

K-means and Expectation Maximization

- latent variable models
 - k-means
- expectation maximization

Unsupervised learning

- so far we've focused on supervised learning
- supervised* machine learning is about **learning functions from data**

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$

- unsupervised* machine learning is about finding interesting patterns in data

$$\mathcal{D} = \{(x_i)\}_{i=1}^n$$

- knowledge discovery*
- pattern recognition*
- under-specified

- learn $p(\mathbf{x})$ instead of $p(\mathbf{x}, y)$ or $p(y|\mathbf{x})$
- learn $p_\theta(\mathbf{x})$

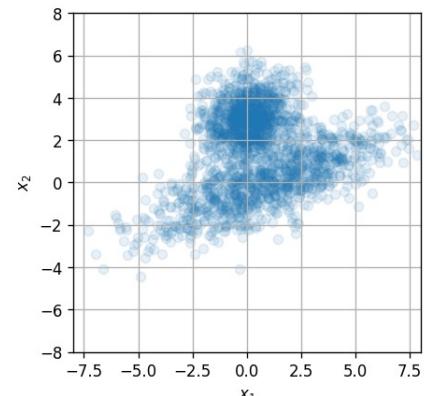
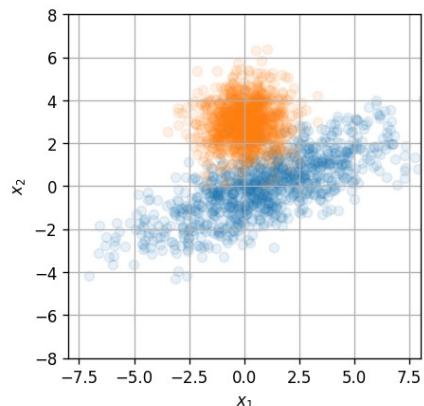
- examples:

density estimation

clustering

community detection in graphs

anomaly detection



Drawing Faces (or MNIST)



Ian Goodfellow @goodfellow_jian · Jan 14, 2019

4.5 years of GAN progress on face generation. arxiv.org/abs/1406.2661
arxiv.org/abs/1511.06434 arxiv.org/abs/1606.07536
arxiv.org/abs/1710.10196 arxiv.org/abs/1812.04948



42

1.4K

3.8K

↑

Variational Auto-Encoder (VAE)

0	4	6	1	6	0	9	4	2	0
8	6	0	7	6	8	3	0	6	8
0	3	6	1	1	9	2	0	9	2
0	2	4	9	3	4	9	6	0	9
3	3	1	2	1	3	8	7	0	2
8	9	0	6	2	3	8	5	6	9
5	0	7	1	6	8	3	3	4	2
6	2	6	0	9	2	5	7	9	6
1	3	7	8	6	1	6	4	1	0
4	0	7	0	0	8	0	0	0	1

Latent Variable Models

- goal: learn $p(\mathbf{x})$ from $\mathcal{D} = \{(\mathbf{x}_i)\}_{i=1}^n$

$$p(\mathbf{x}) = \int p(\mathbf{x}, z) dz$$

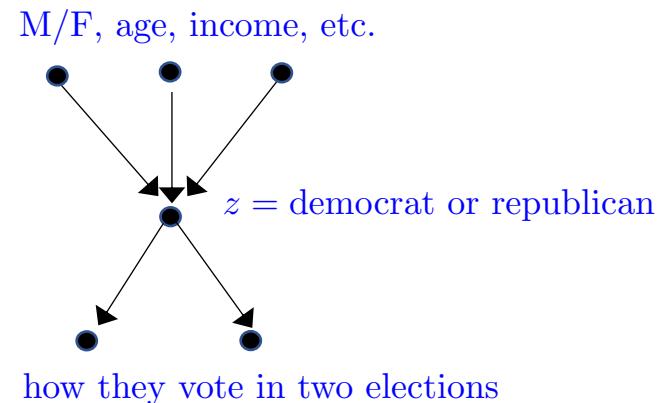
$$p(\mathbf{x}) = \sum_z p(\mathbf{x}, z)$$

$$p(\mathbf{x}) = \int p(z) p(\mathbf{x}|z) dz$$

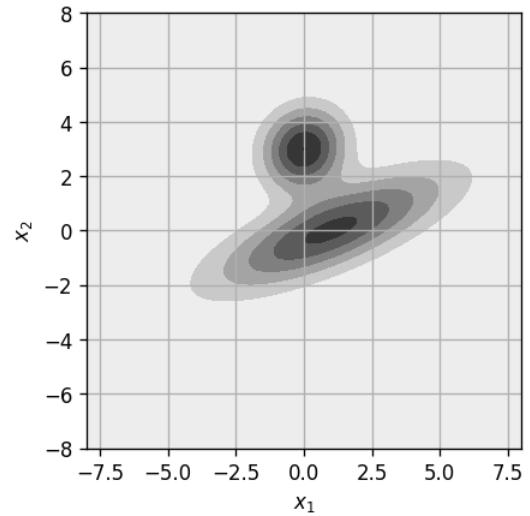
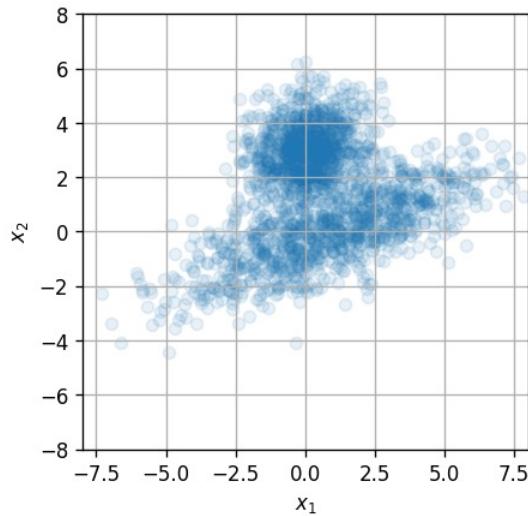
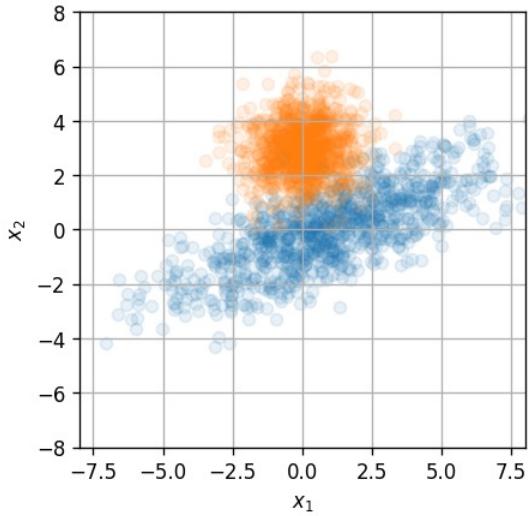
$$p(\mathbf{x}) = \sum_z p(z) p(\mathbf{x}|z)$$

- z is a *latent* variable

- *might* be mathematically convenient
- maybe $p(\mathbf{x}|z)$ has a simple form
- z might be meaningful, or maybe just convenient
- MNIST: z = digit, pen width, slant, etc.
- *latent* means hidden or concealed
 - if z is discrete, then *mixture* model



Gaussian Mixture Models (GMM)



- Gaussian mixture model:

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

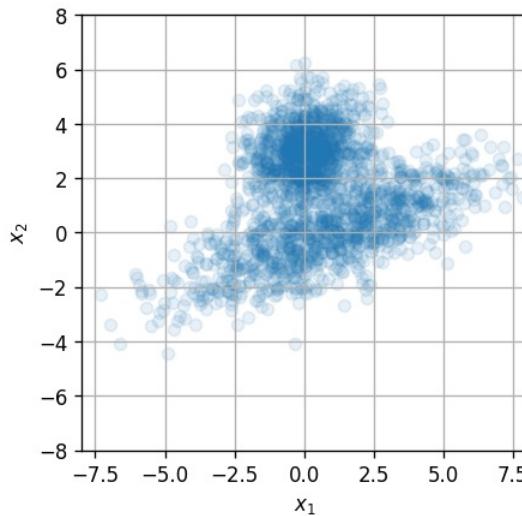
$$\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{e^{-(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}-\boldsymbol{\mu})/2}}{(2\pi)^{n/2} \sqrt{\det \boldsymbol{\Sigma}}}$$

- need to learn:

$$\pi_1, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \pi_2, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2, \dots$$

- uses of GMMs: density estimation, clustering

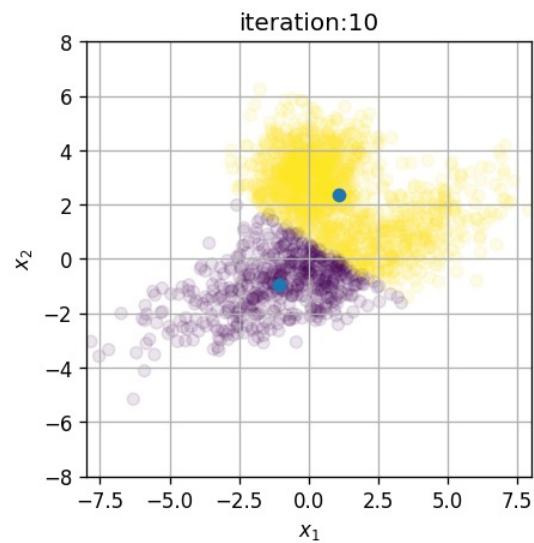
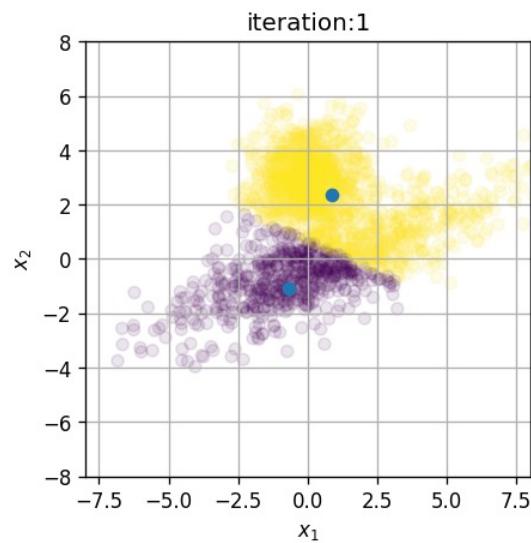
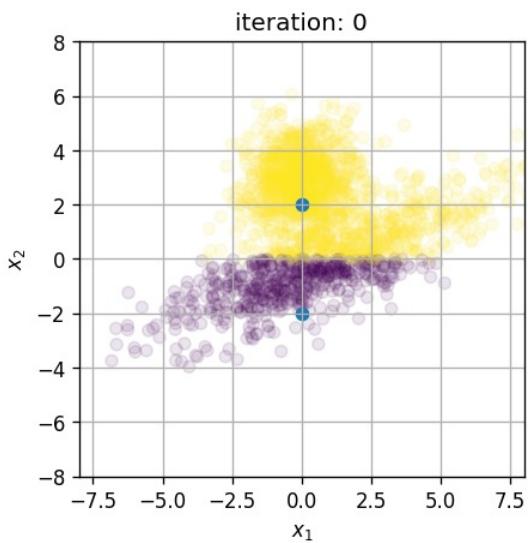
Clustering with k-Means



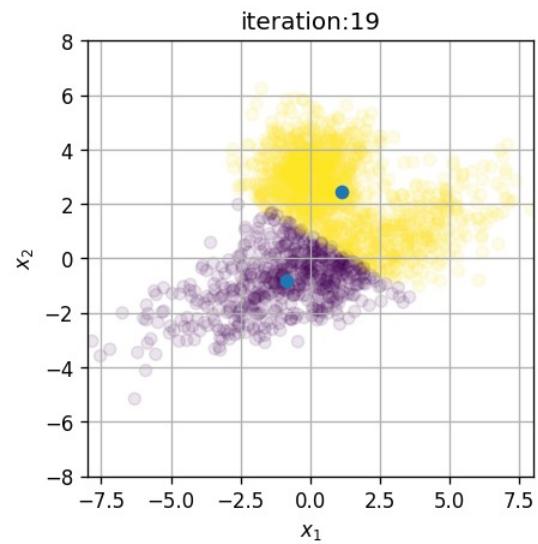
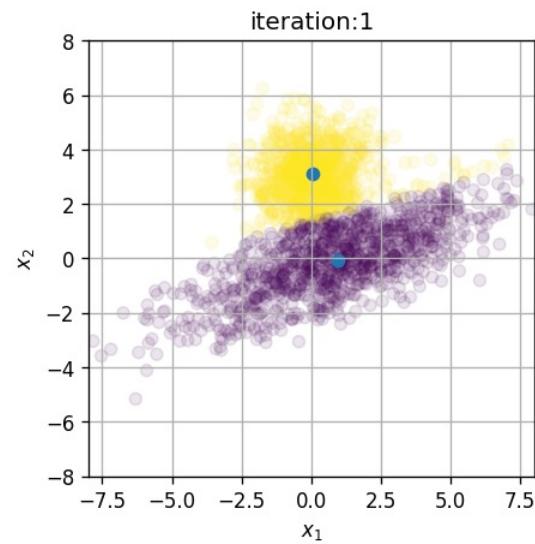
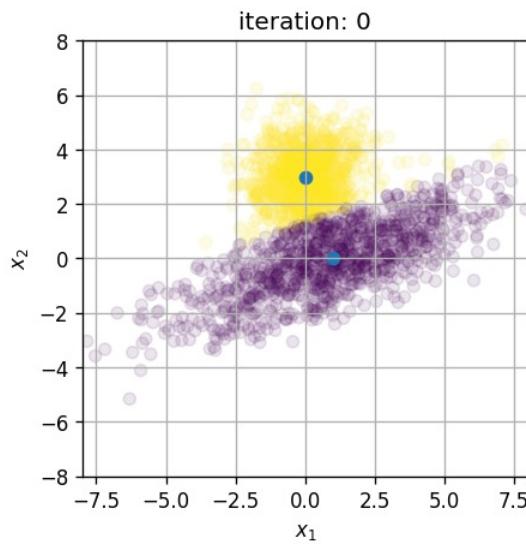
initialize cluster centers $\boldsymbol{\mu}_k$, $k = 1, \dots, K$

1. assign each point to the closest cluster center: $z_i = \operatorname{argmin}_k \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_2^2$
2. update cluster centers: $\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{i:z_i=k} \mathbf{x}_i$
3. repeat 1, 2

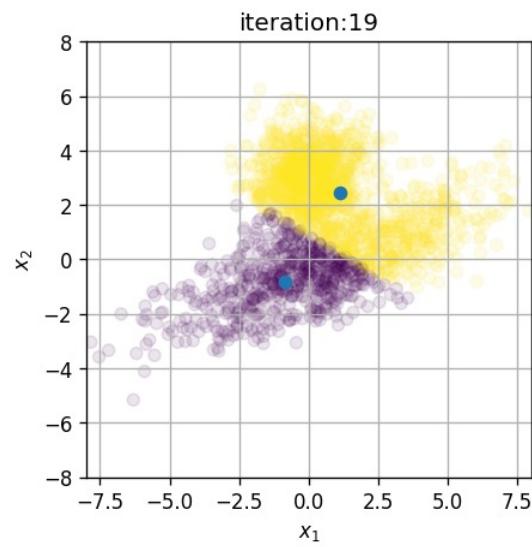
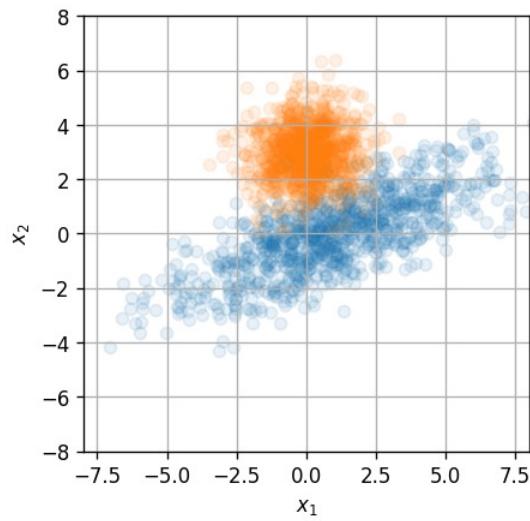
K-means



K-means



K-means



Expectation Maximization (EM) for GMMs

$$\text{Bayes: } p(z|\mathbf{x}) = \frac{p(z)p(\mathbf{x}|z)}{p(\mathbf{x})}$$

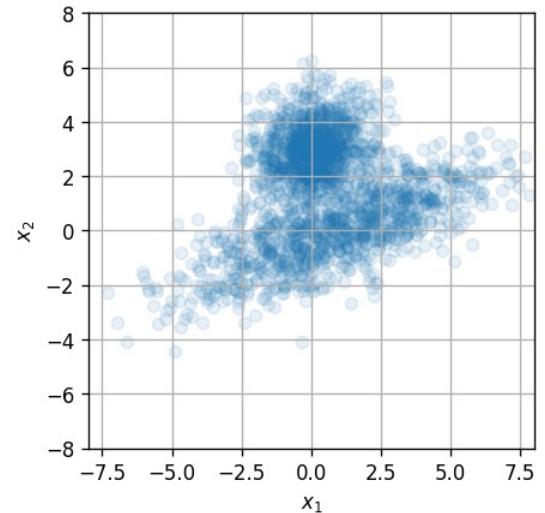
initialize means: $\boldsymbol{\mu}_k, j = 1, \dots, K$

initialize covariances: $\boldsymbol{\Sigma}_k = \mathbf{I}$

initialize priors: $\pi_k = \frac{1}{K}$

E-step: 1. compute the responsibility of each point to each cluster (*soft* assign each point):

$$p(z = k|\mathbf{x}_i) = r_{ik} = \frac{\pi_k \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k'} \pi_{k'} \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'})}$$

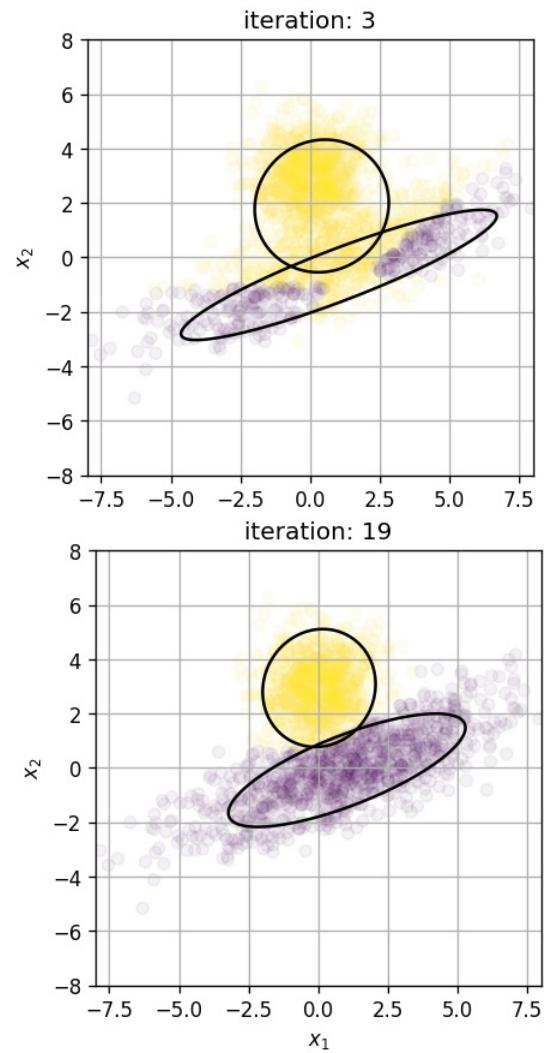
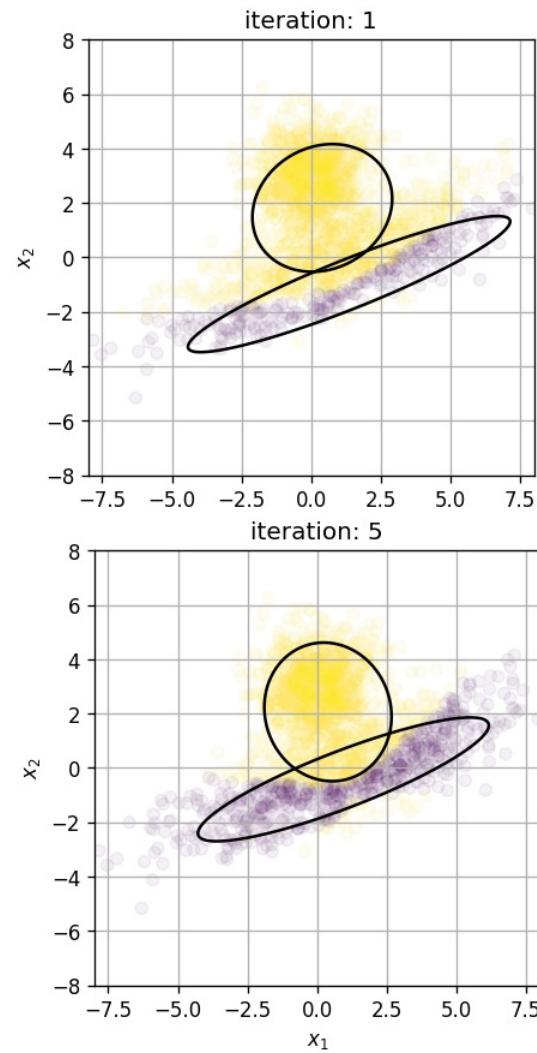
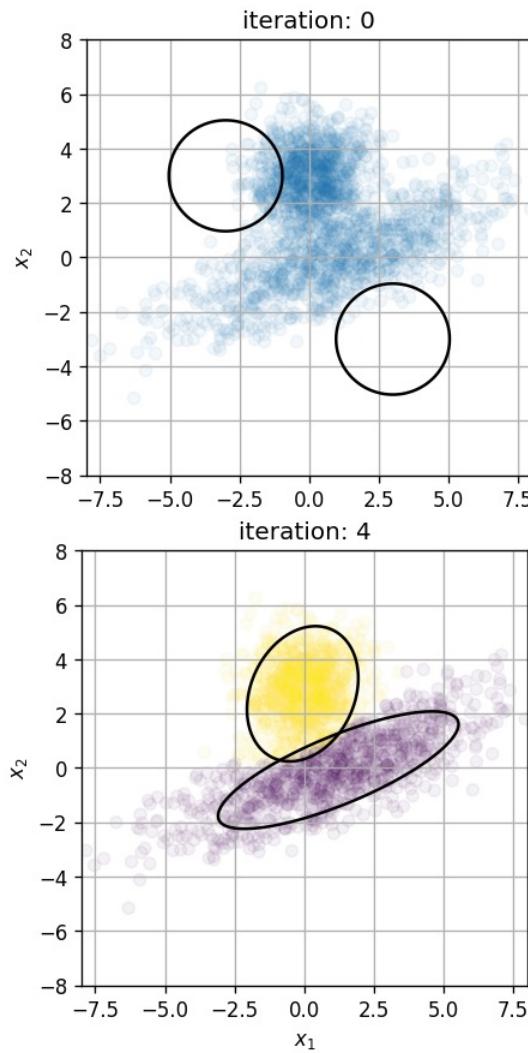


M-step: 2. estimate prior, mean and covariance for each mixture component

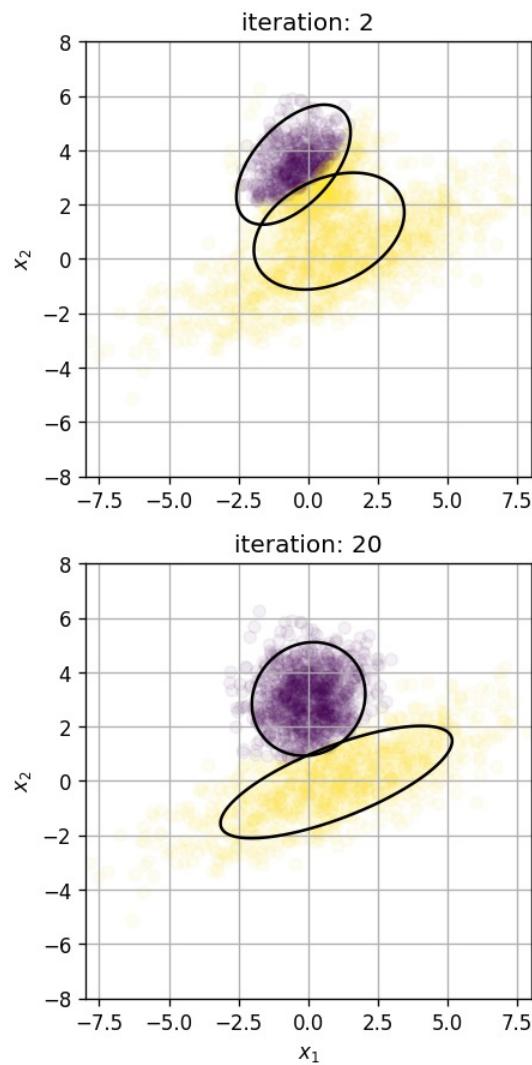
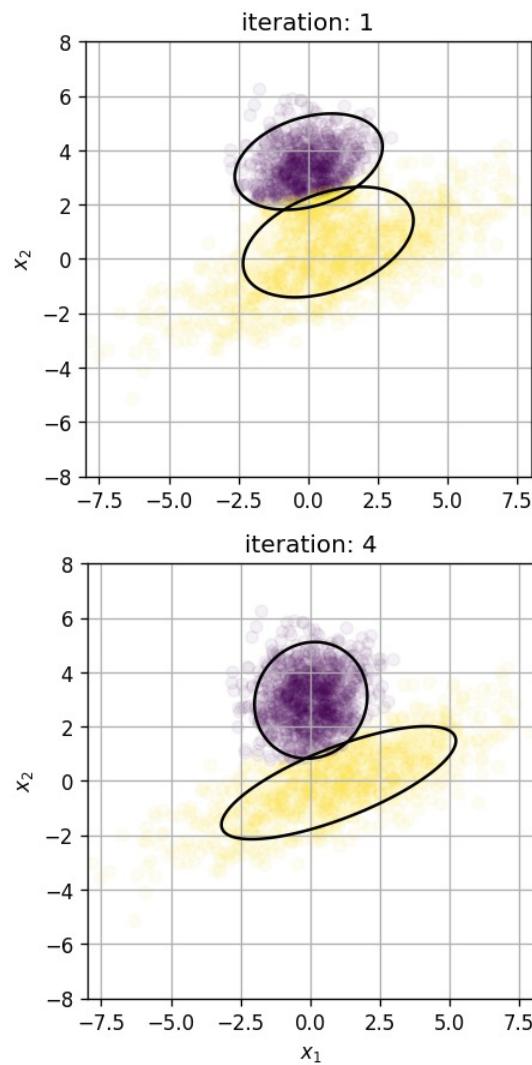
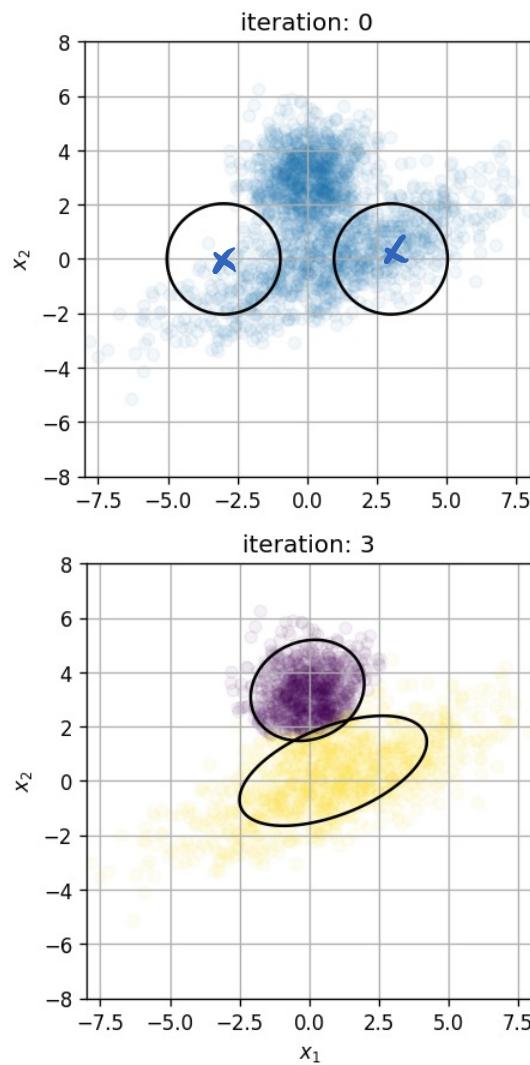
$$\pi_k = \frac{1}{N} \sum_i \mathbb{I}\{z_i = k\} \quad \pi_k = \frac{1}{N} \sum_i E[\mathbb{I}\{z_i = k\}] = \frac{1}{N} \sum_i p(z = k|\mathbf{x}_i)$$

$$\pi_k = \frac{1}{N} \sum_i r_{ik} \quad \boldsymbol{\mu}_k = \frac{1}{\sum_i r_{ik}} \sum_i r_{ik} \mathbf{x}_i \quad \boldsymbol{\Sigma}_k = \frac{1}{\sum_i r_{ik}} \sum_i r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T$$

3. repeat 1, 2



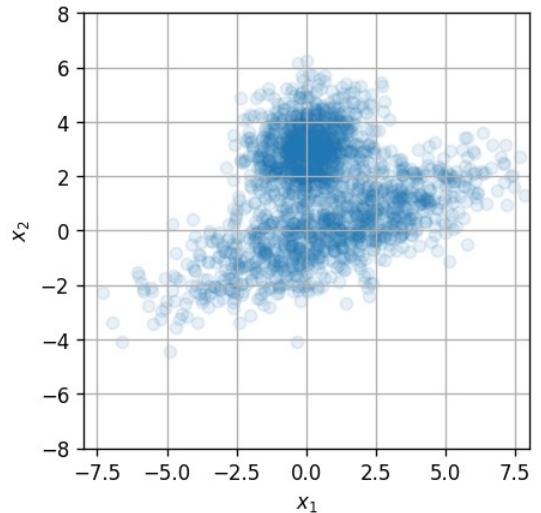
EM



Expectation Maximization (EM)

- observe $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$
- goal: maximize the log-likelihood:

$$\begin{aligned}\log L(\theta) &= \sum_{i=1}^n \log p_\theta(\mathbf{x}_i) = \sum_{i=1}^n \log \sum_{z_i} p_\theta(z_i, \mathbf{x}_i) \\ &= \sum_{i=1}^n \log \sum_{z_i} p(z_i) p_\theta(\mathbf{x}_i | z_i)\end{aligned}$$



E-step:

- for each data point, compute $p_{\theta_{t-1}}(z_i | \mathbf{x}_i)$.
i.e., compute the posterior of the latent variable for fixed θ_{t-1}

M-step:

- maximum likelihood for parameters θ given the expected value of z_i from E-step

$$\theta_t = \arg \max_{\theta} \sum_i E_{z \sim p_{\theta_t}(z | \mathbf{x}_i)} \log p_\theta(\mathbf{x}_i, z)$$

- repeat

K-means and EM

- K means

initialize cluster centers $\boldsymbol{\mu}_k$, $k = 1, \dots, K$

1. assign each point to the closest cluster center: $z_i = \arg \min_k \|\mathbf{x} - \boldsymbol{\mu}_k\|_2^2$
2. update cluster centers: $\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{i:z_i=k} \mathbf{x}_i$
3. repeat 1, 2

- K means and EM

fix covariances: $\boldsymbol{\Sigma}_k = \mathbf{I}$

fix priors: $\pi_k = \frac{1}{K}$

E-step: 1. compute the **most** responsible cluster for each point

$$z_i = \arg \max_k \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k'} \pi_{k'} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'})} \quad \text{hard assignment to each cluster}$$

M-step:

2. estimate mean for each mixture component
3. repeat 1, 2

EM

```
from sklearn import mixture
gmm = mixture.GaussianMixture(n_components=2, covariance_type='full')
gmm.fit(X_rows)

print('covariances: ', gmm.covariances_)
print('means: ', gmm.means_)
```

```
covariances:  
[[[8.22408446 3.03067247]  
 [3.03067247 1.97311564]]  
  
 [[1.04096518 0.073664 ]  
 [0.073664 1.18887948]]]  
  
means:  
[[ 1.02671502 -0.08439335]  
 [ 0.02955104 2.94435719]]
```

$$\mathcal{N} \left(\begin{bmatrix} 0 \\ 3 \end{bmatrix}, \begin{bmatrix} 1 & 0.1 \\ 0.1 & 1 \end{bmatrix} \right)$$
$$\mathcal{N} \left(\begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 8 & 3 \\ 3 & 2 \end{bmatrix} \right)$$

