likelihood is maximized

Solution is the same as the k-means solution! (equivalence)

K-means as a generative model

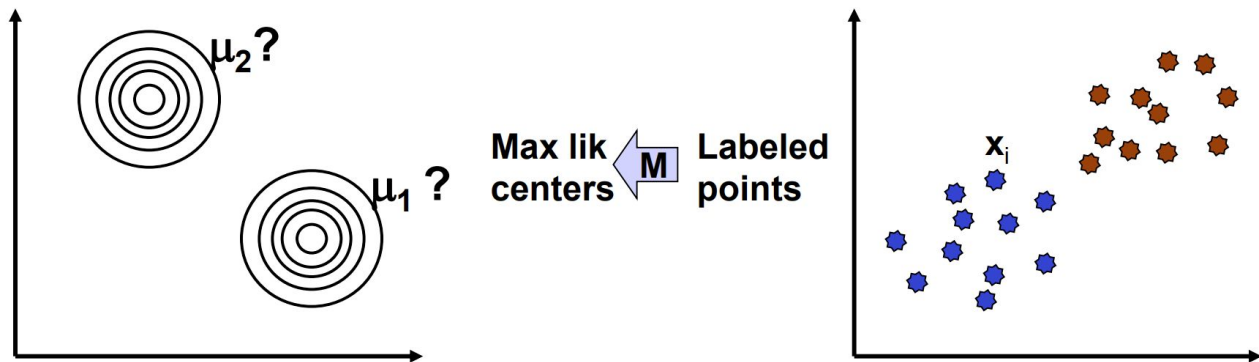


Problem: find cluster centers and labels for data points such that the total likelihood is maximized

E step: If centers are known \rightarrow Estimate memberships

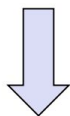
M step: If assignments known \rightarrow Compute centroids

K-means as a generative model



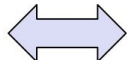
$$\boxed{\arg \max_{\mu}} \left\{ \log \prod_i P(\mathbf{x}_i | \mu) \right\} = \arg \max_{\mu} \sum_i \left\{ -\frac{1}{2}(\mathbf{x}_i - \mu)^2 + \log \left(\frac{1}{\sqrt{2\pi}} \right) \right\}$$

Seeking the **max likelihood** estimate of the cluster mean



EM solution

Equivalent



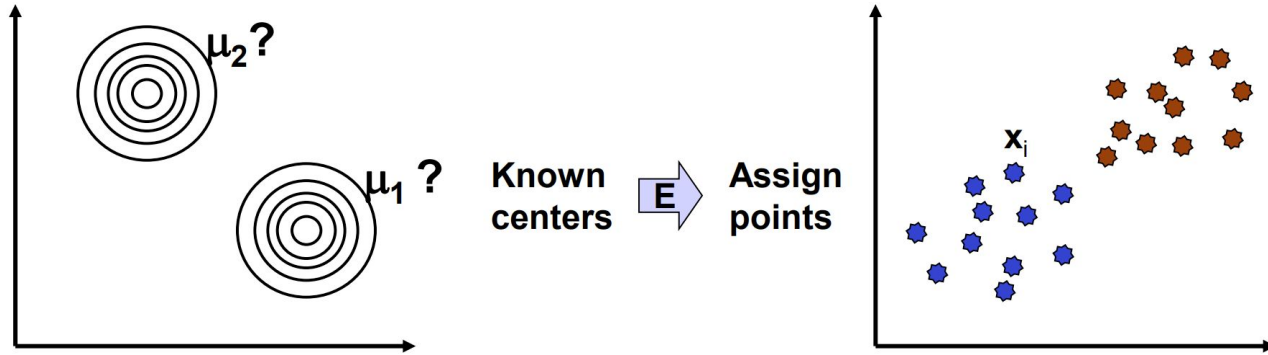
$$= \arg \min_{\mu} \sum_i \boxed{(\mathbf{x}_i - \mu)^2}$$

Solution is the **centroid** of the \mathbf{x}_i



K-means solution

K-means as a generative model



$$\boxed{\arg \max_k} P_k(\mathbf{x}_i | \boldsymbol{\mu}_i) = \arg \max_k \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{(\mathbf{x}_i - \mathbf{u}_k)^2}{2} \right\} = \text{arg min}_k \boxed{(\mathbf{x}_i - \mathbf{u}_k)^2}$$

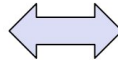
Seeking the label k that maximizes likelihood of point

Solution is the nearest center

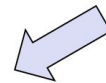


EM solution

Equivalent



K-means solution



Comparison between methods

	K-means		Fuzzy K-means	
	<u>algorithmic</u> formulation	<u>probabilistic</u> interpretation	<u>algorithmic</u> formulation	<u>probabilistic</u> interpretation
Initialization	Initialize K centers μ_k	Initialize model parameters	Initialize K centers μ_k	Initialize model parameters
E-step: Estimate prob of hidden labels (point assignments to classes)	Assign x_i label of <u>nearest center</u> distance $d_{i,k} = (x_i - \mu_k)^2$	Estimate <u>most likely missing label</u> given previous parameters	Calculate <u>probability of membership</u> for each point to each class $P(\text{label } K x_i, \mu_k)$	Estimate <u>probability over missing labels</u> given previous parameters
M-step: Update params to max likelihood estimates given assignments	Move μ_k to <u>centroid</u> of all points with that label	Choose new <u>max likelihood</u> params given points in label	Move μ_k to <u>weighted centroid</u> of all points, each weighted by $P(\text{label})$	Choose new params to maximize <u>expected likelihood</u> given <u>label estimates</u>
Iteration	Iterate	Iterate	Iterate	Iterate

Characteristics of k-means

K-means partitions the points into Voronoi cells

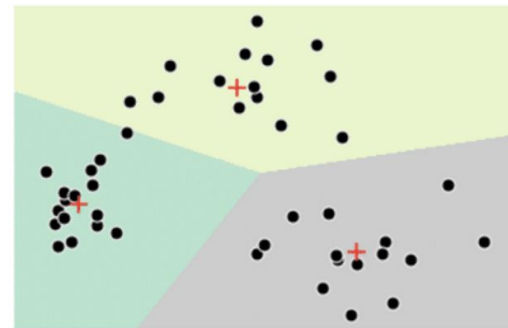
Does it converge?

Yes.

Does it guarantee optimality?

No. (Initialization matters)

Prior dimension reduction helpful



Disadvantages:

- Needs k as input
- Unable to handle noisy data and outliers
- Not suitable to discover clusters with non-convex shapes

Hierarchical clustering

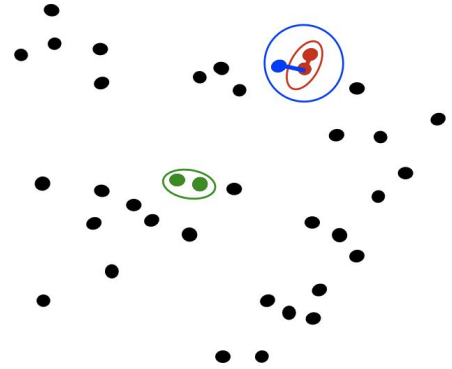
Idea: agglomerative clustering

- First merge very similar instances
- Incrementally build larger clusters out of smaller clusters

Algorithm: Maintain a set of clusters

- Initially, each instance in its own cluster
- Repeat:
 - Pick the two closest clusters
 - Merge them into a new cluster
 - Stop when there's only one cluster left

How to
define?



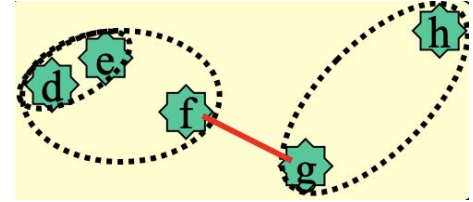
Output:

Produces not one clustering, but a family of clusterings represented by a dendrogram

Hierarchical clustering: define “closest” for clusters

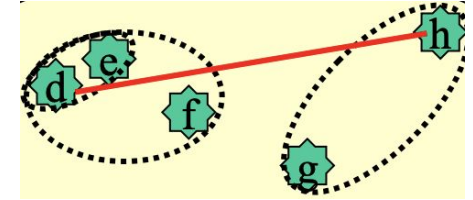
Closest pair
(single-link)

$$CD(X,Y)=\min_{x \in X, y \in Y} D(x,y)$$



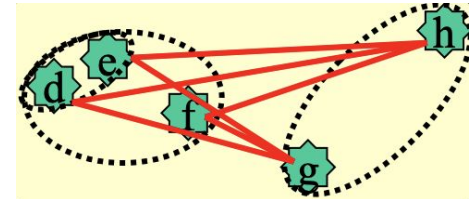
Furthest pair
(complete-link)

$$CD(X,Y)=\max_{x \in X, y \in Y} D(x,y)$$



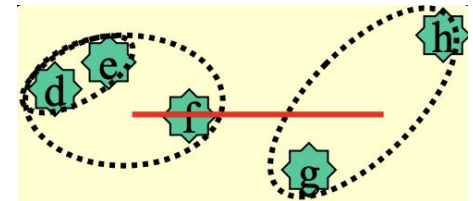
average-link

$$CD(X,Y)=\text{avg}_{x \in X, y \in Y} D(x,y)$$



centroids method

$$CD(X,Y)=D(\text{avg}(X), \text{avg}(Y))$$

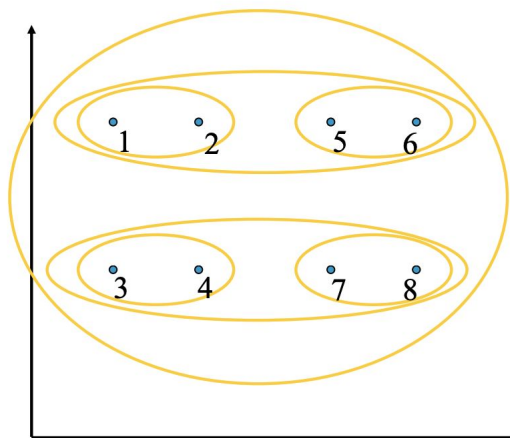


Hierarchical clustering: define “closest” for clusters

Cluster distance affects both results and runtime

Closest pair

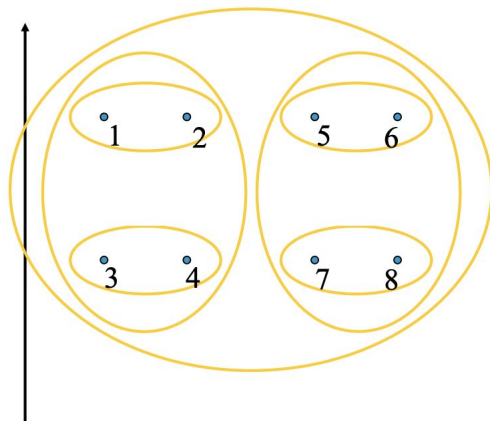
(single-link clustering)



Potentially long and skinny clusters

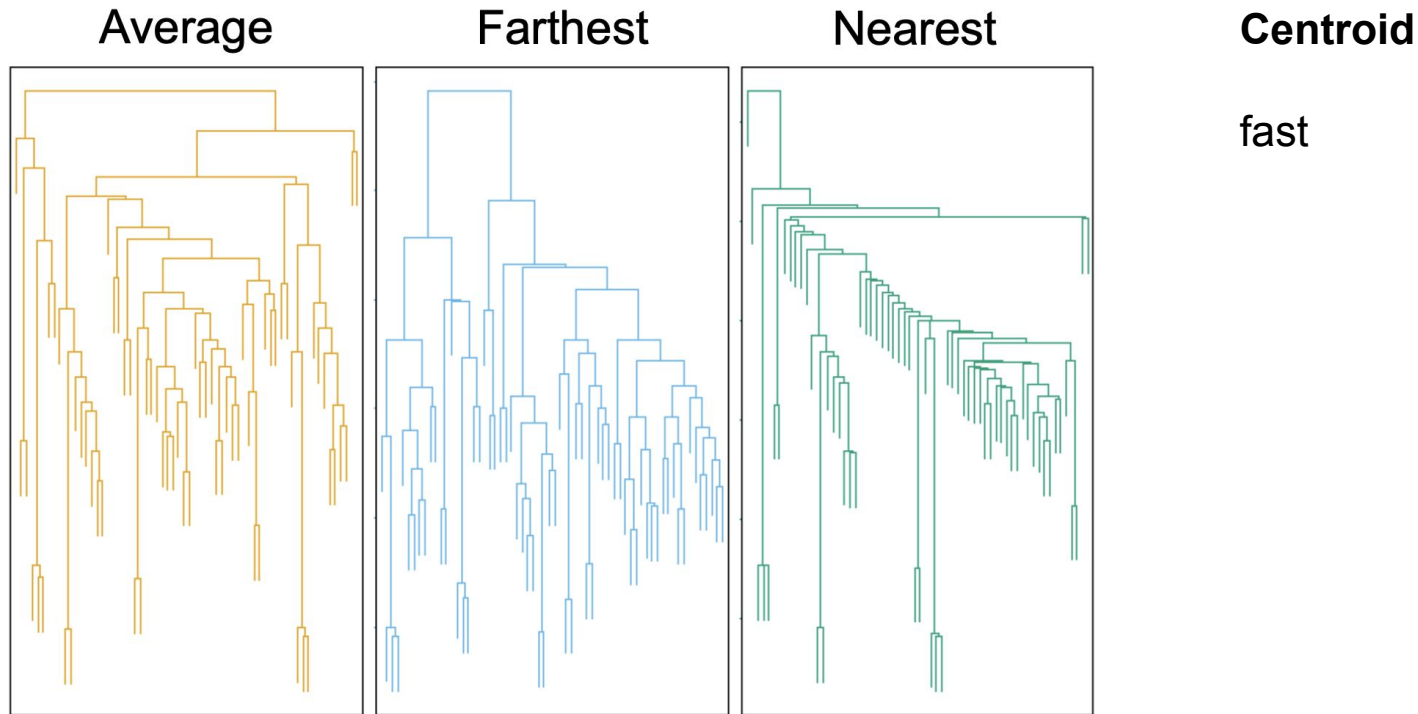
Farthest pair

(complete-link clustering)



tight clusters

Hierarchical clustering: define “closest” for clusters



fast

Robust to noise
Widely used

Figure credit: David Sontag

Hierarchical clustering: distance measures

Table 1 Gene expression similarity measures

Manhattan distance
(city-block distance, L1 norm)

$$d_{fg} = \sum_c |e_{fc} - e_{gc}|$$

Euclidean distance
(L2 norm)

$$d_{fg} = \sqrt{\sum_c (e_{fc} - e_{gc})^2}$$

Mahalanobis distance

$$d_{fg} = (\mathbf{e}_f - \mathbf{e}_g)' \Sigma^{-1} (\mathbf{e}_f - \mathbf{e}_g), \text{ where } \Sigma \text{ is the (full or within-cluster) covariance matrix of the data}$$

Pearson correlation
(centered correlation)

$$d_{fg} = 1 - r_{fg}, \text{ with } r_{fg} = \frac{\sum_c (e_{fc} - \bar{e}_f)(e_{gc} - \bar{e}_g)}{\sqrt{\sum_c (e_{fc} - \bar{e}_f)^2 \sum_c (e_{gc} - \bar{e}_g)^2}}$$

Uncentered correlation
(angular separation, cosine angle)

$$d_{fg} = 1 - r_{fg}, \text{ with } r_{fg} = \frac{\sum_c e_{fc} e_{gc}}{\sqrt{\sum_c e_{fc}^2 \sum_c e_{gc}^2}}$$

Spellman rank correlation

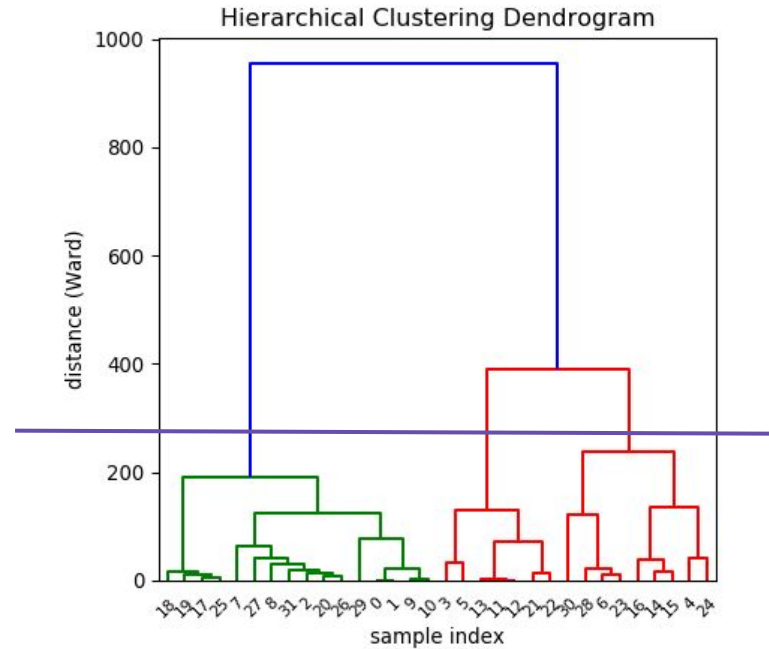
As Pearson correlation, but replace \mathbf{e}_{gc} with the rank of \mathbf{e}_{gc} within the expression values of gene g across all conditions $\mathbf{c} = 1 \dots C$

Absolute or squared correlation

$$d_{fg} = 1 - |r_{fg}| \text{ or } d_{fg} = 1 - r_{fg}^2$$

d_{fg} , distance between expression patterns for genes f and g . e_{gc} , expression level of gene g under condition c .

Hierarchical clustering: obtaining the clusters



Hierarchical clustering: Summary

- No need to specify the number of clusters in advance.
- Hierarchical structure maps nicely onto human intuition for some domains
- They do not scale well: where n is the number of total objects.
- Like any heuristic search algorithms, local optima are a problem.
- Interpretation of results can be subjective