

CSE 564  
VISUALIZATION & VISUAL ANALYTICS  
CLUSTER ANALYSIS & DIMENSION  
REDUCTION

**KLAUS MUELLER**

COMPUTER SCIENCE DEPARTMENT  
STONY BROOK UNIVERSITY AND SUNY KOREA

Lecture	Topic	Projects
1	Intro and logistics	
2	Basic visualizations and tasks, data types, examples, ethical considerations	
3	Data preparation (cleaning, imputation, data set integration)	
4	AI-assisted coding for VIS applications (design, debugging, refactoring)	Project #1 out
5	Big data and data reduction (distance/sim metrics, intro to clustering)	
6	High-D data: concept, subspaces, dimension reduction, PCA	
7	Cluster analysis: hierarchical, density, model, embedding, temporal	
8	Perception and cognition (human visual system, color, contrast, bias)	Project #2(a) out
9	Visual design and aesthetics	
10	Visualization of multivariate and high-dimensional data: direct methods	
11	Visualization of multivariate and high-D data: projections & embeddings	
12	Visualization and AI: mutual support and capabilities (VIS4AI, AI4VIS)	Project #2(b) out
13	Principles of interaction: drive what is visualized, analyzed & how (HCI4VIS)	
14	Visual analytics (VA), human-centered AI, mixed-initiative system	
15	Midterm #1 (tentative date)	
16	VA system design and evaluation, collaborative VA, uncertainty, provenance	
17	Midterm #1 discussion (tentative date)	Final proj. proposal call out
18	Visualization of hierarchical data	
19	Visualization of maps and data with geo-reference	
20	Visualization of graphs, networks (incl. derivation of causal networks)	Final project proposal due
21	Vis. of time-varying, time-series, streaming data, progressive visualization	
22	Visualization of text, LLMs, and semantic data	
23	Ed Tufte revisited: principles, critiques and limits, responsible visualization	
24	Design of effective infographics	Final proj. prelim report due
25	Foundations scientific and medical visualization, intro to volume rendering	
26	Scientific visualization	Bonus project out (Vol Ren)
27	Story telling with data, data journalism	
28	Midterm #2 (tentative date)	
Final	Final project demo on zoom (public)	All final proj. materials due

# CLUSTERING = MODELING SIMILARITY STRUCTURE

## Clustering Families

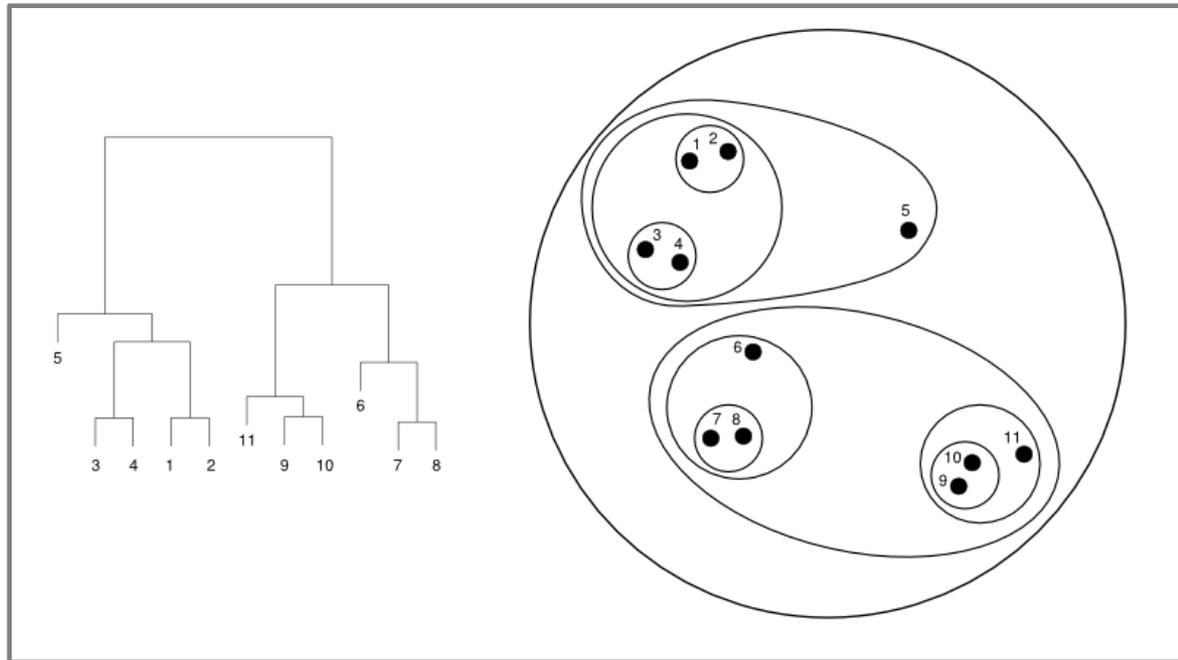
- Partition-based → K-means (see previous lecture)
- Hierarchical → Agglomerative
- Density-based → DBSCAN
- Model-based → Gaussian Mixture (EM)
- Embedding-based → Deep clustering
- Time-series → Tensors

## Often a precursor to other analyses

- preprocessing step for classification and outlier detection
- sampling and data reduction

# HIERARCHICAL (AGGLOMERATIVE) CLUSTERING

# HIERARCHICAL CLUSTERING

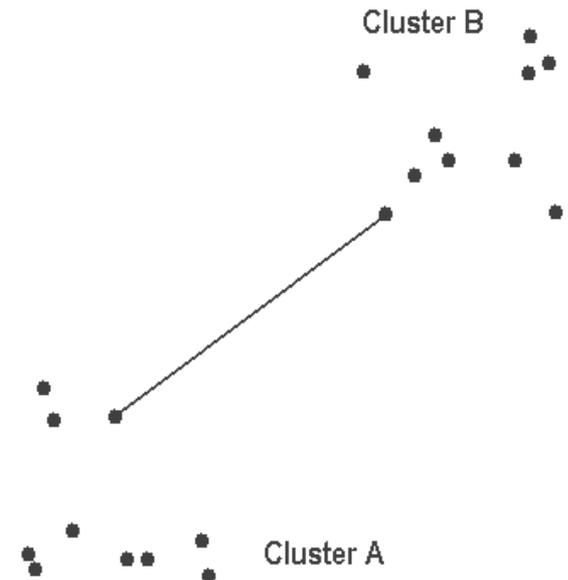
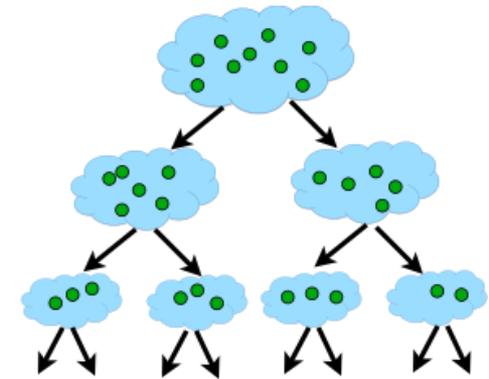


Two options for building the dendrogram on the left

- top down (divisive)
- bottom up (agglomerative)

# BOTTOM-UP AGGLOMERATIVE METHODS

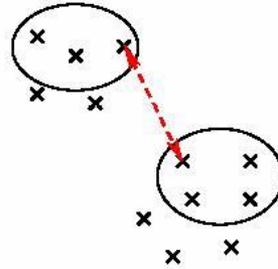
**Algorithm** *AgglomerativeMerge*(Data:  $\mathcal{D}$ )  
**begin**  
  Initialize  $n \times n$  distance matrix  $M$  using  $\mathcal{D}$ ;  
  **repeat**  
    Pick closest pair of clusters  $i$  and  $j$  using  $M$ ;  
    Merge clusters  $i$  and  $j$ ;  
    Delete rows/columns  $i$  and  $j$  from  $M$  and create  
      a new row and column for newly merged cluster;  
    Update the entries of new row and column of  $M$ ;  
  **until** termination criterion;  
  return current merged cluster set;  
**end**



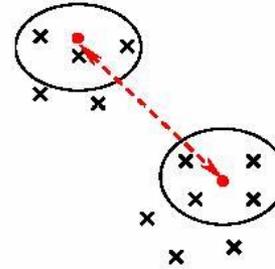
How to merge?

# MERGE CRITERIA

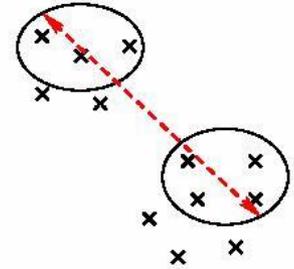
- Simple linkage



- Average linkage



- Complete linkage



## Single (best-case) linkage

- distance = minimum distance between all  $m_i \cdot m_j$  pairs of objects
- joins the closest pair

## Complete (worst-case) linkage

- distance = maximum distance between all  $m_i \cdot m_j$  pairs of objects
- joins the pair furthest apart

## Group-average linkage

- distance = average distance between all object pairs in the groups

## Other methods:

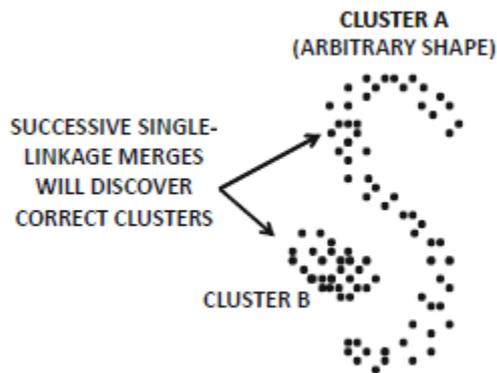
- closest centroid, variance-minimization, Ward's method

# COMPARISON

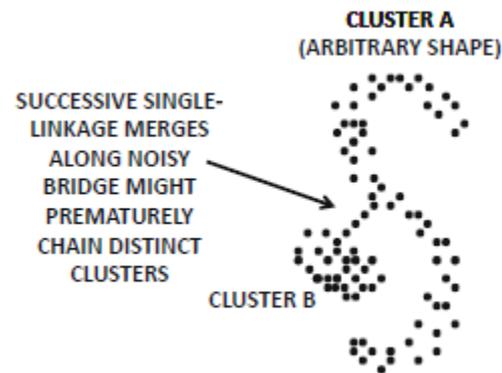
Centroid-based methods tend to merge large clusters

Single linkage method can merge chains of closely related points to discover clusters of arbitrary shape

- but can also (inappropriately) merge two unrelated clusters, when the chaining is caused by noisy points between two clusters



(a) Good case with no noise

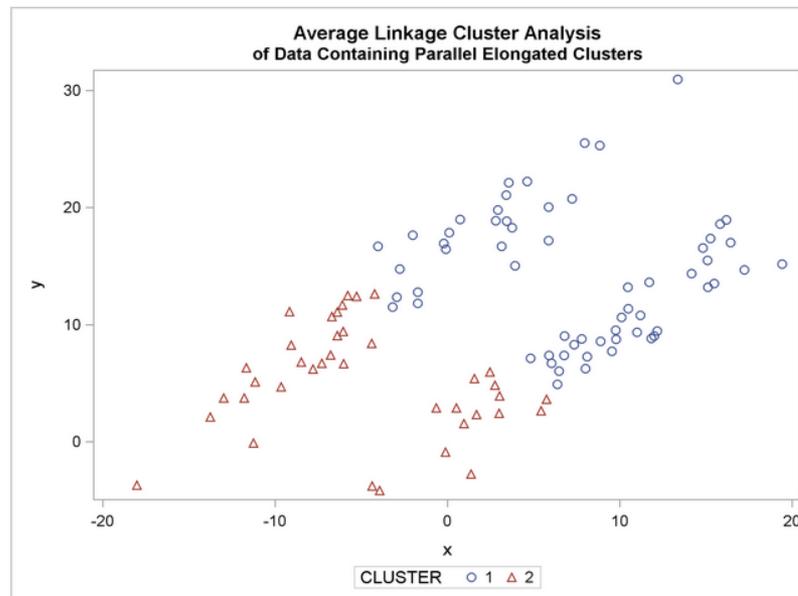


(b) Bad case with noise

# COMPARISON

Complete (worst-case) linkage method tends to create spherical clusters with similar diameter

- will break up the larger odd-shaped clusters into smaller spheres
- also gives too much importance to data points at the noisy fringes of a cluster

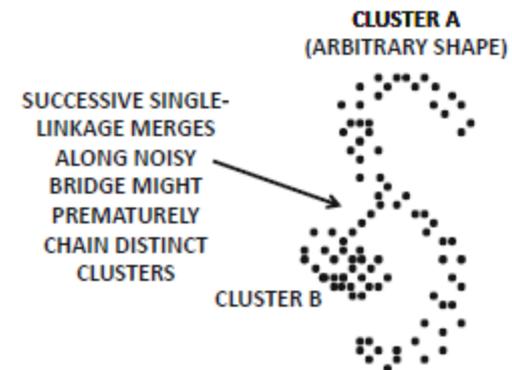


# COMPARISON

The group average, variance, and Ward's methods are more robust to noise due to the use of multiple linkages in the distance computation

Hierarchical methods are sensitive to a small number of mistakes made during the merging process

- can be due to noise
- no way to undo these mistakes



(b) Bad case with noise

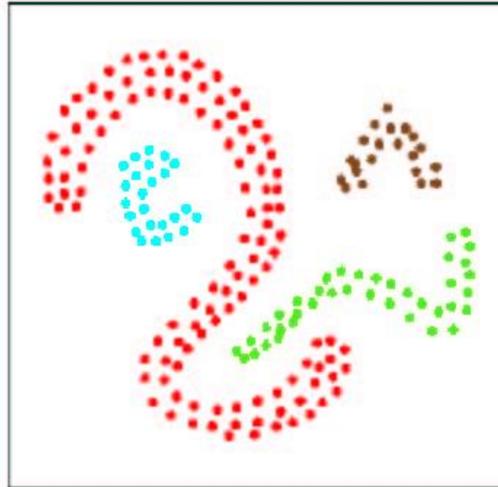
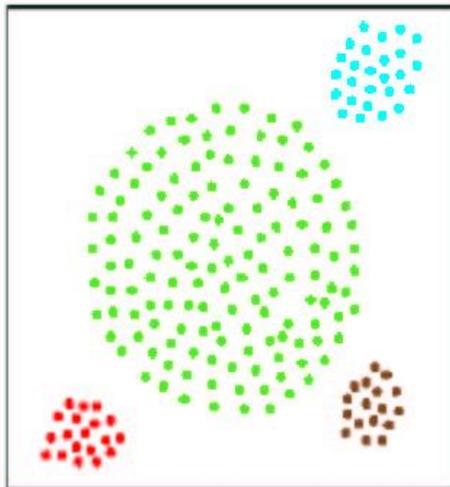
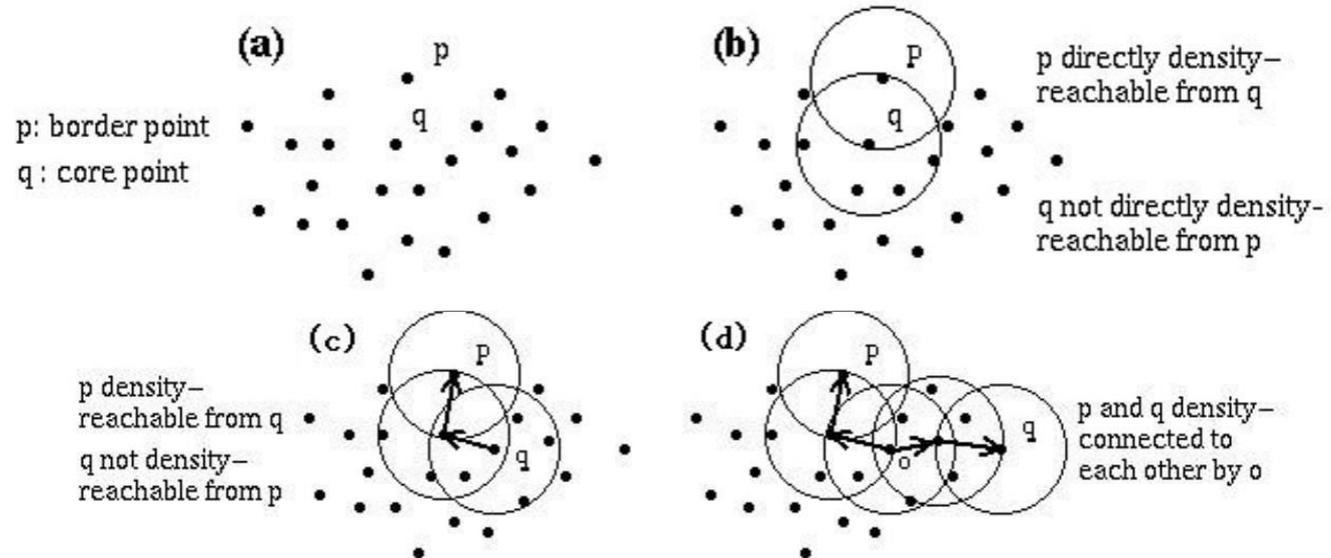
# DENSITY-BASED CLUSTERING

# DBSCAN

Highly-cited density-based hierarchical clustering algorithm (Ester et al. 1996)

- clusters are defined as density-connected sets
- epsilon-distance neighbor criterion (Eps)  
$$N_{Eps}(p) = \{q \in D \mid \text{dist}(p,q) \leq Eps\}$$
- minimum point cluster membership and core point (MinPts)  
$$|N_{Eps}(q)| \geq \text{MinPts}$$
- notions of density-connected & density-reachable (direct, indirect)
- a point  $p$  is directly density-reachable from a point  $q$  wrt. Eps, MinPts if  
$$p \in N_{Eps}(q) \text{ and}$$
  
$$|N_{Eps}(q)| \geq \text{MinPts} \text{ (core point condition)}$$

# DBSCAN

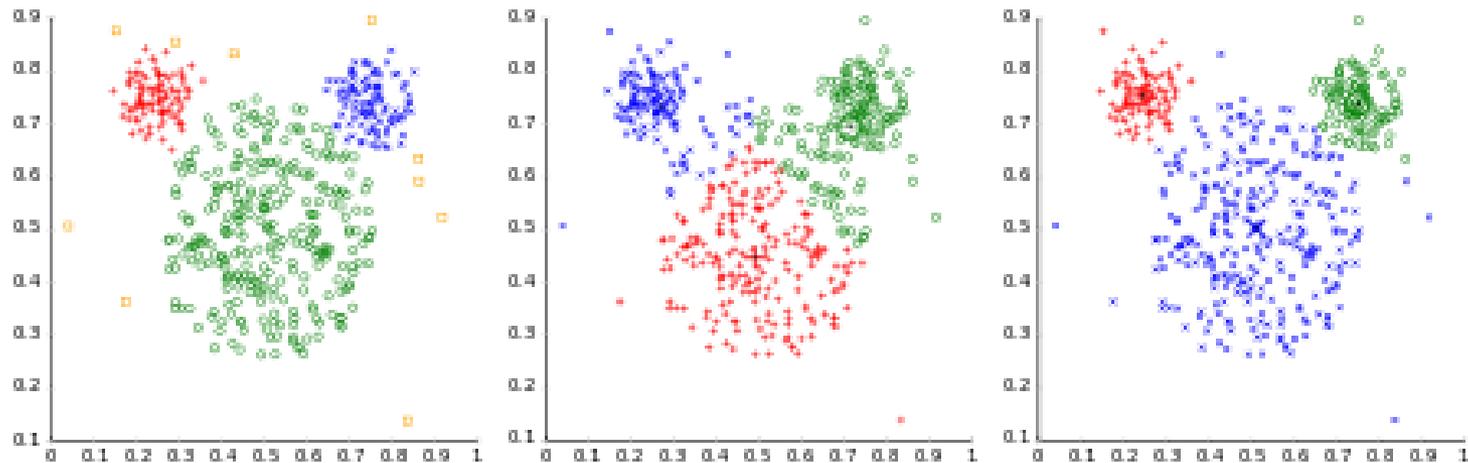


# MODEL-BASED CLUSTERING

# PROBABILISTIC EXTENSION TO K-MEANS

First a comparison:

Different cluster analysis results on "mouse" data set:  
Original Data      k-Means Clustering      EM Clustering



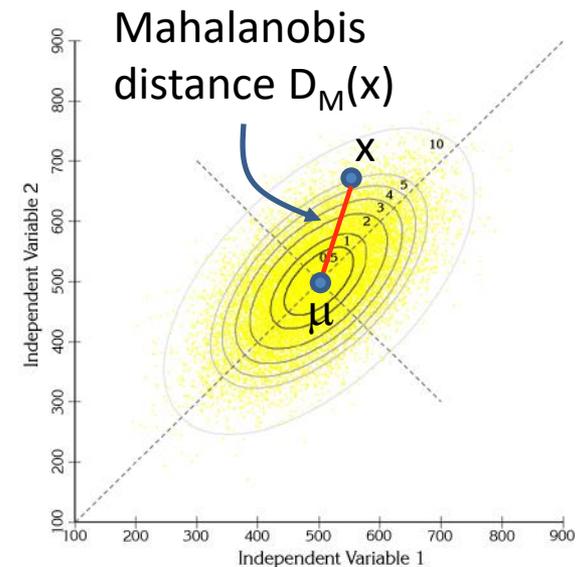
# MAHALANOBIS DISTANCE

The distance between a point  $X$  and a distribution  $D$

- measures how many standard deviations  $X$  is away from the mean  $\mu$  of  $D$
- $S$  is the covariance matrix of the distribution  $D$
- the Mahalanobis distance  $D_M$  of a point  $x$  to a cluster center  $\mu$  is

$$D_M(x) = \sqrt{(x - \mu)^T S^{-1} (x - \mu)}.$$

- $x$  and  $\mu$  are  $N$ -dimensional vectors
- $S$  is the  $N \times N$  covariance matrix
- the outcome  $D_M(x)$  is a single-dimensional number



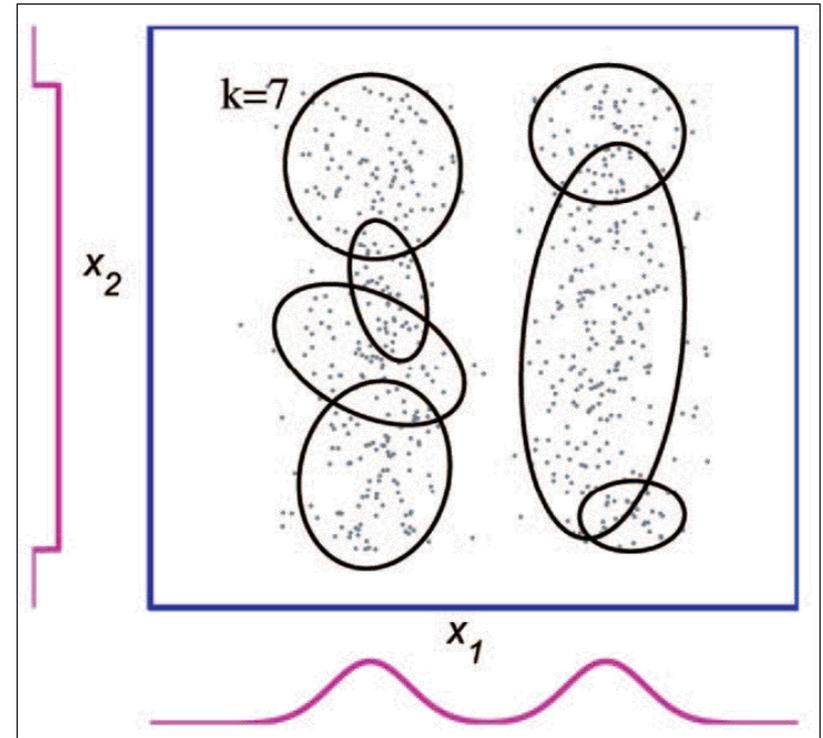
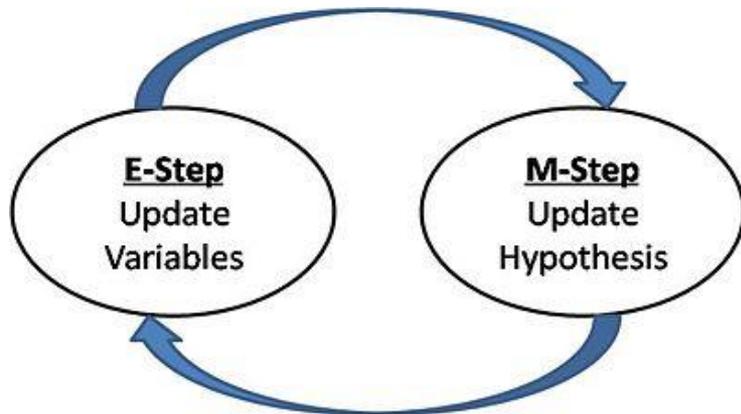
# PROBABILISTIC CLUSTERING

Is a better match for point distributions

- overlapping clusters are now possible
- better match with real world?
- Gaussian mixtures

Need a probabilistic algorithm

- Expectation-Maximization



# EM Algorithm (Mixture Model)

probability that data point  $d_i$  is in class  $c_j$   
(= Mahalanobis distance of  $d_i$  to  $c_j$ )

- Initialize K cluster centers
- Iterate between two steps
  - **E**xpectation step: assign points to  $m$  clusters/classes

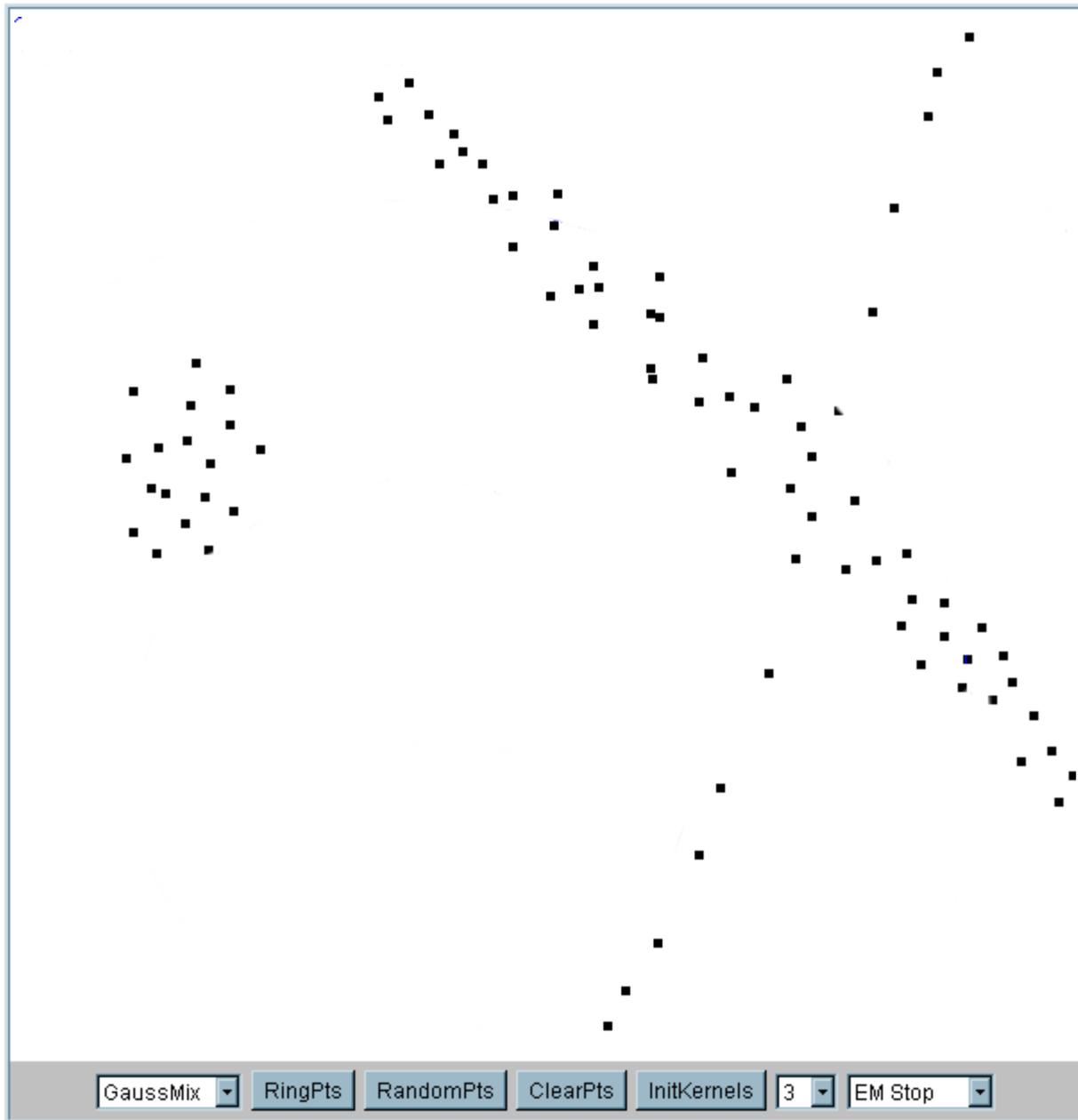
$$P(d_i \in c_k) = \frac{w_k \Pr(d_i | c_k)}{\sum_j w_j \Pr(d_i | c_j)}$$

$$w_k = \frac{\sum_i \Pr(d_i \in c_k)}{N} = \text{probability of class } c_k$$

- **M**aximation step: estimate model parameters

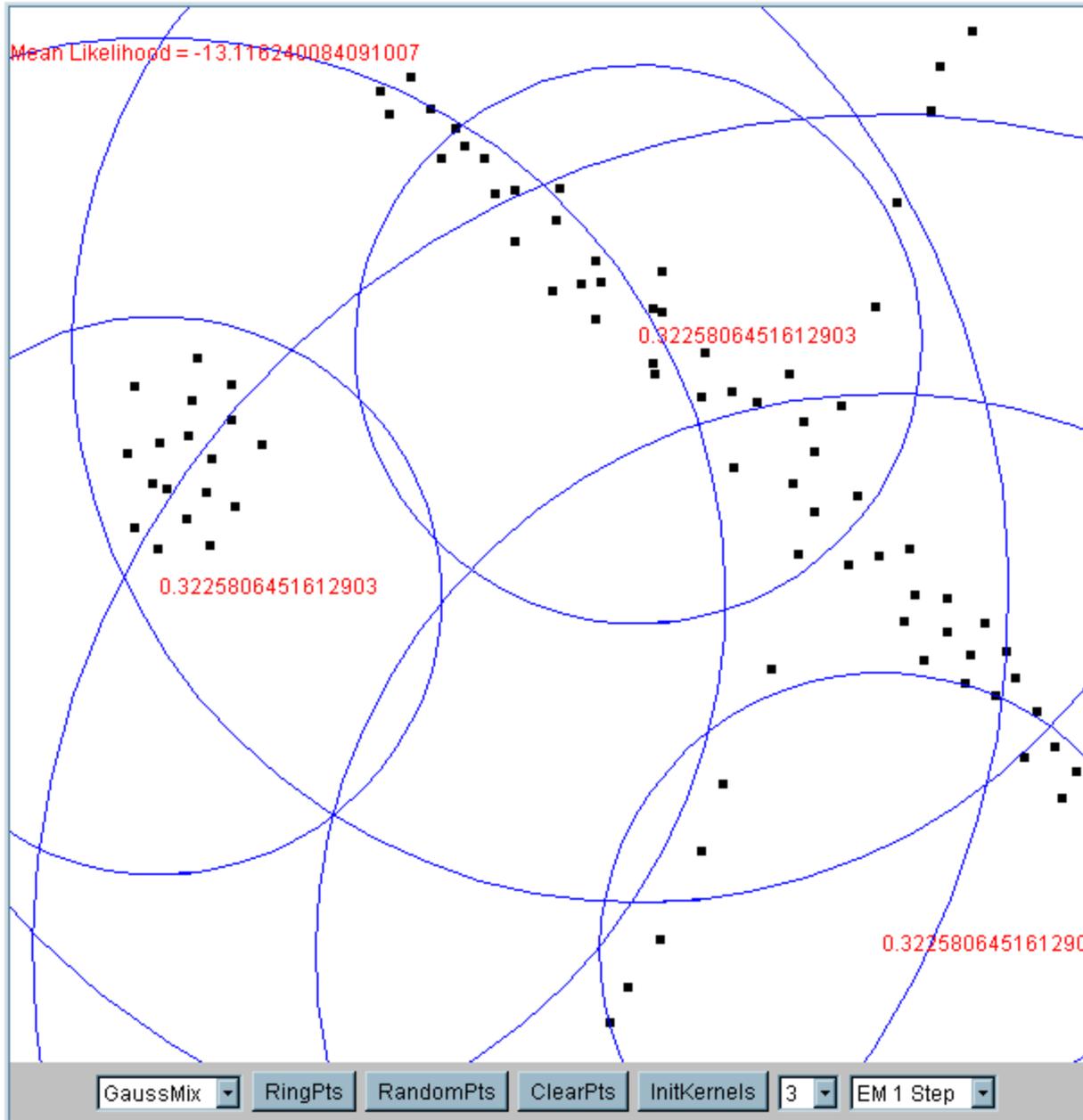
do similar also for  
covariance matrix S

$$\mu_k = \frac{1}{m} \sum_{i=1}^m \frac{d_i P(d_i \in c_k)}{\sum_k P(d_i \in c_k)}$$

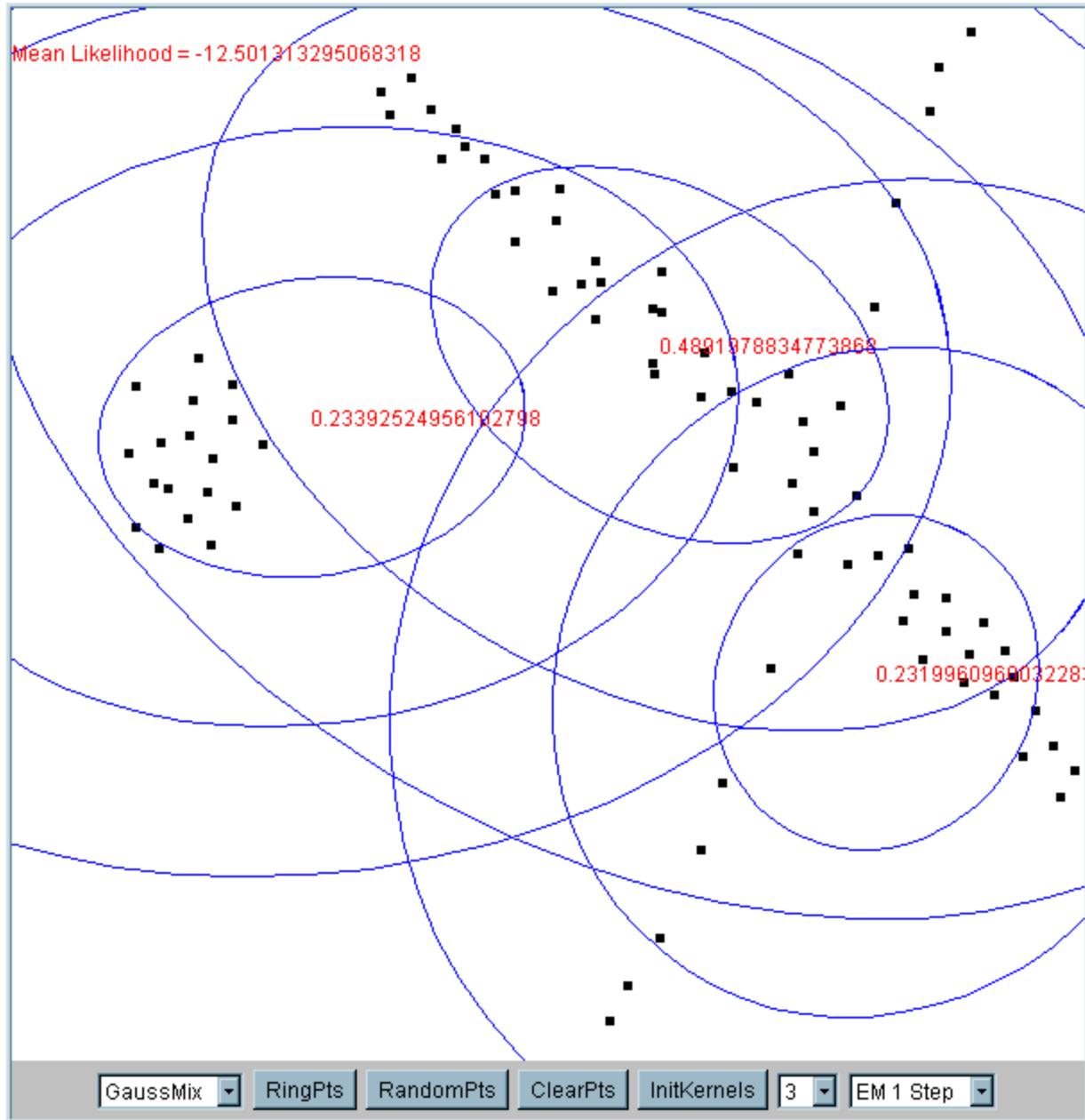


# Iteration 1

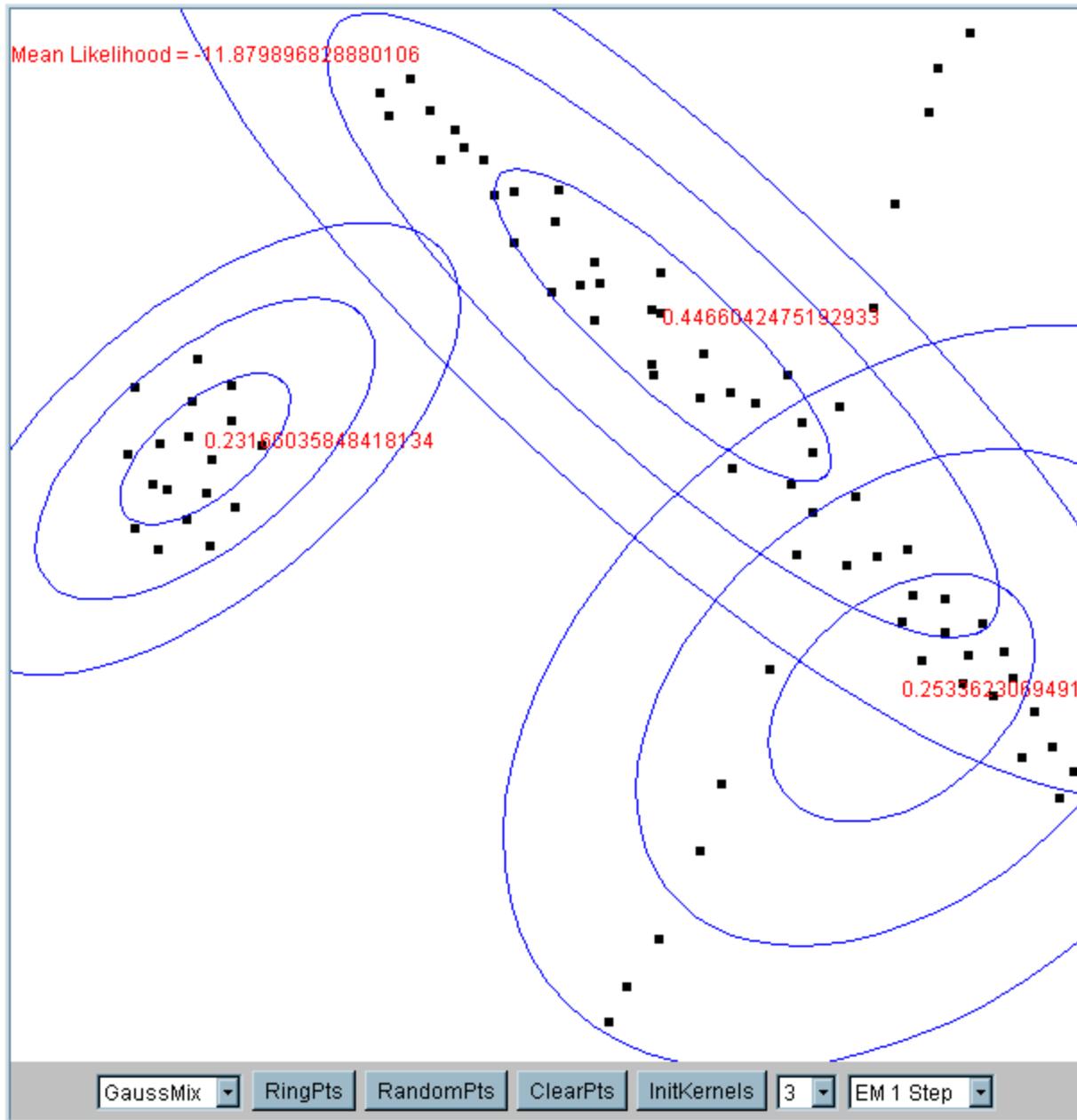
The cluster means are randomly assigned



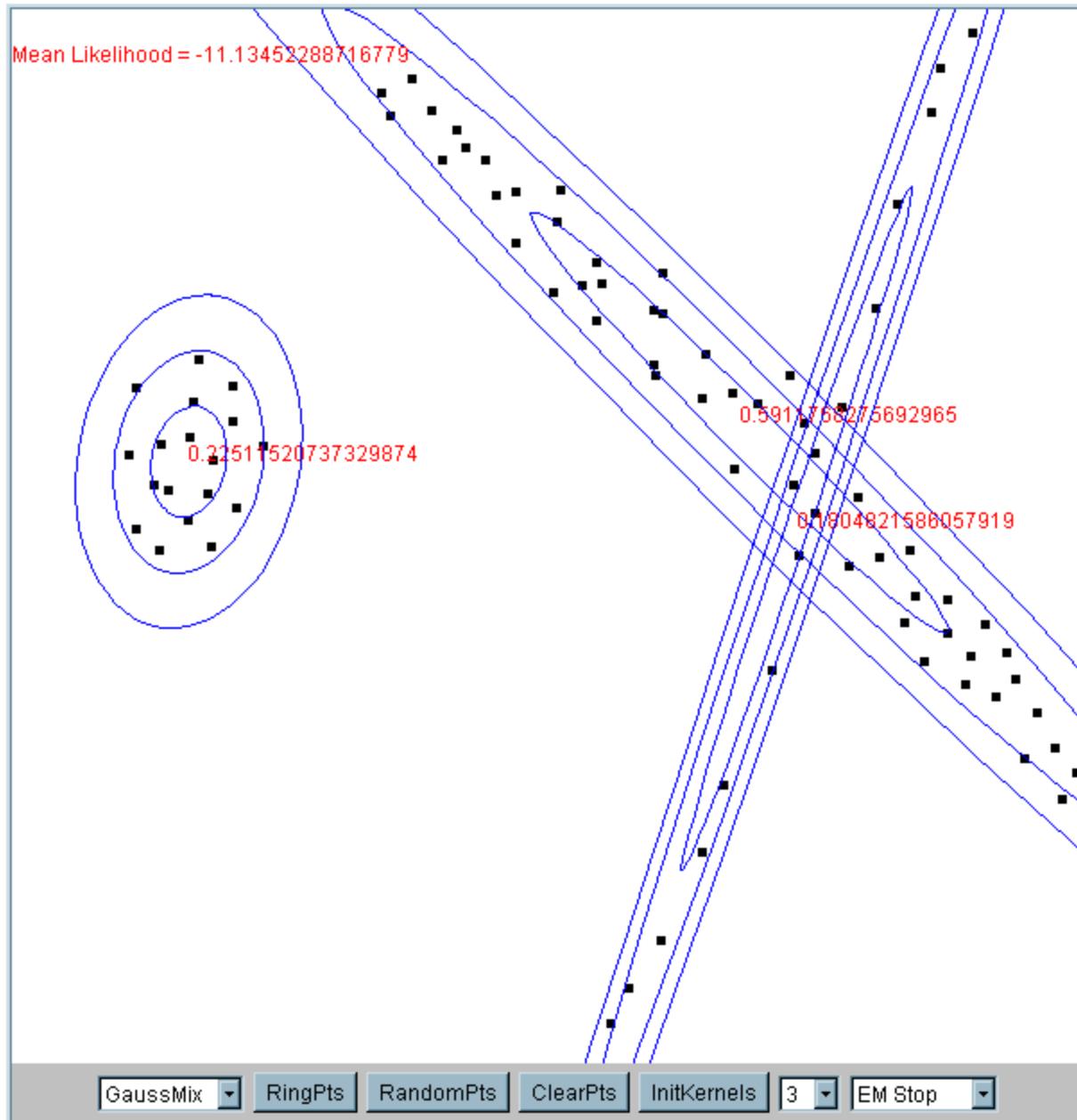
Iteration 2



Iteration 5



Iteration 25

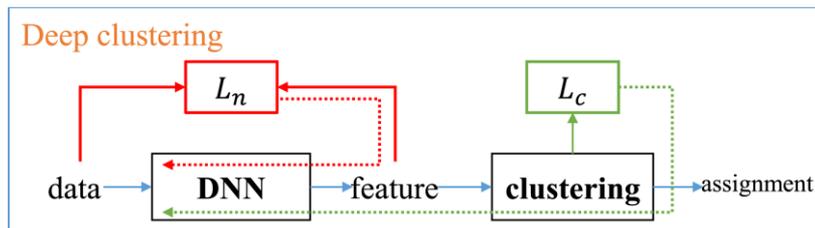


# DEEP CLUSTERING

# STRATEGY

## Embed data into a latent space

- Reduce dimensionality via variational auto-encoder
- Perform clustering in latent space



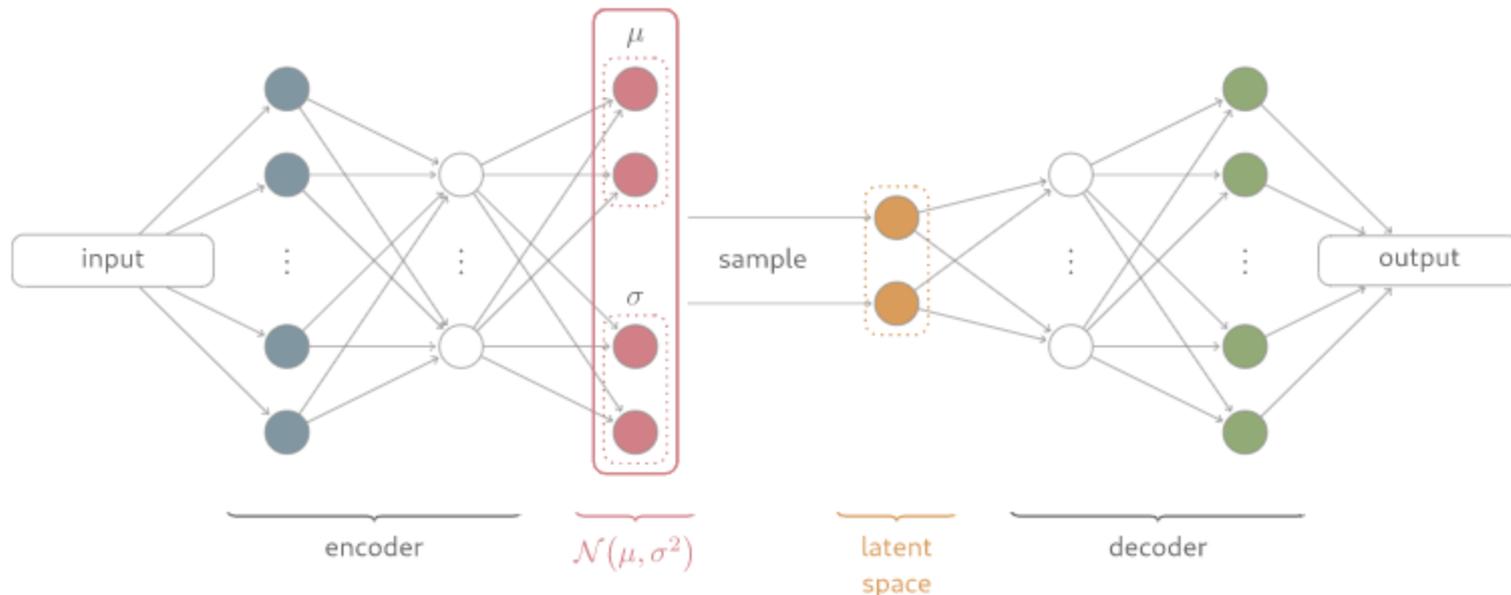
## Key advantages

- Traditional clustering (e.g., k-means, DBSCAN) assumes data lies in a flat Euclidean space
- Deep clustering
  - Learns a non-linear latent manifold
  - “Unfolds” complex structures in high-D data
  - Makes clusters more separable than in raw feature space

# REDUCTION VIA NEURAL NETWORK

## Train a Variational Autoencoder (VAE)

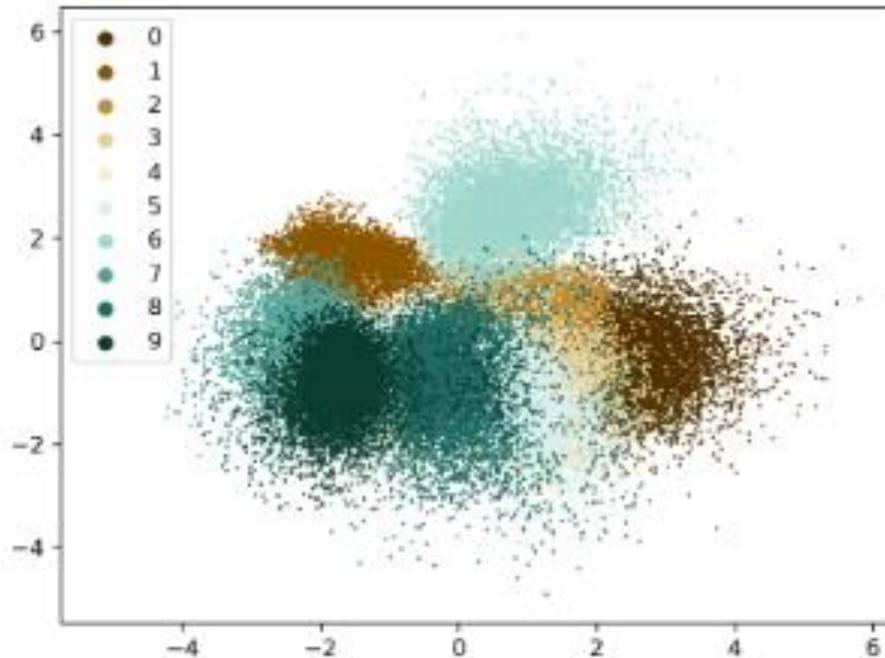
- optimize the output reconstruction loss of the input
- also optimize the latent distribution to be standard normal



# REDUCTION VIA NN: RESULTS

Dataset: 60,000 images of handwritten digits (MINST)

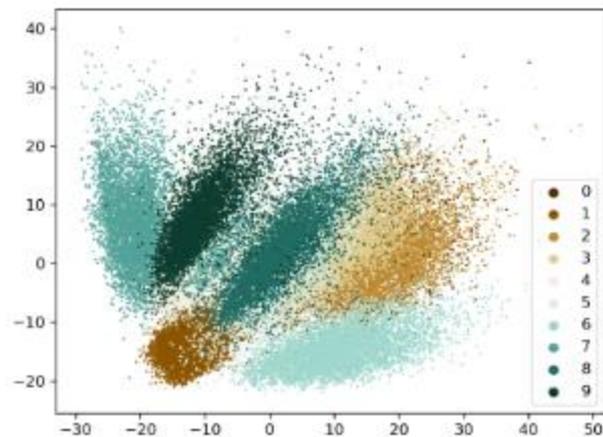
- each image is  $28 \times 28 \rightarrow 784$  D space



PCA projection of its 4D latent space

# REDUCTION VIA NN: RESULTS

Result when not assuring a standard normal distribution in the latent space

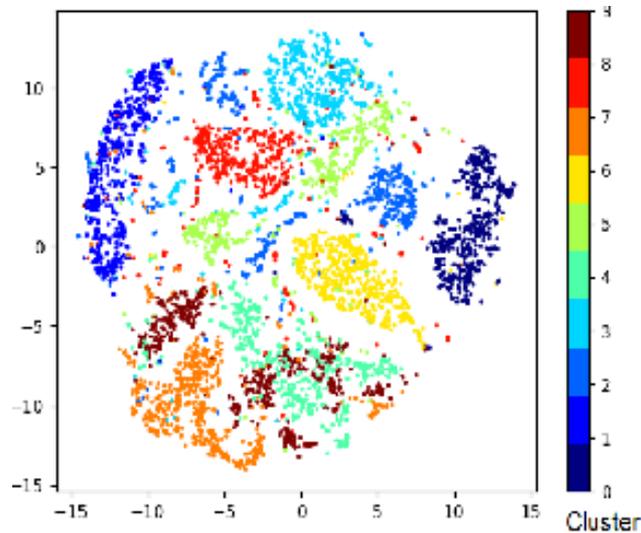


$$l_i(\theta, \phi) = -\mathbb{E}_{z \sim q_\theta(z|x_i)}[\log p_\phi(x_i | z)] + \mathbb{KL}(q_\theta(z | x_i) || p(z))$$

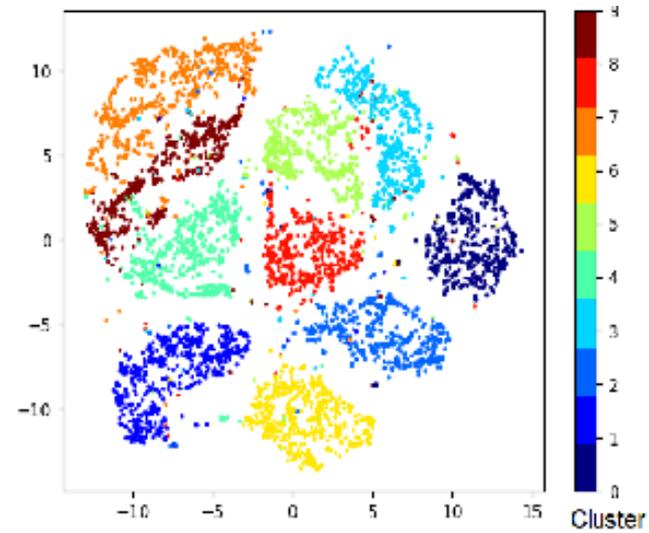
Reconstruction loss

Kullback-Leibler divergence

# EXAMPLE



(a) k-Means



(b) Autoencoder + k-Means

t-SNE visualizations for clustering on MNIST dataset

By: Aljalbout et al. 'Clustering with Deep Learning: Taxonomy and New Methods'

# ALLOWS DIMENSIONALITY REDUCTION WITH SEMANTICS

High-D feature spaces suffer from:

- Distance concentration
- Noise dominance
- Irrelevant dimensions

VAE latent spaces:

- Compress data into a small, information-dense representation
- Preserve semantic similarity, not just numeric similarity
- Suppress noise and redundant dimensions
- This improves cluster stability and interpretability

# SMOOTH, CONTINUOUS LATENT SPACE

Because VAEs enforce a structured latent distribution (typically Gaussian):

- Nearby latent points correspond to similar data instances
- Clusters become compact and well-behaved
- Boundary artifacts caused by raw feature sparsity are reduced

This is especially useful for:

- Image data
- Text embeddings
- Sensor data

# BETTER CLUSTER SHAPE ASSUMPTIONS

Classic clustering often assumes:

- Spherical clusters (k-means)
- Uniform density (DBSCAN)

Deep clustering:

- Allows arbitrary cluster shapes in original space
- Converts them into more regular shapes in latent space
- Enables simple algorithms (e.g., k-means) to work better

Key idea:

- Complexity is absorbed by the encoder, not the clustering algorithm.

# LIMITATIONS

## Limitation of deep clustering

- Less interpretable latent dimensions
- Requires more data
- Harder to tune and debug
- Training instability
- More computationally expensive

# EXTRA: INTERPOLATION IN LATENT SPACE

What's the advantage of it?

- latent space allows easy interpolation
- move between samples in latent space and reconstruct novel instances by the decoder
- not easily possible using other non-linear layouts like MDS, T-SNE

See an example [here](#)

# CLUSTERING OF TIME-SERIES DATA

# TIME SERIES DATA

Rectangular data set with a temporal component

	A	K	L	M	N
1	State	Burglary	Larceny-theft	Motor Vehicle Theft	Arson2
2	CALIFORNIA	2,616	6,298	3,344	71
3	MICHIGAN	1,049	979	154	72
4	MICHIGAN	5,638	8,451	5,828	296
5	TENNESSEE	5,604	12,141	1,373	177
6	MISSOURI	1,960	6,432	1,542	79
7	MARYLAND	3,372	8,761	1,936	137
8	ALABAMA	1,942	3,964	451	
9	OHIO	3,759	5,118	2,008	148
10	ILLINOIS	875	2,242	196	18
11	ARKANSAS	1,852	5,012	589	51
12	CALIFORNIA	2,212	4,450	1,100	43
13	WISCONSIN	2,819	7,622	1,719	120
14	GEORGIA	3,007	8,209	2,328	37

- assume you have these data for each year
- how to handle that, you might ask?

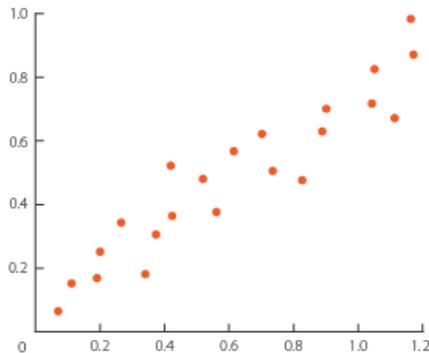
# TIME CUBE

Assume for now we have

- two attributes (burglary, theft)
- both observed over time

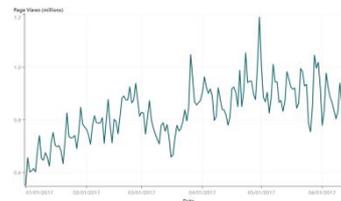
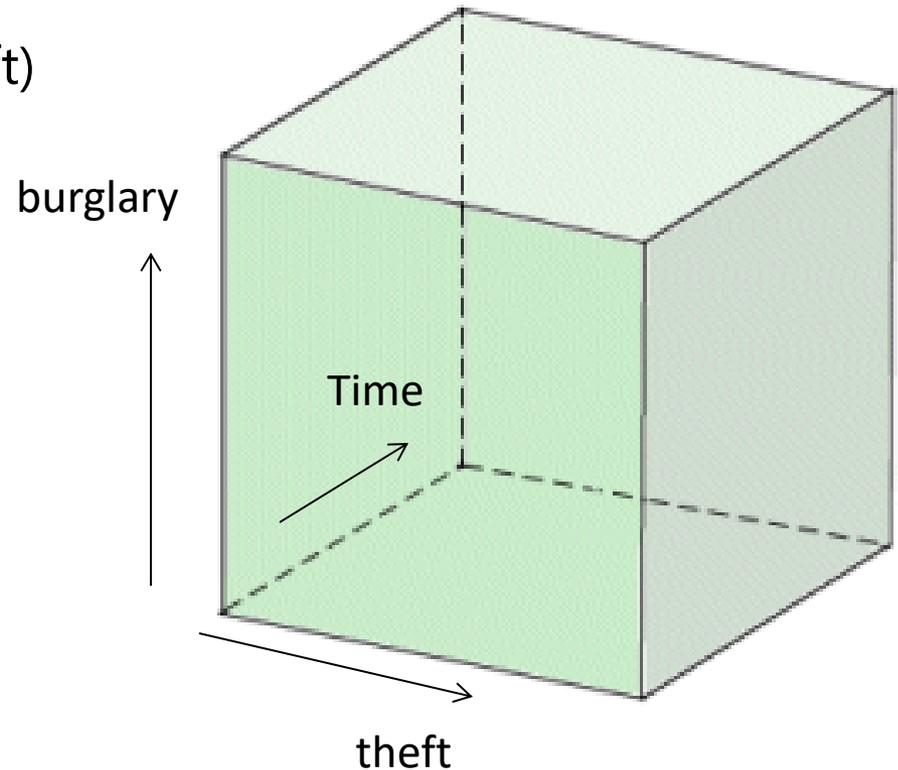
Can visualize

burglary



theft

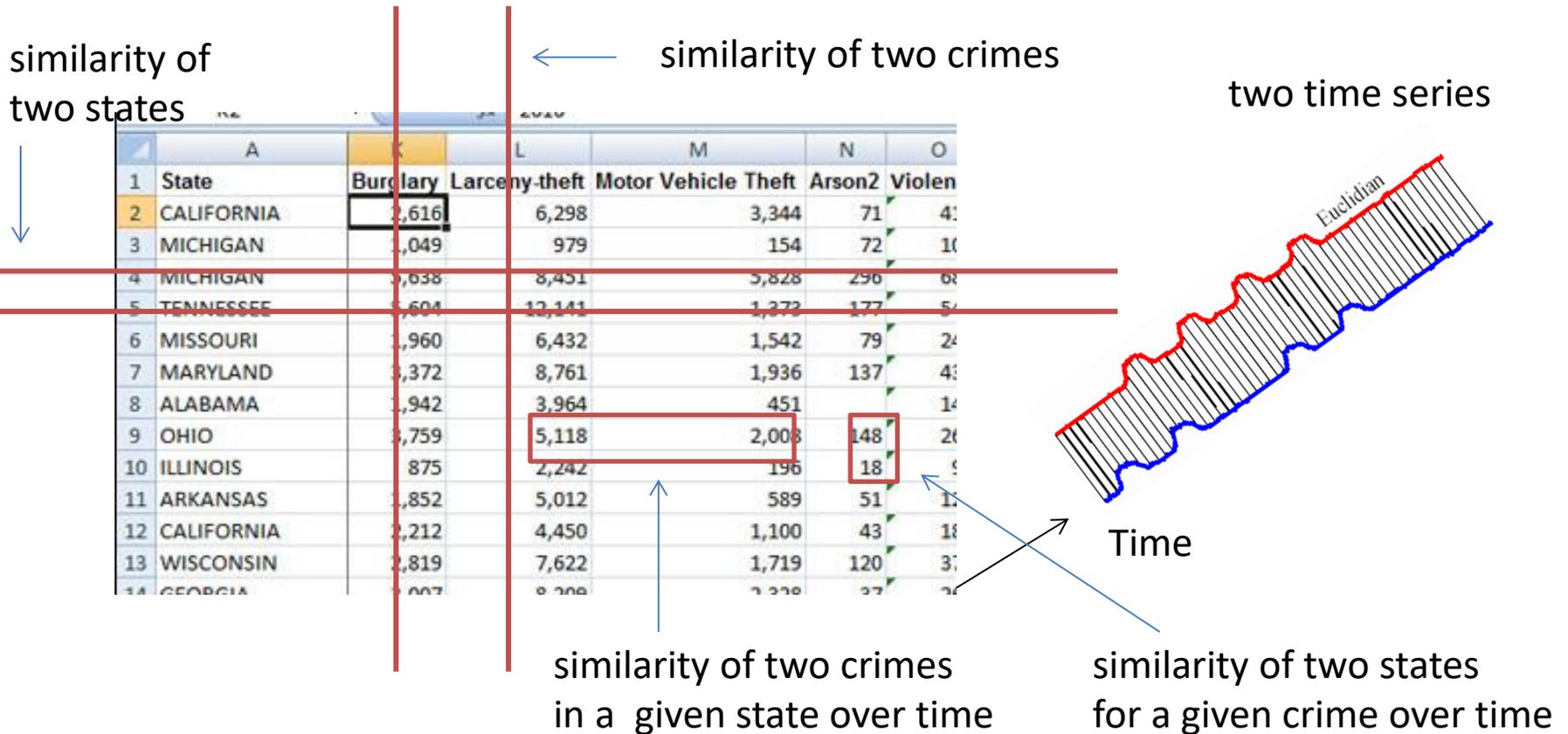
- but each point is a time series!



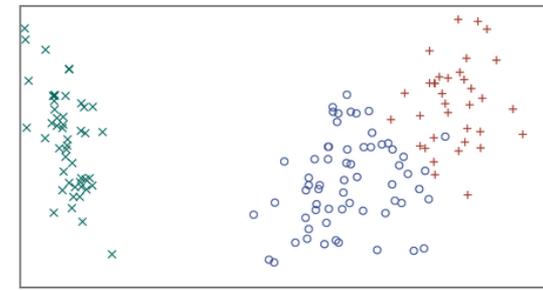
# SIMILARITY MEASURES

Needed it for clustering

- recall Euclidean, correlation, cosine distances

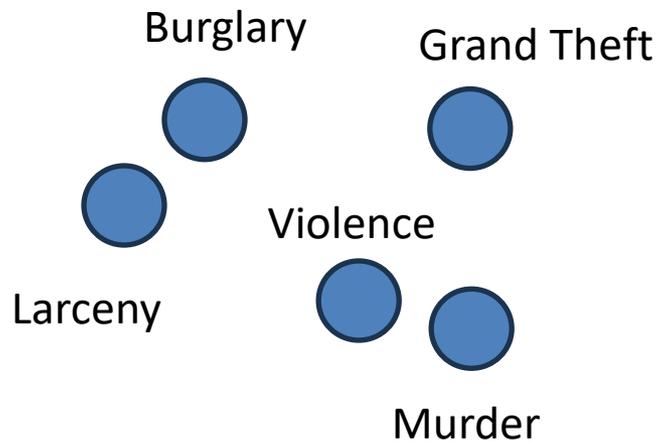


# CLUSTERING ROWS OR COLUMNS

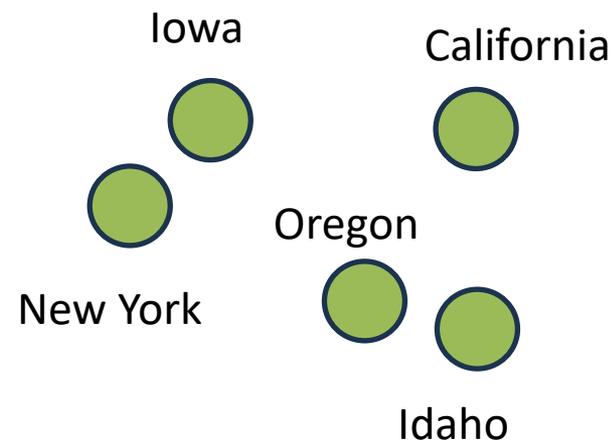


What can be clustered with these measures?

- crimes (averaged over time)
- states (averaged over time)
- crimes in a given state (taking time series into account)
- states for a given crime (taking time series into account)



(Crimes for New York)

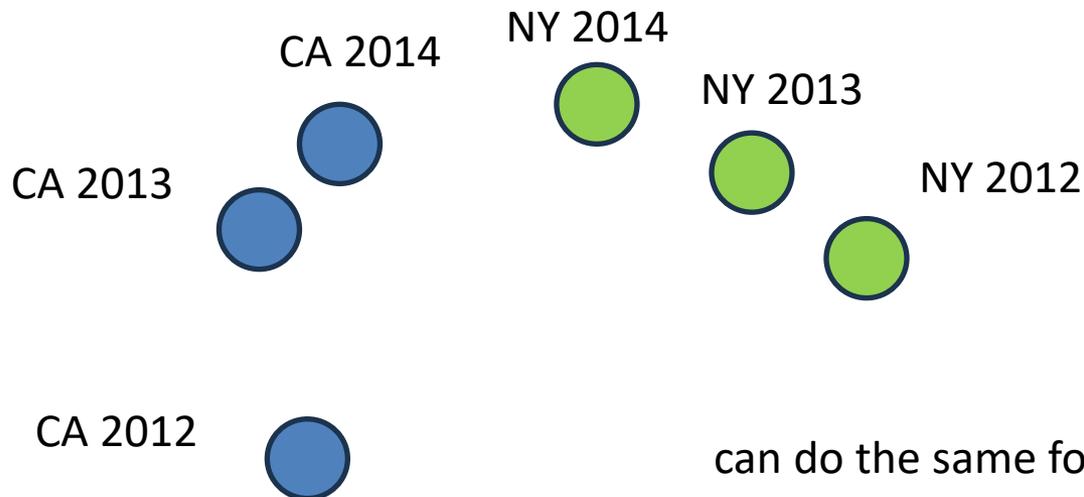


(States for Murder)

# SEPARATE TIME

You may want to just keep time instances as separate entities

- that will work too
- then you might discover clusters that are sensitive to time
- or you can see how the years relate to another along a trajectory



# TIME SERIES ALIGNMENT

The time series might not be aligned

- one crime might cause another
- can apply dynamic time warping (see next)

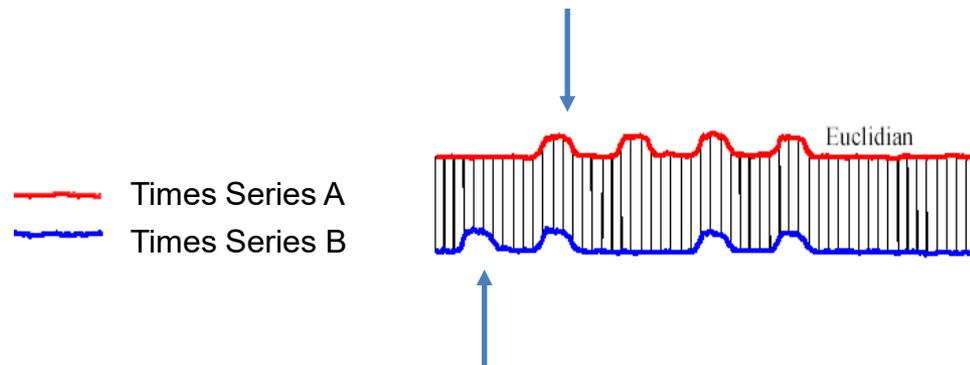
# $L_p$ NORM AND ITS SHORTCOMINGS

Standard pairwise distance

$$Dist(\bar{X}, \bar{Y}) = \left( \sum_{i=1}^n |x_i - y_i|^p \right)^{1/p}$$

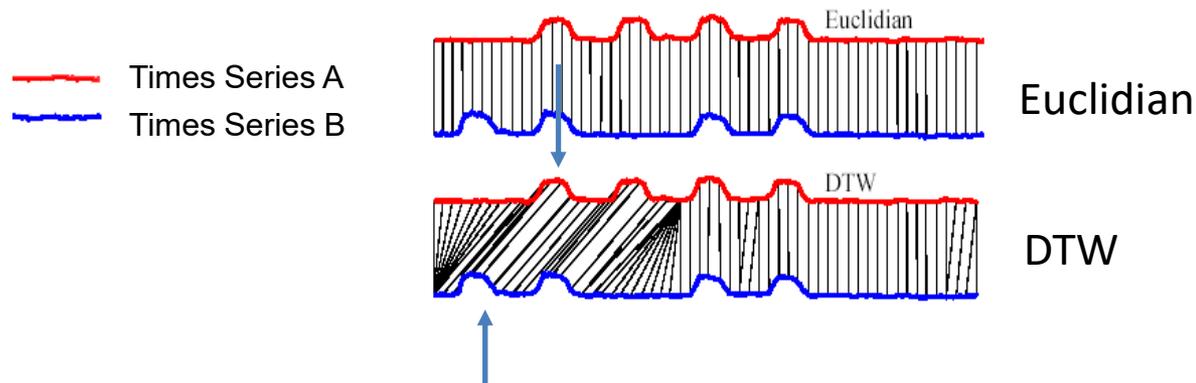
Shortcomings:

- designed for time series of equal length
- cannot address distortions on the temporal (contextual) attributes



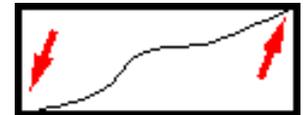
# DYNAMIC TIME WARPING DISTANCE

Can better accommodate local mismatches

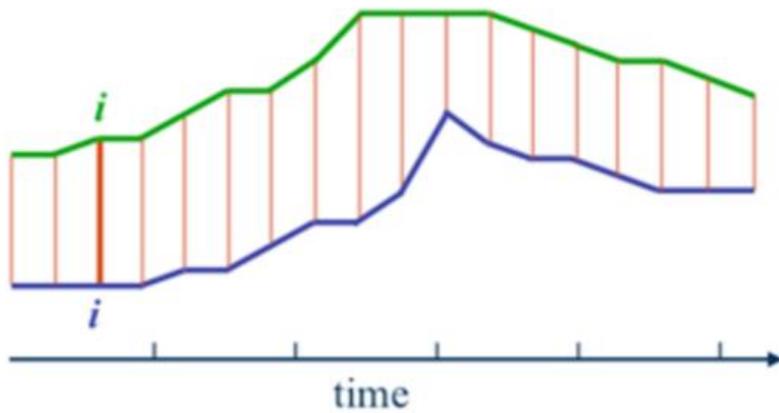


Three constraints

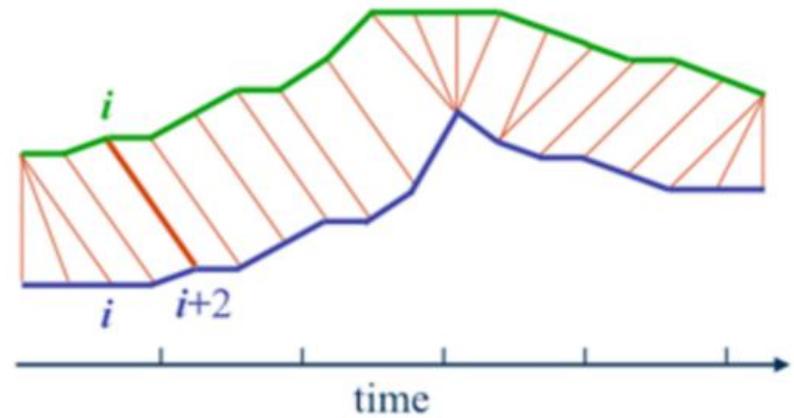
- no skipping of beginning or ends of either sequence
- continuity – no jumps
- monotonicity – can't go back in time



# DTW – FIND THE MINIMUM COST PATH

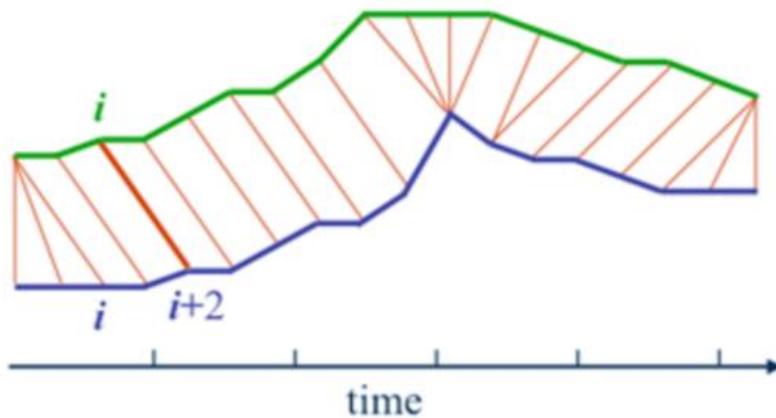


Euclidian



DTW

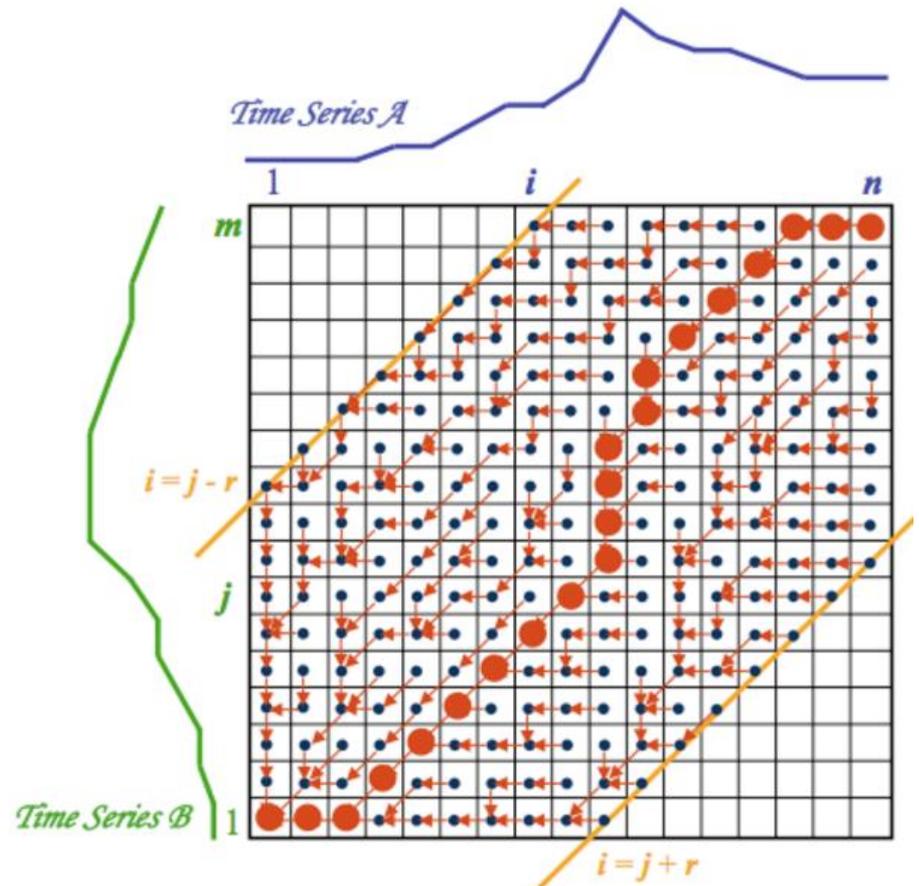
# DTW – FIND THE MINIMUM COST PATH



DTW

Compute using dynamic programming

Available in [python](#)



# SUMMARY

## Structure $\neq$ visualization

- Visual layouts may *suggest* clusters
- Clustering must still be defined algorithmically
- Representation learning (e.g., deep clustering) can make structure clearer — but it does not replace critical judgment
- Subsequent visualization is key, possibly accompanied by user interaction -> subject of a future lecture
- Clusters are hypotheses about structure —their usefulness depends on assumptions, representation, and context and can be optimized -> subject of a future lecture

# SUMMARY VISUALIZED

## 1 Same similarity, different structure

Given similarity, different clustering philosophies impose different structure.

Numeric similarity

Temporal similarity

Algorithm



Temporal similarity



## 3 Structure ≠ visualization

Visual layouts may suggest clusters; clustering must be defined algorithmically.

Similarity matrix

Clustering = modeling similarity structure

Hierarchical

Density-Based

Model-Based (GMM)

Deep Clustering (VAE)

No single "best" clustering—only choices

Same tag: ...

## 2 What you cluster depends on data representation

The same data answers different questions depending on representation.

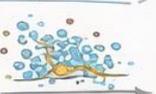
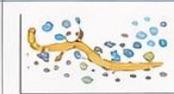
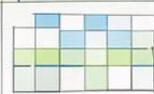
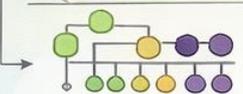
Time series

Rows

Columns

Time Series

...



Cluster crimes by state patterns

States by crime

Cluster states by crime profiles

Cluster time series (alignment)

Before clustering, decide: What are the samples (rows)?  
What are the attributes (columns)?

## 4 Big Picture

Clusters are hypotheses about structure—their usefulness depends on assumptions, representation, and context.



(Auto-clustering & evaluation are covered separately) ...