

Unsupervised clustering

Guillaume TOCHON & Joseph CHAZALON

LRDE

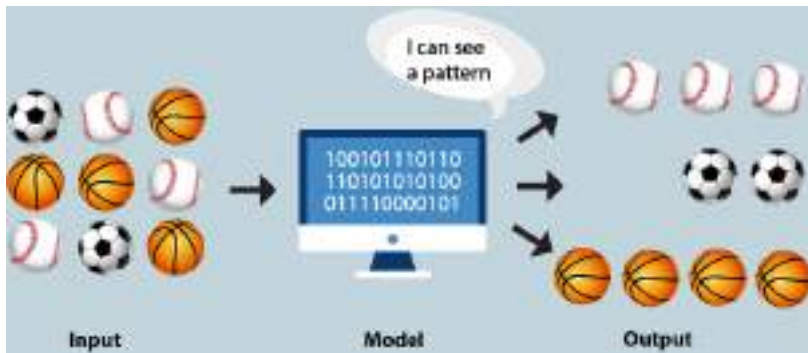


Why do we care?

Clustering: Group the input data into clusters that share some characteristics.

(yup, that's a vague definition)

- Find general patterns in the data (data mining problem)
- Visualize the data (in a simpler way)
- Infer some properties of a given data point based on how it relates to other data points (statistical learning)



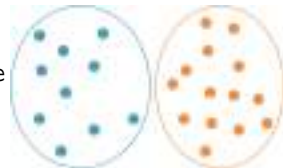
Why is it tricky?

- Clustering belongs to unsupervised learning \Rightarrow no ground truth available to learn/evaluate the quality of a clustering algorithm.



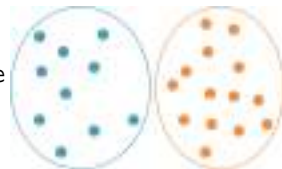
Why is it tricky?

- Clustering belongs to unsupervised learning \Rightarrow no ground truth available to learn/evaluate the quality of a clustering algorithm.



Why is it tricky?

→ Clustering belongs to unsupervised learning \Rightarrow no ground truth available to learn/evaluate the quality of a clustering algorithm.



→ How to assess how much data points are related to each other?

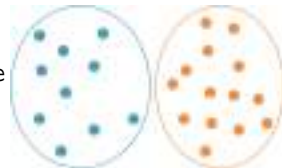
\Rightarrow Which criteria (features) are the more relevant for our problem?

\Rightarrow Which metric makes the most sense?



Why is it tricky?

→ Clustering belongs to unsupervised learning \Rightarrow no ground truth available to learn/evaluate the quality of a clustering algorithm.



→ How to assess how much data points are related to each other?

\Rightarrow Which criteria (features) are the more relevant for our problem?

\Rightarrow Which metric makes the most sense?



→ How to assess the soundness of the created clusters? Is that even a relevant question?

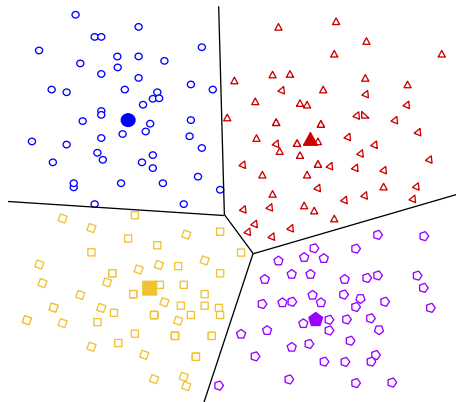


Families of clustering approaches

Clusteringception

One can try to divide existing clustering approaches into several categories:

Centroid-based clustering Clusters are summarized using a single representative point, and points are assigned to clusters based on their distance to this *centroid* (k-means and its direct variants).



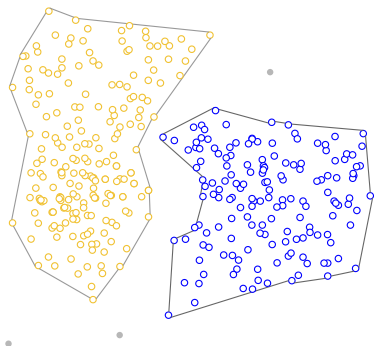
Families of clustering approaches

Clusteringception

One can try to divide existing clustering approaches into several categories:

Centroid-based clustering Clusters are summarized using a single representative point, and points are assigned to clusters based on their distance to this *centroid* (k-means and its direct variants).

Density-based clustering Clusters are defined as set of dense points in the feature space (mean shift clustering, DBSCAN, ...)



Families of clustering approaches

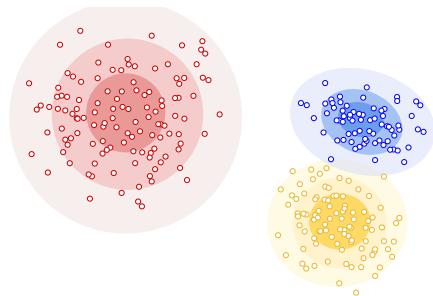
Clusteringception

One can try to divide existing clustering approaches into several categories:

Centroid-based clustering Clusters are summarized using a single representative point, and points are assigned to clusters based on their distance to this *centroid* (k-means and its direct variants).

Density-based clustering Clusters are defined as set of dense points in the feature space (mean shift clustering, DBSCAN, ...)

Distribution-based clustering Clusters are defined based on the likelihood of points to belong to the same probability distribution (Gaussian mixture models, ...)



Families of clustering approaches

Clusteringception

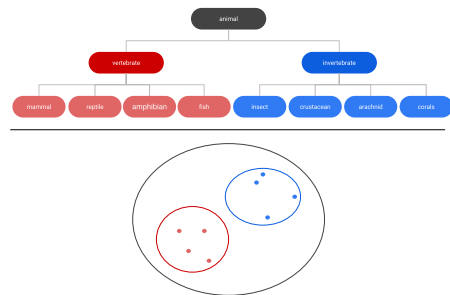
One can try to divide existing clustering approaches into several categories:

Centroid-based clustering Clusters are summarized using a single representative point, and points are assigned to clusters based on their distance to this *centroid* (k-means and its direct variants).

Density-based clustering Clusters are defined as set of dense points in the feature space (mean shift clustering, DB-SCAN, ...)

Distribution-based clustering Clusters are defined based on the likelihood of points to belong to the same probability distribution (Gaussian mixture models, ...)

Hierarchical clustering Clusters are organized in a hierarchical way (Hierarchical Agglomerative Clustering, Recursive K-Means).



Families of clustering approaches

Clusteringception

One can try to divide existing clustering approaches into several categories:

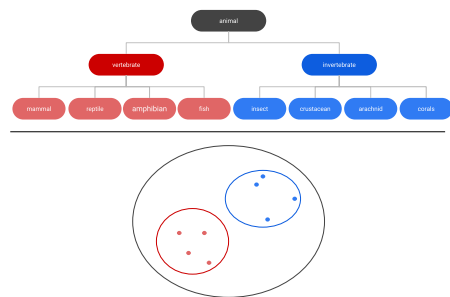
Centroid-based clustering Clusters are summarized using a single representative point, and points are assigned to clusters based on their distance to this *centroid* (k-means and its direct variants).

Density-based clustering Clusters are defined as set of dense points in the feature space (mean shift clustering, DB-SCAN, ...)

Distribution-based clustering Clusters are defined based on the likelihood of points to belong to the same probability distribution (Gaussian mixture models, ...)

Hierarchical clustering Clusters are organized in a hierarchical way (Hierarchical Agglomerative Clustering, Recursive K-Means).

And so on ...

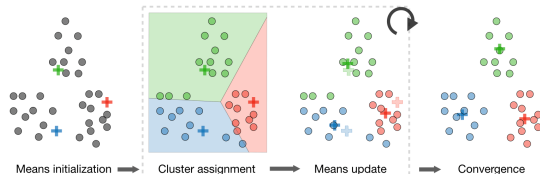


k-means clustering

Partition n observations $\mathbf{x}_1, \dots, \mathbf{x}_n$ into k clusters $\mathbf{C} = \{C_1, \dots, C_k\}$ where each observation \mathbf{x}_i belongs to the clusters C_{j^*} whose mean μ_{j^*} is the closest: $\mathbf{x}_i \in S_{j^*}$ with $j^* = \arg \min_j \|\mathbf{x}_i - \mu_j\|_2$.

Algorithm Basic K-means algorithm

- 1: Select K points as initial centroids.
 - 2: **repeat**
 - 3: Form K clusters by assigning each point to its closest centroid.
 - 4: Recompute the centroid of each cluster.
 - 5: **until** Centroids do not change.
-

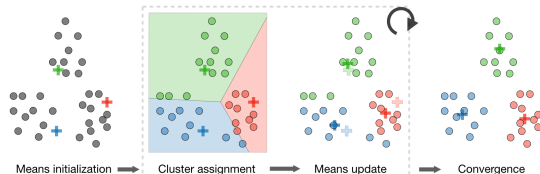


k-means clustering

Partition n observations $\mathbf{x}_1, \dots, \mathbf{x}_n$ into k clusters $\mathbf{C} = \{C_1, \dots, C_k\}$ where each observation \mathbf{x}_i belongs to the clusters C_{j^*} whose mean μ_{j^*} is the closest: $\mathbf{x}_i \in S_{j^*}$ with $j^* = \arg \min_j \|\mathbf{x}_i - \mu_j\|_2$.

Algorithm Basic K-means algorithm

- 1: Select K points as initial centroids
- 2: repeat
- 3: Form K clusters by assigning each point to its closest centroid.
- 4: Recompute the centroid of each cluster.
- 5: until Centroids do not change.



→ Minimizes within-cluster sum of squares (variance)

→ Overall optimization problem:

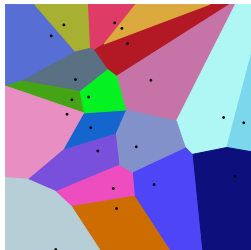
$$\arg \min_{\mathbf{C}=\{C_1, \dots, C_k\}} \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} \|\mathbf{x} - \mu_i\|^2$$

→ NP-hard problem, no guarantee to find the global optimum

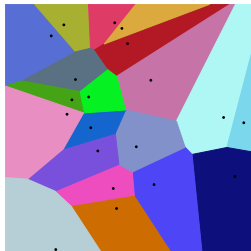
→ Stochastic and very sensitive to initial conditions

→ Sensitive to outliers (thank you, L_2 norm...)

→ Yet it's probably the most used clustering algorithm out there



Voronoi tessellation: partition of the Euclidean space relatively to discrete points/seeds. Each region/Voronoi cell is composed of all the points in the space that are closer to the cell seed than to any other seed



Voronoi tessellation: partition of the Euclidean space relatively to discrete points/seeds. Each region/Voronoi cell is composed of all the points in the space that are closer to the cell seed than to any other seed

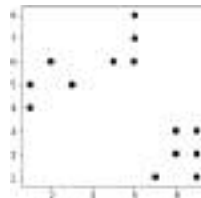
- k-means provides a way to obtain a Voronoi tessellation of the input space, where seeds are the final cluster means.
- Alternatively, one can use some pre-computed Voronoi tessellation seeds as initial clusters for k-means



Determining the optimal number of clusters with the elbow method

Empirical af, but still better than nothing...

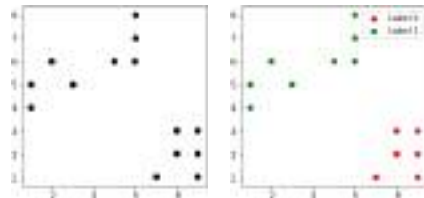
How many clusters for this data set?



Determining the optimal number of clusters with the elbow method

Empirical af, but still better than nothing...

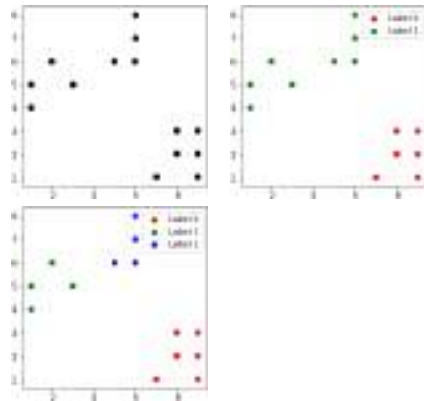
How many clusters for this data set?



Determining the optimal number of clusters with the elbow method

Empirical af, but still better than nothing...

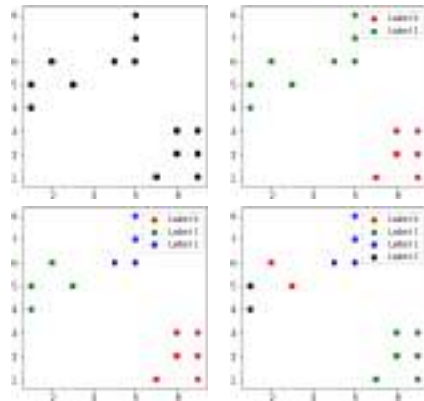
How many clusters for this data set?



Determining the optimal number of clusters with the elbow method

Empirical af, but still better than nothing...

How many clusters for this data set?



Determining the optimal number of clusters with the elbow method

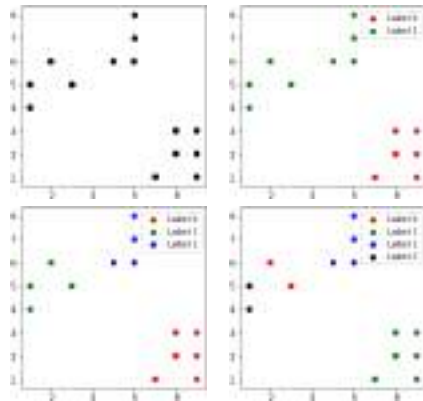
Empirical af, but still better than nothing...

How many clusters for this data set?

- Compute explained variance for an increasing number of clusters k

$$\text{Var}(\mathbf{C} = \{C_1, \dots, C_k\}) = \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2$$

- Plot and find the bend of the elbow



Determining the optimal number of clusters with the elbow method

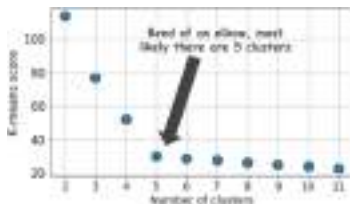
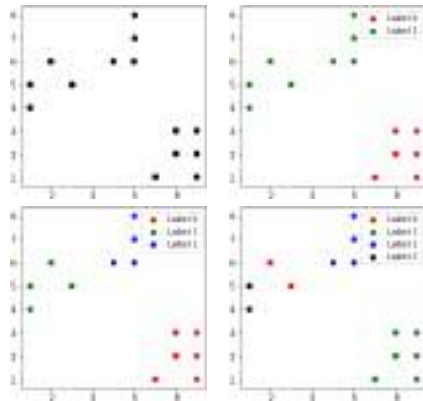
Empirical af, but still better than nothing...

How many clusters for this data set?

- Compute explained variance for an increasing number of clusters k

$$\text{Var}(\mathbf{C} = \{C_1, \dots, C_k\}) = \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2$$

- Plot and find the bend of the elbow



Determining the optimal number of clusters with the elbow method

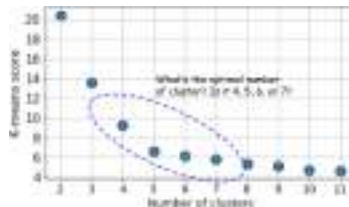
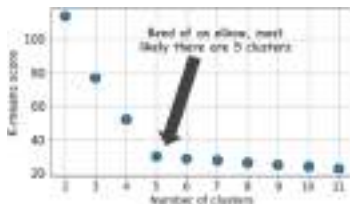
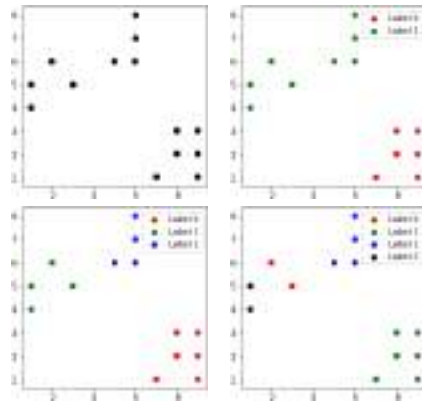
Empirical af, but still better than nothing...

How many clusters for this data set?

- Compute explained variance for an increasing number of clusters k

$$\text{Var}(\mathbf{C} = \{C_1, \dots, C_k\}) = \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2$$

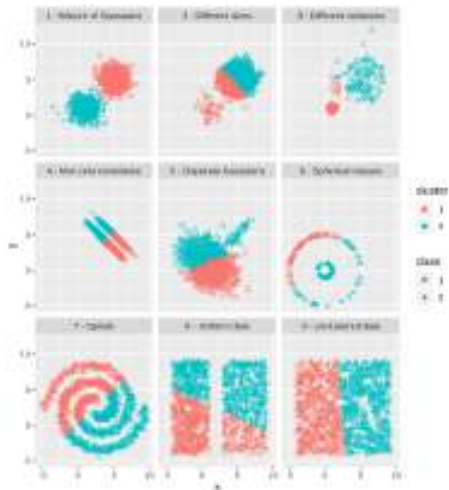
- Plot and find the bend of the elbow



Sometimes, k-means works. . .

But most of the time, not as expected. And it is probably because of the L_2 norm that k-means tries to minimize:

- Sensible to curse of dimensionality
- Form “normalized Gaussian” clusters (spherical)
- Does not adapt to manifold geometry
- Sensible to class imbalance
- Sensible to outliers



Simple Linear Iterative Clustering

A kickass image segmentation algorithm using k-means

SLIC superpixels uses a modified k-means clustering in the *Labxy* space to produce k clusters regularly sampled and perceptually coherent from a color point of view.

Algorithm SLIC superpixel segmentation

- 1: Initialize cluster centers $C_k = [l_k, a_k, b_k, x_k, y_k]^T$ by sampling pixels at regular grid steps S .
- 2: Perturb cluster centers in an $n \times n$ neighborhood, to the lowest gradient position.
- 3: **repeat**
- 4: **for** each cluster center C_k **do**
- 5: Assign the best matching pixels from a $2S \times 2S$ square neighborhood around the cluster center according to the distance measure (Eq. 1).
- 6: **end for**
- 7: Compute new cluster centers and residual error E {L1 distance between previous centers and recomputed centers}
- 8: **until** $E \leq \text{threshold}$
- 9: Enforce connectivity.

$$d_{lab} = \sqrt{(l_k - l_i)^2 + (a_k - a_i)^2 + (b_k - b_i)^2}$$
$$(1) \quad d_{xy} = \sqrt{(x_k - x_i)^2 + (y_k - y_i)^2}$$
$$D_k = d_{lab} + \frac{m}{S} d_{xy}$$



k-medoids clustering

A possible extension to k-means

Clusters centroids are not initial data points \rightarrow can be problematic

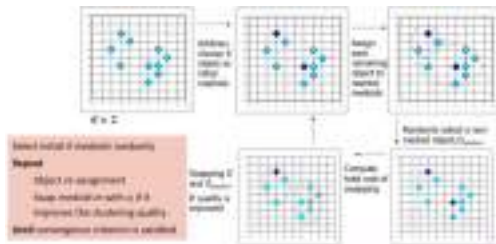
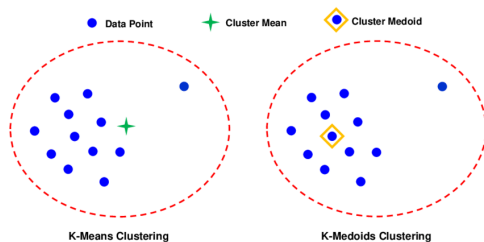
\Rightarrow Replace centroid by medoid (point with the smallest distance to all other points in the cluster)

$$\mathbf{m}_C = \arg \min_{\mathbf{x} \in C} \sum_{\mathbf{x}_i \in C} d(\mathbf{x}, \mathbf{x}_i)$$

\Rightarrow k-medoids algorithm

Overall objective: find k medoids $\mathbf{m}_1, \dots, \mathbf{m}_k$ that minimize the partitioning cost

$$\sum_{i=1}^k \sum_{\mathbf{x} \in C_i} d(\mathbf{x}, \mathbf{m}_i)$$

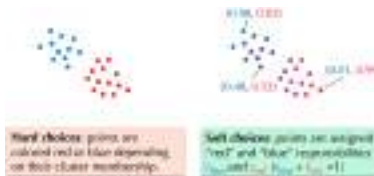
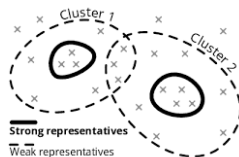


Fuzzy c-means clustering

Let it fuzz

k-means is a *hard* clustering method → each data point 100% belongs to its cluster.

Soft (aka *fuzzy*) clustering methods allow each data point to belong to several clusters with various degrees of membership.

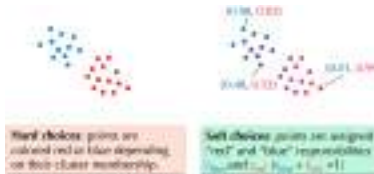
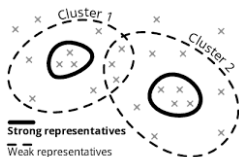


Fuzzy c-means clustering

Let it fuzz

k-means is a *hard* clustering method \rightarrow each data point 100% belongs to its cluster.

Soft (aka *fuzzy*) clustering methods allow each data point to belong to several clusters with various degrees of membership.

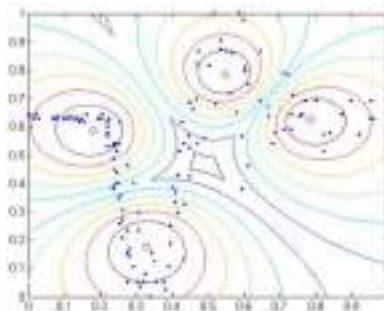


FCM clustering: outputs clusters $\mathcal{C}_1, \dots, \mathcal{C}_k$ and membership matrix $\mathbf{W}_{n \times k}$ ($w_{ij} = \%(x_i \in \mathcal{C}_j)$)

$$\Rightarrow \arg \min_{\mathbf{C}=\{\mathcal{C}_1, \dots, \mathcal{C}_k\}} \sum_{i=1}^n \sum_{j=1}^k w_{ij}^m \|x_i - \mu_j\|^2$$

$$\Rightarrow \text{Alternatively update } \mu_j = \frac{\sum_{x_i} w_{ij}^m x_i}{\sum_{x_i} w_{ij}^m} \text{ and}$$

$$w_{ij} = \frac{1}{\sum_{l=1}^k \left(\frac{\|x_i - \mu_j\|}{\|x_i - \mu_l\|} \right)^{\frac{2}{m-1}}}$$

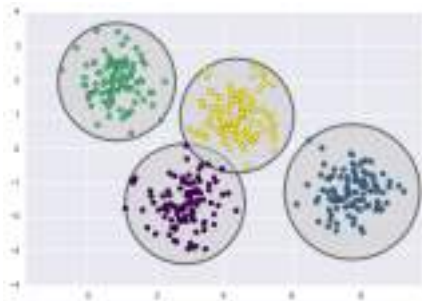


Gaussian mixture models

k-means on steroids

k-means works for spherical clusters, but fails in any other case \Rightarrow try harder

Model probability density function f of data as a mixture of multivariate Gaussian

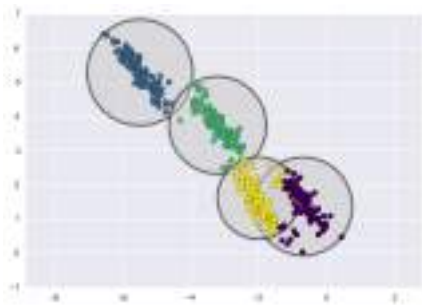


Gaussian mixture models

k-means on steroids

k-means works for spherical clusters, but fails in any other case \Rightarrow try harder

Model probability density function f of data as a mixture of multivariate Gaussian

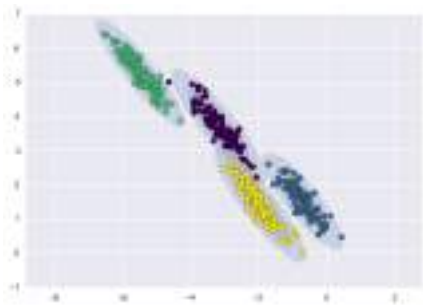


Gaussian mixture models

k-means on steroids

k-means works for spherical clusters, but fails in any other case \Rightarrow try harder

Model probability density function f of data as a mixture of multivariate Gaussian

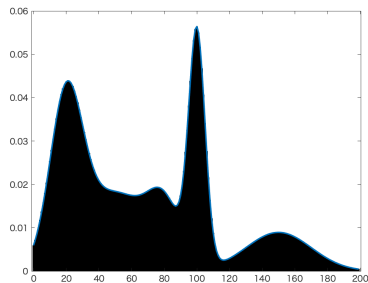
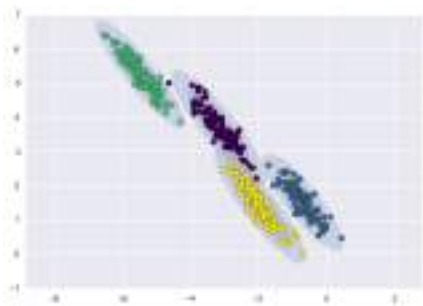


Gaussian mixture models

k-means on steroids

k-means works for spherical clusters, but fails in any other case \Rightarrow try harder

Model probability density function f of data as a mixture of multivariate Gaussian

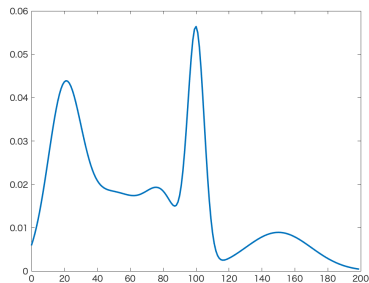
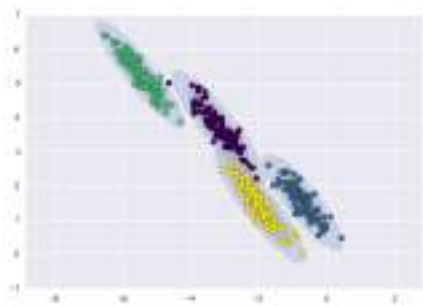


Gaussian mixture models

k-means on steroids

k-means works for spherical clusters, but fails in any other case \Rightarrow try harder

Model probability density function f of data as a mixture of multivariate Gaussian

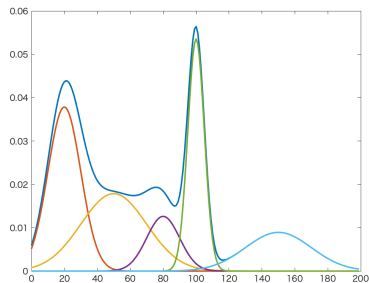
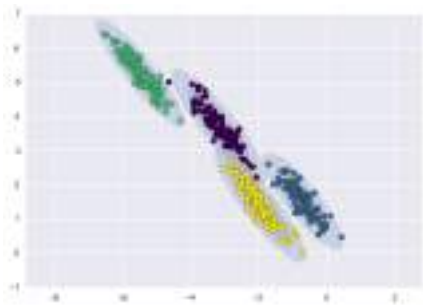


Gaussian mixture models

k-means on steroids

k-means works for spherical clusters, but fails in any other case \Rightarrow try harder

Model probability density function f of data as a mixture of multivariate Gaussian

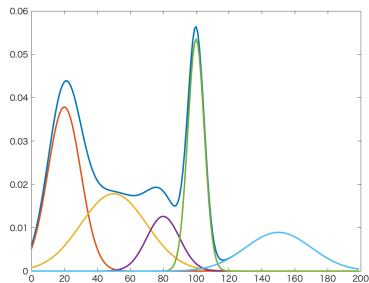
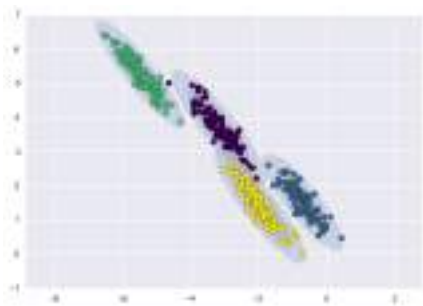


Gaussian mixture models

k-means on steroids

k-means works for spherical clusters, but fails in any other case \Rightarrow try harder

Model probability density function f of data as a mixture of multivariate Gaussian



Gaussian mixture model: $f(\mathbf{x}) = \sum_{i=1}^k \phi_i \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$

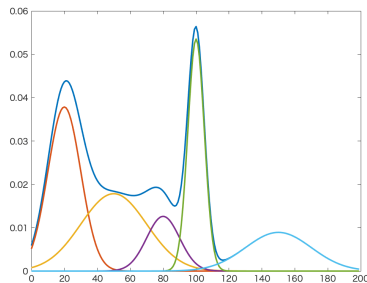
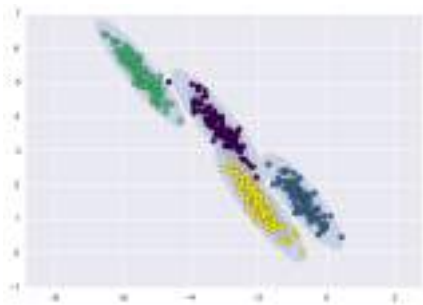
recall that $\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^N \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$ for $\mathbf{x} \in \mathbb{R}^N$

Gaussian mixture models

k-means on steroids

k-means works for spherical clusters, but fails in any other case \Rightarrow try harder

Model probability density function f of data as a mixture of multivariate Gaussian



Gaussian mixture model: $f(\mathbf{x}) = \sum_{i=1}^k \phi_i \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$

recall that $\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^N \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$ for $\mathbf{x} \in \mathbb{R}^N$

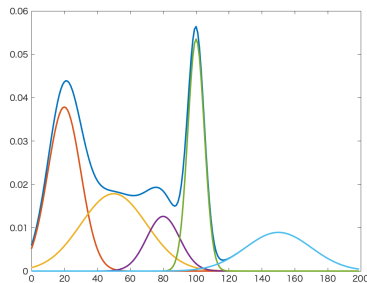
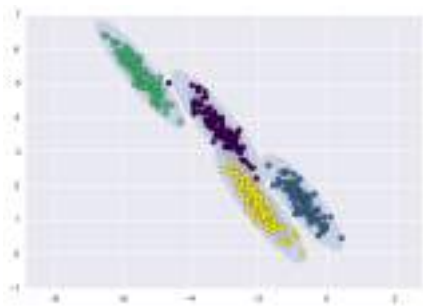
$\rightarrow \phi_i$ are mixture component weights ($\sum_{i=1}^k \phi_i = 1$)

Gaussian mixture models

k-means on steroids

k-means works for spherical clusters, but fails in any other case \Rightarrow try harder

Model probability density function f of data as a mixture of multivariate Gaussian



Gaussian mixture model: $f(\mathbf{x}) = \sum_{i=1}^k \phi_i \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$

recall that $\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^N \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$ for $\mathbf{x} \in \mathbb{R}^N$

$\rightarrow \phi_i$ are mixture component weights ($\sum_{i=1}^k \phi_i = 1$)

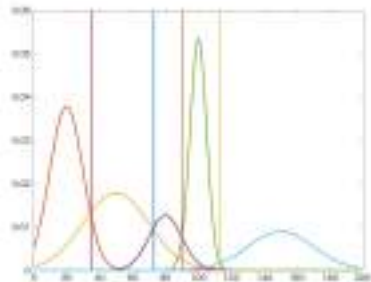
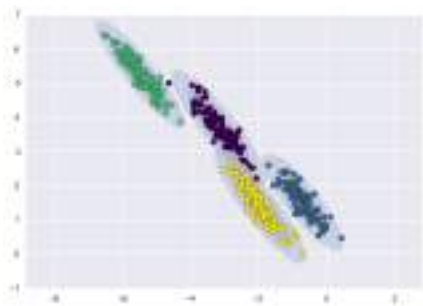
\rightarrow How to estimate $\phi_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i \forall i = 1, \dots, k$?

Gaussian mixture models

k-means on steroids

k-means works for spherical clusters, but fails in any other case \Rightarrow try harder

Model probability density function f of data as a mixture of multivariate Gaussian



Gaussian mixture model: $f(\mathbf{x}) = \sum_{i=1}^k \phi_i \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$

recall that $\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^N \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$ for $\mathbf{x} \in \mathbb{R}^N$

$\rightarrow \phi_i$ are mixture component weights ($\sum_{i=1}^k \phi_i = 1$)

\rightarrow How to estimate $\phi_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i \forall i = 1, \dots, k$?

The EM algorithm

The nightmare of any EECS grad student around the world

Initialization

- Select k random points as initial means $\hat{\mu}_1, \dots, \hat{\mu}_k$
- Init all covariance matrices $\hat{\Sigma}_1, \dots, \hat{\Sigma}_k$ as whole data sample covariance matrix $\hat{\Sigma}$
- Set uniform mixture weights $\hat{\phi}_1, \dots, \hat{\phi}_k = \frac{1}{k}$

Alternate until convergence

Expectation step

- Compute membership weight $\hat{\gamma}_{ij}$ of \mathbf{x}_i with respect to j^{th} component $\mathcal{N}(\mathbf{x}|\mu_j, \Sigma_j)$

$$\hat{\gamma}_{ij} = \frac{\hat{\phi}_j \mathcal{N}(\mathbf{x}_i | \hat{\mu}_j, \hat{\Sigma}_j)}{\sum_{m=1}^k \hat{\phi}_m \mathcal{N}(\mathbf{x}_i | \hat{\mu}_m, \hat{\Sigma}_m)}$$

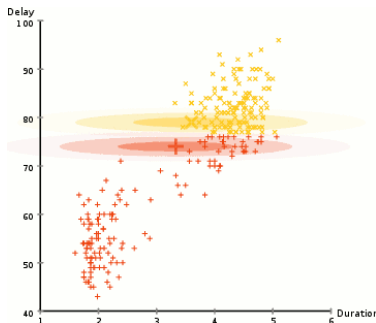
$$\hat{\gamma}_{ij} \equiv \text{posterior probability of } j^{\text{th}} \text{ component given data } \mathbf{x}_i \rightarrow \sum_{j=1}^k \hat{\gamma}_{ij} = 1$$

Maximization step

- Update weights (in that order)

$$\hat{\phi}_j = \frac{N_j}{n} \text{ with } N_j = \sum_i \hat{\gamma}_{ij} ;$$

$$\hat{\mu}_j = \frac{1}{N_j} \sum_i \hat{\gamma}_{ij} \mathbf{x}_i ;$$



$$\hat{\Sigma}_j = \frac{1}{N_j} \sum_i \hat{\gamma}_{ij} (\mathbf{x}_i - \hat{\mu}_j)(\mathbf{x}_i - \hat{\mu}_j)^T$$

The EM algorithm

The nightmare of any EECS grad student around the world

Initialization

- Select k random points as initial means $\hat{\mu}_1, \dots, \hat{\mu}_k$
- Init all covariance matrices $\hat{\Sigma}_1, \dots, \hat{\Sigma}_k$ as whole data sample covariance matrix $\hat{\Sigma}$
- Set uniform mixture weights $\hat{\phi}_1, \dots, \hat{\phi}_k = \frac{1}{k}$

Alternate until convergence

Expectation step

- Compute membership weight $\hat{\gamma}_{ij}$ of \mathbf{x}_i with respect to j^{th} component $\mathcal{N}(\mathbf{x}|\mu_j, \Sigma_j)$

$$\hat{\gamma}_{ij} = \frac{\hat{\phi}_j \mathcal{N}(\mathbf{x}_i | \hat{\mu}_j, \hat{\Sigma}_j)}{\sum_{m=1}^k \hat{\phi}_m \mathcal{N}(\mathbf{x}_i | \hat{\mu}_m, \hat{\Sigma}_m)}$$

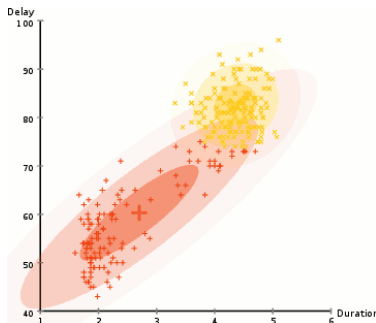
$$\hat{\gamma}_{ij} \equiv \text{posterior probability of } j^{\text{th}} \text{ component given data } \mathbf{x}_i \rightarrow \sum_{j=1}^k \hat{\gamma}_{ij} = 1$$

Maximization step

- Update weights (in that order)

$$\hat{\phi}_j = \frac{N_j}{n} \text{ with } N_j = \sum_i \hat{\gamma}_{ij} ;$$

$$\hat{\mu}_j = \frac{1}{N_j} \sum_i \hat{\gamma}_{ij} \mathbf{x}_i ;$$



The EM algorithm

The nightmare of any EECS grad student around the world

Initialization

- Select k random points as initial means $\hat{\mu}_1, \dots, \hat{\mu}_k$
- Init all covariance matrices $\hat{\Sigma}_1, \dots, \hat{\Sigma}_k$ as whole data sample covariance matrix $\hat{\Sigma}$
- Set uniform mixture weights $\hat{\phi}_1, \dots, \hat{\phi}_k = \frac{1}{k}$

Alternate until convergence

Expectation step

- Compute membership weight $\hat{\gamma}_{ij}$ of \mathbf{x}_i with respect to j^{th} component $\mathcal{N}(\mathbf{x}|\mu_j, \Sigma_j)$

$$\hat{\gamma}_{ij} = \frac{\hat{\phi}_j \mathcal{N}(\mathbf{x}_i | \hat{\mu}_j, \hat{\Sigma}_j)}{\sum_{m=1}^k \hat{\phi}_m \mathcal{N}(\mathbf{x}_i | \hat{\mu}_m, \hat{\Sigma}_m)}$$

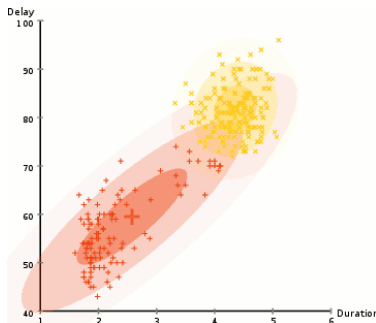
$$\hat{\gamma}_{ij} \equiv \text{posterior probability of } j^{\text{th}} \text{ component given data } \mathbf{x}_i \rightarrow \sum_{j=1}^k \hat{\gamma}_{ij} = 1$$

Maximization step

- Update weights (in that order)

$$\hat{\phi}_j = \frac{N_j}{n} \text{ with } N_j = \sum_i \hat{\gamma}_{ij} ;$$

$$\hat{\mu}_j = \frac{1}{N_j} \sum_i \hat{\gamma}_{ij} \mathbf{x}_i ;$$



The EM algorithm

The nightmare of any EECS grad student around the world

Initialization

- Select k random points as initial means $\hat{\mu}_1, \dots, \hat{\mu}_k$
- Init all covariance matrices $\hat{\Sigma}_1, \dots, \hat{\Sigma}_k$ as whole data sample covariance matrix $\hat{\Sigma}$
- Set uniform mixture weights $\hat{\phi}_1, \dots, \hat{\phi}_k = \frac{1}{k}$

Alternate until convergence

Expectation step

- Compute membership weight $\hat{\gamma}_{ij}$ of \mathbf{x}_i with respect to j^{th} component $\mathcal{N}(\mathbf{x}|\mu_j, \Sigma_j)$

$$\hat{\gamma}_{ij} = \frac{\hat{\phi}_j \mathcal{N}(\mathbf{x}_i | \hat{\mu}_j, \hat{\Sigma}_j)}{\sum_{m=1}^k \hat{\phi}_m \mathcal{N}(\mathbf{x}_i | \hat{\mu}_m, \hat{\Sigma}_m)}$$

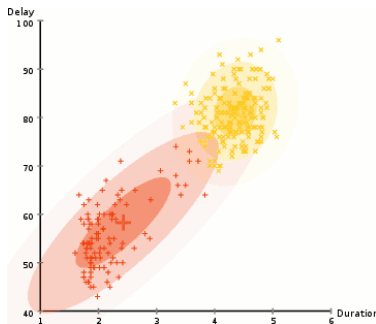
$$\hat{\gamma}_{ij} \equiv \text{posterior probability of } j^{\text{th}} \text{ component given data } \mathbf{x}_i \rightarrow \sum_{j=1}^k \hat{\gamma}_{ij} = 1$$

Maximization step

- Update weights (in that order)

$$\hat{\phi}_j = \frac{N_j}{n} \text{ with } N_j = \sum_i \hat{\gamma}_{ij} ;$$

$$\hat{\mu}_j = \frac{1}{N_j} \sum_i \hat{\gamma}_{ij} \mathbf{x}_i ;$$



The EM algorithm

The nightmare of any EECS grad student around the world

Initialization

- Select k random points as initial means $\hat{\mu}_1, \dots, \hat{\mu}_k$
- Init all covariance matrices $\hat{\Sigma}_1, \dots, \hat{\Sigma}_k$ as whole data sample covariance matrix $\hat{\Sigma}$
- Set uniform mixture weights $\hat{\phi}_1, \dots, \hat{\phi}_k = \frac{1}{k}$

Alternate until convergence

Expectation step

- Compute membership weight $\hat{\gamma}_{ij}$ of \mathbf{x}_i with respect to j^{th} component $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \Sigma_j)$

$$\hat{\gamma}_{ij} = \frac{\hat{\phi}_j \mathcal{N}(\mathbf{x}_i | \hat{\boldsymbol{\mu}}_j, \hat{\Sigma}_j)}{\sum_{m=1}^k \hat{\phi}_m \mathcal{N}(\mathbf{x}_i | \hat{\boldsymbol{\mu}}_m, \hat{\Sigma}_m)}$$

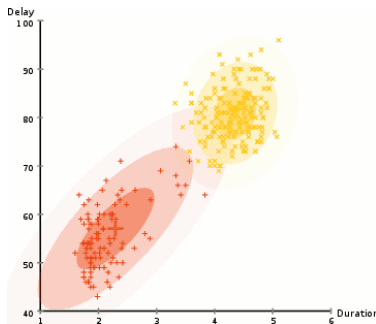
$$\hat{\gamma}_{ij} \equiv \text{posterior probability of } j^{\text{th}} \text{ component given data } \mathbf{x}_i \rightarrow \sum_{j=1}^k \hat{\gamma}_{ij} = 1$$

Maximization step

- Update weights (in that order)

$$\hat{\phi}_j = \frac{N_j}{n} \text{ with } N_j = \sum_i \hat{\gamma}_{ij} ;$$

$$\hat{\boldsymbol{\mu}}_j = \frac{1}{N_j} \sum_i \hat{\gamma}_{ij} \mathbf{x}_i ;$$



$$\hat{\Sigma}_j = \frac{1}{N_j} \sum_i \hat{\gamma}_{ij} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_j)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_j)^T$$

The EM algorithm

The nightmare of any EECS grad student around the world

Initialization

- Select k random points as initial means $\hat{\mu}_1, \dots, \hat{\mu}_k$
- Init all covariance matrices $\hat{\Sigma}_1, \dots, \hat{\Sigma}_k$ as whole data sample covariance matrix $\hat{\Sigma}$
- Set uniform mixture weights $\hat{\phi}_1, \dots, \hat{\phi}_k = \frac{1}{k}$

Alternate until convergence

Expectation step

- Compute membership weight $\hat{\gamma}_{ij}$ of \mathbf{x}_i with respect to j^{th} component $\mathcal{N}(\mathbf{x}|\mu_j, \Sigma_j)$

$$\hat{\gamma}_{ij} = \frac{\hat{\phi}_j \mathcal{N}(\mathbf{x}_i | \hat{\mu}_j, \hat{\Sigma}_j)}{\sum_{m=1}^k \hat{\phi}_m \mathcal{N}(\mathbf{x}_i | \hat{\mu}_m, \hat{\Sigma}_m)}$$

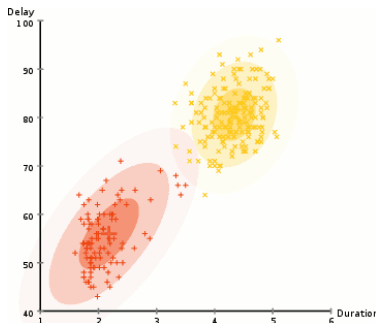
$$\hat{\gamma}_{ij} \equiv \text{posterior probability of } j^{\text{th}} \text{ component given data } \mathbf{x}_i \rightarrow \sum_{j=1}^k \hat{\gamma}_{ij} = 1$$

Maximization step

- Update weights (in that order)

$$\hat{\phi}_j = \frac{N_j}{n} \text{ with } N_j = \sum_i \hat{\gamma}_{ij} ;$$

$$\hat{\mu}_j = \frac{1}{N_j} \sum_i \hat{\gamma}_{ij} \mathbf{x}_i ;$$



$$\hat{\Sigma}_j = \frac{1}{N_j} \sum_i \hat{\gamma}_{ij} (\mathbf{x}_i - \hat{\mu}_j)(\mathbf{x}_i - \hat{\mu}_j)^T$$

The EM algorithm

The nightmare of any EECS grad student around the world

Initialization

- Select k random points as initial means $\hat{\mu}_1, \dots, \hat{\mu}_k$
- Init all covariance matrices $\hat{\Sigma}_1, \dots, \hat{\Sigma}_k$ as whole data sample covariance matrix $\hat{\Sigma}$
- Set uniform mixture weights $\hat{\phi}_1, \dots, \hat{\phi}_k = \frac{1}{k}$

Alternate until convergence

Expectation step

- Compute membership weight $\hat{\gamma}_{ij}$ of \mathbf{x}_i with respect to j^{th} component $\mathcal{N}(\mathbf{x}|\mu_j, \Sigma_j)$

$$\hat{\gamma}_{ij} = \frac{\hat{\phi}_j \mathcal{N}(\mathbf{x}_i | \hat{\mu}_j, \hat{\Sigma}_j)}{\sum_{m=1}^k \hat{\phi}_m \mathcal{N}(\mathbf{x}_i | \hat{\mu}_m, \hat{\Sigma}_m)}$$

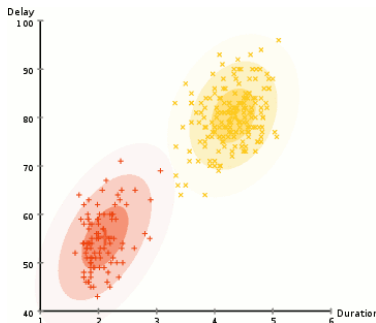
$$\hat{\gamma}_{ij} \equiv \text{posterior probability of } j^{\text{th}} \text{ component given data } \mathbf{x}_i \rightarrow \sum_{j=1}^k \hat{\gamma}_{ij} = 1$$

Maximization step

- Update weights (in that order)

$$\hat{\phi}_j = \frac{N_j}{n} \text{ with } N_j = \sum_i \hat{\gamma}_{ij} ;$$

$$\hat{\mu}_j = \frac{1}{N_j} \sum_i \hat{\gamma}_{ij} \mathbf{x}_i ;$$



The EM algorithm

The nightmare of any EECS grad student around the world

Initialization

- Select k random points as initial means $\hat{\mu}_1, \dots, \hat{\mu}_k$
- Init all covariance matrices $\hat{\Sigma}_1, \dots, \hat{\Sigma}_k$ as whole data sample covariance matrix $\hat{\Sigma}$
- Set uniform mixture weights $\hat{\phi}_1, \dots, \hat{\phi}_k = \frac{1}{k}$

Alternate until convergence

Expectation step

- Compute membership weight $\hat{\gamma}_{ij}$ of \mathbf{x}_i with respect to j^{th} component $\mathcal{N}(\mathbf{x}|\mu_j, \Sigma_j)$

$$\hat{\gamma}_{ij} = \frac{\hat{\phi}_j \mathcal{N}(\mathbf{x}_i | \hat{\mu}_j, \hat{\Sigma}_j)}{\sum_{m=1}^k \hat{\phi}_m \mathcal{N}(\mathbf{x}_i | \hat{\mu}_m, \hat{\Sigma}_m)}$$

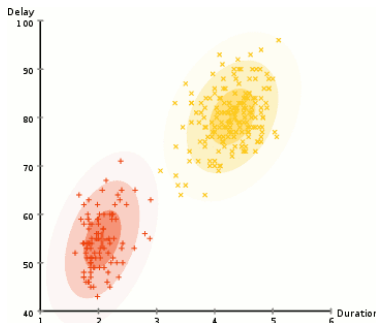
$$\hat{\gamma}_{ij} \equiv \text{posterior probability of } j^{\text{th}} \text{ component given data } \mathbf{x}_i \rightarrow \sum_{j=1}^k \hat{\gamma}_{ij} = 1$$

Maximization step

- Update weights (in that order)

$$\hat{\phi}_j = \frac{N_j}{n} \text{ with } N_j = \sum_i \hat{\gamma}_{ij} ;$$

$$\hat{\mu}_j = \frac{1}{N_j} \sum_i \hat{\gamma}_{ij} \mathbf{x}_i ;$$



$$\hat{\Sigma}_j = \frac{1}{N_j} \sum_i \hat{\gamma}_{ij} (\mathbf{x}_i - \hat{\mu}_j)(\mathbf{x}_i - \hat{\mu}_j)^T$$

The EM algorithm

The nightmare of any EECS grad student around the world

Initialization

- Select k random points as initial means $\hat{\mu}_1, \dots, \hat{\mu}_k$
- Init all covariance matrices $\hat{\Sigma}_1, \dots, \hat{\Sigma}_k$ as whole data sample covariance matrix $\hat{\Sigma}$
- Set uniform mixture weights $\hat{\phi}_1, \dots, \hat{\phi}_k = \frac{1}{k}$

Alternate until convergence

Expectation step

- Compute membership weight $\hat{\gamma}_{ij}$ of \mathbf{x}_i with respect to j^{th} component $\mathcal{N}(\mathbf{x}|\mu_j, \Sigma_j)$

$$\hat{\gamma}_{ij} = \frac{\hat{\phi}_j \mathcal{N}(\mathbf{x}_i | \hat{\mu}_j, \hat{\Sigma}_j)}{\sum_{m=1}^k \hat{\phi}_m \mathcal{N}(\mathbf{x}_i | \hat{\mu}_m, \hat{\Sigma}_m)}$$

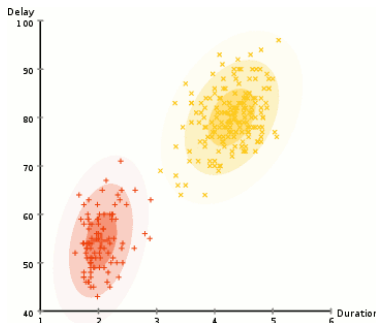
$$\hat{\gamma}_{ij} \equiv \text{posterior probability of } j^{\text{th}} \text{ component given data } \mathbf{x}_i \rightarrow \sum_{j=1}^k \hat{\gamma}_{ij} = 1$$

Maximization step

- Update weights (in that order)

$$\hat{\phi}_j = \frac{N_j}{n} \text{ with } N_j = \sum_i \hat{\gamma}_{ij} ;$$

$$\hat{\mu}_j = \frac{1}{N_j} \sum_i \hat{\gamma}_{ij} \mathbf{x}_i ;$$



The EM algorithm

The nightmare of any EECS grad student around the world

Initialization

- Select k random points as initial means $\hat{\mu}_1, \dots, \hat{\mu}_k$
- Init all covariance matrices $\hat{\Sigma}_1, \dots, \hat{\Sigma}_k$ as whole data sample covariance matrix $\hat{\Sigma}$
- Set uniform mixture weights $\hat{\phi}_1, \dots, \hat{\phi}_k = \frac{1}{k}$

Alternate until convergence

Expectation step

- Compute membership weight $\hat{\gamma}_{ij}$ of \mathbf{x}_i with respect to j^{th} component $\mathcal{N}(\mathbf{x}|\mu_j, \Sigma_j)$

$$\hat{\gamma}_{ij} = \frac{\hat{\phi}_j \mathcal{N}(\mathbf{x}_i | \hat{\mu}_j, \hat{\Sigma}_j)}{\sum_{m=1}^k \hat{\phi}_m \mathcal{N}(\mathbf{x}_i | \hat{\mu}_m, \hat{\Sigma}_m)}$$

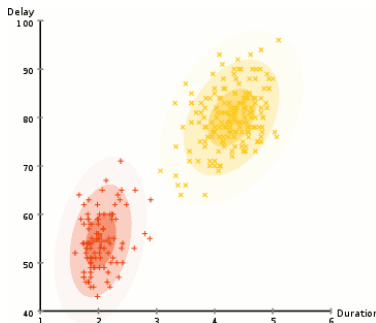
$$\hat{\gamma}_{ij} \equiv \text{posterior probability of } j^{\text{th}} \text{ component given data } \mathbf{x}_i \rightarrow \sum_{j=1}^k \hat{\gamma}_{ij} = 1$$

Maximization step

- Update weights (in that order)

$$\hat{\phi}_j = \frac{N_j}{n} \text{ with } N_j = \sum_i \hat{\gamma}_{ij} ;$$

$$\hat{\mu}_j = \frac{1}{N_j} \sum_i \hat{\gamma}_{ij} \mathbf{x}_i ;$$



$$\hat{\Sigma}_j = \frac{1}{N_j} \sum_i \hat{\gamma}_{ij} (\mathbf{x}_i - \hat{\mu}_j)(\mathbf{x}_i - \hat{\mu}_j)^T$$

The EM algorithm

The nightmare of any EECS grad student around the world

Initialization

- Select k random points as initial means $\hat{\mu}_1, \dots, \hat{\mu}_k$
- Init all covariance matrices $\hat{\Sigma}_1, \dots, \hat{\Sigma}_k$ as whole data sample covariance matrix $\hat{\Sigma}$
- Set uniform mixture weights $\hat{\phi}_1, \dots, \hat{\phi}_k = \frac{1}{k}$

Alternate until convergence

Expectation step

- Compute membership weight $\hat{\gamma}_{ij}$ of \mathbf{x}_i with respect to j^{th} component $\mathcal{N}(\mathbf{x}|\mu_j, \Sigma_j)$

$$\hat{\gamma}_{ij} = \frac{\hat{\phi}_j \mathcal{N}(\mathbf{x}_i | \hat{\mu}_j, \hat{\Sigma}_j)}{\sum_{m=1}^k \hat{\phi}_m \mathcal{N}(\mathbf{x}_i | \hat{\mu}_m, \hat{\Sigma}_m)}$$

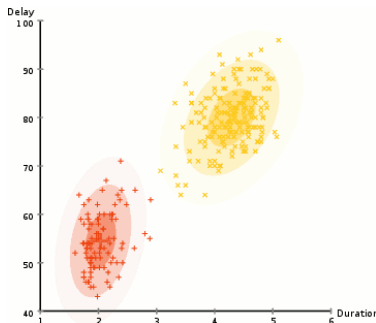
$$\hat{\gamma}_{ij} \equiv \text{posterior probability of } j^{\text{th}} \text{ component given data } \mathbf{x}_i \rightarrow \sum_{j=1}^k \hat{\gamma}_{ij} = 1$$

Maximization step

- Update weights (in that order)

$$\hat{\phi}_j = \frac{N_j}{n} \text{ with } N_j = \sum_i \hat{\gamma}_{ij} ;$$

$$\hat{\mu}_j = \frac{1}{N_j} \sum_i \hat{\gamma}_{ij} \mathbf{x}_i ;$$

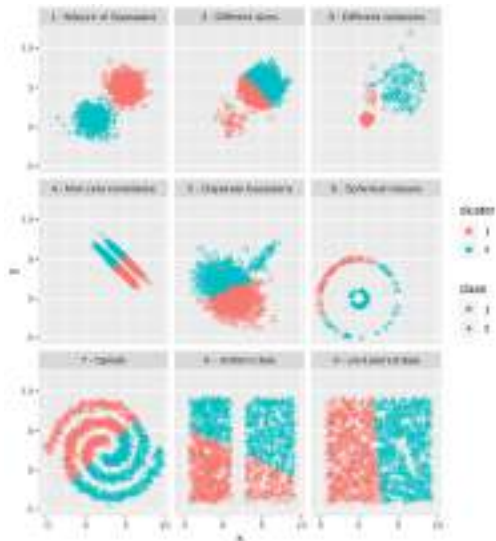


$$\hat{\Sigma}_j = \frac{1}{N_j} \sum_i \hat{\gamma}_{ij} (\mathbf{x}_i - \hat{\mu}_j)(\mathbf{x}_i - \hat{\mu}_j)^T$$

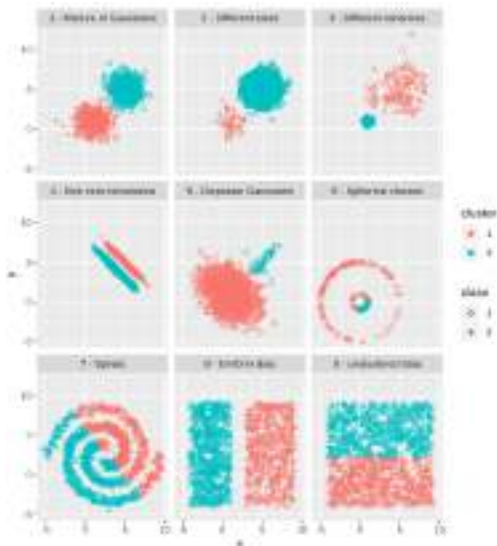
k-means vs GMM

Let the fight begin

k-means



GMM

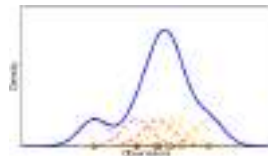
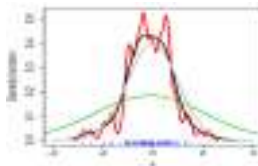
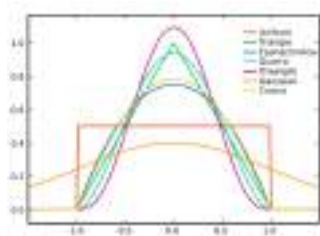


Worth the pain, right?

Kernel density estimation

Nonparametric estimation

Goal: Estimate probability density function f based on observation $x_1 \dots, x_n$ only, assumed to derive from f (otherwise wtf are we doing here?)



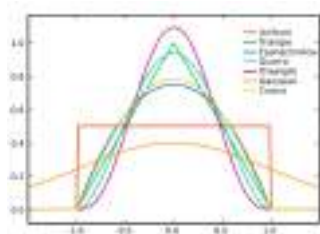
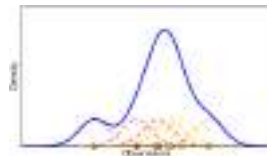
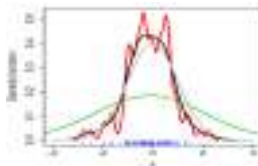
The kernel density estimator with bandwidth h at a given point x is given by

$$\hat{f}_h(x) = \frac{1}{n} \sum_{i=1}^n K_h(x - x_i) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right),$$

Kernel density estimation

Nonparametric estimation

Goal: Estimate probability density function f based on observation $x_1 \dots, x_n$ only, assumed to derive from f (otherwise wtf are we doing here?)



The kernel density estimator with bandwidth h at a given point x is given by

$$\hat{f}_h(x) = \frac{1}{n} \sum_{i=1}^n K_h(x - x_i) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right),$$

