

MLE part 2

Gaussian Mixture Model

- ▶ Suppose data is drawn from k Gaussians, meaning

$$Y = j \sim \text{Discrete}(\pi_1, \dots, \pi_k),$$

$$\mathbf{X} = \mathbf{x} | Y = j \sim \mathcal{N}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j),$$

and the parameters are $\boldsymbol{\theta} = ((\pi_1, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \dots, (\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))$.

(Note: this is a **generative** model, and we have a way to sample.)

Gaussian Mixture Model

- Suppose data is drawn from k Gaussians, meaning

$$Y = j \sim \text{Discrete}(\pi_1, \dots, \pi_k),$$
$$\mathbf{X} = \mathbf{x} | Y = j \sim \mathcal{N}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j),$$

and the parameters are $\boldsymbol{\theta} = ((\pi_1, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \dots, (\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))$.
(Note: this is a **generative** model, and we have a way to sample.)

- The probability density (with parameters $\boldsymbol{\theta} = ((\pi_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j))_{j=1}^k$) at a given \mathbf{x} is

$$p_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{j=1}^k p_{\boldsymbol{\theta}}(\mathbf{x} | y = j) p_{\boldsymbol{\theta}}(y = j) = \sum_{j=1}^k p_{\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j}(\mathbf{x} | Y = j) \pi_j,$$

and the likelihood problem is

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^n \ln \sum_{j=1}^k \frac{\pi_j}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_j|}} \exp\left(-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1}(\mathbf{x}_i - \boldsymbol{\mu}_j)\right).$$

The \ln and the \exp are no longer next to each other; we can't just

Lloyd's method for k -means

Original k -means formulation

$$\phi((\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k)) = \sum_{i=1}^n \min_j \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2.$$

Lloyd's method for k -means

Original k -means formulation

$$\phi((\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k)) = \sum_{i=1}^n \min_j \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2.$$

To make an algorithm, we introduced *assignment matrix* $\mathbf{A} \in \mathcal{A}_{n,k}$:

$$\phi((\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k); \mathbf{A}) = \sum_{i=1}^n \sum_{j=1}^k A_{ij} \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2.$$

Lloyd's method for k -means

Original k -means formulation

$$\phi((\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k)) = \sum_{i=1}^n \min_j \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2.$$

To make an algorithm, we introduced *assignment matrix* $\mathbf{A} \in \mathcal{A}_{n,k}$:

$$\phi((\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k); \mathbf{A}) = \sum_{i=1}^n \sum_{j=1}^k A_{ij} \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2.$$

Let's do the same thing with Gaussians!

Gaussian mixture likelihood with *responsibility matrix* \mathbf{R}

Let's replace $\sum_{i=1}^n \ln \sum_{j=1}^k \pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i)$ with

$$\sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln (\pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i))$$

where $\mathbf{R} \in \mathcal{R}_{n,k} := \{\mathbf{R} \in [0, 1]^{n \times k} : \mathbf{R}\mathbf{1}_k = \mathbf{1}_n\}$ is a **responsibility matrix**.

Gaussian mixture likelihood with *responsibility matrix* \mathbf{R}

Let's replace $\sum_{i=1}^n \ln \sum_{j=1}^k \pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i)$ with

$$\sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln (\pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i))$$

where $\mathbf{R} \in \mathcal{R}_{n,k} := \{\mathbf{R} \in [0, 1]^{n \times k} : \mathbf{R}\mathbf{1}_k = \mathbf{1}_n\}$ is a **responsibility matrix**.

Holding \mathbf{R} fixed and optimizing θ gives

$$\begin{aligned}\pi_j &:= \frac{\sum_{i=1}^n R_{ij}}{\sum_{i=1}^n \sum_{l=1}^k R_{il}} = \frac{\sum_{i=1}^n R_{ij}}{n}, \\ \mu_j &:= \frac{\sum_{i=1}^n R_{ij} \mathbf{x}_i}{\sum_{i=1}^n R_{ij}} = \frac{\sum_{i=1}^n R_{ij} \mathbf{x}_i}{n\pi_j}, \\ \Sigma_j &:= \frac{\sum_{i=1}^n R_{ij} (\mathbf{x}_i - \mu_j)(\mathbf{x}_i - \mu_j)^\top}{n\pi_j}.\end{aligned}$$

(Should use new mean in Σ_j so that all derivatives 0.)

Updating μ_j

Recall our new likelihood with responsibilities R :

$$\sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln \pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i)$$

(In the literature, this quantity is “**expected complete data likelihood**”.)

Updating μ_j

Recall our new likelihood with responsibilities R :

$$\sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln \pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i)$$

(In the literature, this quantity is “**expected complete data likelihood**”.)

Taking derivative and setting to 0:

$$\begin{aligned} 0 &= \sum_{i=1}^n R_{ij} \nabla_{\mu_j} \left(\ln \exp\left(-\frac{1}{2}(\mathbf{x}_i - \mu_j)^\top \Sigma_j^{-1}(\mathbf{x}_i - \mu_j)\right) + \text{terms w/o } \mu_j \right) \\ &= \sum_{i=1}^n R_{ij} \Sigma_j^{-1}(\mathbf{x}_i - \mu_j). \end{aligned}$$

$$\text{Rearranging, } \mu_j = \frac{\sum_{i=1}^n R_{ij} \mathbf{x}_i}{n \pi_j}.$$

Updating π

Recall our new likelihood with responsibilities R :

$$\sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln \pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i)$$

Updating π

Recall our new likelihood with responsibilities R :

$$\sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln \pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i)$$

Taking derivative and setting to 0:

$$0 = \sum_{i=1}^n \frac{R_{ij}}{\pi_j};$$

oops?

Updating π

Recall our new likelihood with responsibilities R :

$$\sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln \pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i)$$

Taking derivative and setting to 0:

$$0 = \sum_{i=1}^n \frac{R_{ij}}{\pi_j};$$

oops?

Fix: we forgot the constraints on π !

Updating π

Include constraint $\sum_{j=1}^k \pi_j = 1$ with a **Lagrangian**:

$$\sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln \pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i) + \lambda \left(1 - \sum_{j=1}^k \pi_j \right)$$

Updating π

Include constraint $\sum_{j=1}^k \pi_j = 1$ with a **Lagrangian**:

$$\sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln \pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i) + \lambda \left(1 - \sum_{j=1}^k \pi_j \right)$$

Differentiating and setting this Lagrangian to 0, we get

$$\lambda = \sum_{i=1}^n \frac{R_{ij}}{\pi_j}, \quad \text{and} \quad \sum_j \pi_j = 1.$$

Updating π

Include constraint $\sum_{j=1}^k \pi_j = 1$ with a **Lagrangian**:

$$\sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln \pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i) + \lambda \left(1 - \sum_{j=1}^k \pi_j \right)$$

Differentiating and setting this Lagrangian to 0, we get

$$\lambda = \sum_{i=1}^n \frac{R_{ij}}{\pi_j}, \quad \text{and} \quad \sum_j \pi_j = 1.$$

Together, $\pi_j = \sum_{i=1}^n R_{ij} / \lambda$, and

$$1 = \sum_{j=1}^k \pi_j = \sum_{j=1}^k \sum_{i=1}^n \frac{R_{ij}}{\lambda} = \frac{n}{\lambda},$$

so $\lambda = n$ and $\pi_j = \sum_{i=1}^n R_{ij} / n$.

Starting again from likelihood with responsibilities R :

$$\sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln \pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i).$$

Updating Σ_j

Starting again from likelihood with responsibilities R :

$$\sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln \pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i).$$

Taking derivative and setting to 0,

$$0 = \sum_{i=1}^n R_{ij} \nabla_{\Sigma_j} \left(-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_j)^\top \Sigma_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j) - \frac{1}{2} \ln |\Sigma_j| + \text{other stuff} \right).$$

Updating Σ_j

Starting again from likelihood with responsibilities R :

$$\sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln \pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i).$$

Taking derivative and setting to 0,

$$0 = \sum_{i=1}^n R_{ij} \nabla_{\Sigma_j} \left(-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_j)^\top \Sigma_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j) - \frac{1}{2} \ln |\Sigma_j| + \text{other stuff} \right).$$

By magic matrix derivative rules,

$$\Sigma_j^{-1} = \sum_{i=1}^n R_{ij} (\mathbf{x}_i - \boldsymbol{\mu}_j) (\mathbf{x}_i - \boldsymbol{\mu}_j)^\top / (n\pi_j).$$

Summary of θ optimization

Replace $\sum_{i=1}^n \ln \sum_{j=1}^k \pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i)$ with

$$\sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln (\pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i)).$$

Hold R fixed and optimize θ :

$$\begin{aligned}\pi_j &:= \frac{\sum_{i=1}^n R_{ij}}{\sum_{i=1}^n \sum_{l=1}^k R_{il}} = \frac{\sum_{i=1}^n R_{ij}}{n}; \\ \boldsymbol{\mu}_j &:= \frac{\sum_{i=1}^n R_{ij} \mathbf{x}_i}{\sum_{i=1}^n R_{ij}} = \frac{\sum_{i=1}^n R_{ij} \mathbf{x}_i}{n\pi_j}, \\ \boldsymbol{\Sigma}_j &:= \frac{\sum_{i=1}^n R_{ij} (\mathbf{x}_i - \boldsymbol{\mu}_j)(\mathbf{x}_i - \boldsymbol{\mu}_j)^\top}{n\pi_j}.\end{aligned}$$

Summary of θ optimization

Replace $\sum_{i=1}^n \ln \sum_{j=1}^k \pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i)$ with

$$\sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln (\pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i)).$$

Hold R fixed and optimize θ :

$$\begin{aligned}\pi_j &:= \frac{\sum_{i=1}^n R_{ij}}{\sum_{i=1}^n \sum_{l=1}^k R_{il}} = \frac{\sum_{i=1}^n R_{ij}}{n}; \\ \mu_j &:= \frac{\sum_{i=1}^n R_{ij} \mathbf{x}_i}{\sum_{i=1}^n R_{ij}} = \frac{\sum_{i=1}^n R_{ij} \mathbf{x}_i}{n\pi_j}, \\ \Sigma_j &:= \frac{\sum_{i=1}^n R_{ij} (\mathbf{x}_i - \mu_j)(\mathbf{x}_i - \mu_j)^\top}{n\pi_j}.\end{aligned}$$

How to optimize R_{ij} ?

- ▶ Likelihood lacks the \min_j from the k -means cost.
- ▶ We'll now develop the E-M method, which picks R in a way that guarantees **likelihood** increases.

E-M (Expectation-Maximization)

Generalizing the assignment matrix to GMMs

We introduced an **assignment matrix** $A \in \{0, 1\}^{n \times k}$:

- ▶ For each x_i , define $\mu(x_i)$ to be a closest center:

$$\|x_i - \mu(x_i)\| = \min_j \|x_i - \mu_j\|.$$

- ▶ For each i , set $A_{ij} = \mathbb{1}[\mu(x_i) = \mu_j]$.

Generalizing the assignment matrix to GMMs

We introduced an **assignment matrix** $\mathbf{A} \in \{0, 1\}^{n \times k}$:

- ▶ For each \mathbf{x}_i , define $\boldsymbol{\mu}(\mathbf{x}_i)$ to be a closest center:

$$\|\mathbf{x}_i - \boldsymbol{\mu}(\mathbf{x}_i)\| = \min_j \|\mathbf{x}_i - \boldsymbol{\mu}_j\|.$$

- ▶ For each i , set $A_{ij} = \mathbb{1}[\boldsymbol{\mu}(\mathbf{x}_i) = \boldsymbol{\mu}_j]$.

- ▶ **Key property:** by this choice,

$$\phi(\mathbf{C}; \mathbf{A}) = \sum_{i=1}^n \sum_{j=1}^k A_{ij} \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2 = \sum_{i=1}^n \min_j \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2 = \phi(\mathbf{C});$$

therefore we can decrease $\phi(\mathbf{C}) = \phi(\mathbf{C}; \mathbf{A})$

first by optimizing \mathbf{C} to get $\phi(\mathbf{C}'; \mathbf{A}) \leq \phi(\mathbf{C}; \mathbf{A})$,

then setting \mathbf{A} as above to get

$$\phi(\mathbf{C}') = \phi(\mathbf{C}'; \mathbf{A}') \leq \phi(\mathbf{C}'; \mathbf{A}) \leq \phi(\mathbf{C}; \mathbf{A}) = \phi(\mathbf{C}).$$

In other words: we minimize $\phi(\mathbf{C})$ via $\phi(\mathbf{C}; \mathbf{A})$.

Generalizing the assignment matrix to GMMs

We introduced an **assignment matrix** $\mathbf{A} \in \{0, 1\}^{n \times k}$:

- ▶ For each \mathbf{x}_i , define $\boldsymbol{\mu}(\mathbf{x}_i)$ to be a closest center:

$$\|\mathbf{x}_i - \boldsymbol{\mu}(\mathbf{x}_i)\| = \min_j \|\mathbf{x}_i - \boldsymbol{\mu}_j\|.$$

- ▶ For each i , set $A_{ij} = \mathbb{1}[\boldsymbol{\mu}(\mathbf{x}_i) = \boldsymbol{\mu}_j]$.

- ▶ **Key property:** by this choice,

$$\phi(\mathbf{C}; \mathbf{A}) = \sum_{i=1}^n \sum_{j=1}^k A_{ij} \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2 = \sum_{i=1}^n \min_j \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2 = \phi(\mathbf{C});$$

therefore we can decrease $\phi(\mathbf{C}) = \phi(\mathbf{C}; \mathbf{A})$

first by optimizing \mathbf{C} to get $\phi(\mathbf{C}'; \mathbf{A}) \leq \phi(\mathbf{C}; \mathbf{A})$,

then setting \mathbf{A} as above to get

$$\phi(\mathbf{C}') = \phi(\mathbf{C}'; \mathbf{A}') \leq \phi(\mathbf{C}'; \mathbf{A}) \leq \phi(\mathbf{C}; \mathbf{A}) = \phi(\mathbf{C}).$$

In other words: we minimize $\phi(\mathbf{C})$ via $\phi(\mathbf{C}; \mathbf{A})$.

What fulfills the same role for \mathcal{L} ?

Latent variable models.

Since $1 = \sum_{j=1}^k p_{\theta}(y_i = j | \mathbf{x}_i)$ and $p_{\theta}(y_i = j | \mathbf{x}_i) = \frac{p_{\theta}(y_i = j, \mathbf{x}_i)}{p_{\theta}(\mathbf{x}_i)}$, then

$$\begin{aligned}\mathcal{L}(\theta) &= \sum_{i=1}^n \ln p_{\theta}(\mathbf{x}_i) = \sum_{i=1}^n 1 \cdot \ln p_{\theta}(\mathbf{x}_i) = \sum_{i=1}^n \sum_{j=1}^k p_{\theta}(y_i = j | \mathbf{x}_i) \ln p_{\theta}(\mathbf{x}_i) \\ &= \sum_{i=1}^n \sum_{j=1}^k p_{\theta}(y_i = j | \mathbf{x}_i) \ln \frac{p_{\theta}(\mathbf{x}_i, y_i = j)}{p_{\theta}(y_i = j | \mathbf{x}_i)}.\end{aligned}$$

Latent variable models.

Since $1 = \sum_{j=1}^k p_{\theta}(y_i = j | \mathbf{x}_i)$ and $p_{\theta}(y_i = j | \mathbf{x}_i) = \frac{p_{\theta}(\mathbf{x}_i, y_i = j)}{p_{\theta}(\mathbf{x}_i)}$, then

$$\begin{aligned}\mathcal{L}(\theta) &= \sum_{i=1}^n \ln p_{\theta}(\mathbf{x}_i) = \sum_{i=1}^n 1 \cdot \ln p_{\theta}(\mathbf{x}_i) = \sum_{i=1}^n \sum_{j=1}^k p_{\theta}(y_i = j | \mathbf{x}_i) \ln p_{\theta}(\mathbf{x}_i) \\ &= \sum_{i=1}^n \sum_{j=1}^k p_{\theta}(y_i = j | \mathbf{x}_i) \ln \frac{p_{\theta}(\mathbf{x}_i, y_i = j)}{p_{\theta}(y_i = j | \mathbf{x}_i)}.\end{aligned}$$

Therefore: define **augmented likelihood**

$$\mathcal{L}(\theta; \mathbf{R}) := \sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln \frac{p_{\theta}(\mathbf{x}_i, y_i = j)}{R_{ij}};$$

note that $R_{ij} := p_{\theta}(y_i = j | \mathbf{x}_i)$ **implies** $\mathcal{L}(\theta; \mathbf{R}) = \mathcal{L}(\theta)$.

Define **augmented likelihood**

$$\mathcal{L}(\boldsymbol{\theta}; \mathbf{R}) := \sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_i, y_i = j)}{R_{ij}},$$

with **responsibility matrix** $\mathbf{R} \in \mathcal{R}_{n,k} := \{\mathbf{R} \in [0, 1]^{n \times k} : \mathbf{R}\mathbf{1}_k = \mathbf{1}_n\}$.

Alternate two steps:

- ▶ **E-step**: set $(\mathbf{R}_t)_{ij} := p_{\boldsymbol{\theta}_{t-1}}(y_i = j | \mathbf{x}_i)$.
- ▶ **M-step**: set $\boldsymbol{\theta}_t = \arg \max_{\boldsymbol{\theta} \in \Theta} \mathcal{L}(\boldsymbol{\theta}; \mathbf{R}_t)$.

E-M method for latent variable models

Define **augmented likelihood**

$$\mathcal{L}(\boldsymbol{\theta}; \mathbf{R}) := \sum_{i=1}^n \sum_{j=1}^k R_{ij} \ln \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_i, y_i = j)}{R_{ij}},$$

with **responsibility matrix** $\mathbf{R} \in \mathcal{R}_{n,k} := \{\mathbf{R} \in [0, 1]^{n \times k} : \mathbf{R}\mathbf{1}_k = \mathbf{1}_n\}$.

Alternate two steps:

- ▶ **E-step**: set $(\mathbf{R}_t)_{ij} := p_{\boldsymbol{\theta}_{t-1}}(y_i = j | \mathbf{x}_i)$.
- ▶ **M-step**: set $\boldsymbol{\theta}_t = \arg \max_{\boldsymbol{\theta} \in \Theta} \mathcal{L}(\boldsymbol{\theta}; \mathbf{R}_t)$.

Soon: we'll see this gives nondecreasing likelihood!

E-M for Gaussian mixtures

Initialization: a standard choice is $\pi_j = 1/k$, $\Sigma_j = \mathbf{I}$, and $(\mu_j)_{j=1}^k$ given by k -means.

- ▶ **E-step:** Set $R_{ij} = p_{\theta}(y_i = j | \mathbf{x}_i)$, meaning

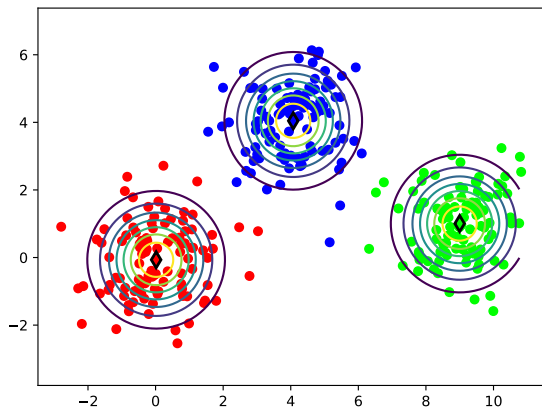
$$R_{ij} = p_{\theta}(y_i = j | \mathbf{x}_i) = \frac{p_{\theta}(y_i = j, \mathbf{x}_i)}{p_{\theta}(\mathbf{x}_i)} = \frac{\pi_j p_{\mu_j, \Sigma_j}(\mathbf{x}_i)}{\sum_{l=1}^k \pi_l p_{\mu_l, \Sigma_l}(\mathbf{x}_i)}.$$

- ▶ **M-step:** solve $\arg \max_{\theta \in \Theta} \mathcal{L}(\theta; \mathbf{R})$, meaning

$$\begin{aligned}\pi_j &:= \frac{\sum_{i=1}^n R_{ij}}{\sum_{i=1}^n \sum_{l=1}^k R_{il}} = \frac{\sum_{i=1}^n R_{ij}}{n}, \\ \mu_j &:= \frac{\sum_{i=1}^n R_{ij} \mathbf{x}_i}{\sum_{i=1}^n R_{ij}} = \frac{\sum_{i=1}^n R_{ij} \mathbf{x}_i}{n\pi_j}, \\ \Sigma_j &:= \frac{\sum_{i=1}^n R_{ij} (\mathbf{x}_i - \mu_j)(\mathbf{x}_i - \mu_j)^{\top}}{n\pi_j}.\end{aligned}$$

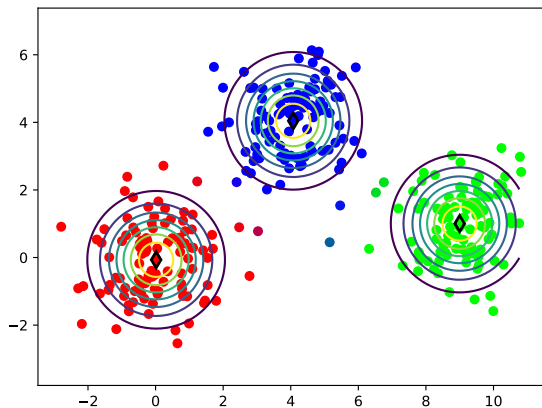
(These are as before.)

Demo: spherical clusters



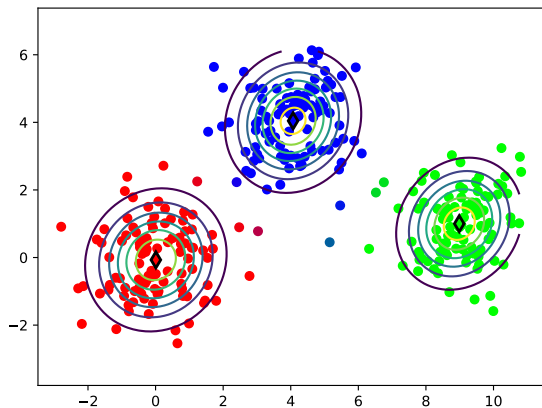
(Initialized with k -means, thus not so dramatic.)

Demo: spherical clusters



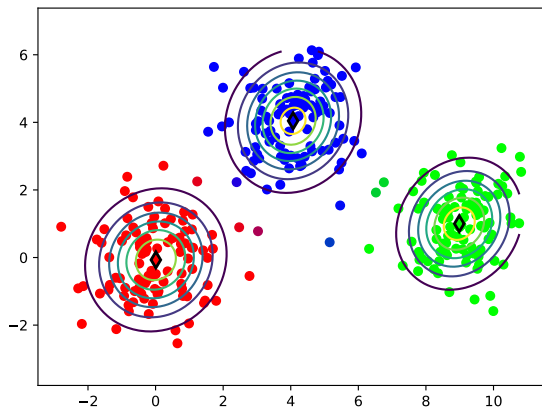
(Initialized with k -means, thus not so dramatic.)

Demo: spherical clusters



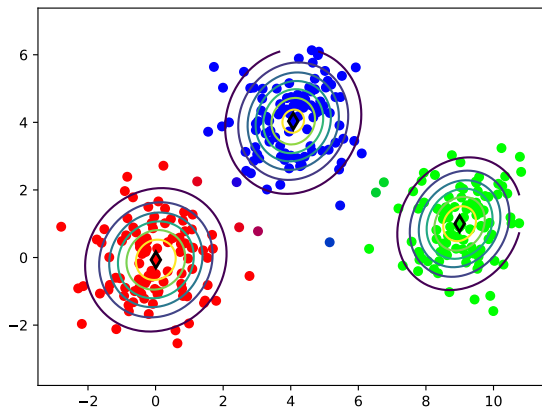
(Initialized with k -means, thus not so dramatic.)

Demo: spherical clusters



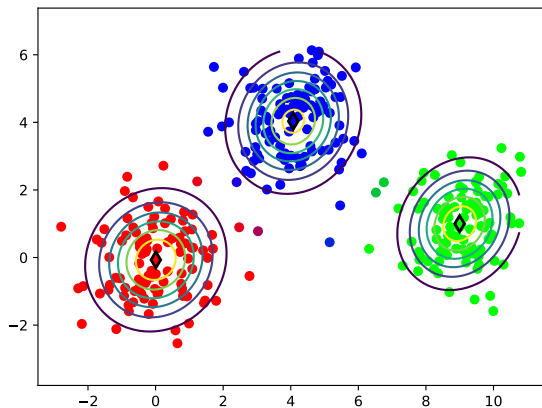
(Initialized with k -means, thus not so dramatic.)

Demo: spherical clusters



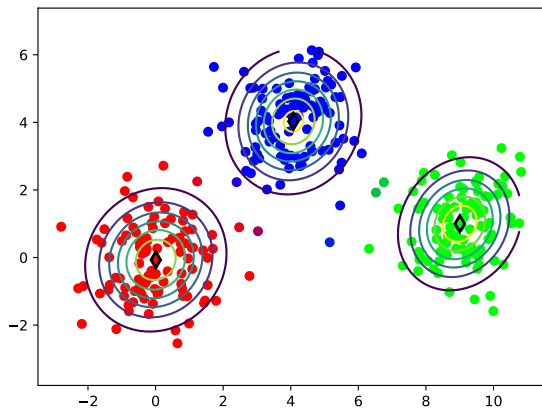
(Initialized with k -means, thus not so dramatic.)

Demo: spherical clusters



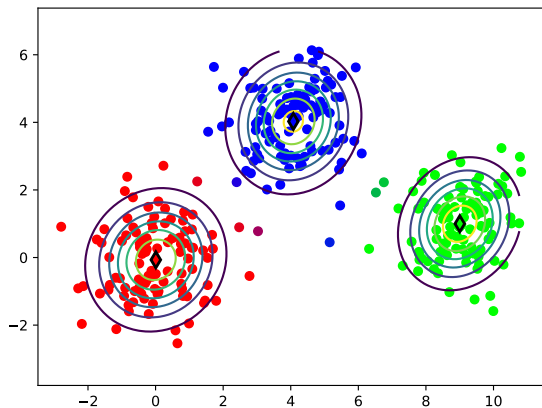
(Initialized with k -means, thus not so dramatic.)

Demo: spherical clusters



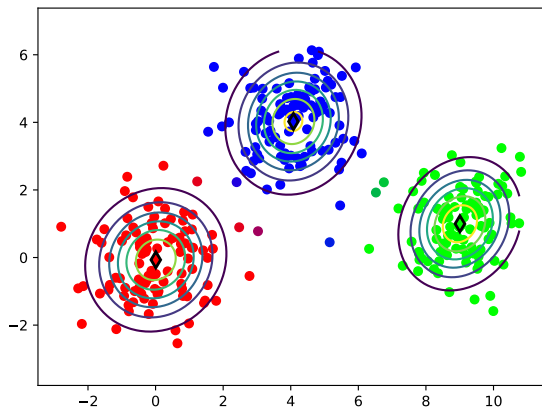
(Initialized with k -means, thus not so dramatic.)

Demo: spherical clusters



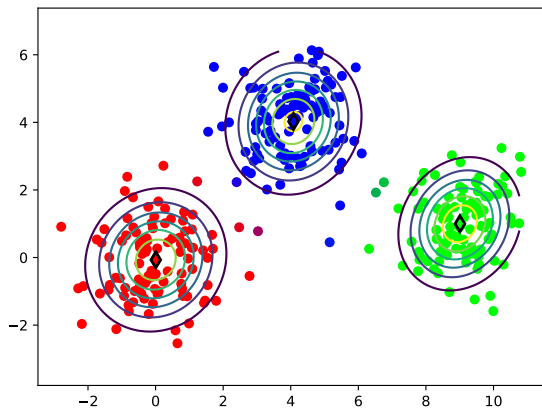
(Initialized with k -means, thus not so dramatic.)

Demo: spherical clusters



(Initialized with k -means, thus not so dramatic.)

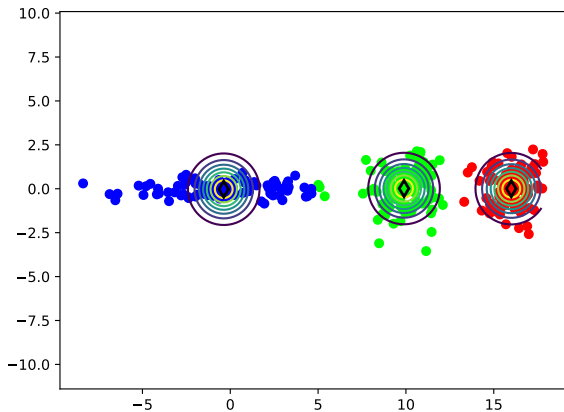
Demo: spherical clusters



(Initialized with k -means, thus not so dramatic.)

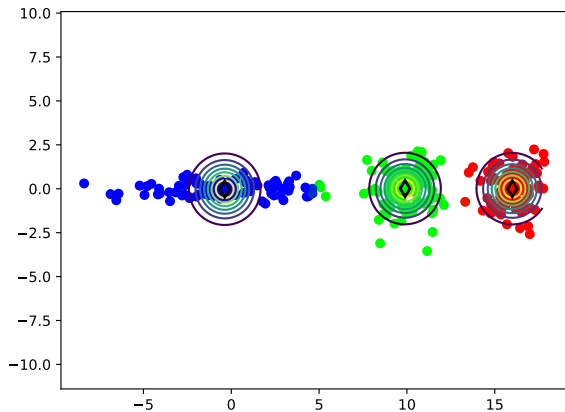
Demo: elliptical clusters

E...



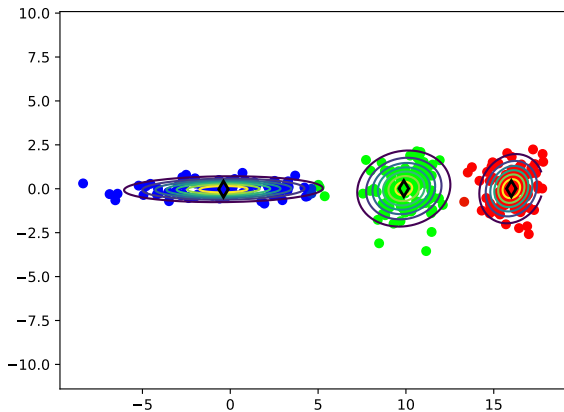
Demo: elliptical clusters

E... M...



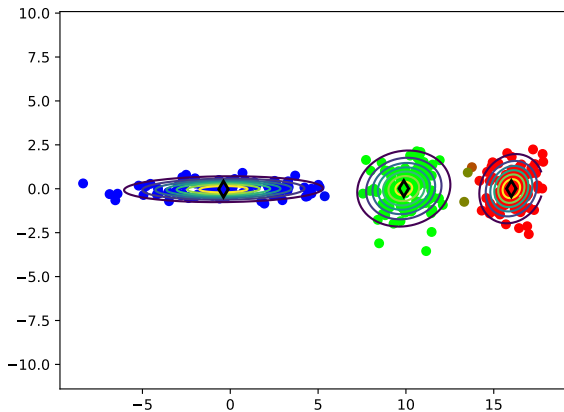
Demo: elliptical clusters

E... M... E...



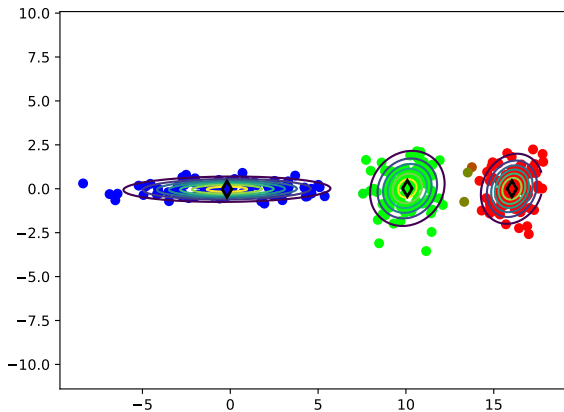
Demo: elliptical clusters

E... M... E... M...



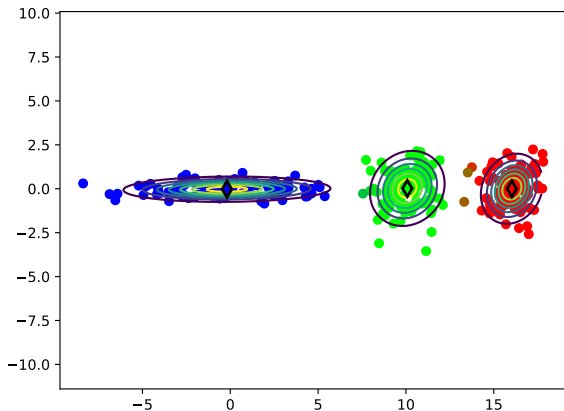
Demo: elliptical clusters

E... M... E... M... E...



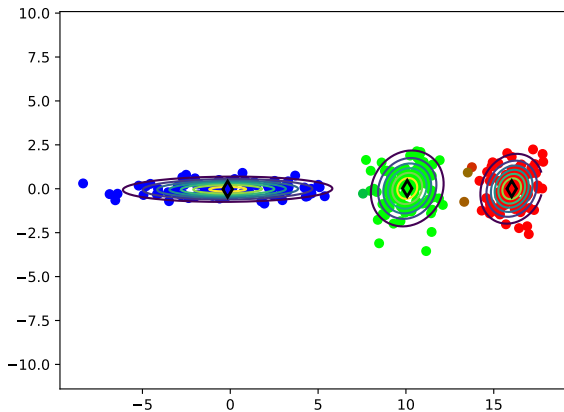
Demo: elliptical clusters

E... M... E... M... E... M...



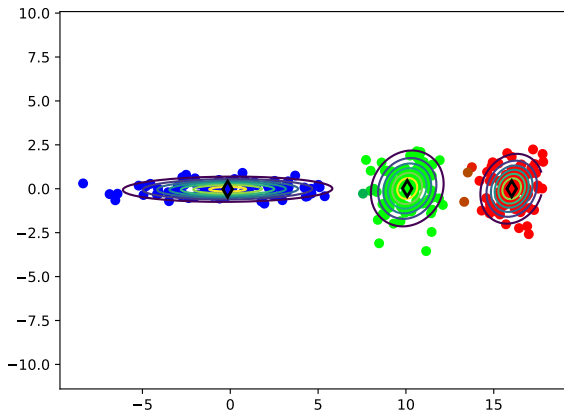
Demo: elliptical clusters

E... M... E... M... E... M... E...



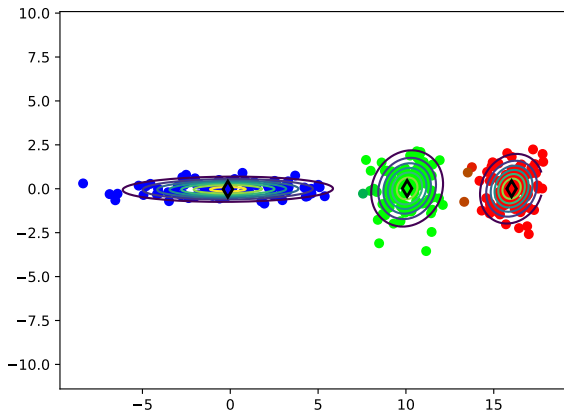
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



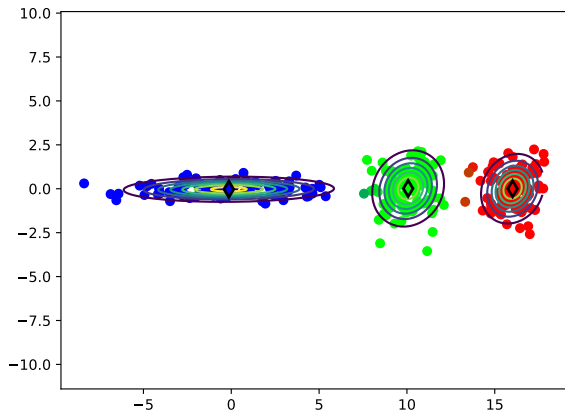
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



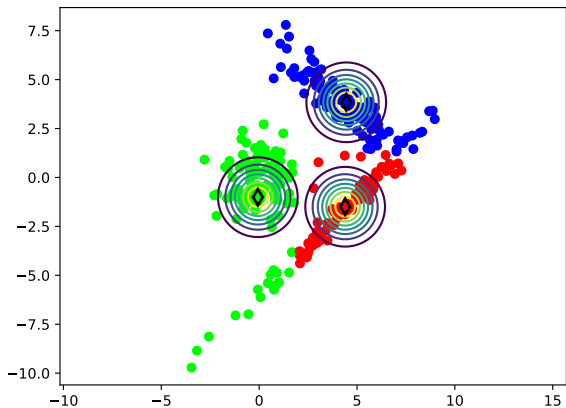
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



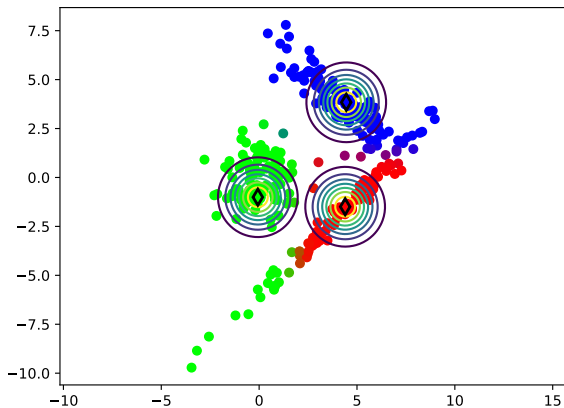
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



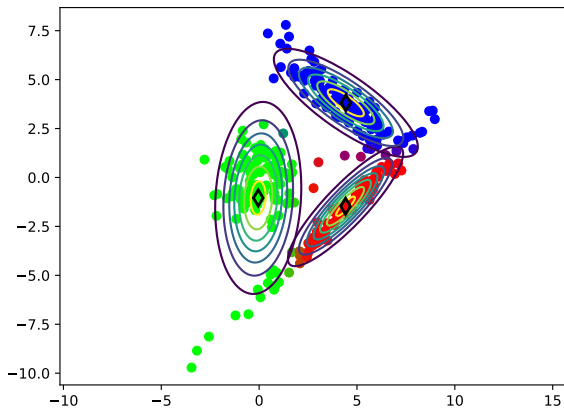
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



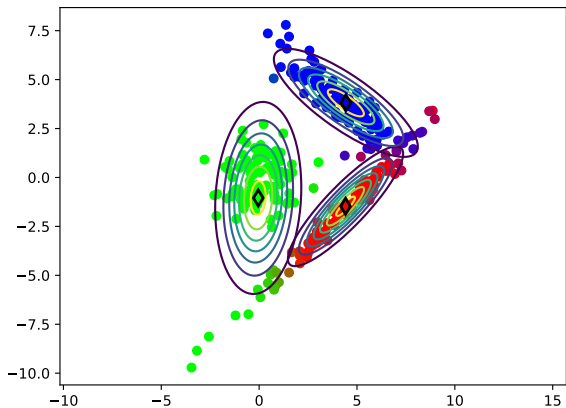
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



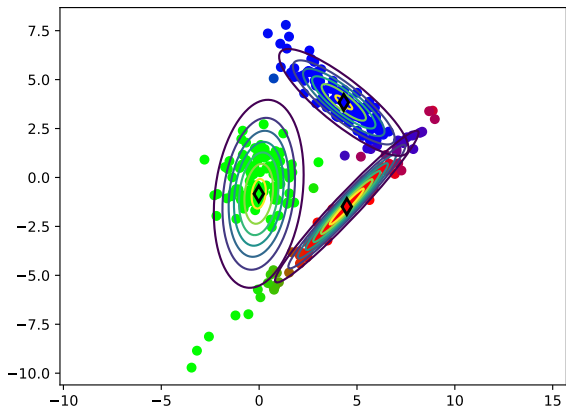
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



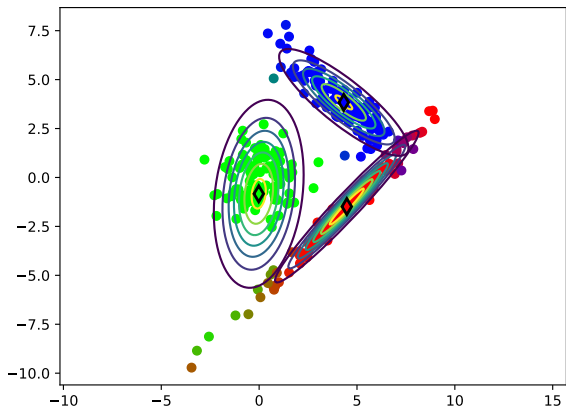
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



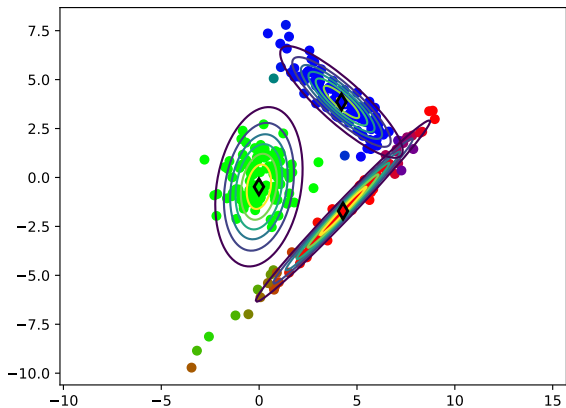
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



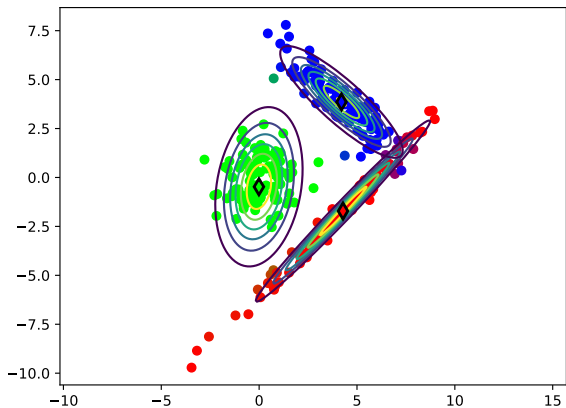
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



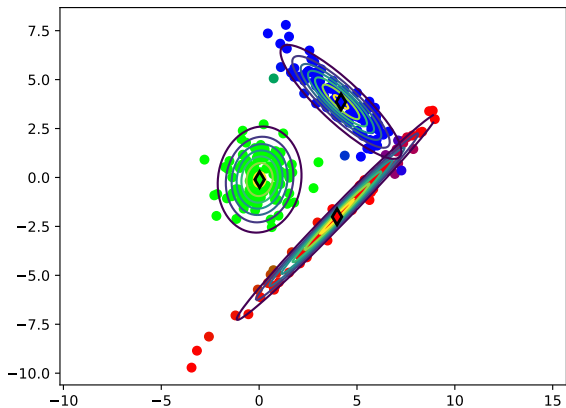
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



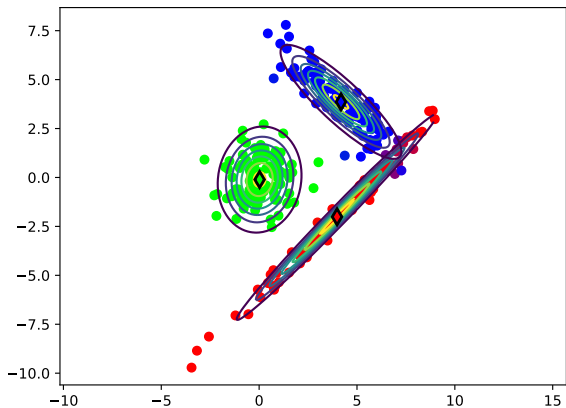
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



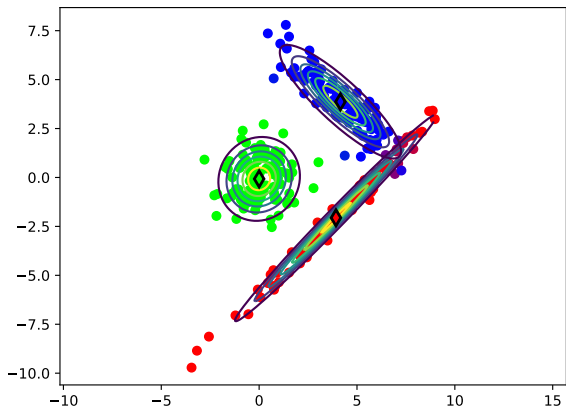
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



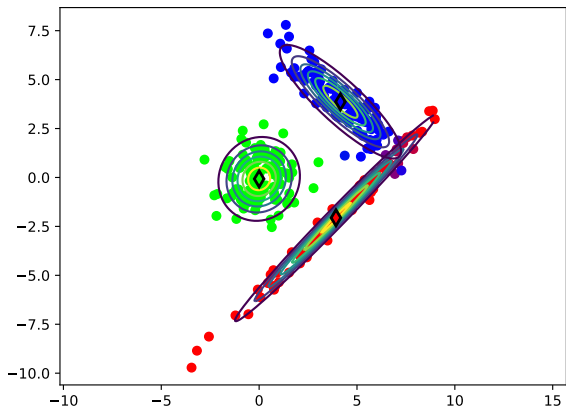
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



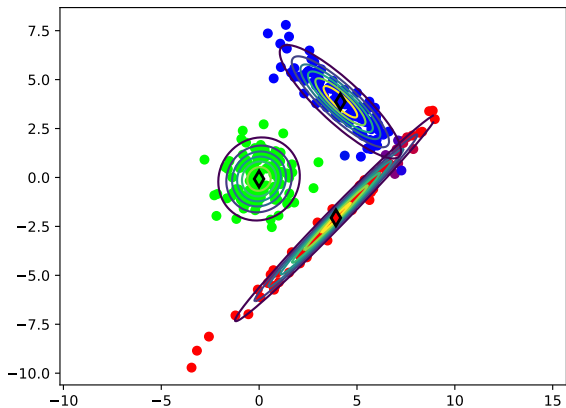
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



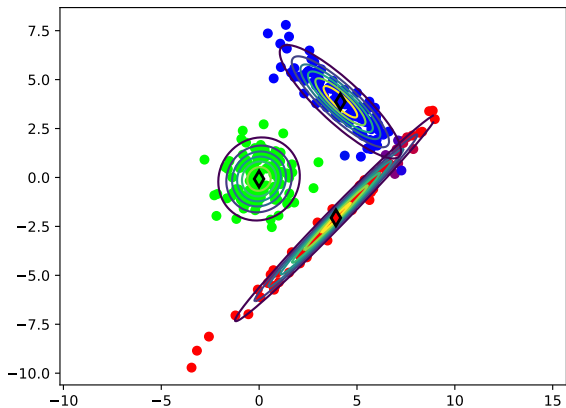
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



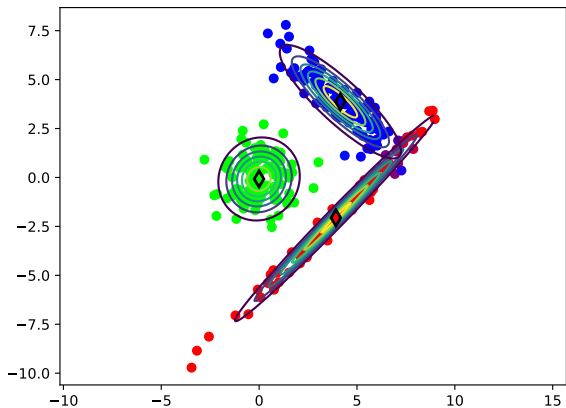
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



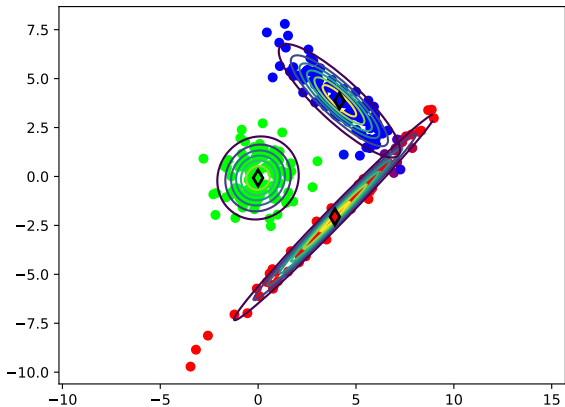
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



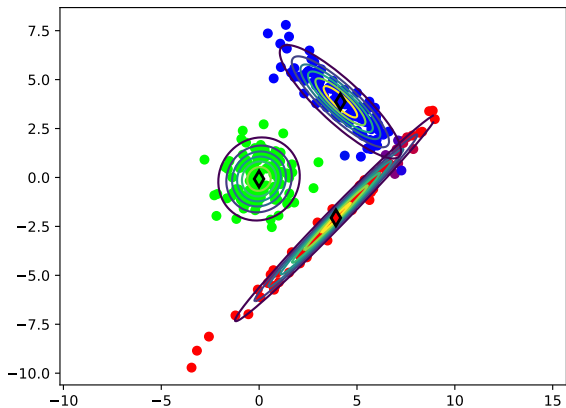
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



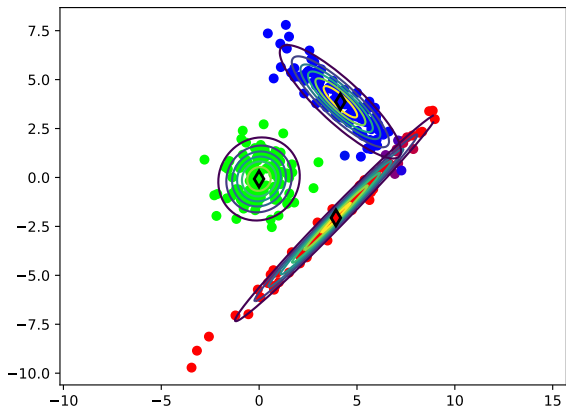
Demo: elliptical clusters

E... M... E... M... E... M... E... M...



Demo: elliptical clusters

E... M... E... M... E... M... E... M...



Theorem.

Suppose $(\mathbf{R}_0, \boldsymbol{\theta}_0) \in \mathcal{R}_{n,k} \times \Theta$ arbitrary,
thereafter $(\mathbf{R}_t, \boldsymbol{\theta}_t)$ given by E-M:

$$(\mathbf{R}_t)_{ij} := p_{\boldsymbol{\theta}_{t-1}}(y = j | \mathbf{x}_i). \quad \text{and} \quad \boldsymbol{\theta}_t := \arg \max_{\boldsymbol{\theta} \in \Theta} \mathcal{L}(\boldsymbol{\theta}; \mathbf{R}_t)$$

Then

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}_t; \mathbf{R}_t) &\leq \max_{\mathbf{R} \in \mathcal{R}_{n \times k}} \mathcal{L}(\boldsymbol{\theta}_t; \mathbf{R}) = \mathcal{L}(\boldsymbol{\theta}_t; \mathbf{R}_{t+1}) = \mathcal{L}(\boldsymbol{\theta}_t) \\ &\leq \mathcal{L}(\boldsymbol{\theta}_{t+1}; \mathbf{R}_{t+1}). \end{aligned}$$

In particular, $\mathcal{L}(\boldsymbol{\theta}_t) \leq \mathcal{L}(\boldsymbol{\theta}_{t+1})$.

Theorem.

Suppose $(\mathbf{R}_0, \boldsymbol{\theta}_0) \in \mathcal{R}_{n,k} \times \Theta$ arbitrary, thereafter $(\mathbf{R}_t, \boldsymbol{\theta}_t)$ given by E-M:

$$(\mathbf{R}_t)_{ij} := p_{\boldsymbol{\theta}_{t-1}}(y = j | \mathbf{x}_i). \quad \text{and} \quad \boldsymbol{\theta}_t := \arg \max_{\boldsymbol{\theta} \in \Theta} \mathcal{L}(\boldsymbol{\theta}; \mathbf{R}_t)$$

Then

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}_t; \mathbf{R}_t) &\leq \max_{\mathbf{R} \in \mathcal{R}_{n \times k}} \mathcal{L}(\boldsymbol{\theta}_t; \mathbf{R}) = \mathcal{L}(\boldsymbol{\theta}_t; \mathbf{R}_{t+1}) = \mathcal{L}(\boldsymbol{\theta}_t) \\ &\leq \mathcal{L}(\boldsymbol{\theta}_{t+1}; \mathbf{R}_{t+1}). \end{aligned}$$

In particular, $\mathcal{L}(\boldsymbol{\theta}_t) \leq \mathcal{L}(\boldsymbol{\theta}_{t+1})$.

Remarks.

- ▶ We proved a similar guarantee for k -means, which is also an *alternating minimization* scheme.
- ▶ Similarly, MLE for Gaussian mixtures is NP-hard; it is also known to need exponentially many samples in k to information-theoretically recover the parameters.

Proof. We've already shown:

- ▶ $\mathcal{L}(\boldsymbol{\theta}_t; \mathbf{R}_{t+1}) = \mathcal{L}(\boldsymbol{\theta}_t)$;
- ▶ $\mathcal{L}(\boldsymbol{\theta}_t; \mathbf{R}_{t+1}) = \max_{\boldsymbol{\theta} \in \Theta} \mathcal{L}(\boldsymbol{\theta}; \mathbf{R}_{t+1}) \leq \mathcal{L}(\boldsymbol{\theta}_{t+1}; \mathbf{R}_{t+1})$ by definition of $\boldsymbol{\theta}_{t+1}$.

We still need to show: $\mathcal{L}(\boldsymbol{\theta}_t; \mathbf{R}_{t+1}) = \max_{\mathbf{R} \in \mathcal{R}_{n,k}} \mathcal{L}(\boldsymbol{\theta}_{t+1}; \mathbf{R})$.

We'll give two proofs.

Proof. We've already shown:

- ▶ $\mathcal{L}(\boldsymbol{\theta}_t; \mathbf{R}_{t+1}) = \mathcal{L}(\boldsymbol{\theta}_t)$;
- ▶ $\mathcal{L}(\boldsymbol{\theta}_t; \mathbf{R}_{t+1}) = \max_{\boldsymbol{\theta} \in \Theta} \mathcal{L}(\boldsymbol{\theta}; \mathbf{R}_{t+1}) \leq \mathcal{L}(\boldsymbol{\theta}_{t+1}; \mathbf{R}_{t+1})$ by definition of $\boldsymbol{\theta}_{t+1}$.

We still need to show: $\mathcal{L}(\boldsymbol{\theta}_t; \mathbf{R}_{t+1}) = \max_{\mathbf{R} \in \mathcal{R}_{n,k}} \mathcal{L}(\boldsymbol{\theta}_{t+1}; \mathbf{R})$.

We'll give two proofs.

By concavity of \ln ("Jensen's inequality" in convexity lectures), for any $\mathbf{R} \in \mathcal{R}_{n,k}$,

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}_t; \mathbf{R}) &= \sum_{i=1}^n \sum_{j=1}^k \mathbf{R}_{ij} \ln \frac{p_{\boldsymbol{\theta}_t}(\mathbf{x}_i, y_i = j)}{\mathbf{R}_{ij}} \\ &\leq \sum_{i=1}^n \ln \left(\sum_{j=1}^k \mathbf{R}_{ij} \frac{p_{\boldsymbol{\theta}_t}(\mathbf{x}_i, y_i = j)}{\mathbf{R}_{ij}} \right) \\ &= \sum_{i=1}^n \ln p_{\boldsymbol{\theta}_t}(\mathbf{x}_i) = \mathcal{L}(\boldsymbol{\theta}_t) = \mathcal{L}(\boldsymbol{\theta}_t; \mathbf{R}_{t+1}). \end{aligned}$$

Since \mathbf{R} was arbitrary, $\max_{\mathbf{R} \in \mathcal{R}} \mathcal{L}(\boldsymbol{\theta}_t; \mathbf{R}) = \mathcal{L}(\boldsymbol{\theta}_t; \mathbf{R}_{t+1})$.

Proof (continued). Here's a second proof of that missing fact. To evaluate $\arg \max_{\mathbf{R} \in \mathcal{R}_{n,k}} \mathcal{L}(\boldsymbol{\theta}; \mathbf{R})$, consider Lagrangian

$$\sum_{i=1}^n \left(\sum_{j=1}^k R_{ij} \ln p_{\boldsymbol{\theta}}(\mathbf{x}_i, y = j) - \sum_{j=1}^k R_{ij} \ln R_{ij} + \lambda_i \left(\sum_{j=1}^k R_{ij} - 1 \right) \right).$$

Fixing i and taking the gradient with respect to R_{ij} for any j ,

$$0 = \ln p_{\boldsymbol{\theta}}(\mathbf{x}_i, y_i = j) - \ln R_{ij} - 1 + \lambda_i,$$

giving $R_{ij} = p_{\boldsymbol{\theta}}(\mathbf{x}_i, y = j) \exp(\lambda_i - 1)$. Since moreover

$$1 = \sum_j R_{ij} = \exp(\lambda_i - 1) \sum_j p_{\boldsymbol{\theta}}(\mathbf{x}_i, y = j) = \exp(\lambda_i - 1) p_{\boldsymbol{\theta}}(\mathbf{x}_i),$$

it follows that $\exp(\lambda_i - 1) = 1/p_{\boldsymbol{\theta}}(\mathbf{x}_i)$,

and the optimal \mathbf{R} satisfies $R_{ij} = p_{\boldsymbol{\theta}}(\mathbf{x}_i, y=j)/p_{\boldsymbol{\theta}}(\mathbf{x}_i) = p_{\boldsymbol{\theta}}(y = j | \mathbf{x}_i)$.



Related issues.

Parameter constraints.

E-M for GMMs **still works** if we freeze or constrain some parameters.

Parameter constraints.

E-M for GMMs **still works** if we freeze or constrain some parameters.

Examples:

- ▶ **No weights:** initialize $\boldsymbol{\pi} = (1/k, \dots, 1/k)$ and never update it.
- ▶ **Diagonal covariance matrices:** update everything as before, except $\boldsymbol{\Sigma}_j := \text{diag}((\boldsymbol{\sigma}_j)_1^2, \dots, (\boldsymbol{\sigma}_j)_d^2)$ where

$$(\boldsymbol{\sigma}_j)_l^2 := \frac{\sum_{i=1}^n R_{ij} (\mathbf{x}_i - \boldsymbol{\mu}_j)_l^2}{n\pi_j};$$

that is: we use coordinate-wise sample variances weighted by \mathbf{R} .

Why is this a good idea?

Parameter constraints.

E-M for GMMs **still works** if we freeze or constrain some parameters.

Examples:

- ▶ **No weights:** initialize $\boldsymbol{\pi} = (1/k, \dots, 1/k)$ and never update it.
- ▶ **Diagonal covariance matrices:** update everything as before, except $\boldsymbol{\Sigma}_j := \text{diag}((\boldsymbol{\sigma}_j)_1^2, \dots, (\boldsymbol{\sigma}_j)_d^2)$ where

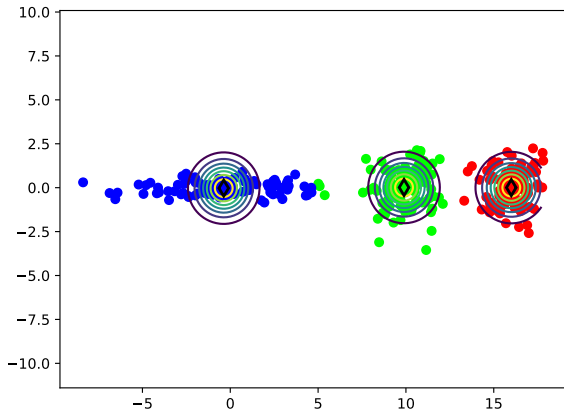
$$(\boldsymbol{\sigma}_j)_l^2 := \frac{\sum_{i=1}^n R_{ij} (\mathbf{x}_i - \boldsymbol{\mu}_j)_l^2}{n\pi_j};$$

that is: we use coordinate-wise sample variances weighted by \mathbf{R} .

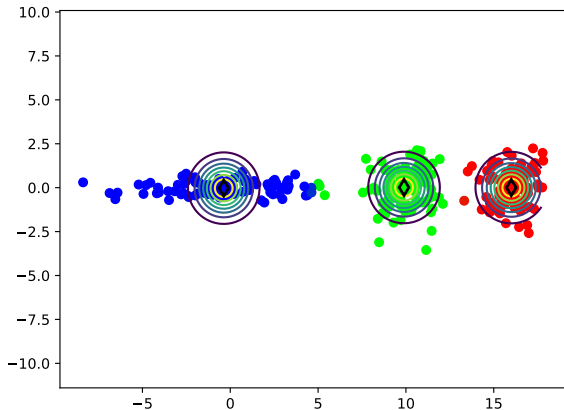
Why is this a good idea?

Computation (of inverse), sample complexity, ...

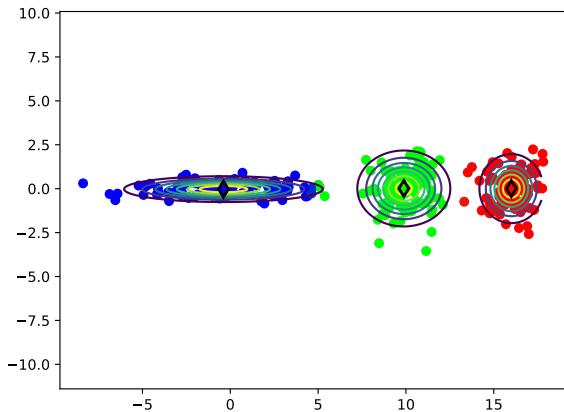
Gaussian Mixture Model with diagonal covariances.



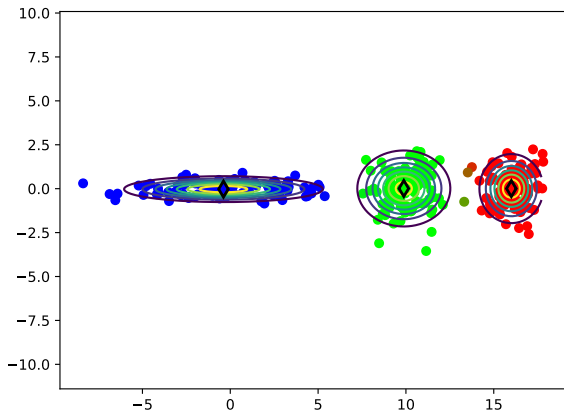
Gaussian Mixture Model with diagonal covariances.



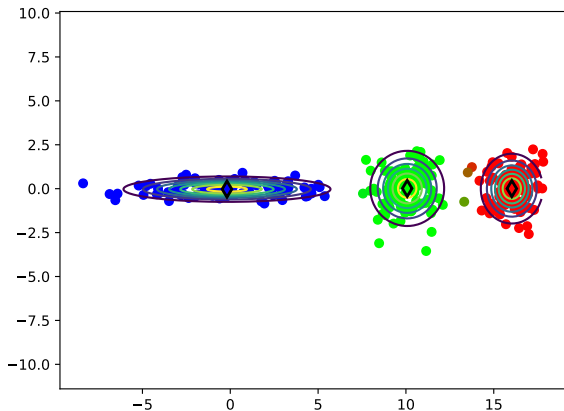
Gaussian Mixture Model with diagonal covariances.



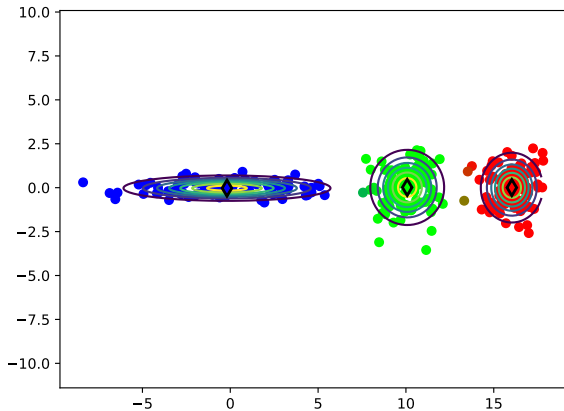
Gaussian Mixture Model with diagonal covariances.



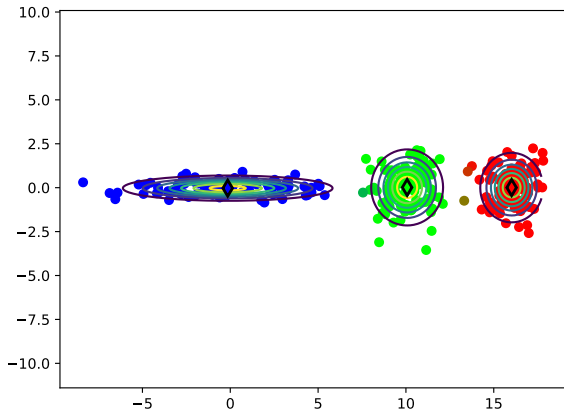
Gaussian Mixture Model with diagonal covariances.



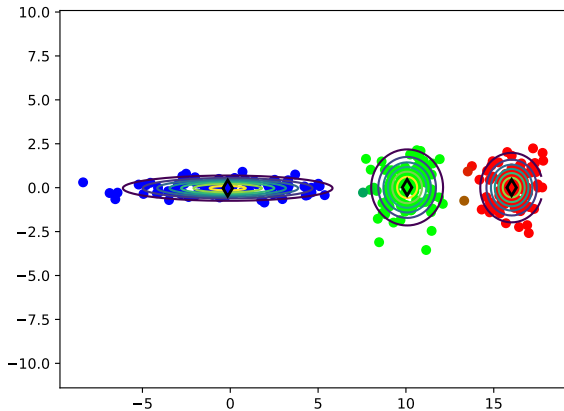
Gaussian Mixture Model with diagonal covariances.



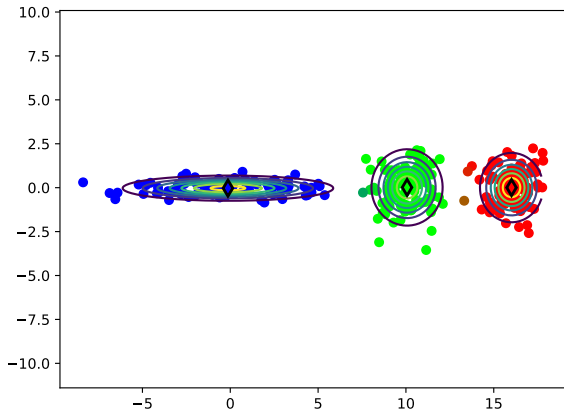
Gaussian Mixture Model with diagonal covariances.



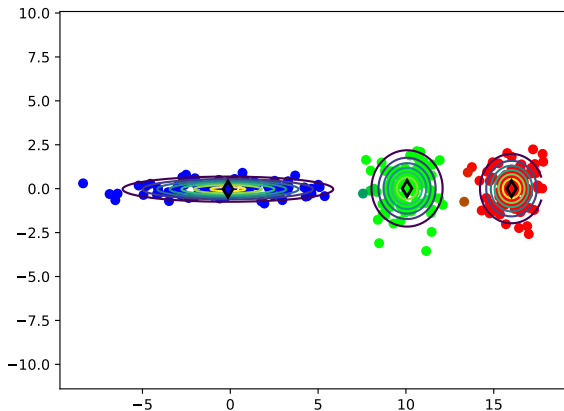
Gaussian Mixture Model with diagonal covariances.



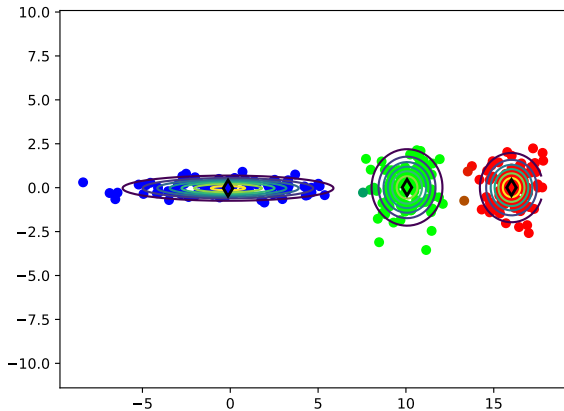
Gaussian Mixture Model with diagonal covariances.



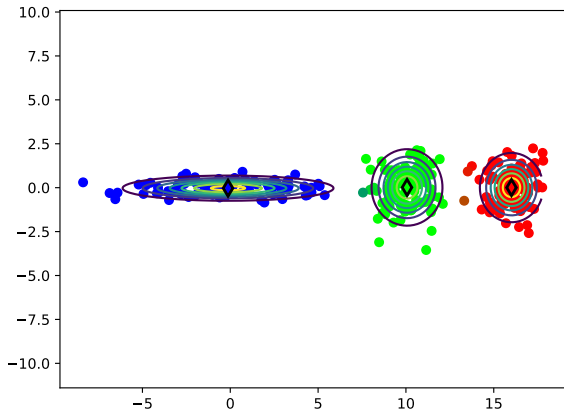
Gaussian Mixture Model with diagonal covariances.



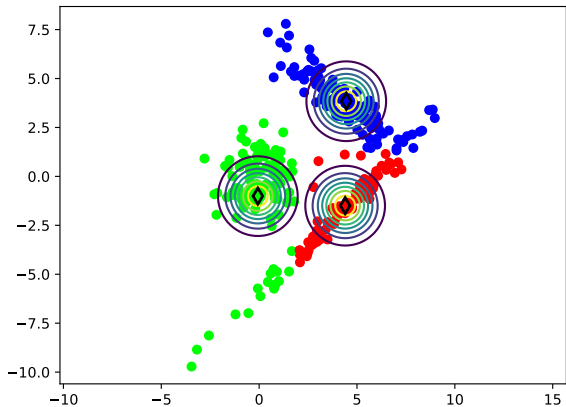
Gaussian Mixture Model with diagonal covariances.



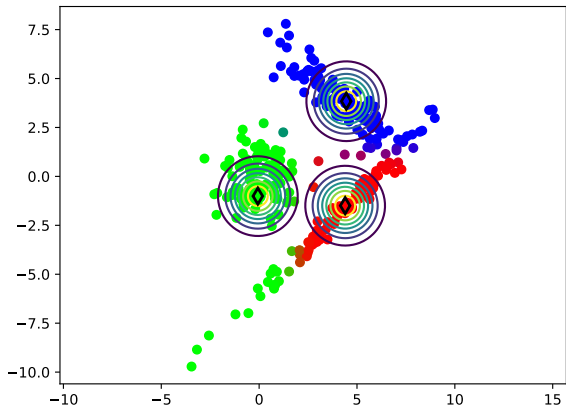
Gaussian Mixture Model with diagonal covariances.



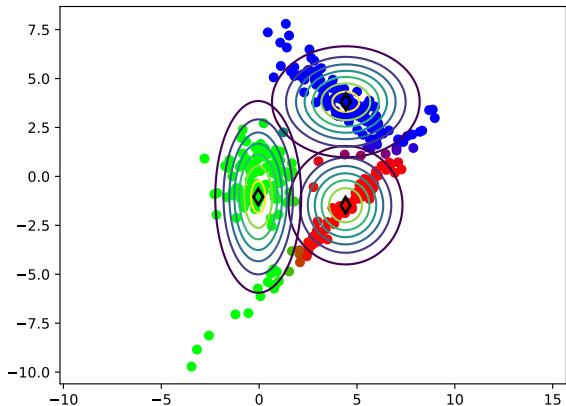
Gaussian Mixture Model with diagonal covariances.



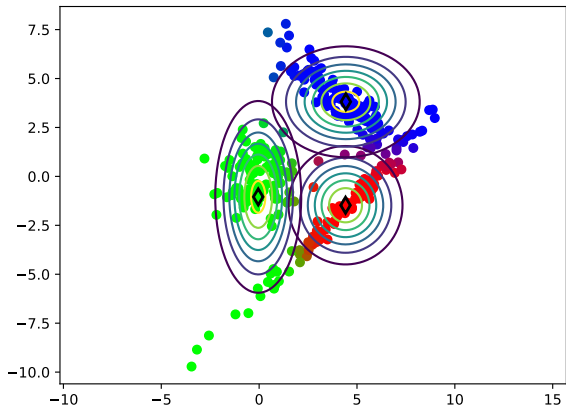
Gaussian Mixture Model with diagonal covariances.



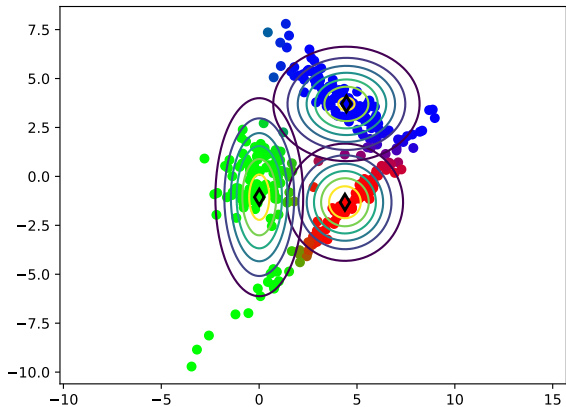
Gaussian Mixture Model with diagonal covariances.



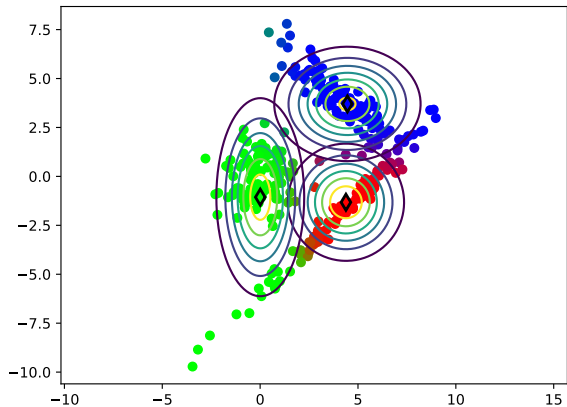
Gaussian Mixture Model with diagonal covariances.



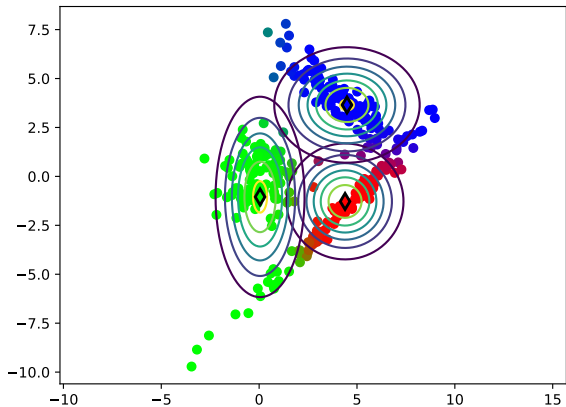
Gaussian Mixture Model with diagonal covariances.



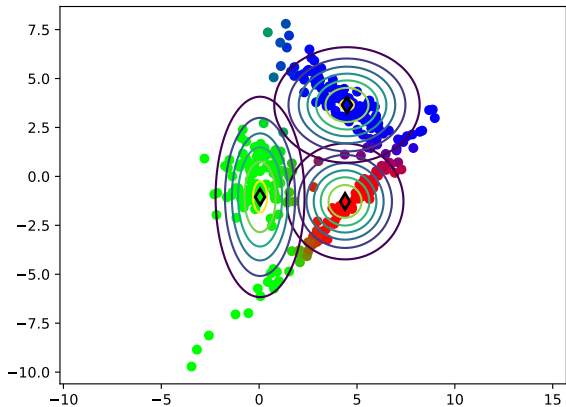
Gaussian Mixture Model with diagonal covariances.



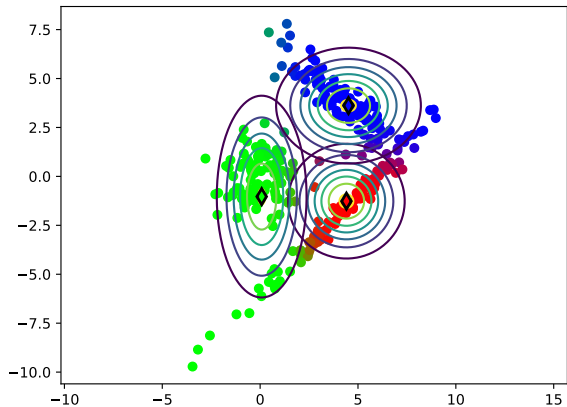
Gaussian Mixture Model with diagonal covariances.



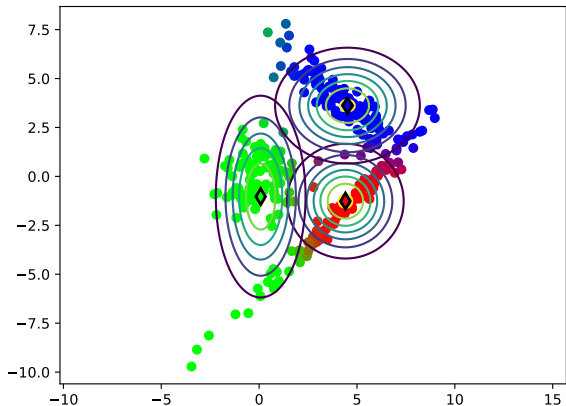
Gaussian Mixture Model with diagonal covariances.



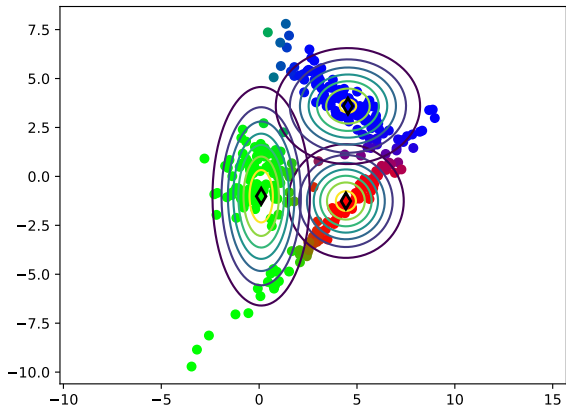
Gaussian Mixture Model with diagonal covariances.



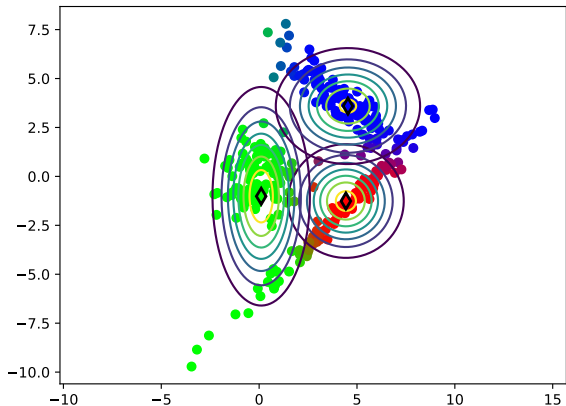
Gaussian Mixture Model with diagonal covariances.



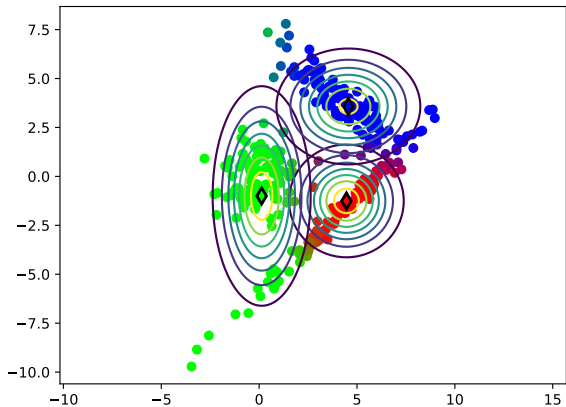
Gaussian Mixture Model with diagonal covariances.



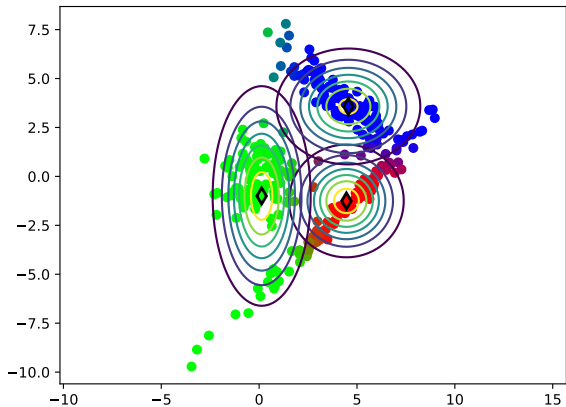
Gaussian Mixture Model with diagonal covariances.



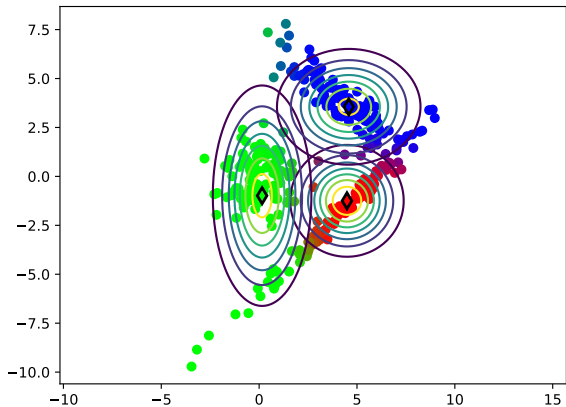
Gaussian Mixture Model with diagonal covariances.



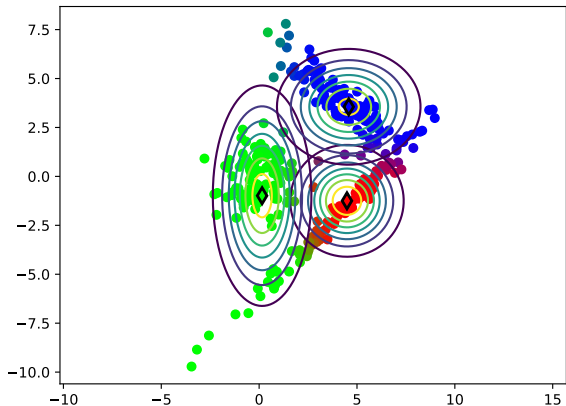
Gaussian Mixture Model with diagonal covariances.



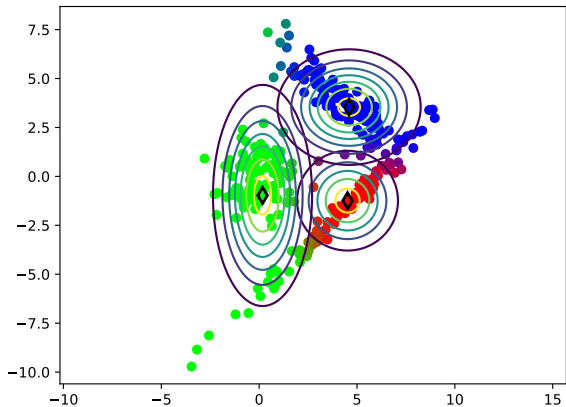
Gaussian Mixture Model with diagonal covariances.



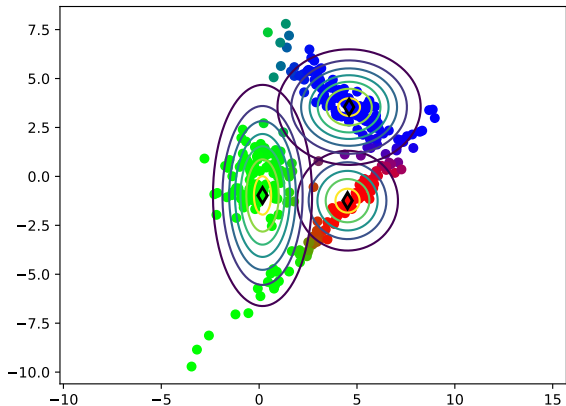
Gaussian Mixture Model with diagonal covariances.



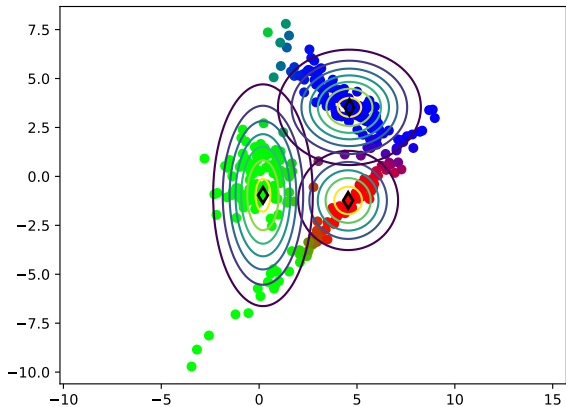
Gaussian Mixture Model with diagonal covariances.



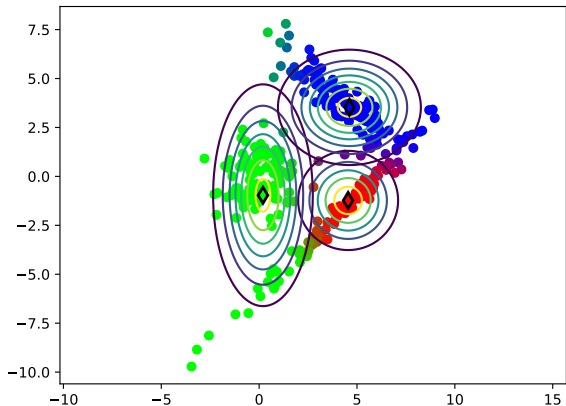
Gaussian Mixture Model with diagonal covariances.



Gaussian Mixture Model with diagonal covariances.



Gaussian Mixture Model with diagonal covariances.



E-M with GMMs suffers from **singularities**:
trivial situations where the likelihood goes to ∞ but the solution is bad.

► Suppose:

$$d = 1, k = 2, \pi_j = 1/2,$$

$$n = 3 \text{ with } x_1 = -1 \text{ and } x_2 = +1 \text{ and } x_3 = +3.$$

$$\text{Initialize with } \mu_1 = 0 \text{ and } \sigma_1 = 1,$$

$$\text{but } \mu_2 = +3 = x_3 \text{ and } \sigma_2 = 1/100.$$

$$\text{Then } \sigma_2 \rightarrow 0 \text{ and } \mathcal{L} \uparrow \infty.$$

Interpolating between k -means and GMM E-M

Same M-step: fix $\boldsymbol{\pi} = (1/k, \dots, 1/k)$ and $\boldsymbol{\Sigma}_j = c\mathbf{I}$ for a fixed $c > 0$.

Interpolating between k -means and GMM E-M

Same M-step: fix $\boldsymbol{\pi} = (1/k, \dots, 1/k)$ and $\boldsymbol{\Sigma}_j = c\mathbf{I}$ for a fixed $c > 0$.

Same E-step: define $q_{ij} := \frac{1}{2}\|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2$; the E-step chooses

$$\begin{aligned} R_{ij} &:= p_{\boldsymbol{\theta}}(y_i = j | \mathbf{x}_i) = \frac{p_{\boldsymbol{\theta}}(y_i = j, \mathbf{x}_i)}{p_{\boldsymbol{\theta}}(\mathbf{x}_i)} = \frac{p_{\boldsymbol{\theta}}(y_i = j, \mathbf{x}_i)}{\sum_{l=1}^k p_{\boldsymbol{\theta}}(y_i = l, \mathbf{x}_i)} \\ &= \frac{\pi_j p_{\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j}(\mathbf{x}_i)}{\sum_{l=1}^k \pi_l p_{\boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l}(\mathbf{x}_i)} = \frac{\exp(-q_{ij}/c)}{\sum_{l=1}^k \exp(-q_{il}/c)} \end{aligned}$$

Fix $i \in \{1, \dots, n\}$ and suppose minimum $q_i := \min_j q_{ij}$ is unique:

Interpolating between k -means and GMM E-M

Same M-step: fix $\pi = (1/k, \dots, 1/k)$ and $\Sigma_j = c\mathbf{I}$ for a fixed $c > 0$.

Same E-step: define $q_{ij} := \frac{1}{2} \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2$; the E-step chooses

$$\begin{aligned} R_{ij} &:= p_{\boldsymbol{\theta}}(y_i = j | \mathbf{x}_i) = \frac{p_{\boldsymbol{\theta}}(y_i = j, \mathbf{x}_i)}{p_{\boldsymbol{\theta}}(\mathbf{x}_i)} = \frac{p_{\boldsymbol{\theta}}(y_i = j, \mathbf{x}_i)}{\sum_{l=1}^k p_{\boldsymbol{\theta}}(y_i = l, \mathbf{x}_i)} \\ &= \frac{\pi_j p_{\boldsymbol{\mu}_j, \Sigma_j}(\mathbf{x}_i)}{\sum_{l=1}^k \pi_l p_{\boldsymbol{\mu}_l, \Sigma_l}(\mathbf{x}_i)} = \frac{\exp(-q_{ij}/c)}{\sum_{l=1}^k \exp(-q_{il}/c)} \end{aligned}$$

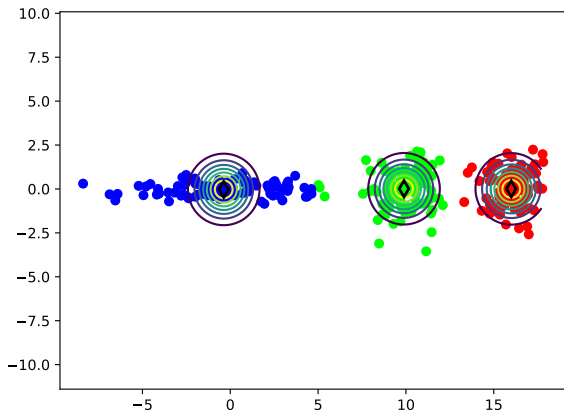
Fix $i \in \{1, \dots, n\}$ and suppose minimum $q_i := \min_j q_{ij}$ is unique:

$$\begin{aligned} \lim_{c \downarrow 0} R_{ij} &= \lim_{c \downarrow 0} \frac{\exp(-q_{ij}/c)}{\sum_{l=1}^k \exp(-q_{il}/c)} = \lim_{c \downarrow 0} \frac{\exp(q_i - q_{ij}/c)}{\sum_{l=1}^k \exp(q_i - q_{il}/c)} \\ &= \begin{cases} 1 & q_{ij} = q_i, \\ 0 & q_{ij} \neq q_i. \end{cases} \end{aligned}$$

That is, R becomes hard assignment A as $c \downarrow 0$.

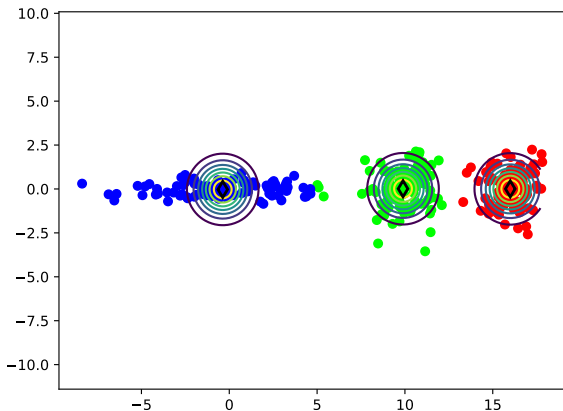
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



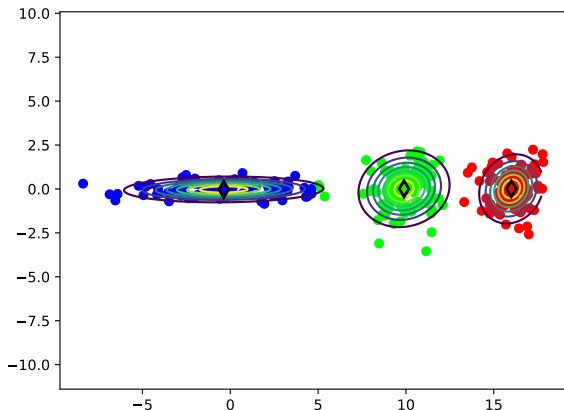
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



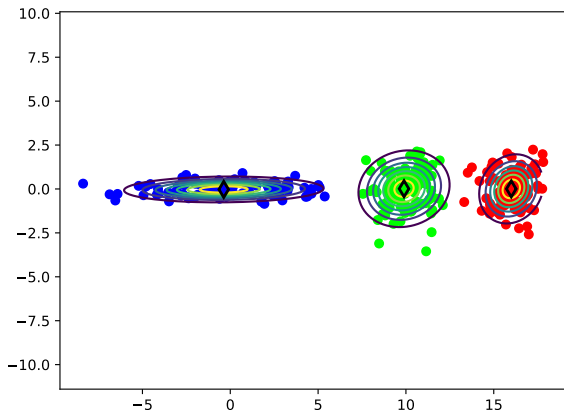
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



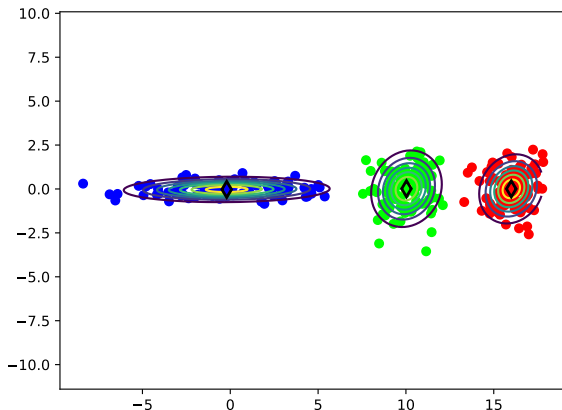
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**,
meaning we can create algorithms that have elements of both.
Here's something like k -means but with weights and covariances.



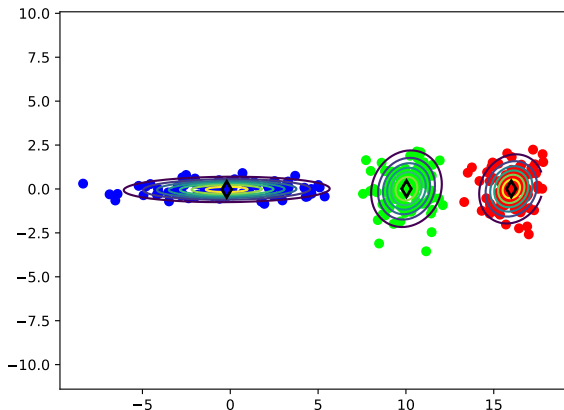
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



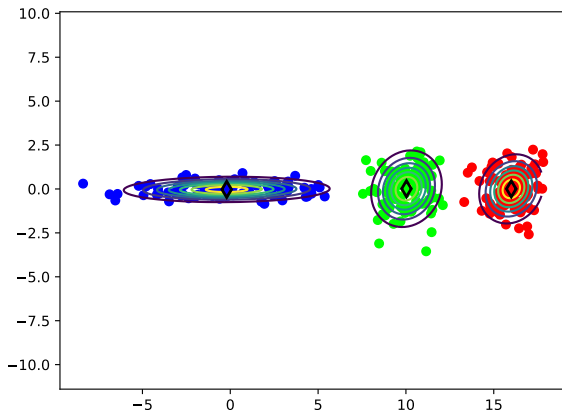
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



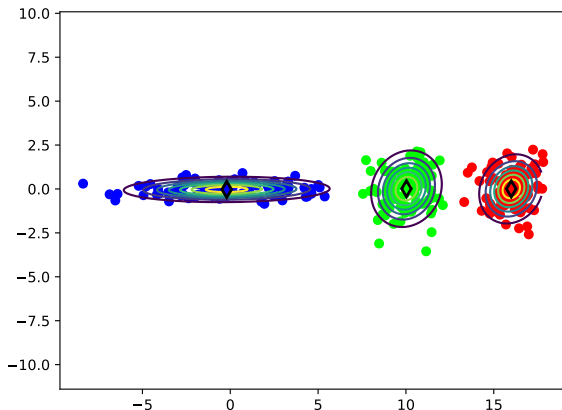
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



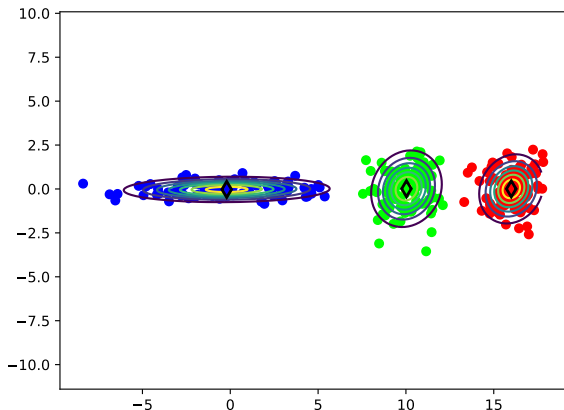
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



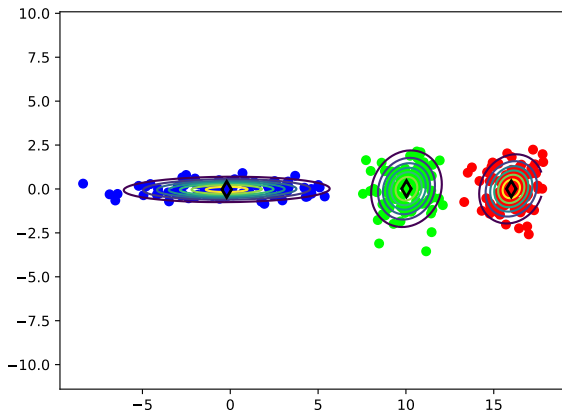
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



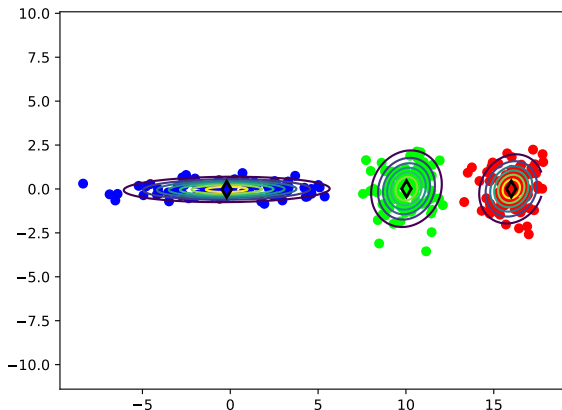
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



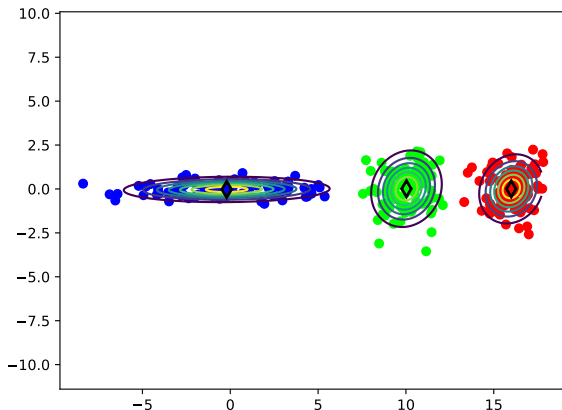
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



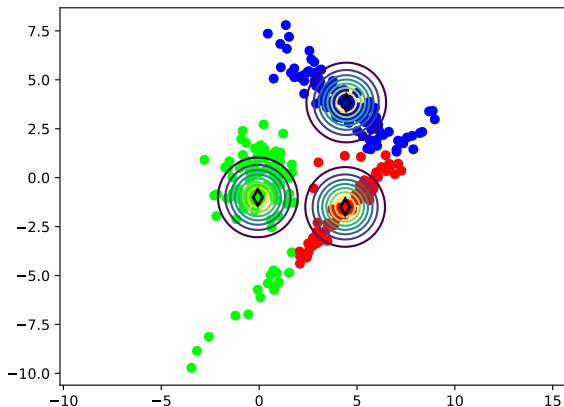
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



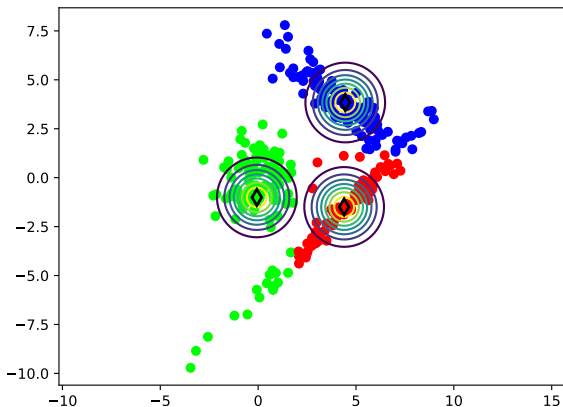
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



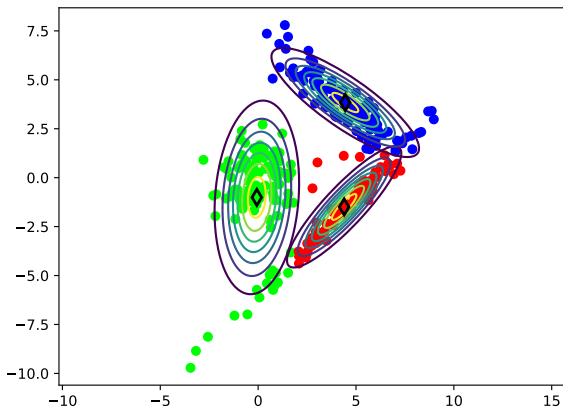
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



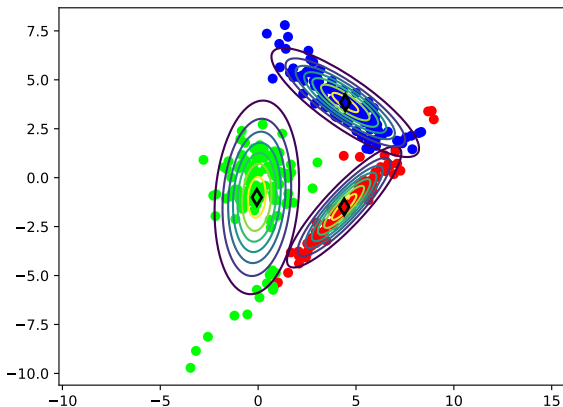
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



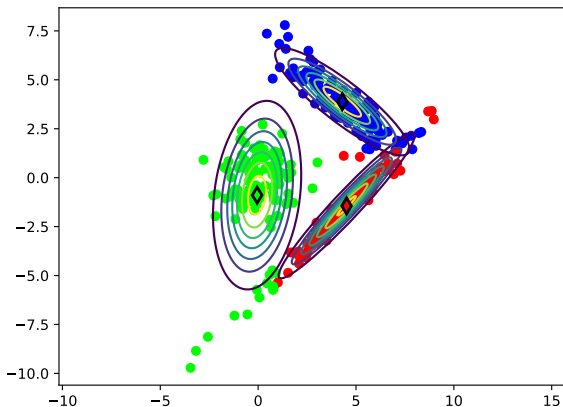
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



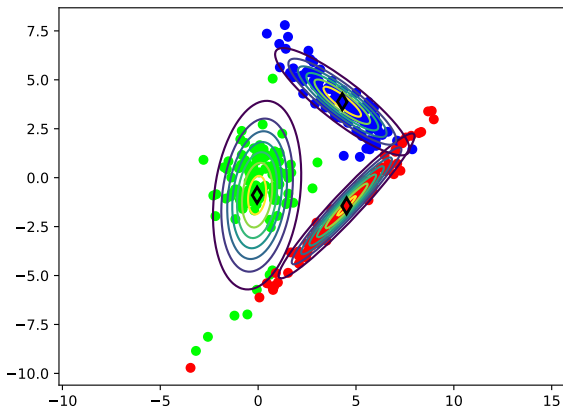
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



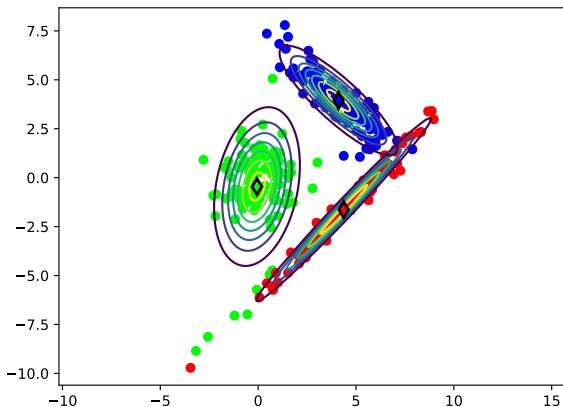
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



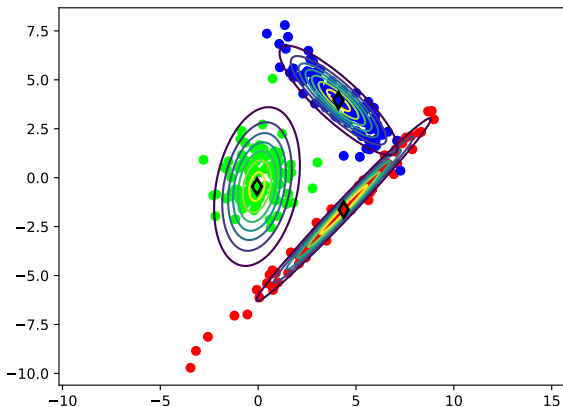
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



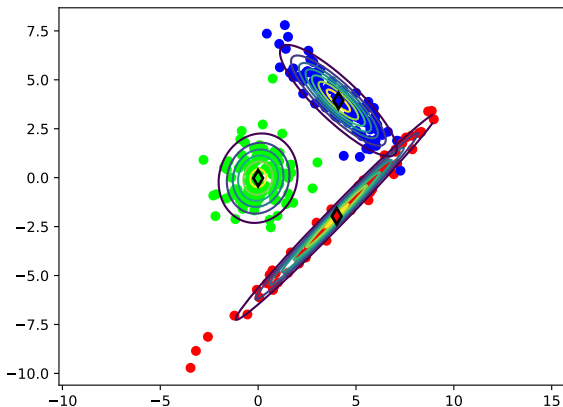
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



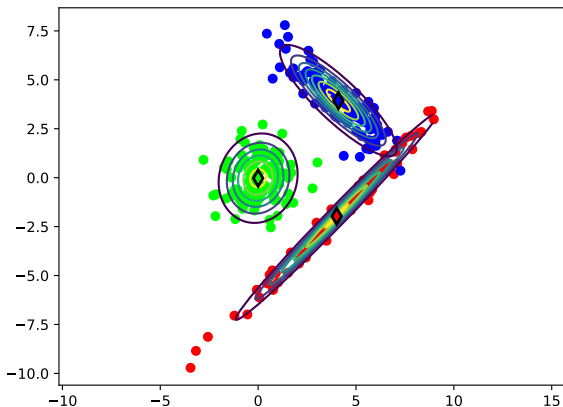
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



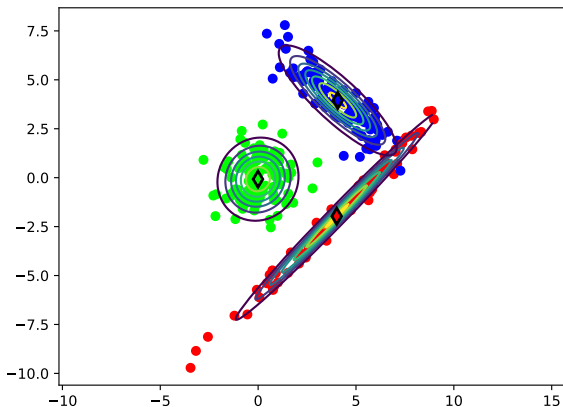
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



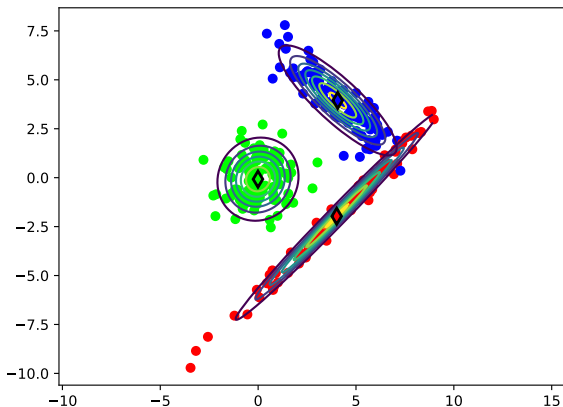
Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



Interpolating between k -means and GMM E-M (part 2)

We can interpolate **algorithmically**, meaning we can create algorithms that have elements of both. Here's something like k -means but with weights and covariances.



Summary of MLE part 2