

# Reinforced EM Algorithm through Clever Initialization for Clustering with Gaussian Mixture Models



Joshua Tobin  
TCD

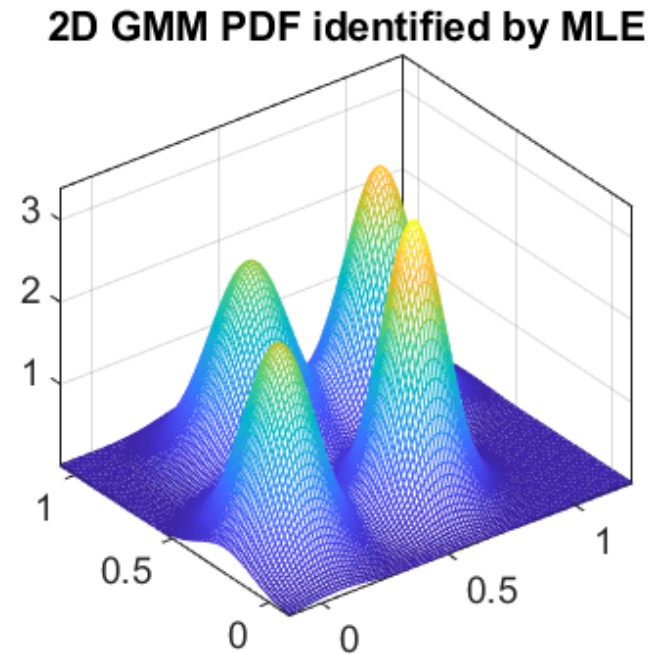
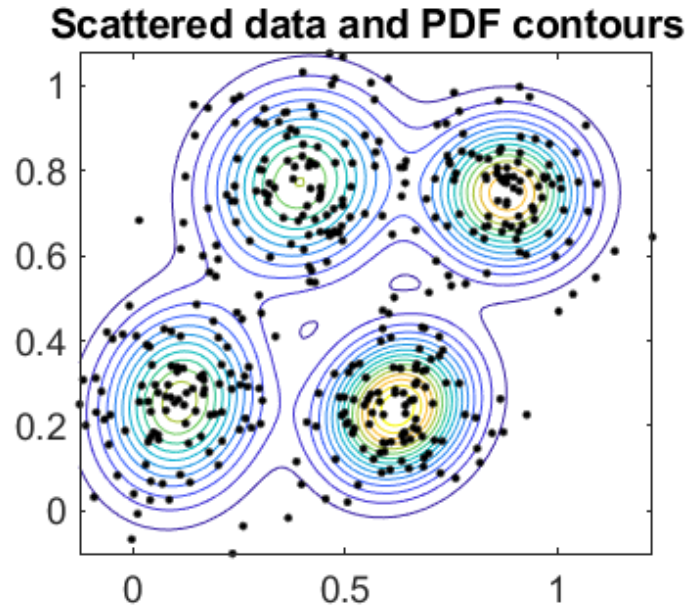


Chin Pang Ho  
CityU HK



Mimi Zhang  
TCD

# GMM for Clustering



A GMM density has the form

$$f(\mathbf{x}) = \sum_{j=1}^m \pi_j \phi(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j).$$

Clustering is done by assigning each  $\mathbf{x}_i$  to the mixture component (i.e., cluster) to which it is most likely to belong a posteriori.

# EM Algorithm

1. Initialize the parameters:  $\{\pi_1, \dots, \pi_m\}$ ,  $\{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_m\}$  and  $\{\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_m\}$ .
2. Compute the responsibilities: for  $i = 1, \dots, n$  and  $j = 1, \dots, m$ ,

$$r_{ij} = \frac{\pi_j \phi(\mathbf{x}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{v=1}^m \pi_v \phi(\mathbf{x}_i; \boldsymbol{\mu}_v, \boldsymbol{\Sigma}_v)}.$$

3. Update the estimates: for  $j = 1, \dots, m$ ,

$$\pi_j = \frac{\sum_{i=1}^n r_{ij}}{n}, \quad \boldsymbol{\mu}_j = \frac{\sum_{i=1}^n r_{ij} \mathbf{x}_i}{\sum_{i=1}^n r_{ij}}, \quad \boldsymbol{\Sigma}_j = \frac{\sum_{i=1}^n r_{ij} (\mathbf{x}_i - \boldsymbol{\mu}_j)(\mathbf{x}_i - \boldsymbol{\mu}_j)^T}{\sum_{i=1}^n r_{ij}}.$$

4. Iterate steps 2 and 3 until convergence.

# EM Algorithm

With random initialization, converge to bad local maxima with probability  $1 - e^{-\mathcal{O}(m)}$ .

1. **Initialize** the parameters:  $\{\pi_1, \dots, \pi_m\}$ ,  $\{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_m\}$  and  $\{\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_m\}$ .
2. Compute the responsibilities: for  $i = 1, \dots, n$  and  $j = 1, \dots, m$ ,

$$r_{ij} = \frac{\pi_j \phi(\mathbf{x}_i; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{v=1}^m \pi_v \phi(\mathbf{x}_i; \boldsymbol{\mu}_v, \boldsymbol{\Sigma}_v)}.$$

3. Update the estimates: for  $j = 1, \dots, m$ ,

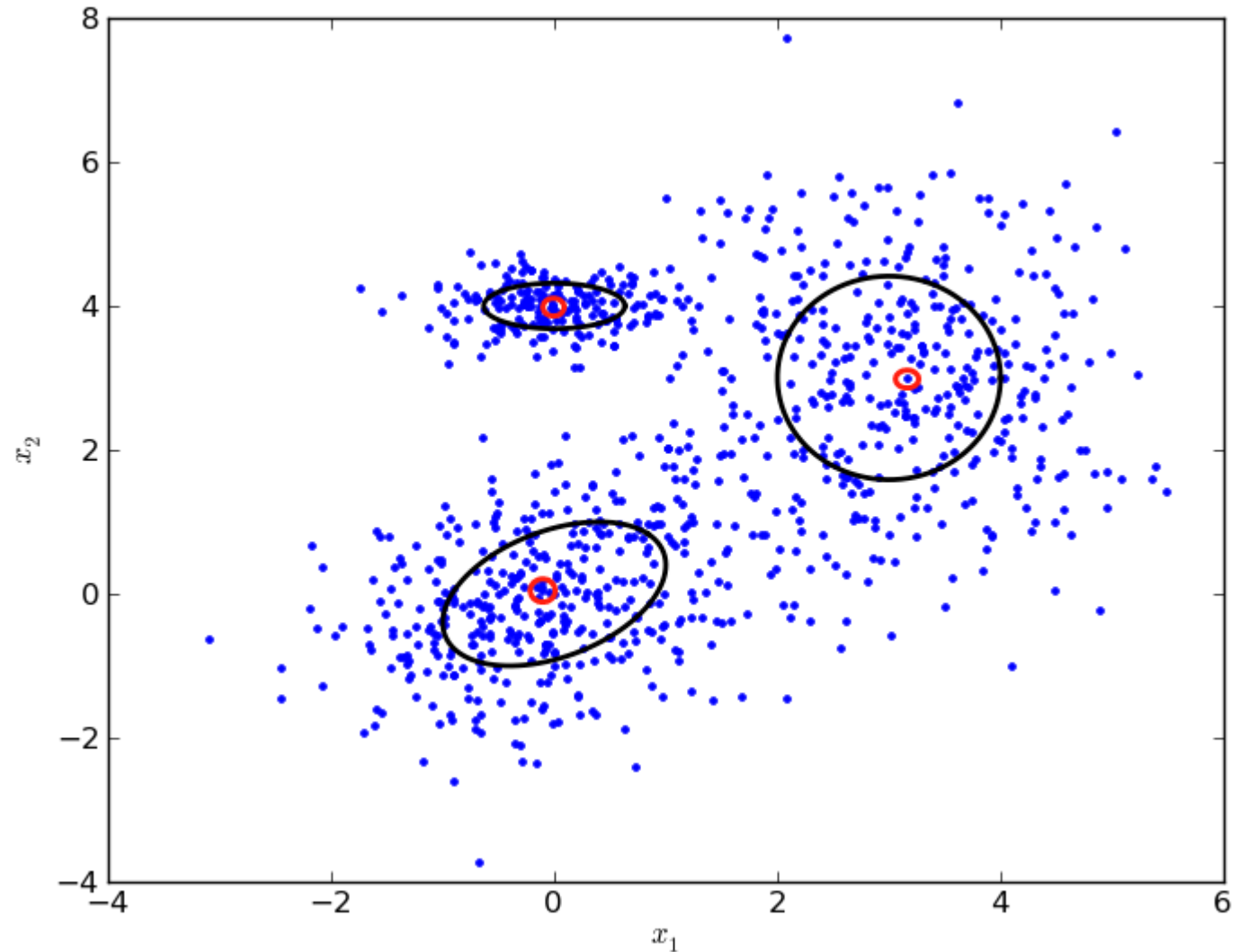
$$\pi_j = \frac{\sum_{i=1}^n r_{ij}}{n}, \quad \boldsymbol{\mu}_j = \frac{\sum_{i=1}^n r_{ij} \mathbf{x}_i}{\sum_{i=1}^n r_{ij}}, \quad \boldsymbol{\Sigma}_j = \frac{\sum_{i=1}^n r_{ij} (\mathbf{x}_i - \boldsymbol{\mu}_j)(\mathbf{x}_i - \boldsymbol{\mu}_j)^T}{\sum_{i=1}^n r_{ij}}.$$

4. Iterate steps 2 and 3 until **convergence**.

May converge to a singularity.

# Exemplar Means

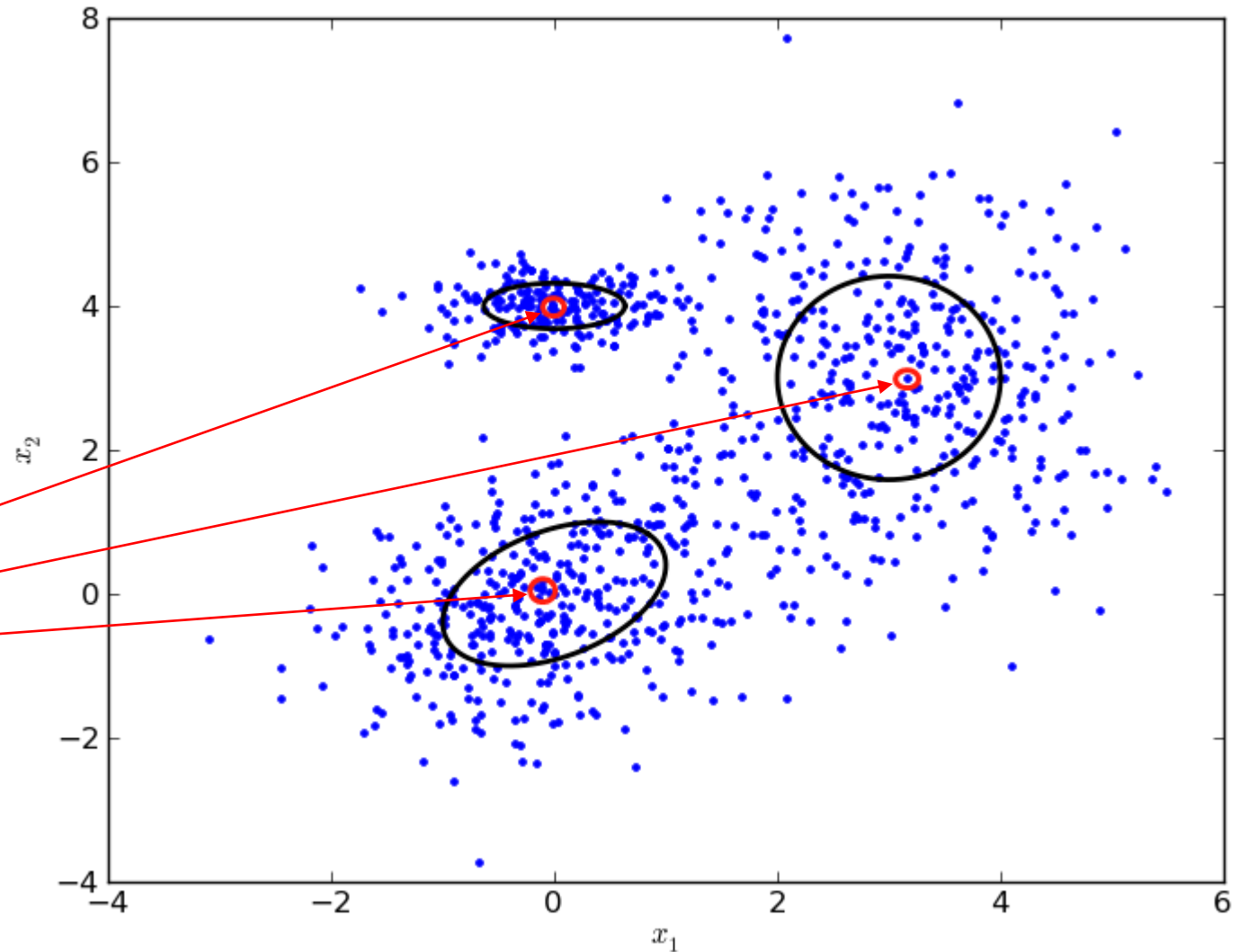
Assume that the clusters are dense enough, such that there is always a data point very close to the real cluster centre.



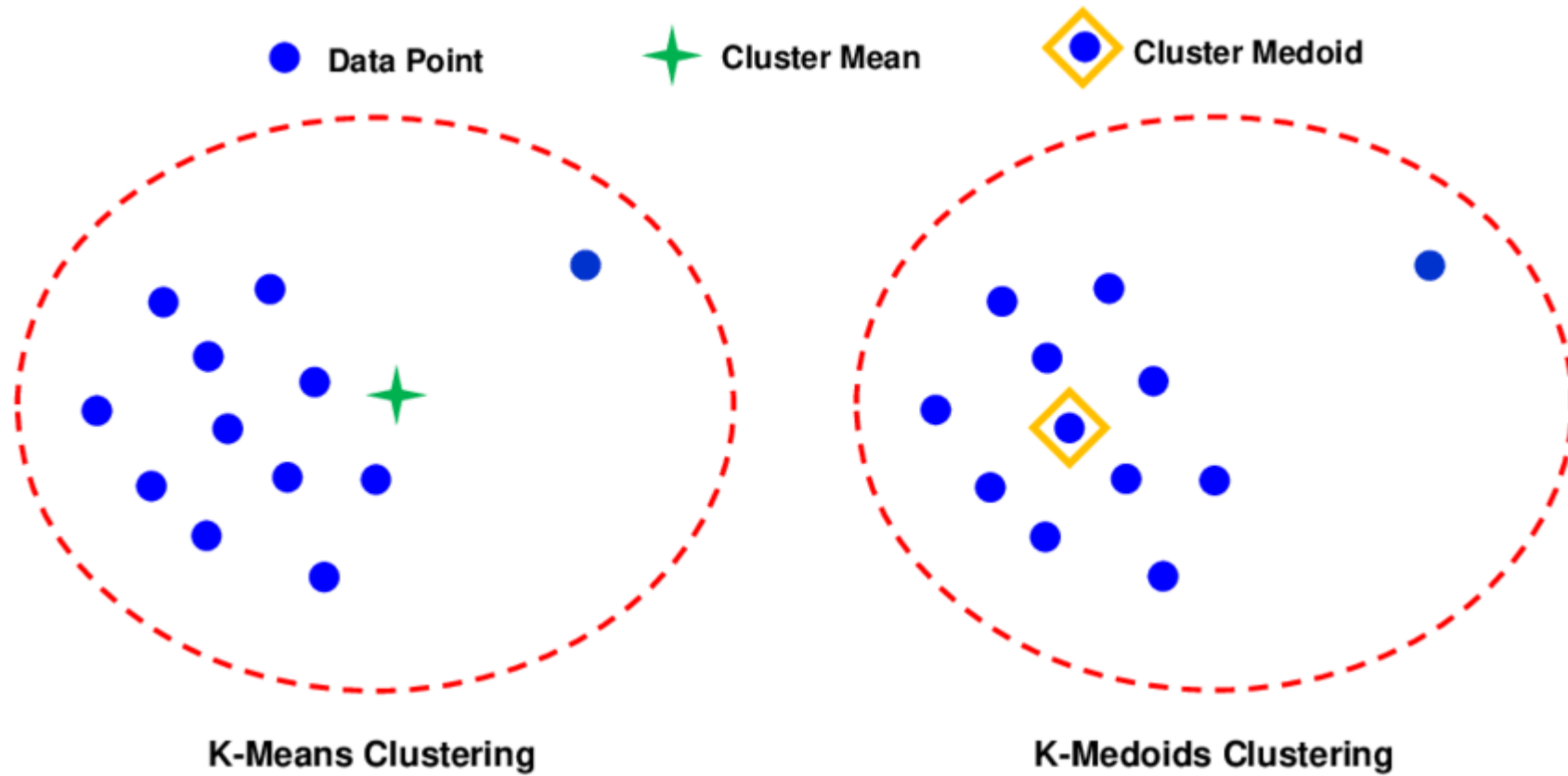
# Exemplar Means

Assume that the clusters are dense enough, such that there is always a data point very close to the real cluster centre.

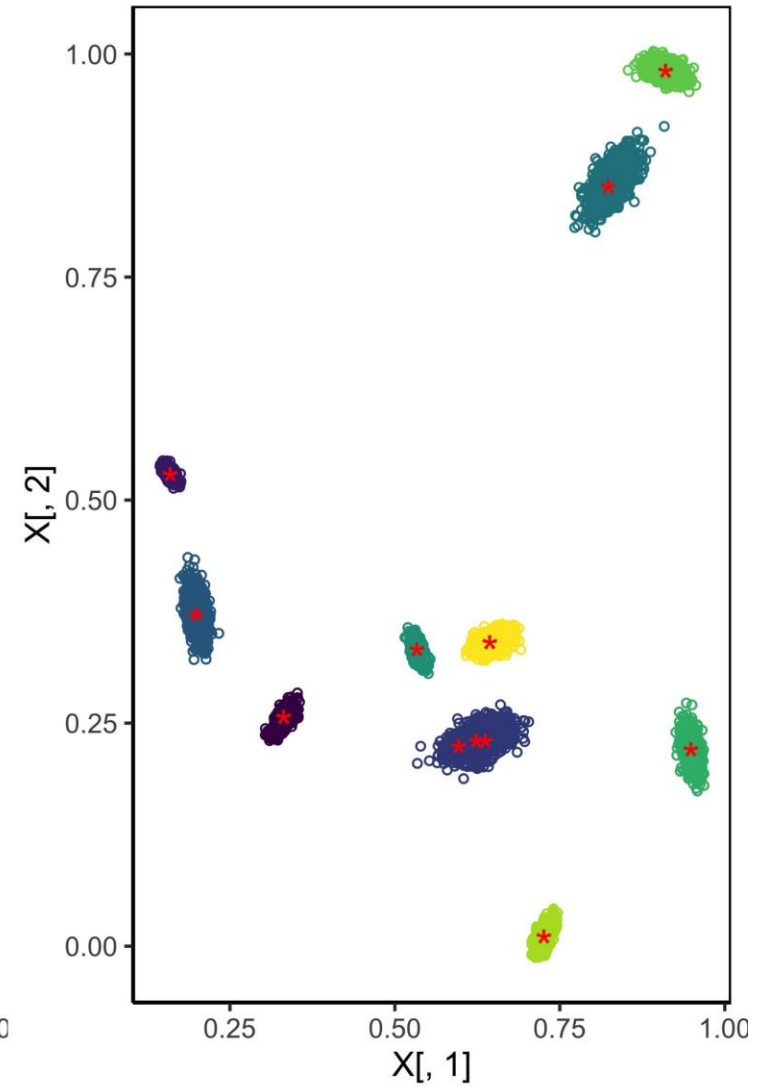
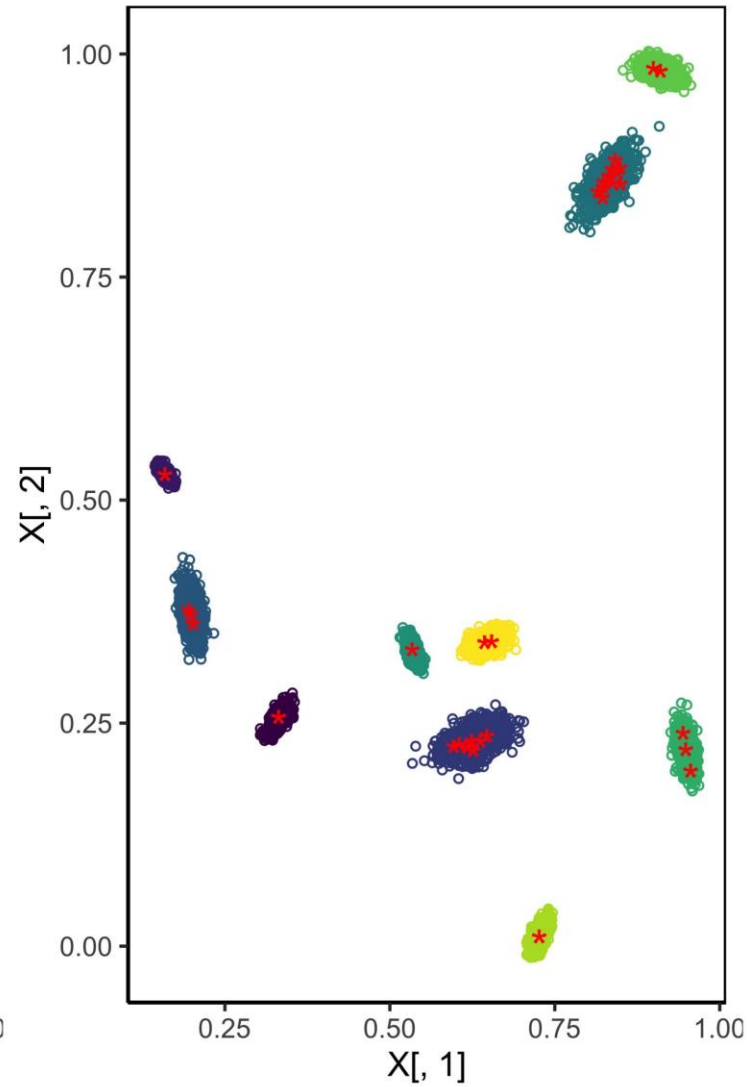
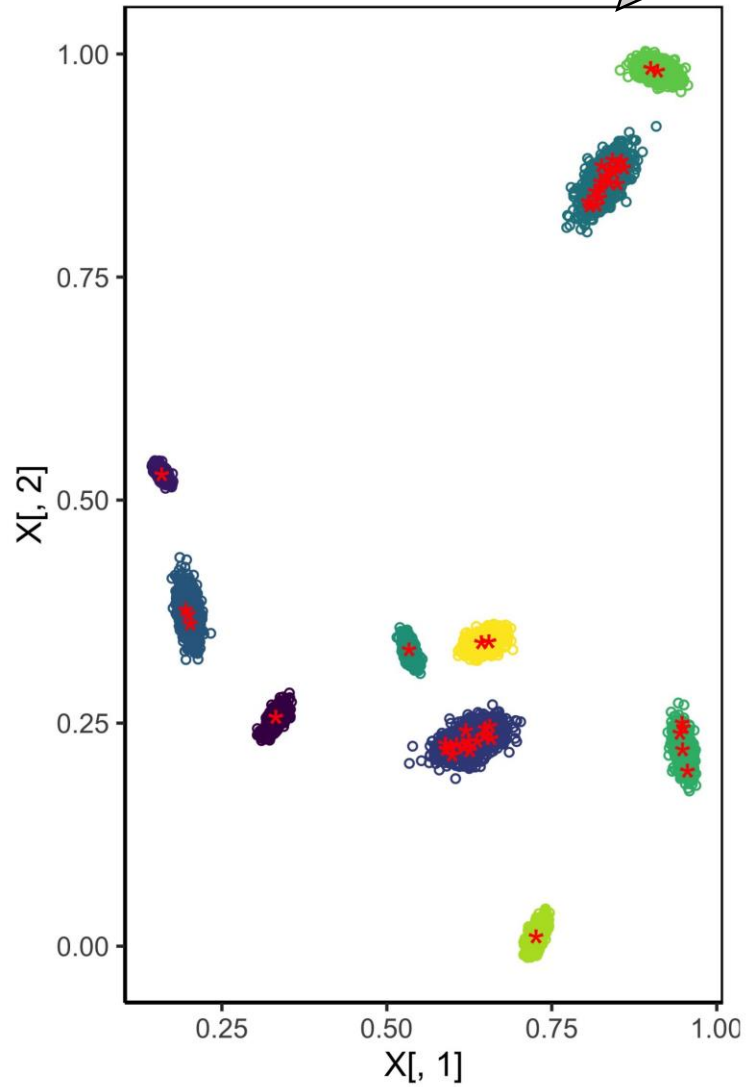
**Fix** Gaussian means at “exemplars” in the dataset.



# Exemplar Means

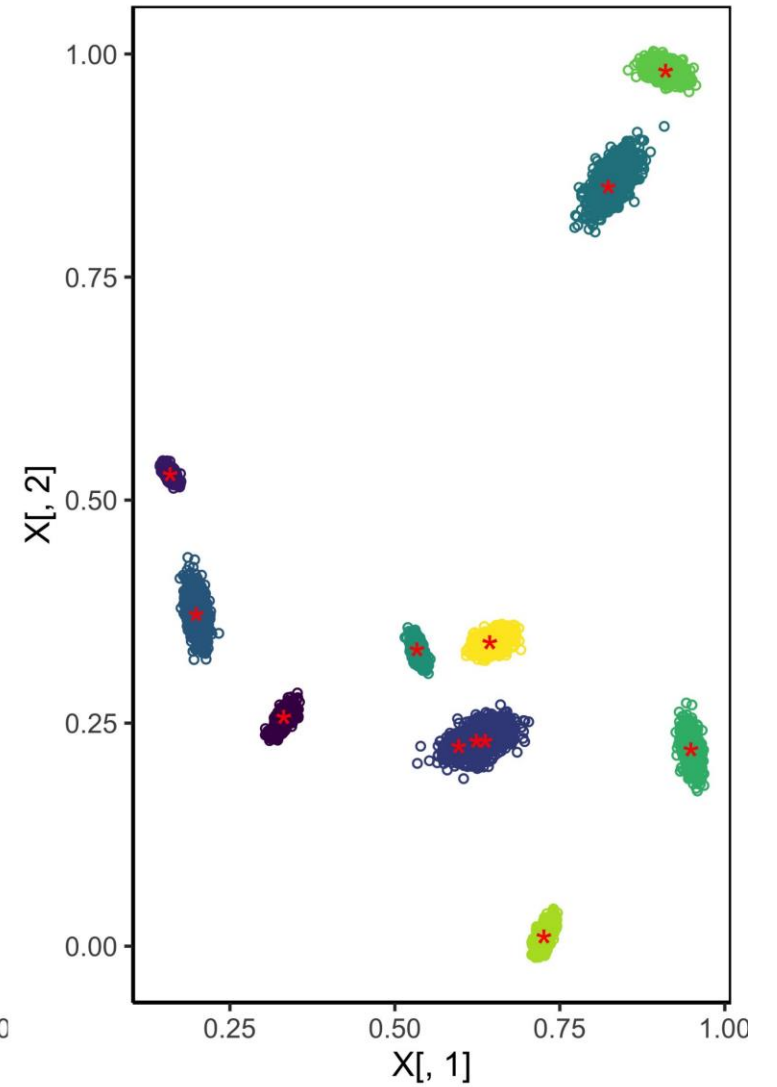
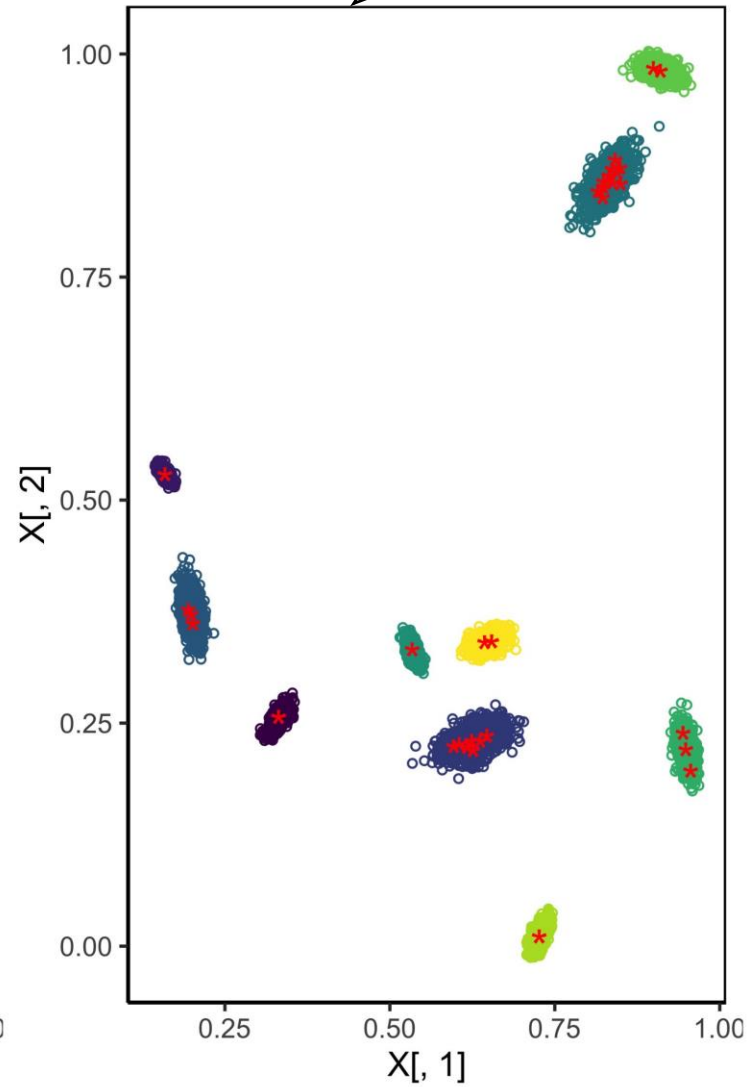
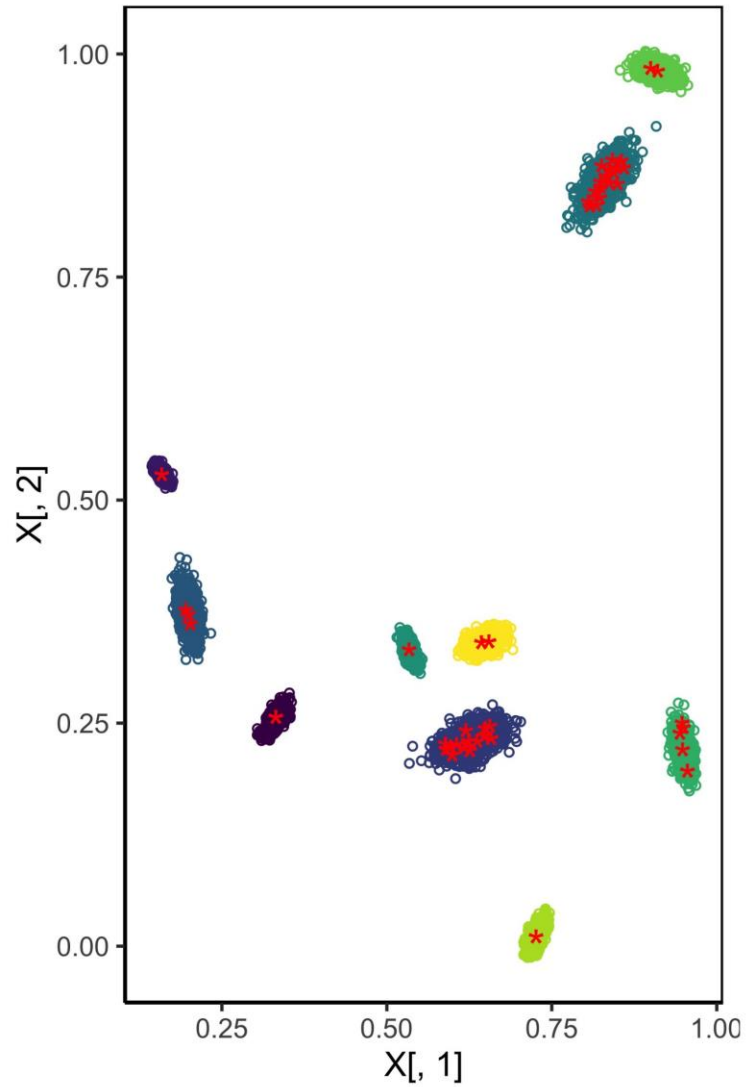


**Initialize** the EM algorithm with an inclusive set of exemplars (i.e., density peaks), identified from the data.

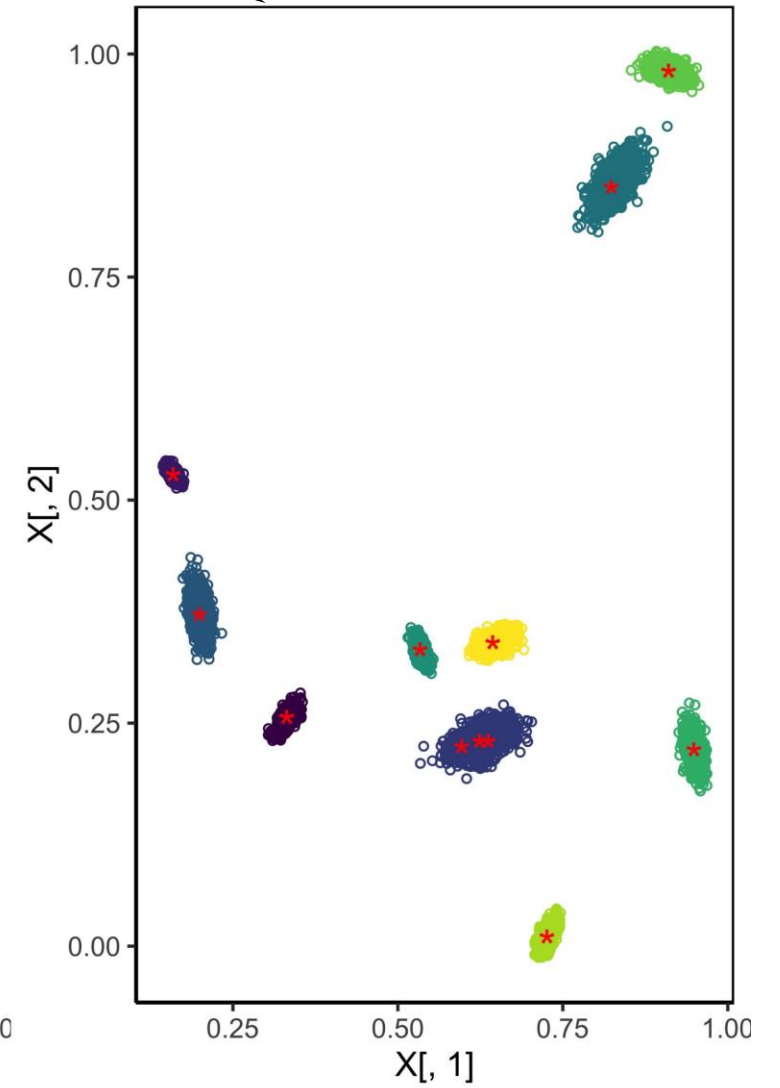
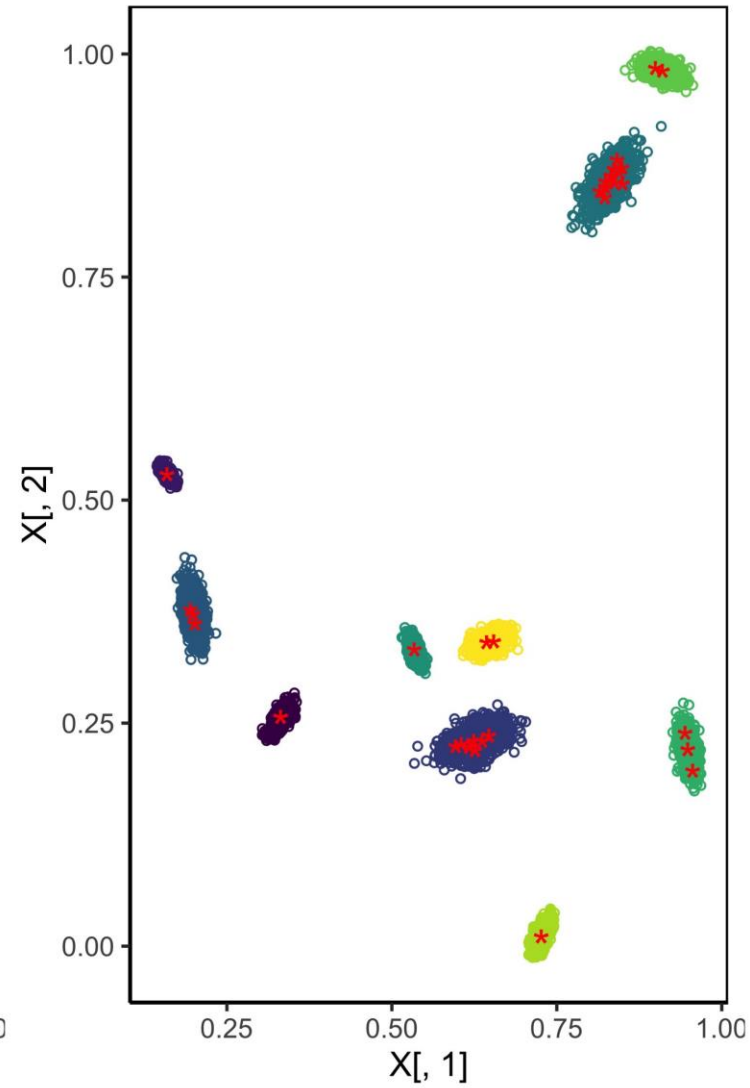
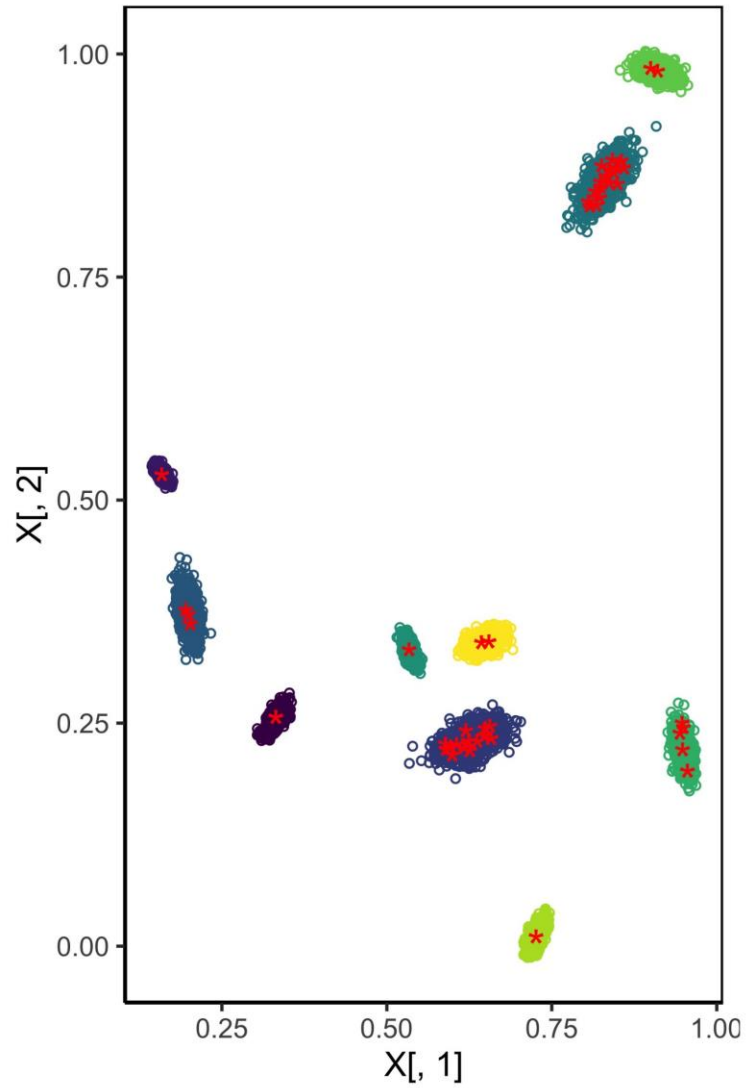


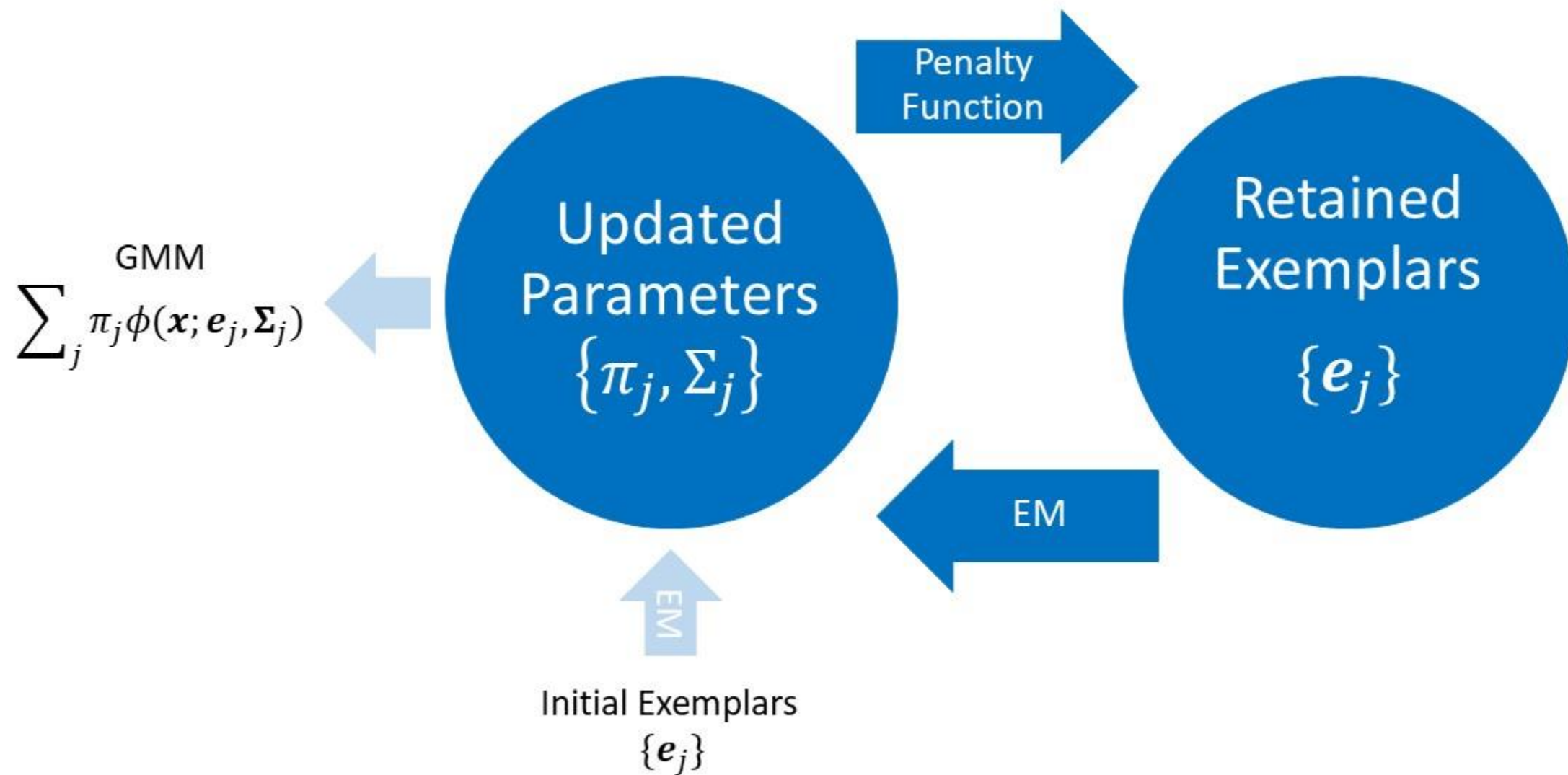


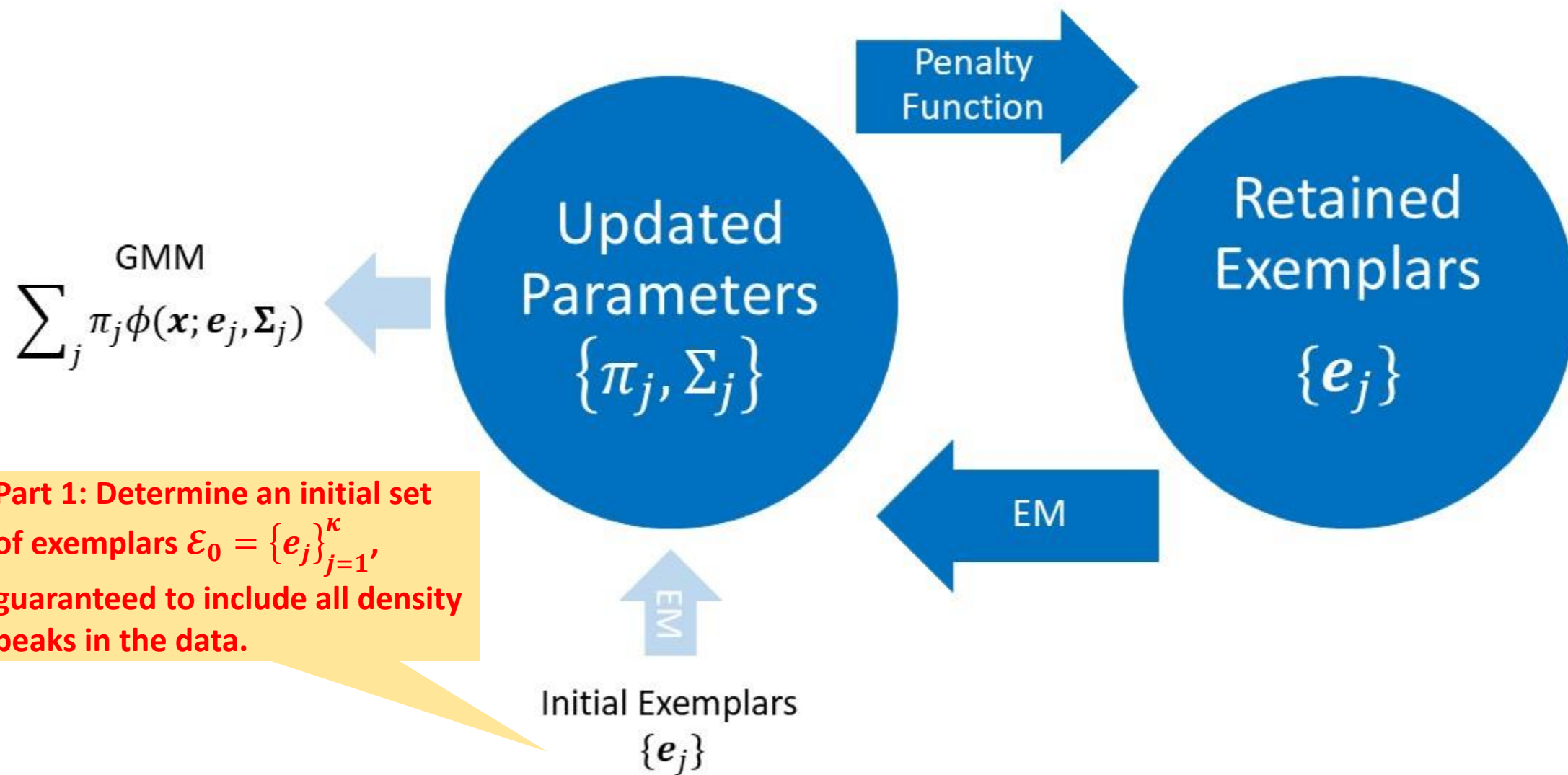
**Estimate** the parameters  $\{\pi_j, \Sigma_j\}$  using the EM algorithm, and then **prune** spurious exemplars.



**Iterate** the two steps (parameter estimation & pruning) until highest AIC or BIC.







Part 2: Prune redundant exemplars.

Penalty  
Function

GMM

$$\sum_j \pi_j \phi(\mathbf{x}; \mathbf{e}_j, \Sigma_j)$$

Updated  
Parameters  
 $\{\pi_j, \Sigma_j\}$

Retained  
Exemplars  
 $\{\mathbf{e}_j\}$

Part 1: Determine an initial set  
of exemplars  $\mathcal{E}_0 = \{\mathbf{e}_j\}_{j=1}^K$ ,  
guaranteed to include all density  
peaks in the data.

EM

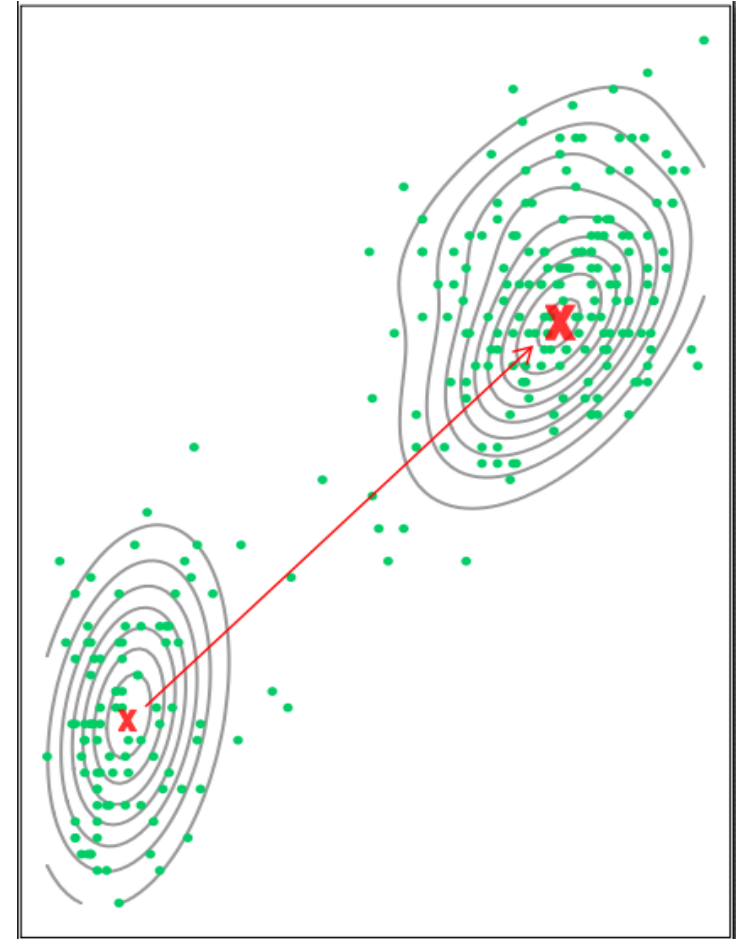
Initial Exemplars  
 $\{\mathbf{e}_j\}$

EM

# Part 1 – Density-Peak Finding

Density peaks are characterized by:

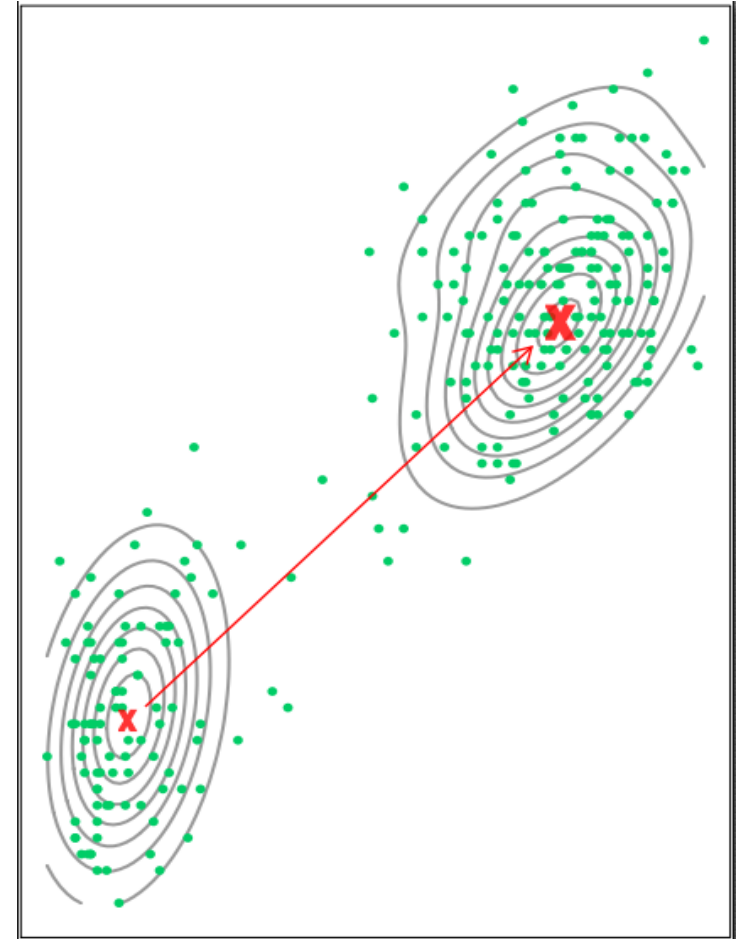
- (1) a higher density than their neighbours
- (2) a relatively large distance from points with higher densities



# Part 1 – Density-Peak Finding

Density peaks are characterized by:

- (1) a higher **density** than their neighbours
- (2) a relatively large **distance** from points with higher densities



# Part 1 – Density-Peak Finding

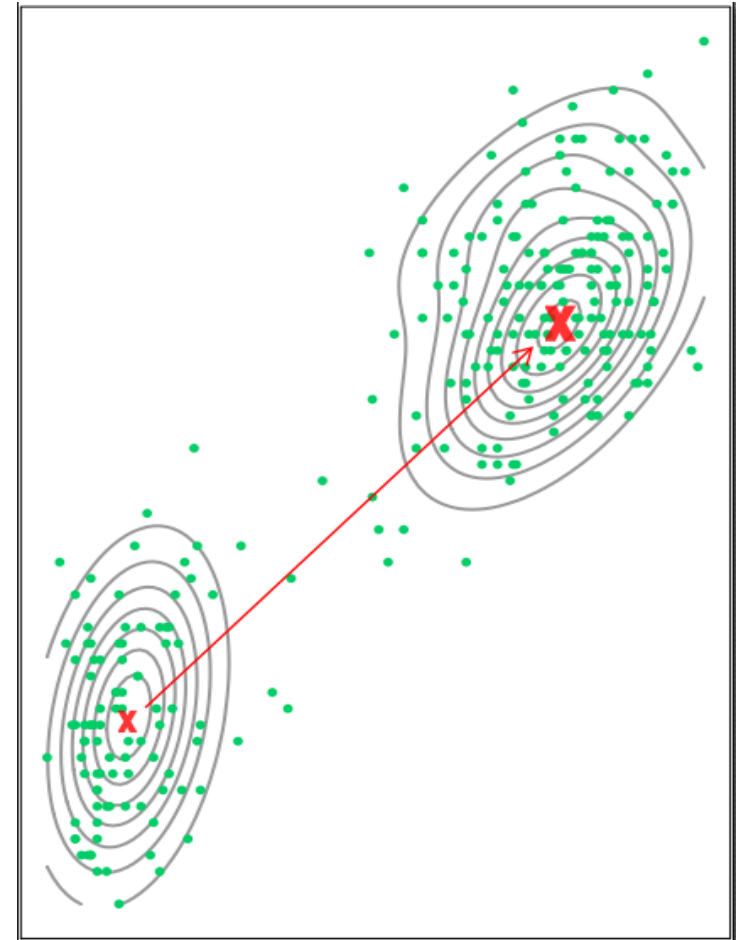
Density peaks are characterized by:

(1) a higher **density** than their neighbours

$\rho(x)$  – Gaussian kernel density estimate

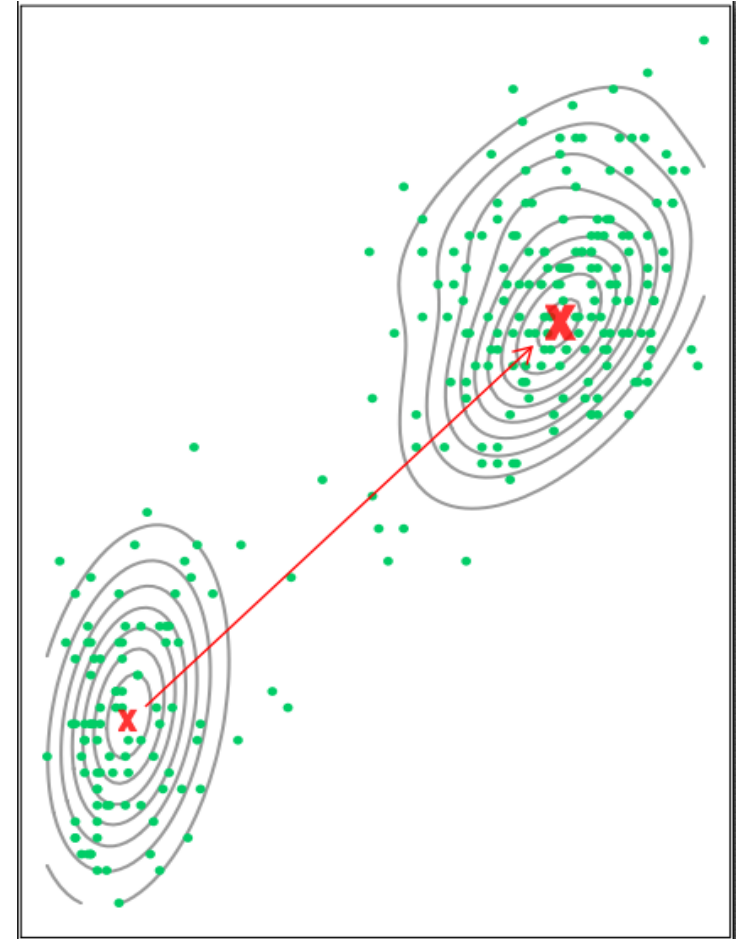
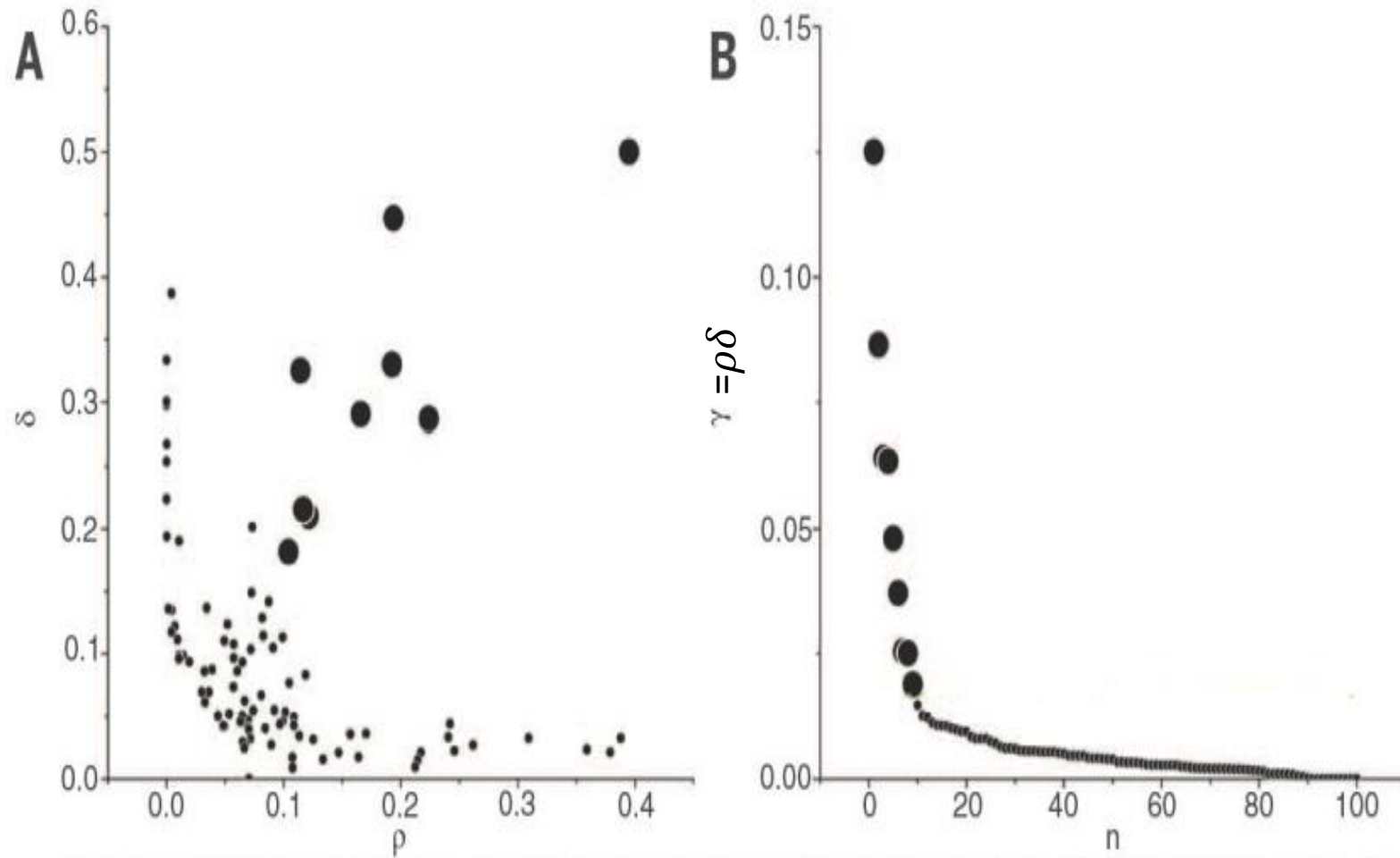
(2) a relatively large **distance** from points with higher densities

$\delta(x)$  -- distance to the nearest neighbour of higher local density

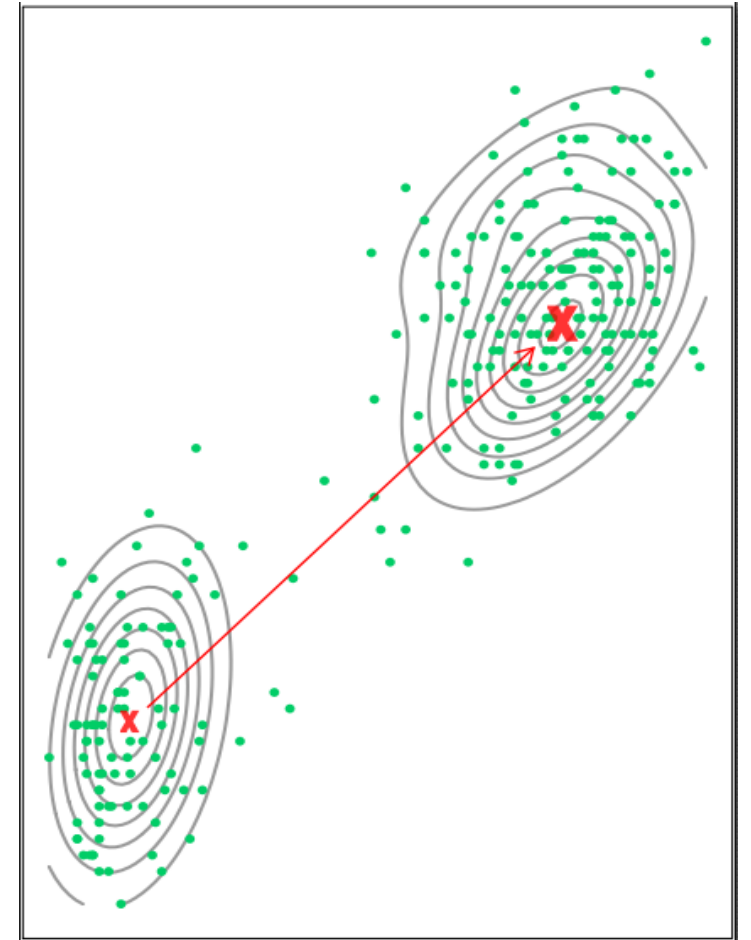
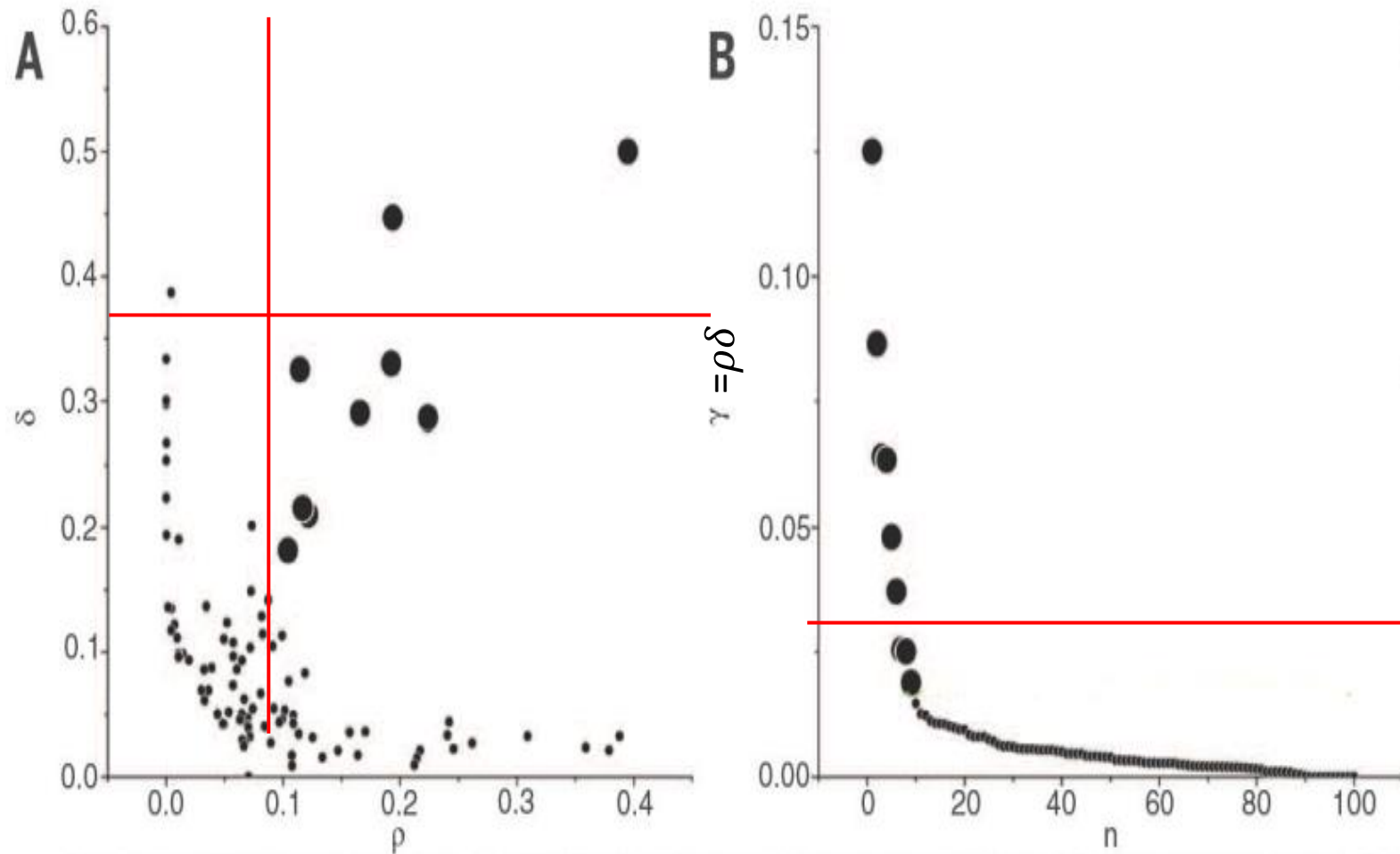




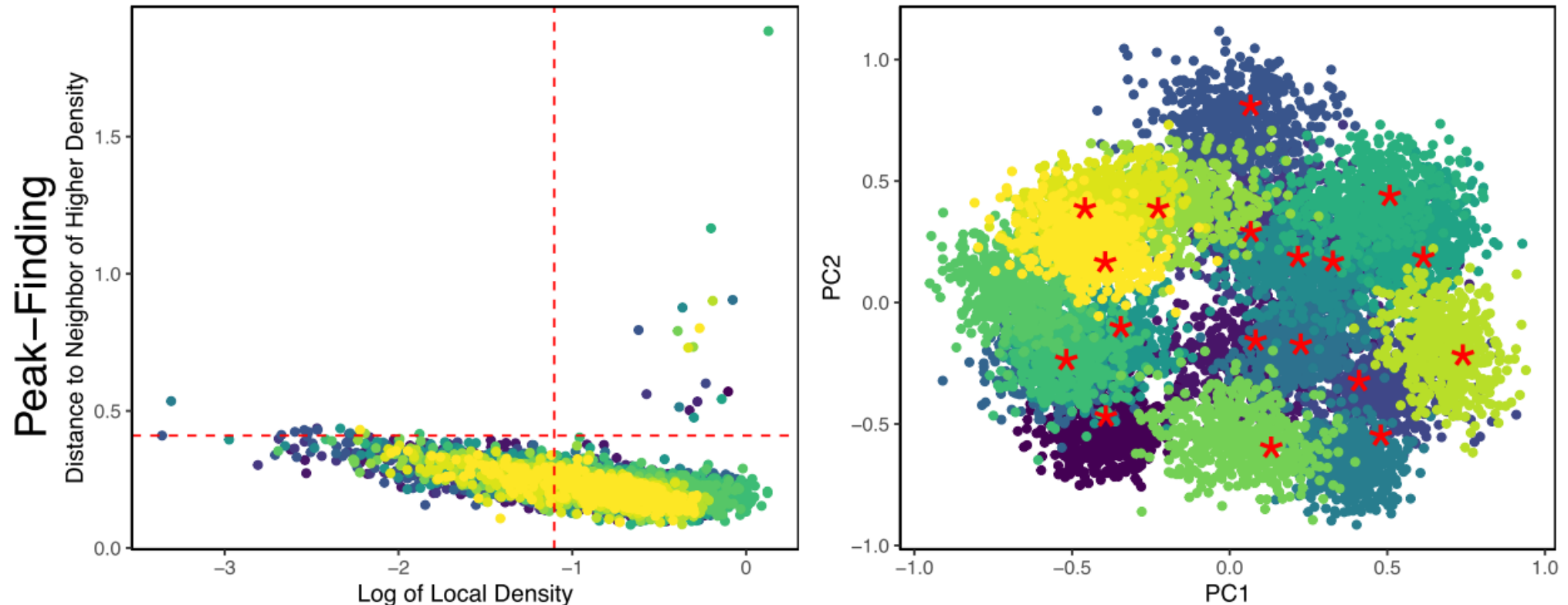
# Part 1 – Density-Peak Finding



# Part 1 – Density-Peak Finding



# Part 1 – Density-Peak Finding



A 10-dimensional dataset contains 20 Gaussian components. Clusters are indicated by different colours. The right figure shows the locations of the selected peaks, projected onto the first two principal components.

# Part 1 – Density-Peak Finding

**Theorem:** For  $n$  large enough, with high probability,  $\mathcal{E}_0 = \{\mathbf{e}_j\}_{j=1}^k$  contains unique estimates for all the true modes of the GMM.

# Part 1 – Density-Peak Finding

**Theorem:** For  $n$  large enough, with high probability,  $\mathcal{E}_0 = \{\mathbf{e}_j\}_{j=1}^k$  contains unique estimates for all the true modes of the GMM.

Relaxing the cut-off levels on the density  $\rho(\mathbf{x})$  and distance  $\delta(\mathbf{x})$ ,  $\mathcal{E}_0$  is guaranteed to include all density peaks in the data.

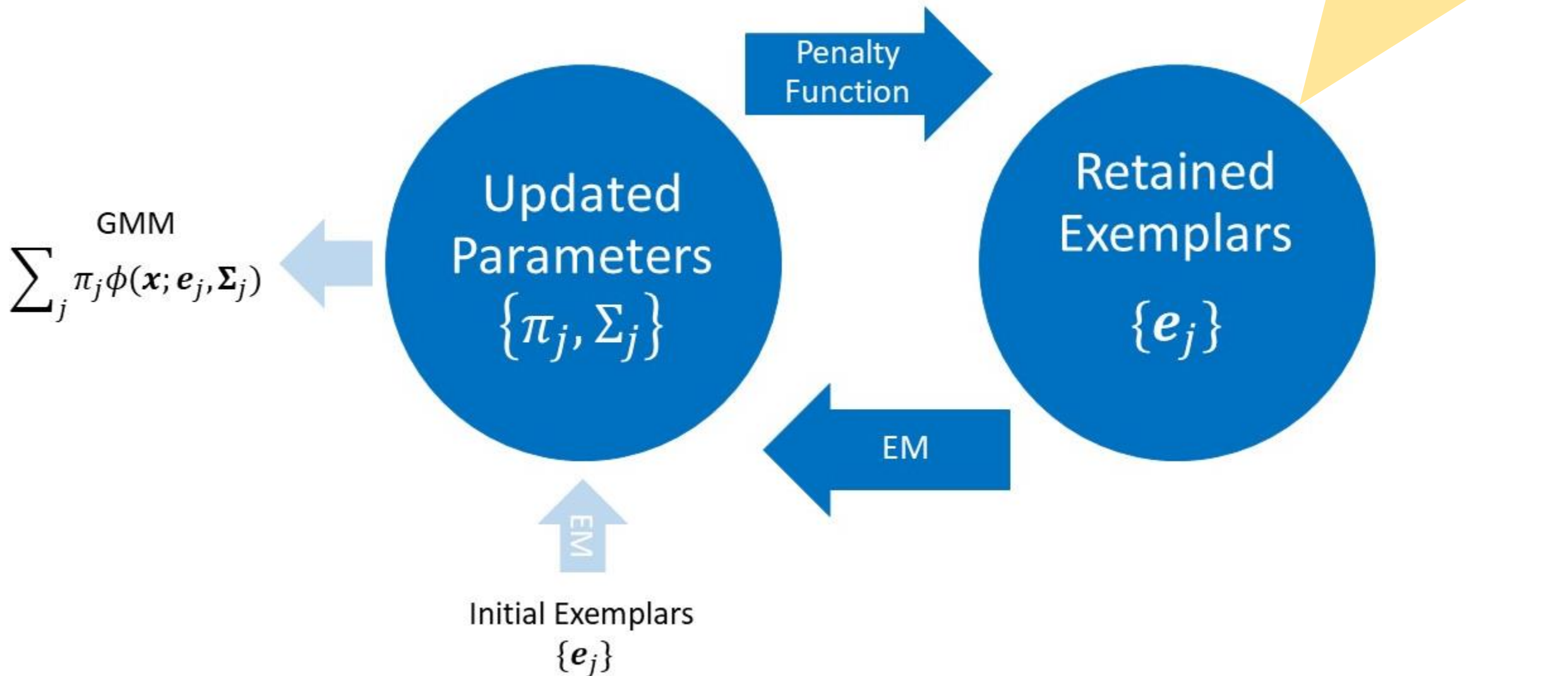
# Part 1 – Density-Peak Finding

**Theorem:** For  $n$  large enough, with high probability,  $\mathcal{E}_0 = \{\mathbf{e}_j\}_{j=1}^k$  contains unique estimates for all the true modes of the GMM.

Relaxing the cut-off levels on the density  $\rho(\mathbf{x})$  and distance  $\delta(\mathbf{x})$ ,  $\mathcal{E}_0$  is guaranteed to include all density peaks in the data.

Apply a pruning strategy to retain only instances that well represent their associated Gaussian components.

## Part 2 – Exemplar Pruning



## Part 2 – Exemplar Pruning

Given  $\mathcal{E}_0 = \{\mathbf{e}_j\}_{j=1}^{\kappa}$ , the log-likelihood function is

$$\sum_{i=1}^n \log \left( \sum_{j=1}^{\kappa} \pi_j \phi(\mathbf{x}_i; \mathbf{e}_j, \mathbf{\Sigma}_j) \right).$$

Introduce sparsity into  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_{\kappa})$ ; if  $\pi_j = 0$ , then the exemplar  $\mathbf{e}_j$  is dismissed as cluster centre.

Objective function:

simplified log-likelihood + cardinality penalty of  $\boldsymbol{\pi}$



## Part 2 – Exemplar Pruning

Given  $\mathcal{E}_0 = \{\mathbf{e}_j\}_{j=1}^{\kappa}$ , the log-likelihood function is

$$\sum_{i=1}^n \log \left( \sum_{j=1}^{\kappa} \pi_j \phi(\mathbf{x}_i; \mathbf{e}_j, \Sigma_j) \right).$$

Introduce sparsity into  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_{\kappa})$ ; if  $\pi_j = 0$ , then the exemplar  $\mathbf{e}_j$  is dismissed as cluster centre.

Objective function:

simplified log-likelihood + cardinality penalty of  $\boldsymbol{\pi}$

## Part 2 – Exemplar Pruning

Given  $\mathcal{E}_0 = \{\mathbf{e}_j\}_{j=1}^{\kappa}$ , the log-likelihood function is

$$\sum_{i=1}^n \log \left( \sum_{j=1}^{\kappa} \pi_j \phi(\mathbf{x}_i; \mathbf{e}_j, \mathbf{\Sigma}_j) \right).$$

By Jensen's inequality we have

$$-\log \left( \sum_{j=1}^{\kappa} \pi_j \phi(\mathbf{x}_i; \mathbf{e}_j, \mathbf{\Sigma}_j) \right) \leq \sum_{j=1}^{\kappa} r_{ij} \log \left( \frac{r_{ij}}{\pi_j \phi(\mathbf{x}_i; \mathbf{e}_j, \mathbf{\Sigma}_j)} \right).$$

$r_{ij}$ 's are the responsibilities in the EM algorithm:  $\pi_j = \frac{1}{n} \sum_{i=1}^n r_{ij}$  and  $\sum_{j=1}^{\kappa} r_{ij} = 1$ .

## Part 2 – Exemplar Pruning

Given  $\mathcal{E}_0 = \{\mathbf{e}_j\}_{j=1}^{\kappa}$ , the log-likelihood function is

$$\sum_{i=1}^n \log \left( \sum_{j=1}^{\kappa} \pi_j \phi(\mathbf{x}_i; \mathbf{e}_j, \mathbf{\Sigma}_j) \right).$$

Therefore, the negative log-likelihood is

$$\begin{aligned} \sum_{i=1}^n -\log \left( \sum_{j=1}^{\kappa} \pi_j \phi(\mathbf{x}_i; \mathbf{e}_j, \mathbf{\Sigma}_j) \right) &\leq \sum_{i=1}^n \sum_{j=1}^{\kappa} r_{ij} \log \left( \frac{r_{ij}}{\pi_j \phi(\mathbf{x}_i; \mathbf{e}_j, \mathbf{\Sigma}_j)} \right) \\ &= \min_{\{\mathbf{r}_{i \in \Delta}\}_{i=1}^n} \sum_{i=1}^n \sum_{j=1}^{\kappa} r_{ij} \log \left( \frac{r_{ij}}{\pi_j \phi(\mathbf{x}_i; \mathbf{e}_j, \mathbf{\Sigma}_j)} \right). \end{aligned}$$

## Part 2 – Exemplar Pruning

Given  $\mathcal{E}_0 = \{\mathbf{e}_j\}_{j=1}^{\kappa}$ , the log-likelihood function is

$$\sum_{i=1}^n \log \left( \sum_{j=1}^{\kappa} \pi_j \phi(\mathbf{x}_i; \mathbf{e}_j, \Sigma_j) \right).$$

Minimizing the negative log-likelihood is equivalent to

$$\min_{\{\Sigma_j > 0\}_{j=1}^{\kappa}} \min_{\{\mathbf{r}_{i \in \Delta}\}_{i=1}^n} \sum_{i=1}^n \sum_{j=1}^{\kappa} r_{ij} \log \left( \frac{r_{ij}}{\pi_j \phi(\mathbf{x}_i; \mathbf{e}_j, \Sigma_j)} \right),$$

which is

$$\min_{\{\Sigma_j > 0\}_{j=1}^{\kappa}} \min_{\{\mathbf{r}_{i \in \Delta}\}_{i=1}^n} \sum_{i=1}^n \sum_{j=1}^{\kappa} r_{ij} \left[ \log \left( \frac{r_{ij}}{\pi_j} \right) + \frac{1}{2} \log(|\Sigma_j|) + \frac{1}{2} (\mathbf{x}_i - \mathbf{e}_j)^T \Sigma_j^{-1} (\mathbf{x}_i - \mathbf{e}_j) \right]$$

## Part 2 – Exemplar Pruning

$$\min_{\{\Sigma_j \succ 0\}_{j=1}^k} \min_{\{\mathbf{r}_i \in \Delta\}_{i=1}^n} \sum_{i=1}^n \sum_{j=1}^k r_{ij} \left[ \log \left( \frac{r_{ij}}{\pi_j} \right) + \frac{1}{2} \log(|\Sigma_j|) + \frac{1}{2} (\mathbf{x}_i - \mathbf{e}_j)^T \Sigma_j^{-1} (\mathbf{x}_i - \mathbf{e}_j) \right]$$

- $\min_{\{\Sigma_j \succ 0\}_{j=1}^k}$  : the covariance estimates  $\{\Sigma_j\}_{j=1}^k$  are fixed when pruning.
- $\log \left( \frac{r_{ij}}{\pi_j} \right)$ : numerical algorithms will behave erratically for any  $\pi_j \rightarrow 0$ .

Our simplified log-likelihood is

$$\min_{\{\mathbf{r}_i \in \Delta\}_{i=1}^n} \sum_{i=1}^n \sum_{j=1}^k r_{ij} \left[ \frac{1}{2} \log(|\Sigma_j|) + \frac{1}{2} (\mathbf{x}_i - \mathbf{e}_j)^T \Sigma_j^{-1} (\mathbf{x}_i - \mathbf{e}_j) \right].$$

## Part 2 – Exemplar Pruning

Given  $\mathcal{E}_0 = \{\mathbf{e}_j\}_{j=1}^{\kappa}$ , the log-likelihood function is

$$\sum_{i=1}^n \log \left( \sum_{j=1}^{\kappa} \pi_j \phi(\mathbf{x}_i; \mathbf{e}_j, \Sigma_j) \right).$$

Introduce sparsity into  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_{\kappa})$ ; if  $\pi_j = 0$ , then the exemplar  $\mathbf{e}_j$  is dismissed as cluster centre.

Objective function:

simplified log-likelihood + cardinality penalty of  $\boldsymbol{\pi}$

## Part 2 – Exemplar Pruning

The classical  $\ell_1$ -norm penalty is not suitable here:  $\|\boldsymbol{\pi}\|_1 = 1$  is constant on the simplex.

Our penalty is in the form of  $\|\boldsymbol{\omega} \circ \boldsymbol{\pi}\|_1$ , where  $\circ$  is the element-wise multiplication operator.

The weight vector  $\boldsymbol{\omega}$  should be data-driven and has the desirable property that gives more penalty to closer exemplars.

## Part 2 – Exemplar Pruning

The weight vector  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_k)$  is computed as

$$\omega_i = \max_{j=1, \dots, k} \{ \Pr(\pi_i \phi(X; \mathbf{e}_i, \boldsymbol{\Sigma}_i) < \pi_j \phi(X; \mathbf{e}_j, \boldsymbol{\Sigma}_j) | X \sim N(\mathbf{e}_i, \boldsymbol{\Sigma}_i)) \}.$$

Interpretation: the weight  $\omega_i$

(1) measures the likelihood that an instance from the  $i$ th mixture component is misclassified (into another mixture component),

(2) reflects the overlapping degree between the component distribution of  $\mathbf{e}_i$  and the other component distributions.



## Part 2 – Exemplar Pruning

The weight vector  $\boldsymbol{\omega} = (\omega_1, \dots, \omega_k)$  is computed as

$$\omega_i = \max_{j=1, \dots, k} \left\{ \Pr(\pi_i \phi(X; \mathbf{e}_i, \boldsymbol{\Sigma}_i) < \pi_j \phi(X; \mathbf{e}_j, \boldsymbol{\Sigma}_j) | X \sim N(\mathbf{e}_i, \boldsymbol{\Sigma}_i)) \right\}.$$

Analytic calculation of  $\omega_i$  is impractical, but numerical computation is readily done (Maitra and Melnykov, 2010).

## Part 2 – Exemplar Pruning

Objective function:

simplified log-likelihood + cardinality penalty of  $\pi$

$$\min_{\{r_i \in \Delta\}_{i=1}^n} \sum_{i=1}^n \sum_{j=1}^{\kappa} r_{ij} \left[ \frac{1}{2} \log(|\Sigma_j|) + \frac{1}{2} (\mathbf{x}_i - \mathbf{e}_j)^T \Sigma_j^{-1} (\mathbf{x}_i - \mathbf{e}_j) \right] + \theta \sum_{i=1}^n r_i^T \boldsymbol{\omega}$$

Equivalent to

$$\min_{\{r_i \in \Delta\}_{i=1}^n} \sum_{i=1}^n r_i^T \left( \frac{1}{2} \boldsymbol{\xi} + \frac{1}{2} \mathbf{d}_i + \theta \boldsymbol{\omega} \right),$$

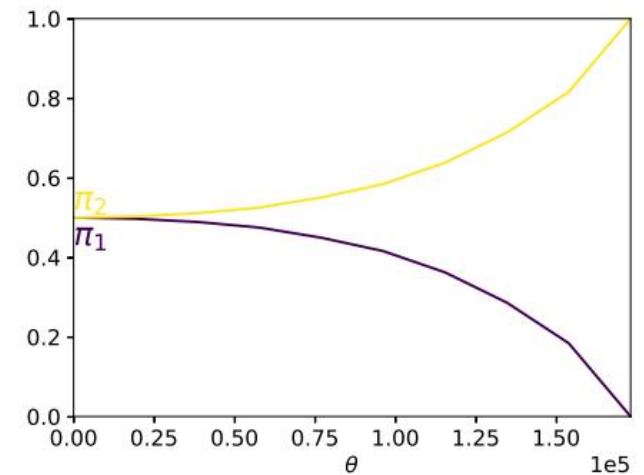
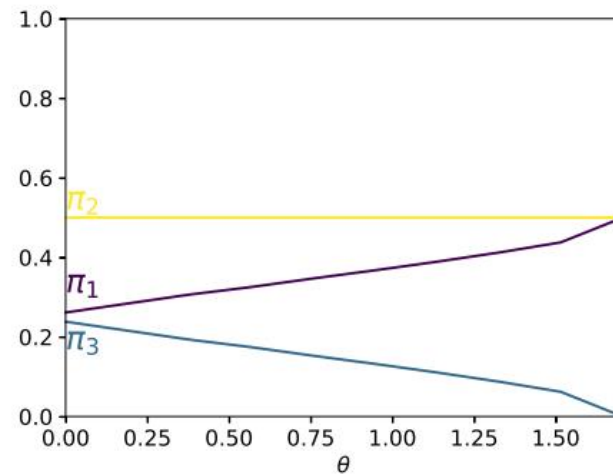
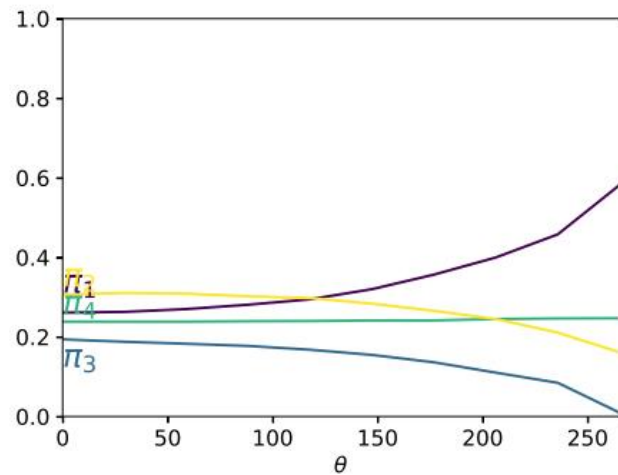
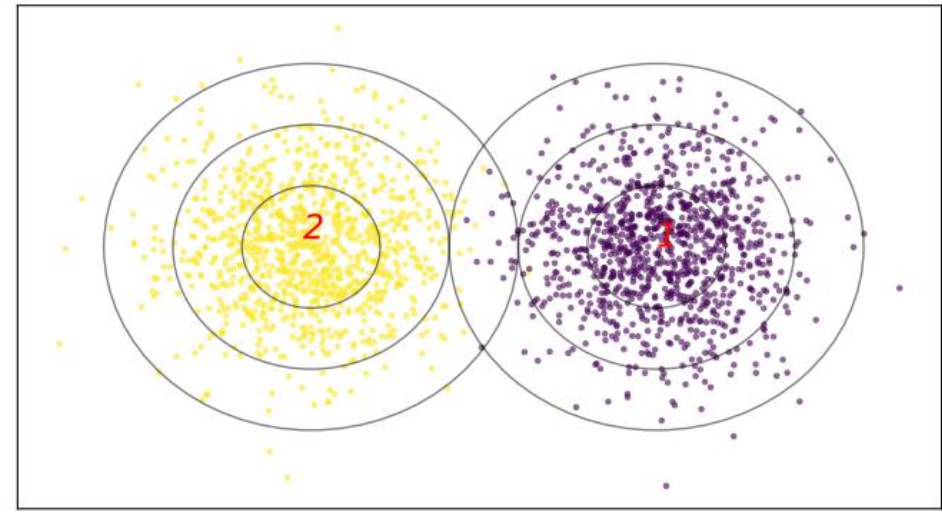
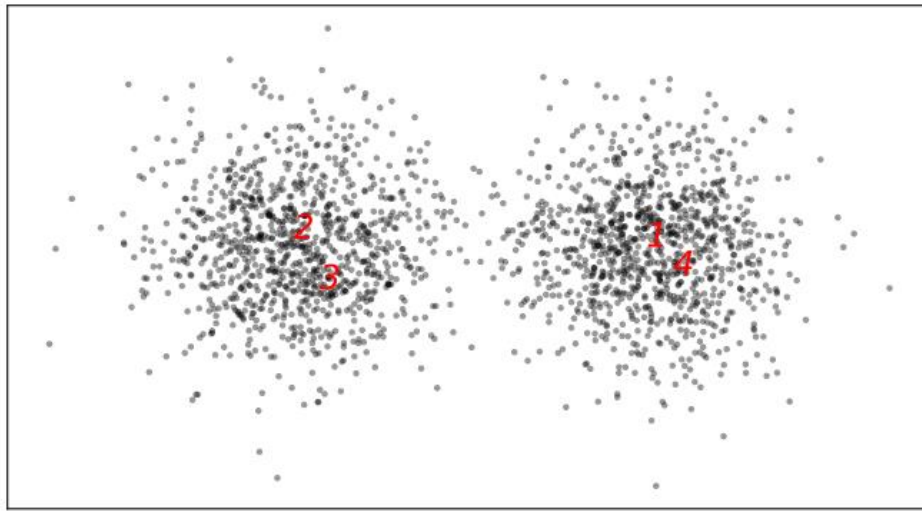
extremely simple (linear and separable).

## Part 2 – Exemplar Pruning

Objective function:

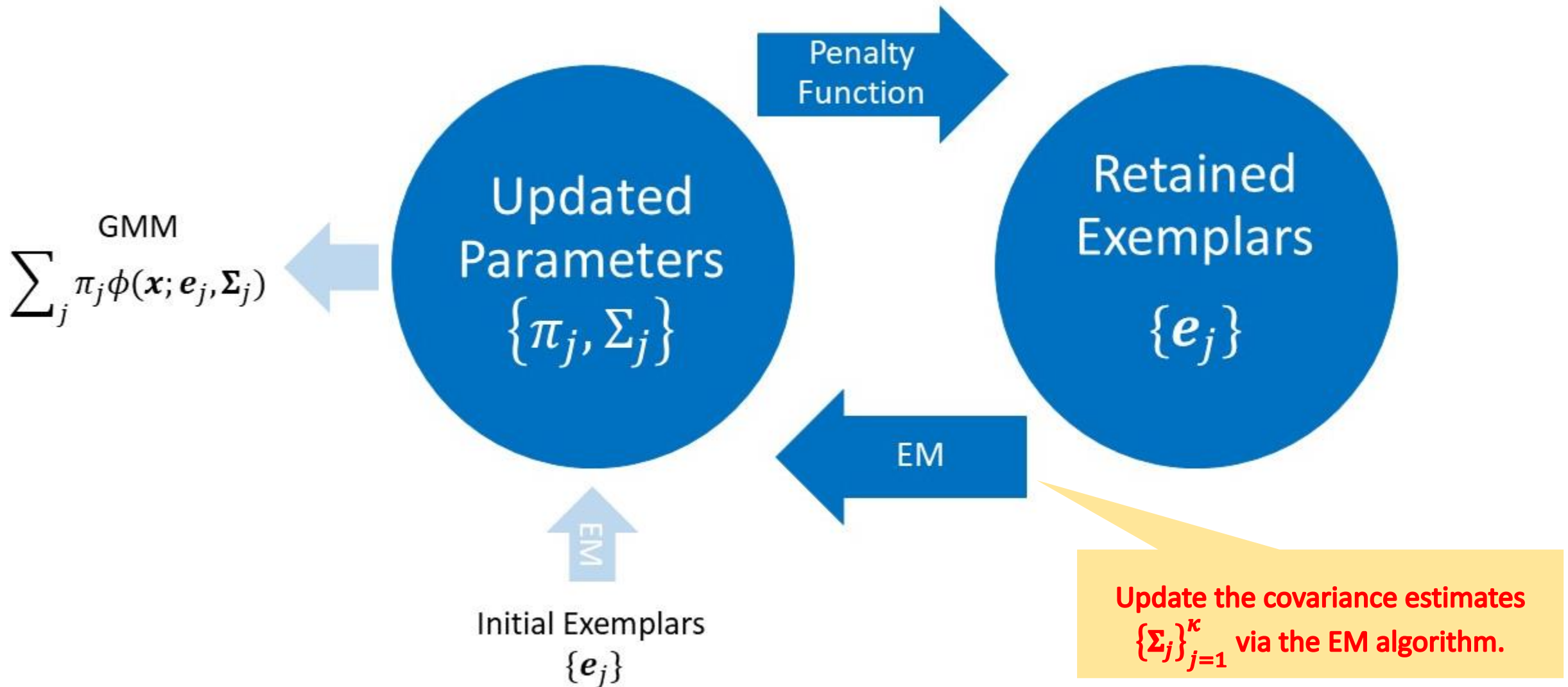
$$\min_{\{\mathbf{r}_i \in \Delta\}_{i=1}^n} \sum_{i=1}^n \mathbf{r}_i^T \left( \frac{1}{2} \boldsymbol{\xi} + \frac{1}{2} \mathbf{d}_i + \theta \boldsymbol{\omega} \right).$$

- (1) The parameter  $\theta$  controls the amount of shrinkage on  $\boldsymbol{\pi}$  ( $= \frac{1}{n} \sum_{i=1}^n \mathbf{r}_i$ ).
- (2) The trajectory of  $\boldsymbol{\pi}$ , as a function of  $\theta$ , can be easily computed by piecewise-linear homotopy methods.



**Top Left:** The data and the four selected exemplars, labelled in decreasing order of  $\rho\delta$ . **Top Right:** The final clustering obtained by REM. **Bottom:** The whole trajectory of  $\pi$ , as a function of  $\theta$ , in each REM iteration. After the first iteration, exemplar  $e_3$  is pruned; after the second iteration, exemplar  $e_4$  is pruned. The bottom right panel shows that  $\theta$  needs to be very large to merge two true cluster centres.

## Part 2 – Exemplar Pruning



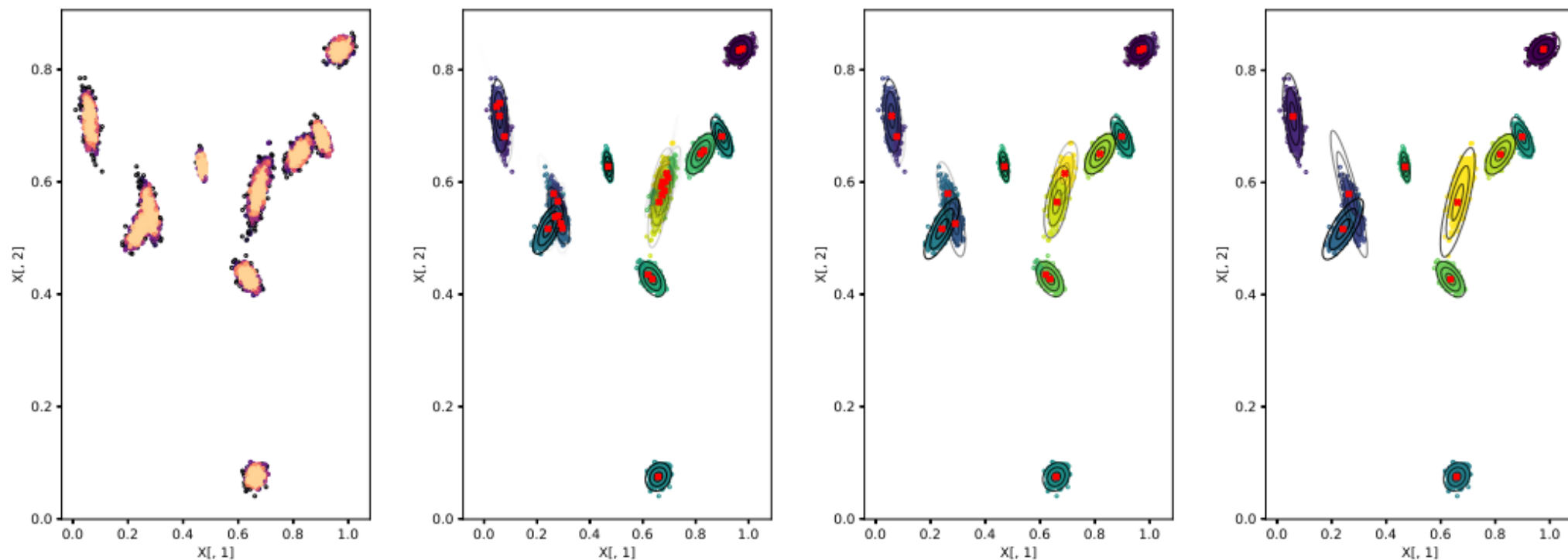
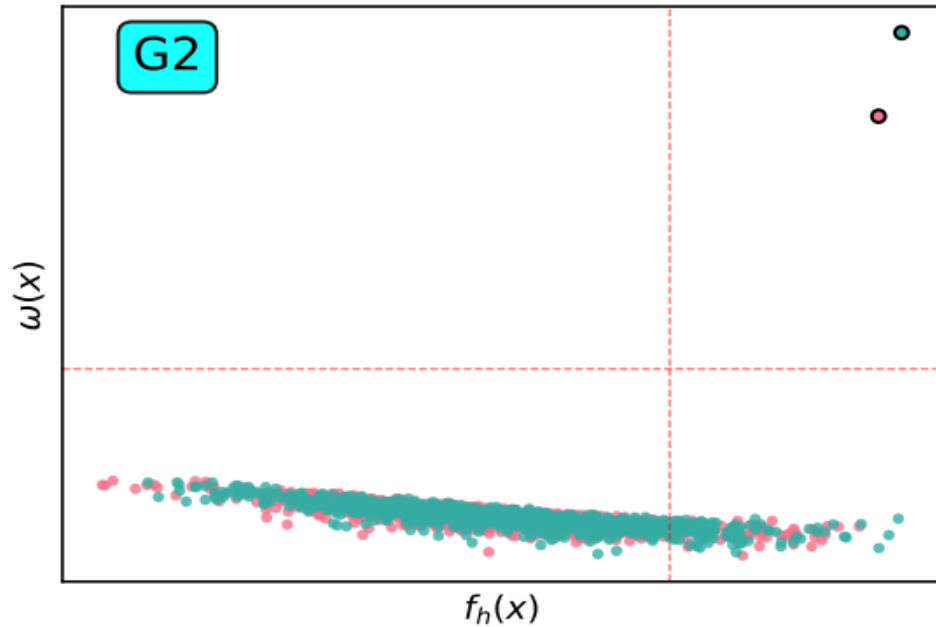


Fig. 4: A worked example of the clustering process for a 2-dimensional dataset with 10 components. (1) The leftmost figure shows the kernel density estimate for each instance, with lighter colors representing instances of higher density. (2) The second figure shows the initial exemplars (in red) with confidence ellipses representing the initial covariance matrix for each exemplar. (3) The third figure shows an intermediate clustering step, when multiple exemplars have been pruned from the initial set. (4) The rightmost figure shows the optimal clustering selected from the sequence using the ICL criterion.

# Evaluation

	$n$	$p$	$m$
G2	2048	128	2

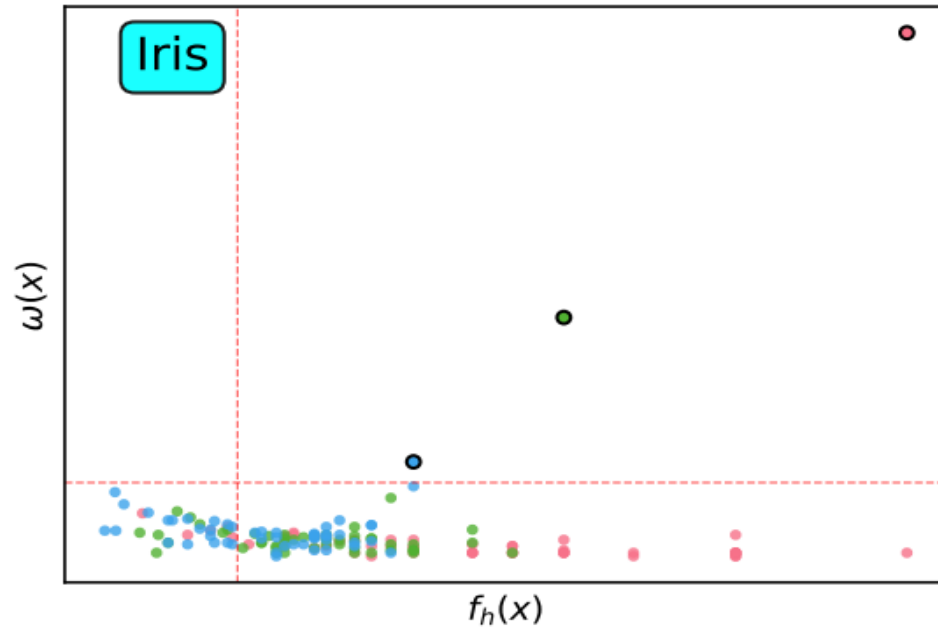


## G2

		ARI	NMI	Time
REM	AIC	1.00	1.00	32.0
	BIC	1.00	1.00	
	ICL	1.00	1.00	
riEM	AIC	0.15	0.19	2025.9
	BIC	0.00	0.00	
kmEM	AIC	0.70	0.72	52.1
	BIC	0.00	0.05	
emEM	AIC	0.75	0.76	181.1
	BIC	0.00	0.00	
rndEM	AIC	1.00	1.00	2433.0
	BIC	0.00	0.00	
Mclust	BIC	0.00	0.00	1153.9

# Evaluation

	$n$	$p$	$m$
Iris	150	4	3



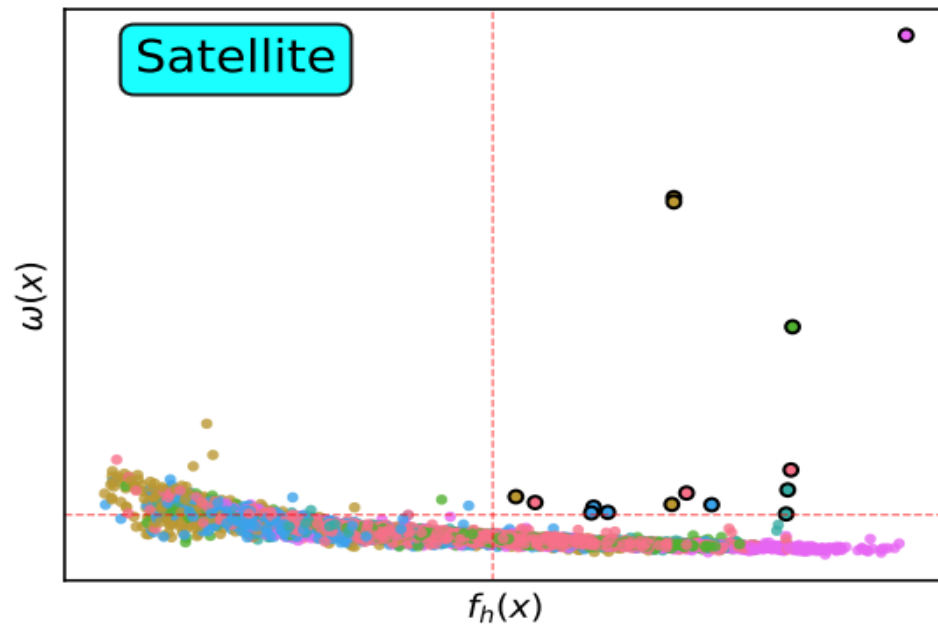
## Iris

		ARI	NMI	Time
REM	AIC	0.90	0.90	0.3
	BIC	0.90	0.90	
	ICL	0.90	0.90	
riEM	AIC	0.74	0.79	90.9
	BIC	0.44	0.65	
kmEM	AIC	0.69	0.77	4.7
	BIC	0.69	0.77	
emEM	AIC	0.86	0.86	2.3
	BIC	0.57	0.73	
rndEM	AIC	0.90	0.90	26.3
	BIC	0.57	0.73	
Mclust	BIC	0.56	0.76	0.7



# Evaluation

	$n$	$p$	$m$
Satellite	4435	36	6

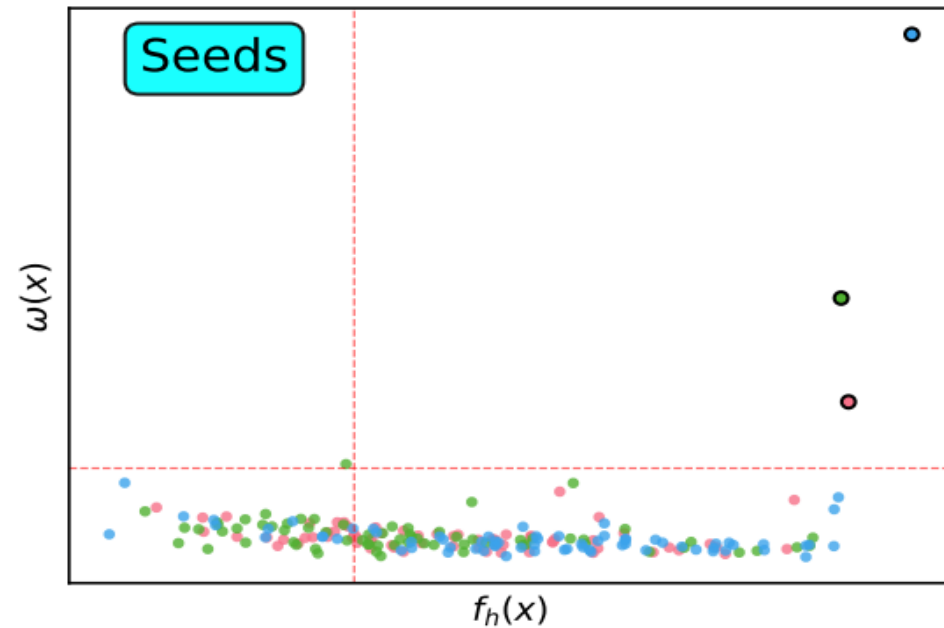


## Satellite

		ARI	NMI	Time
REM	AIC	0.52	0.58	295.0
	BIC	0.52	0.58	
	ICL	0.52	0.58	
riEM	AIC	0.44	0.50	7302.2
	BIC	0.42	0.49	
kmEM	AIC	0.44	0.57	1161.9
	BIC	0.47	0.56	
emEM	AIC	0.42	0.55	1020.4
	BIC	0.44	0.52	
rndEM	AIC	0.46	0.55	7242.5
	BIC	0.48	0.56	
Mclust	BIC	0.47	0.56	437.7

# Evaluation

	$n$	$p$	$m$
Seeds	210	7	3



## Seeds

		ARI	NMI	Time
REM	AIC	0.77	<b>0.74</b>	<b>0.9</b>
	BIC	0.77	<b>0.74</b>	
	ICL	0.77	<b>0.74</b>	
riEM	AIC	0.48	0.56	158.2
	BIC	0.01	0.07	
kmEM	AIC	0.50	0.51	2.3
	BIC	0.51	0.58	
emEM	AIC	0.51	0.58	4.0
	BIC	0.51	0.57	
rndEM	AIC	0.64	0.67	41.2
	BIC	0.66	0.62	
Mclust	BIC	<b>0.79</b>	<b>0.74</b>	1.2

# Reference

- [**Exemplar**] Lashkari, D., Golland, P.: Convex clustering with exemplar-based models. In: Advances in Neural Information Processing Systems 20 (NIPS 2007), pp. 825–832
- [**Exemplar**] Pilanci, M., Ghaoui, L.E., Chandrasekaran, V.: Recovery of sparse probability measures via convex programming. In: Advances in Neural Information Processing Systems 25 (NIPS 2012), pp. 2420–2428
- [**Peak Finding**] Rodriguez, A., Laio, A.: Clustering by fast search and find of density peaks. Science (New York, N.Y.) 344(6191), 1492–1496 (2014)
- [**Overlapping Degree**] Maitra, R., Melnykov, V.: Simulating data to study performance of finite mixture modeling and clustering algorithms. Journal of Computational and Graphical Statistics, vol. 19, no. 2, 2010, pp. 354–76.