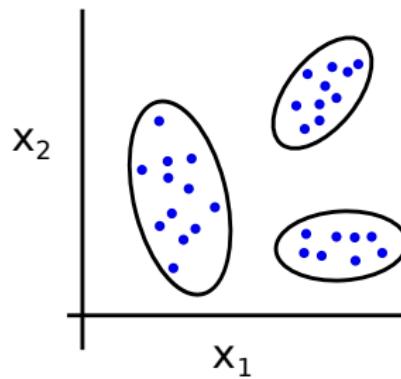


Clustering and the EM algorithm

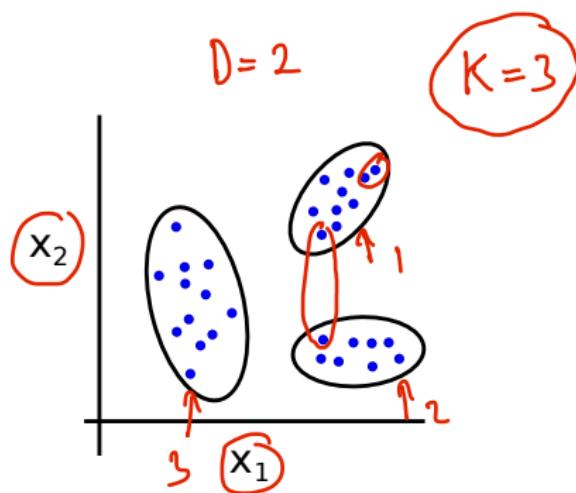
Rich Turner



What is clustering?

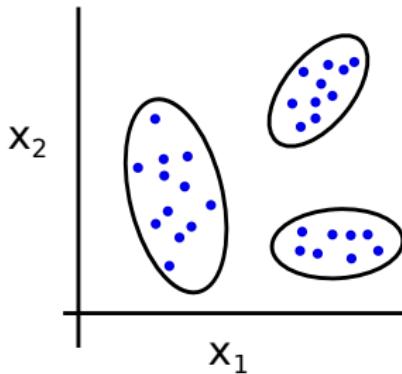
- ▶ **Roughly speaking:** two points belonging to the same cluster are generally more similar to each other or closer to each other than two points belonging to different clusters.

$$\dim(\underline{x}) = D$$



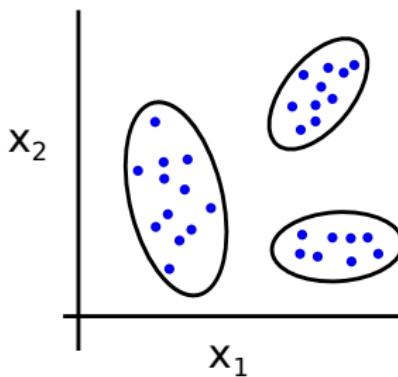
What is clustering?

- ▶ **Roughly speaking:** two points belonging to the same cluster are generally more similar to each other or closer to each other than two points belonging to different clusters.
- ▶ **Notation:** $\mathcal{D} = \{x_1 \dots x_N\} \rightarrow s = \{s_1 \dots s_N\}$
where $D = \dim(x_n)$ and $s_n \in \{1 \dots K\}$

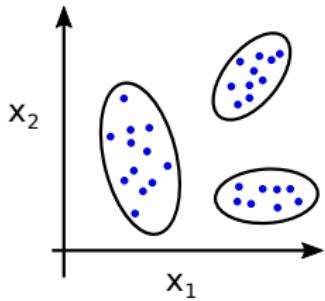


What is clustering?

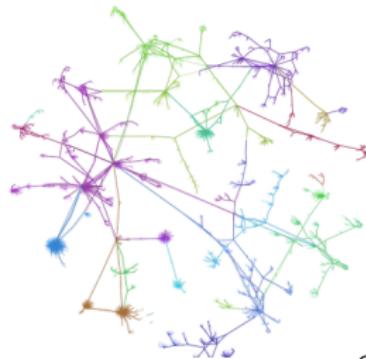
- ▶ **Roughly speaking:** two points belonging to the same cluster are generally more similar to each other or closer to each other than two points belonging to different clusters.
- ▶ **Notation:** $\mathcal{D} = \{x_1 \dots x_N\} \rightarrow s = \{s_1 \dots s_N\}$
where $D = \dim(x_n)$ and $s_n \in \{1 \dots K\}$
- ▶ Unsupervised learning problem (no labels or rewards)



Where is clustering used?

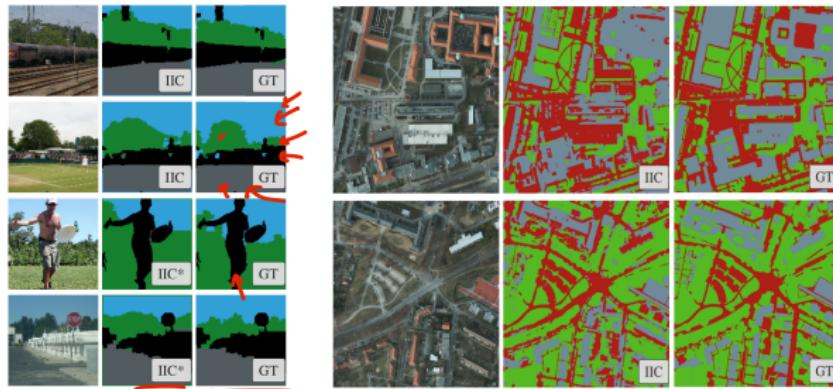


network community detection



Campbell et al
Social Network Analysis

image segmentation



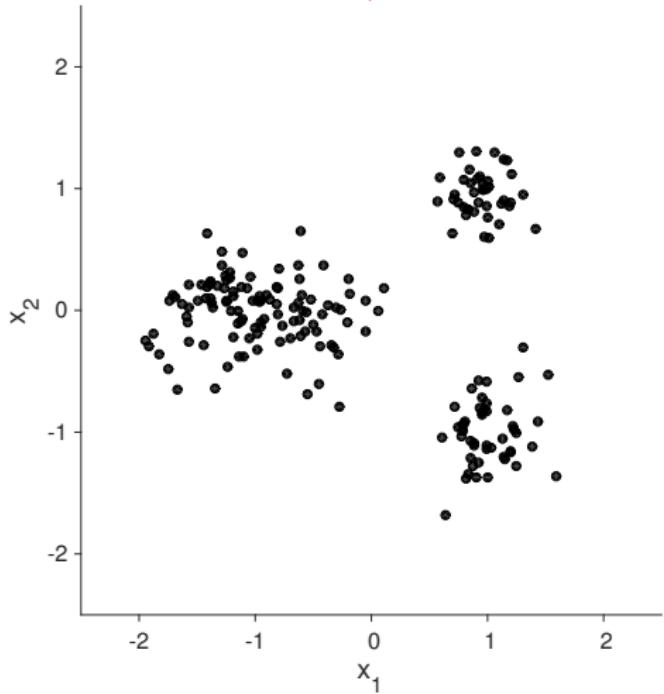
Ji et al Invariant Information Clustering

vector quantisation
genetic clustering ↙
anomaly detection ↙
crime analysis ↘

...

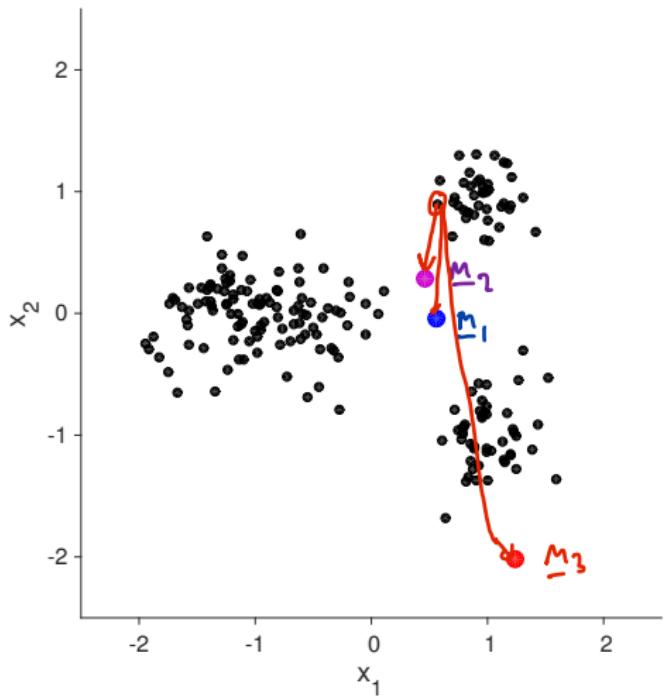
A first clustering algorithm: k-means

$$D=2, N=100$$



cluster centres / means
input: $\mathcal{D} = \{x_1, \dots, x_N\}, x_n \in \mathbb{R}^D$
initialise: $\underline{\underline{m_k}} \in \mathbb{R}^D$ for $k = 1 \dots K$
repeat
 for $n = 1 \dots N$
 $s_n = \arg \min_k ||x_n - m_k||$
 endfor
 for $k = 1 \dots K$
 $m_k = \text{mean}(x_n : s_n = k)$
 endfor
until convergence (s_n fixed)

A first clustering algorithm: k-means



input: $\mathcal{D} = \{x_1, \dots, x_N\}$, $x_n \in \mathbb{R}^D$

initialise: $m_k \in \mathbb{R}^D$ for $k = 1 \dots K$

repeat

for $n = 1 \dots N$

$s_n = \arg \min_k \|x_n - m_k\|$

endfor

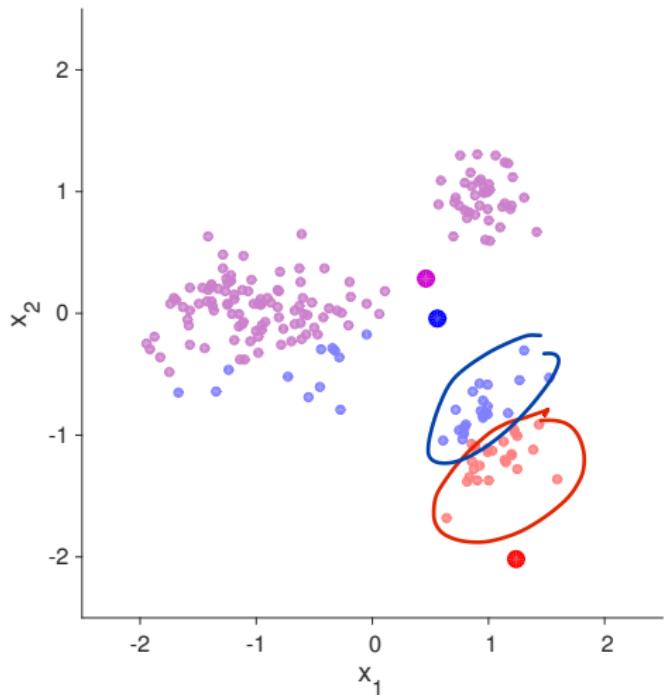
for $k = 1 \dots K$ update

$m_k = \text{mean}(x_n : s_n = k)$

endfor

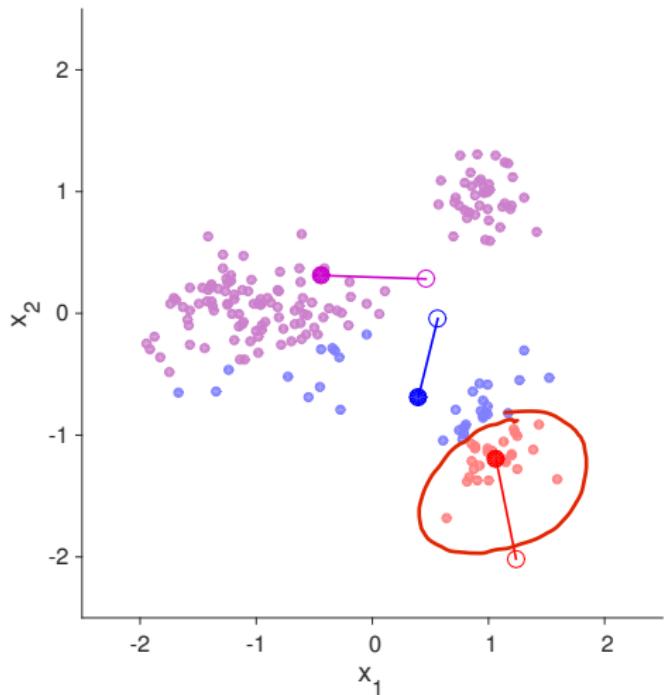
until convergence (s_n fixed)

A first clustering algorithm: k-means



input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$
initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$
repeat
 for $n = 1 \dots N$
 $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$
 endfor
 for $k = 1 \dots K$
 $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$
 endfor
until convergence (s_n fixed)

A first clustering algorithm: k-means



input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$

initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$

repeat

for $n = 1 \dots N$

$s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$

endfor

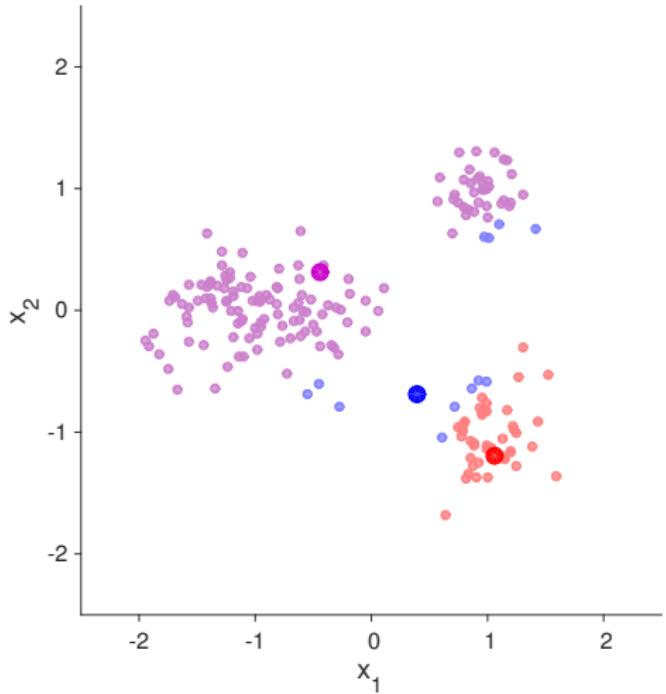
for $k = 1 \dots K$

$\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$

endfor

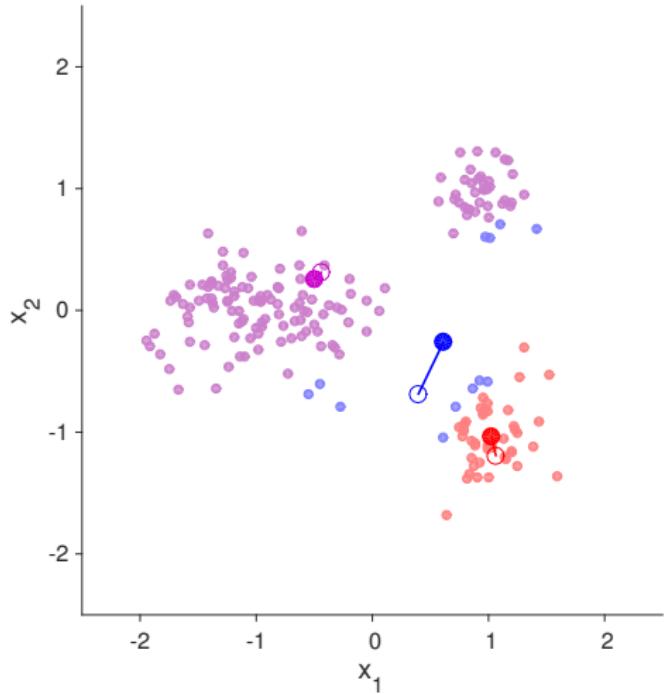
until convergence (s_n fixed)

A first clustering algorithm: k-means



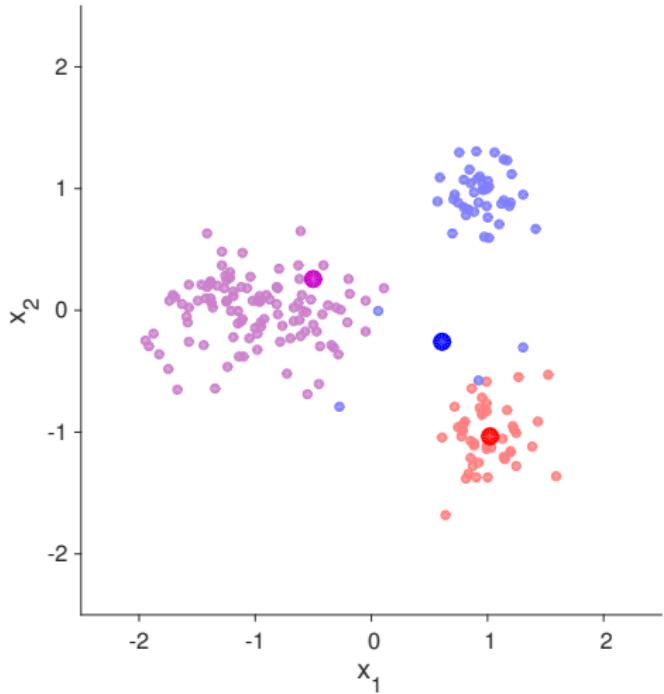
```
input:  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ ,  $\mathbf{x}_n \in \mathbb{R}^D$ 
initialise:  $\mathbf{m}_k \in \mathbb{R}^D$  for  $k = 1 \dots K$ 
repeat
    for  $n = 1 \dots N$ 
         $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$ 
    endfor
    for  $k = 1 \dots K$ 
         $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$ 
    endfor
until convergence ( $s_n$  fixed)
```

A first clustering algorithm: k-means



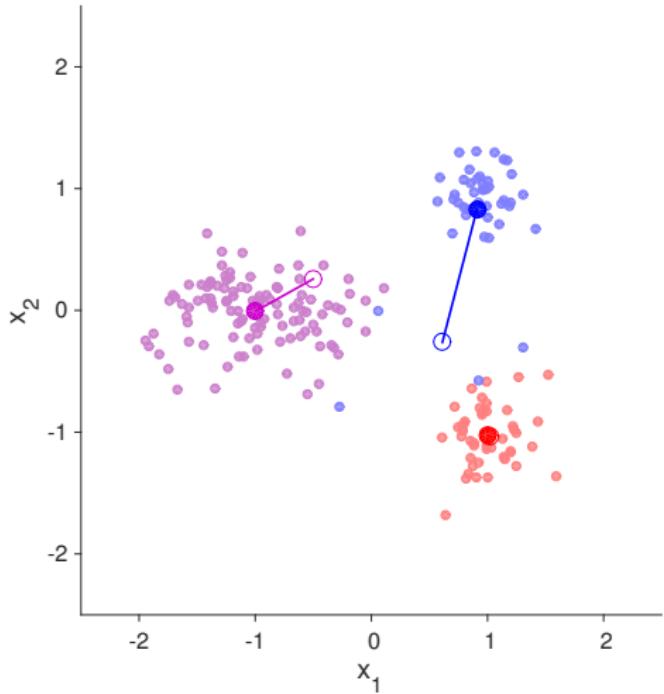
input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$
initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$
repeat
 for $n = 1 \dots N$
 $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$
 endfor
 for $k = 1 \dots K$
 $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$
 endfor
until convergence (s_n fixed)

A first clustering algorithm: k-means



input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$
initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$
repeat
 for $n = 1 \dots N$
 $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$
 endfor
 for $k = 1 \dots K$
 $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$
 endfor
until convergence (s_n fixed)

A first clustering algorithm: k-means



input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$

initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$

repeat

for $n = 1 \dots N$

$s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$

endfor

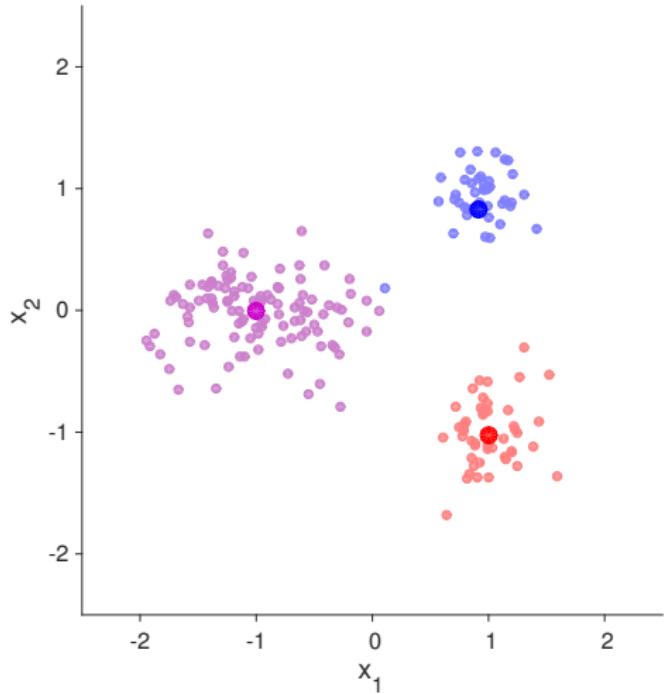
for $k = 1 \dots K$

$\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$

endfor

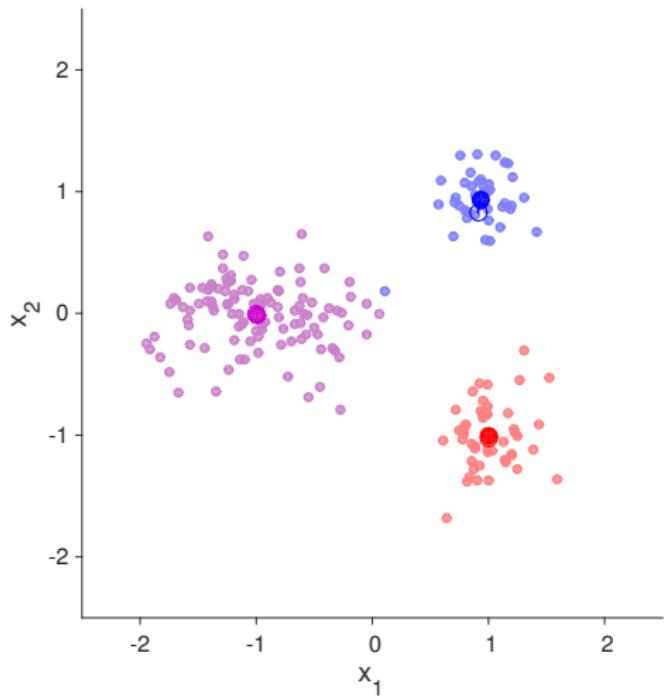
until convergence (s_n fixed)

A first clustering algorithm: k-means



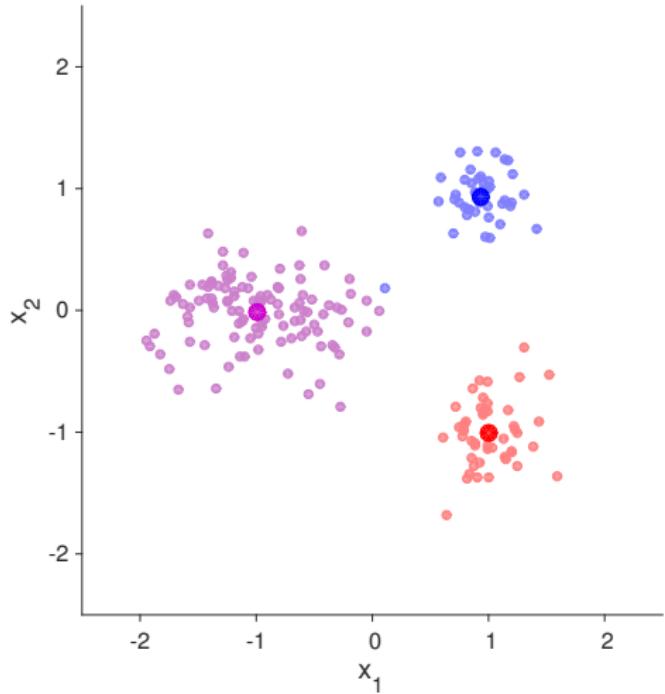
input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$
initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$
repeat
 for $n = 1 \dots N$
 $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$
 endfor
 for $k = 1 \dots K$
 $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$
 endfor
until convergence (s_n fixed)

A first clustering algorithm: k-means



input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$
initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$
repeat
 for $n = 1 \dots N$
 $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$
 endfor
 for $k = 1 \dots K$
 $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$
 endfor
until convergence (s_n fixed)

A first clustering algorithm: k-means



input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$

initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$

repeat

for $n = 1 \dots N$

$$s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$$

endfor

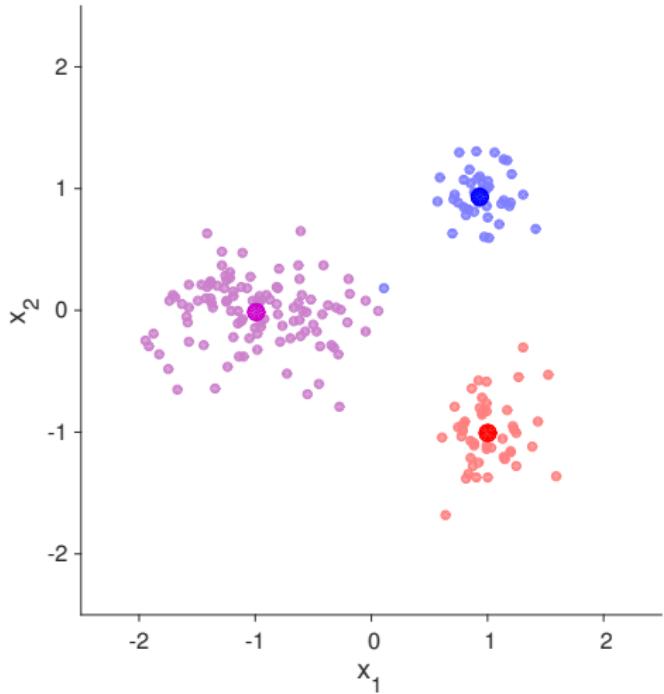
for $k = 1 \dots K$

$$\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$$

endfor

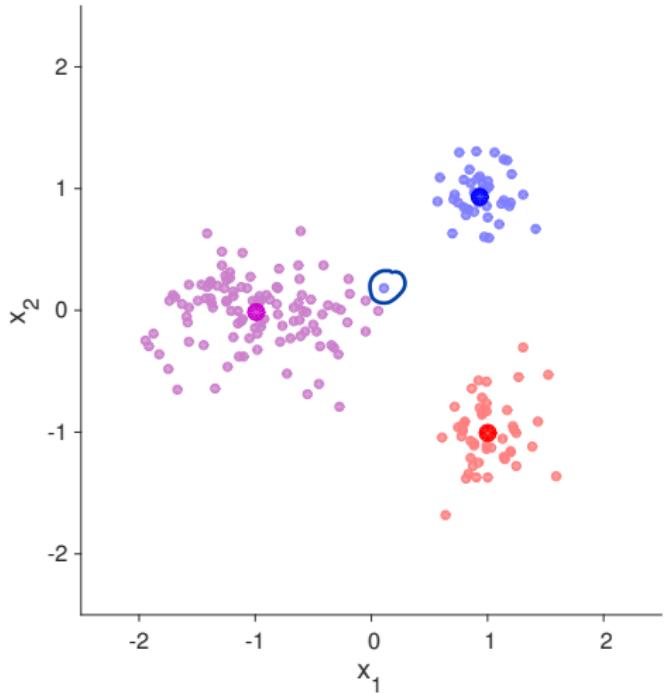
until convergence (s_n fixed)

A first clustering algorithm: k-means



input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$
initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$
repeat
 for $n = 1 \dots N$
 $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$
 endfor
 for $k = 1 \dots K$
 $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$
 endfor
until convergence (s_n fixed)

A first clustering algorithm: k-means



input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$

initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$

repeat

for $n = 1 \dots N$

$s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$

endfor

for $k = 1 \dots K$

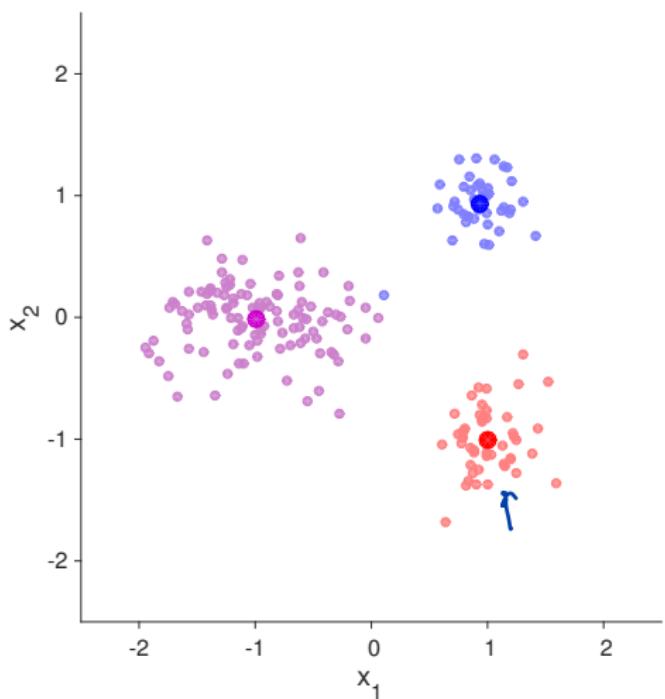
$\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$

endfor

until convergence (s_n fixed)

Question: is K-means guaranteed to converge for any dataset \mathcal{D} ?

could one or more of the cluster centres diverge or oscillate?



input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$

initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$

repeat

for $n = 1 \dots N$

$s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$

endfor

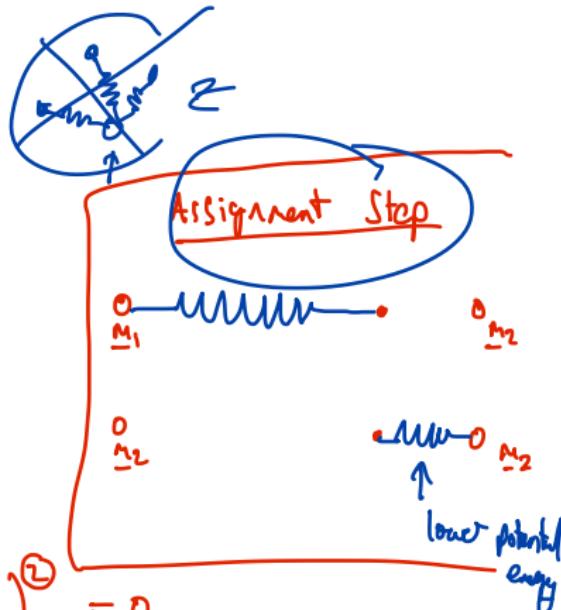
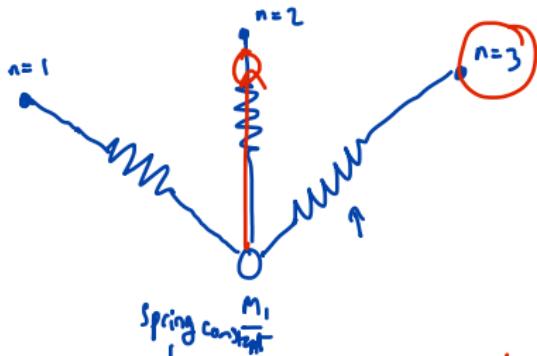
for $k = 1 \dots K$

$\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$

endfor

until convergence (s_n fixed)





$$\frac{d}{dm_1} \text{P.E.} = \frac{1}{2} p \sum_n \left(\| m_1 - z_n \|^2 \right)$$

$$\frac{d}{dm_{1,e}} \frac{1}{2} p \sum_{n=1}^N \sum_{d=1}^D (m_{1,d} - z_{n,d})^2 = 0$$

$$\frac{1}{2} \times p \sum_n (m_{1,e} - z_{n,e})^2 = 0$$

$$m_{1,e} = \frac{1}{N_{\text{cluster}}} \sum_{n=1}^{N_{\text{cluster}}} z_{n,e}$$

$$m_1 = \frac{1}{N_{\text{cluster}}} \sum_{n=1}^{N_{\text{cluster}}} z_n$$

K-means as optimisation

$$s_n \in \{1, \dots, K\}$$

Let $s_{n,k} = 1$ if data point n is assigned to cluster k and zero otherwise

Note: $\sum_{k=1}^K s_{n,k} = 1$

— clusters →

	1	2 ..	$k=3$
1	0	1	0
2	1	0	0
3	1	0	0
$N=4$	0	0	1

\Downarrow

$s_1 = 2, s_2 = 1, s_3 = 1, s_4 = 3$

of datapoints

K-means as optimisation

Let $s_{n,k} = 1$ if data point n is assigned to cluster k and zero otherwise

Note: $\sum_{k=1}^K s_{n,k} = 1$

Cost:

$$\mathcal{C}(\{\underline{s}_{n,k}\}, \{\underline{\boldsymbol{m}}_k\}) = \sum_{n=1}^N \sum_{k=1}^K s_{n,k} \|\boldsymbol{x}_n - \boldsymbol{m}_k\|^2$$

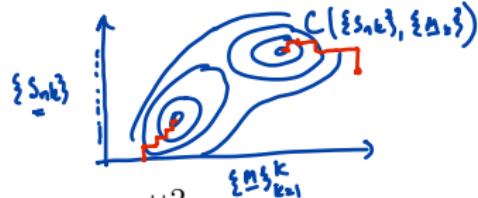
K-means as optimisation

Let $s_{n,k} = 1$ if data point n is assigned to cluster k and zero otherwise

Note: $\sum_{k=1}^K s_{n,k} = 1$

Cost:

$$\mathcal{C}(\{s_{n,k}\}, \{\mathbf{m}_k\}) = \sum_{n=1}^N \sum_{k=1}^K s_{n,k} \|\mathbf{x}_n - \mathbf{m}_k\|^2$$



K-means tries to minimise the cost function \mathcal{C} with respect to $\{s_{n,k}\}$ and $\{\mathbf{m}_k\}$, subject to $\sum_k s_{n,k} = 1$ and $s_{n,k} \in \{0, 1\}$

K-means sequentially:

- ▶ minimises \mathcal{C} with respect to $\{s_{n,k}\}$, holding $\{\mathbf{m}_k\}$ fixed. Assignment
- ▶ minimises \mathcal{C} with respect to $\{\mathbf{m}_k\}$, holding $\{s_{n,k}\}$ fixed. update

K-means as optimisation

Let $s_{n,k} = 1$ if data point n is assigned to cluster k and zero otherwise

Note: $\sum_{k=1}^K s_{n,k} = 1$

Cost:

$$\mathcal{C}(\{s_{n,k}\}, \{\mathbf{m}_k\}) = \sum_{n=1}^N \sum_{k=1}^K s_{n,k} \|\mathbf{x}_n - \mathbf{m}_k\|^2$$

K-means tries to minimise the cost function \mathcal{C} with respect to $\{s_{n,k}\}$ and $\{\mathbf{m}_k\}$, subject to $\sum_k s_{n,k} = 1$ and $s_{n,k} \in \{0, 1\}$

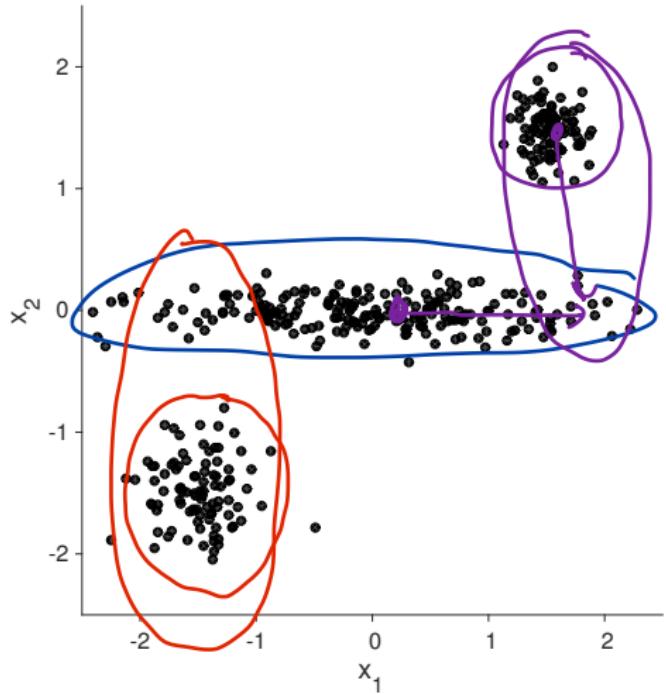
K-means sequentially:

- ▶ minimises \mathcal{C} with respect to $\{s_{n,k}\}$, holding $\{\mathbf{m}_k\}$ fixed.
- ▶ minimises \mathcal{C} with respect to $\{\mathbf{m}_k\}$, holding $\{s_{n,k}\}$ fixed.

\mathcal{C} is a Lyupanov function \implies K-means converges, but...

Finding the global optimum of \mathcal{C} is a hard problem.

Where will K-means converge to when run on these data?



input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$

initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$

repeat

for $n = 1 \dots N$

$s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$

endfor

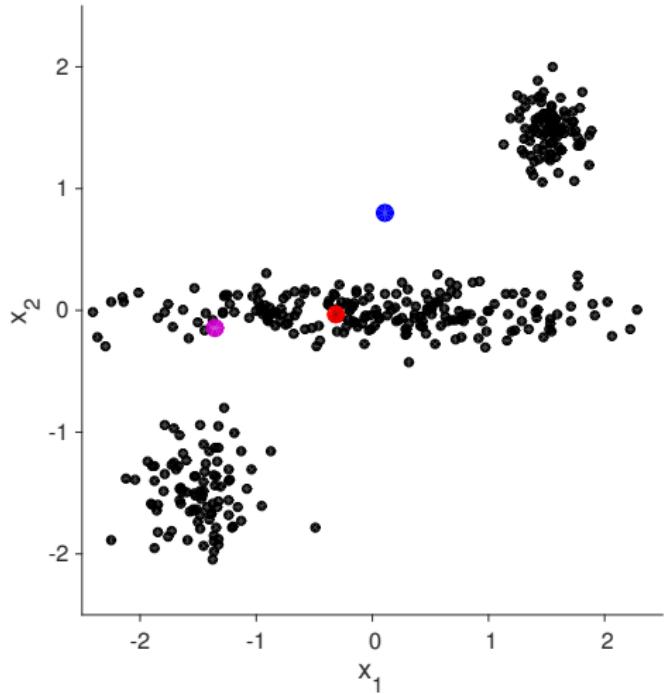
for $k = 1 \dots K$

$\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$

endfor

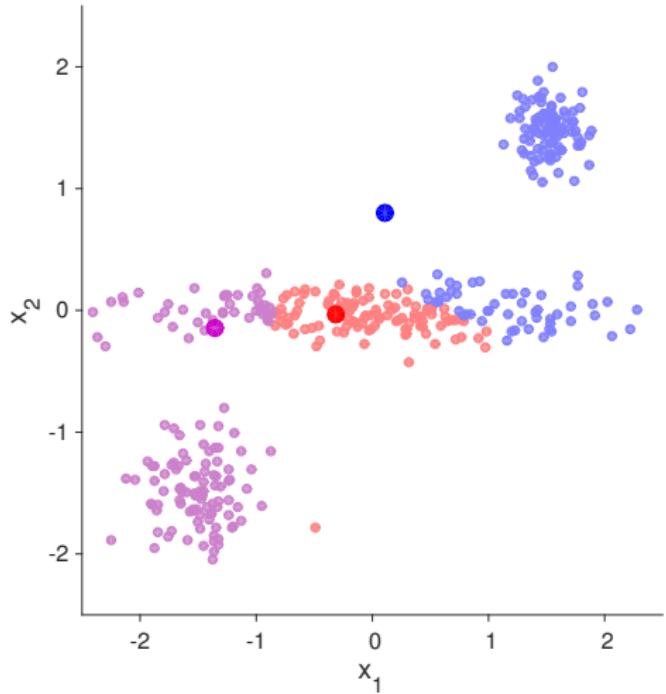
until convergence (s_n fixed)

Where will K-means converge to when run on these data?



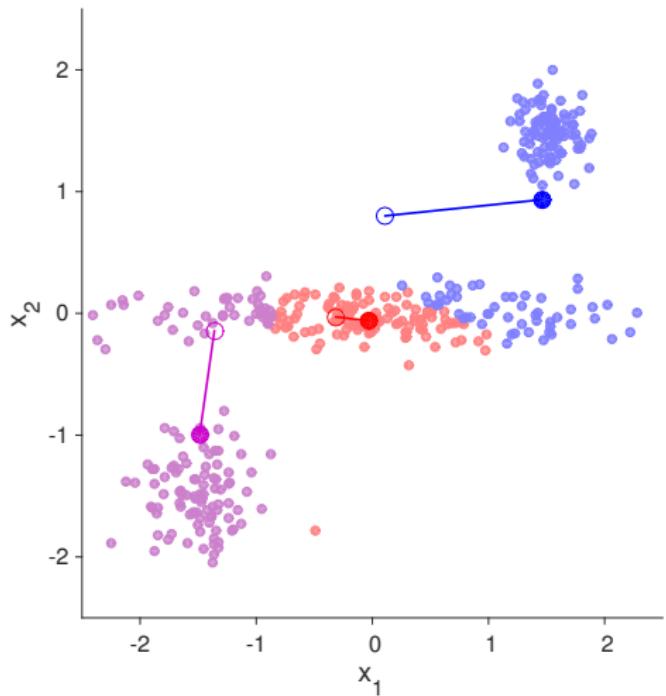
input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$
initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$
repeat
 for $n = 1 \dots N$
 $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$
 endfor
 for $k = 1 \dots K$
 $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$
 endfor
until convergence (s_n fixed)

Where will K-means converge to when run on these data?



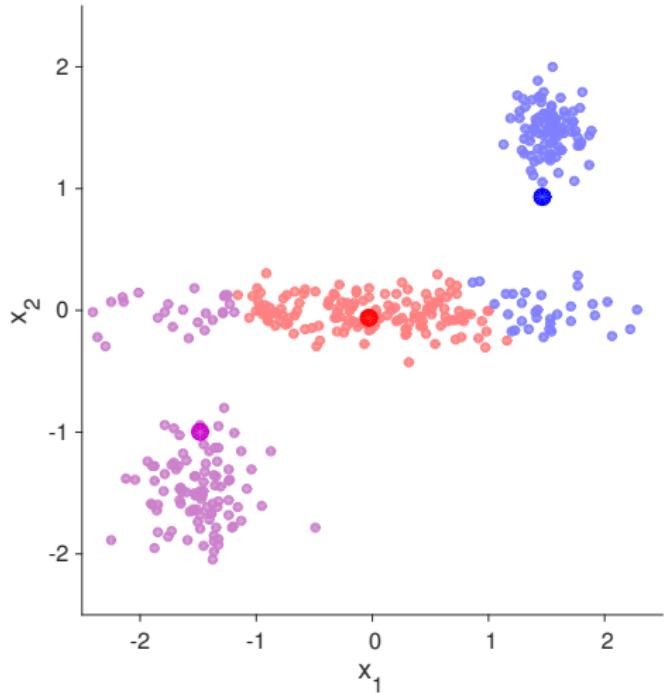
input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$
initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$
repeat
 for $n = 1 \dots N$
 $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$
 endfor
 for $k = 1 \dots K$
 $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$
 endfor
until convergence (s_n fixed)

Where will K-means converge to when run on these data?



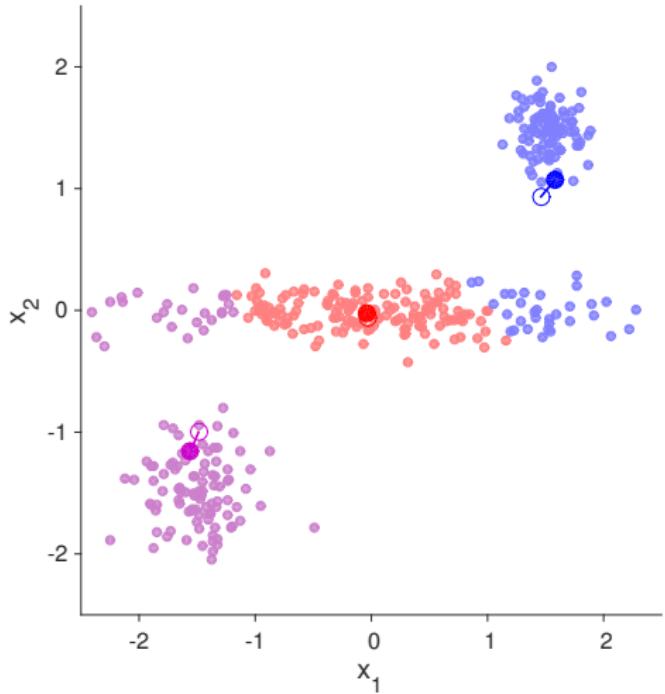
input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$
initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$
repeat
 for $n = 1 \dots N$
 $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$
 endfor
 for $k = 1 \dots K$
 $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$
 endfor
until convergence (s_n fixed)

Where will K-means converge to when run on these data?



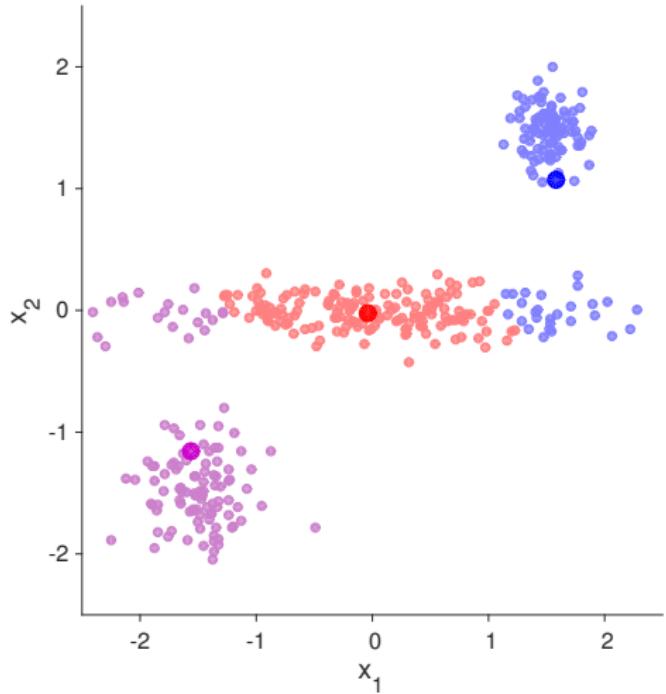
input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$
initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$
repeat
 for $n = 1 \dots N$
 $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$
 endfor
 for $k = 1 \dots K$
 $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$
 endfor
until convergence (s_n fixed)

Where will K-means converge to when run on these data?



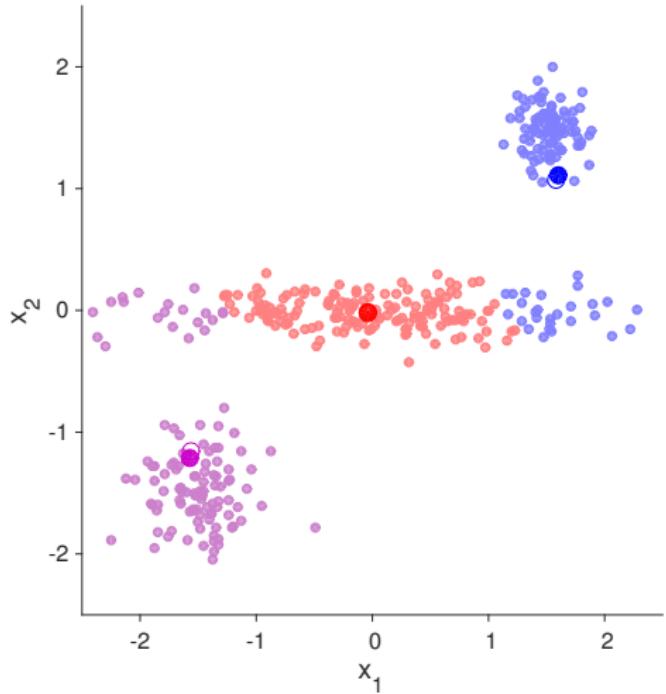
input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$
initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$
repeat
 for $n = 1 \dots N$
 $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$
 endfor
 for $k = 1 \dots K$
 $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$
 endfor
until convergence (s_n fixed)

Where will K-means converge to when run on these data?



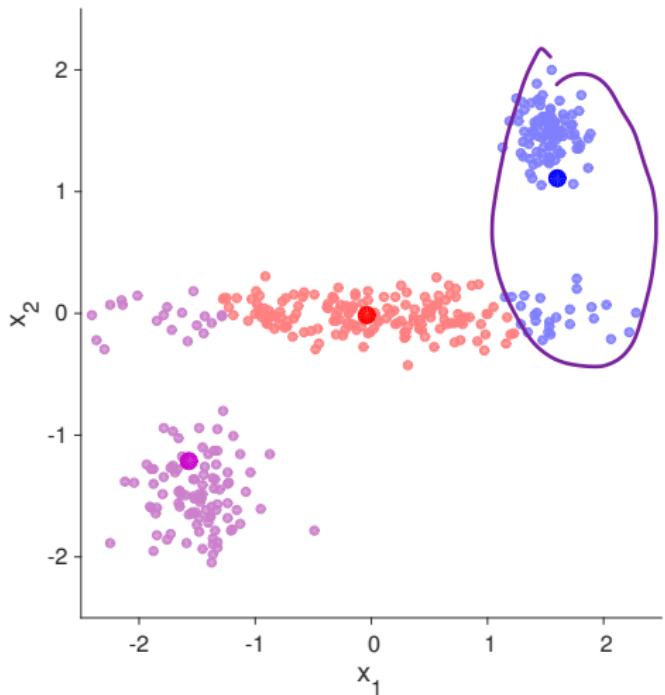
input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$
initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$
repeat
 for $n = 1 \dots N$
 $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$
 endfor
 for $k = 1 \dots K$
 $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$
 endfor
until convergence (s_n fixed)

Where will K-means converge to when run on these data?



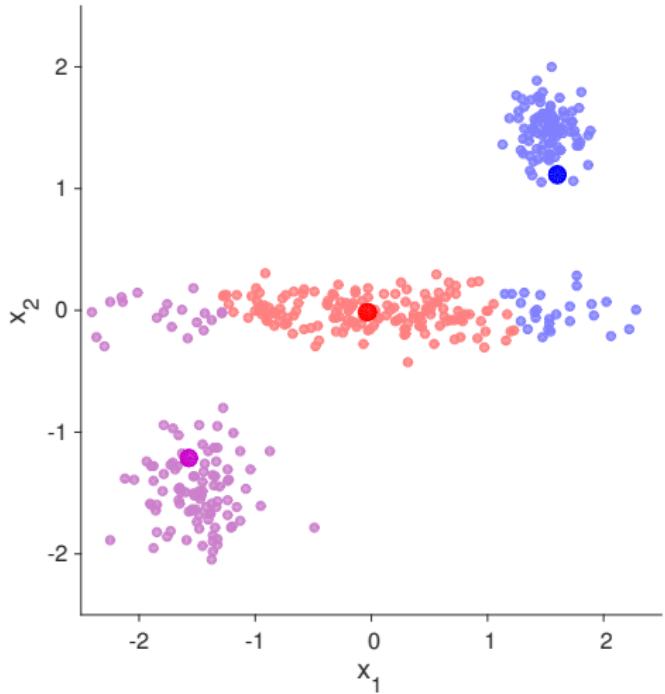
```
input:  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ ,  $\mathbf{x}_n \in \mathbb{R}^D$ 
initialise:  $\mathbf{m}_k \in \mathbb{R}^D$  for  $k = 1 \dots K$ 
repeat
    for  $n = 1 \dots N$ 
         $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$ 
    endfor
    for  $k = 1 \dots K$ 
         $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$ 
    endfor
until convergence ( $s_n$  fixed)
```

Where will K-means converge to when run on these data?



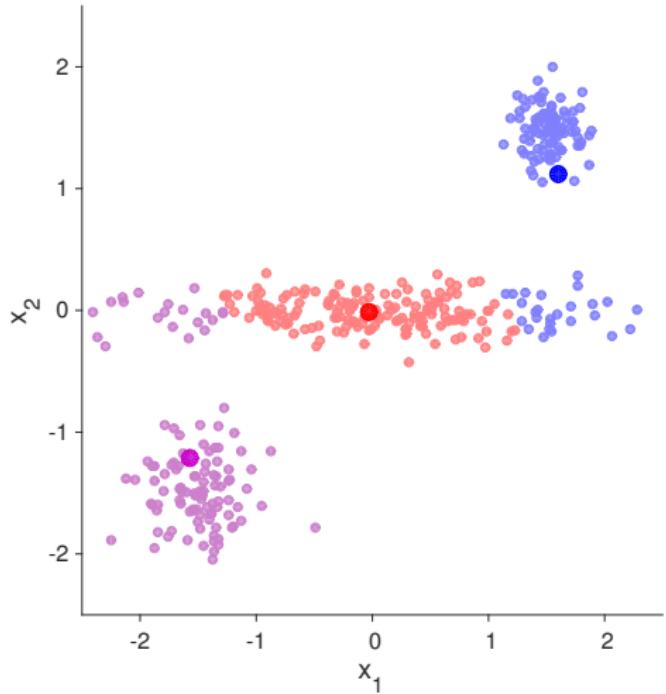
input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$
initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$
repeat
 for $n = 1 \dots N$
 $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$
 endfor
 for $k = 1 \dots K$
 $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$
 endfor
until convergence (s_n fixed)

Where will K-means converge to when run on these data?



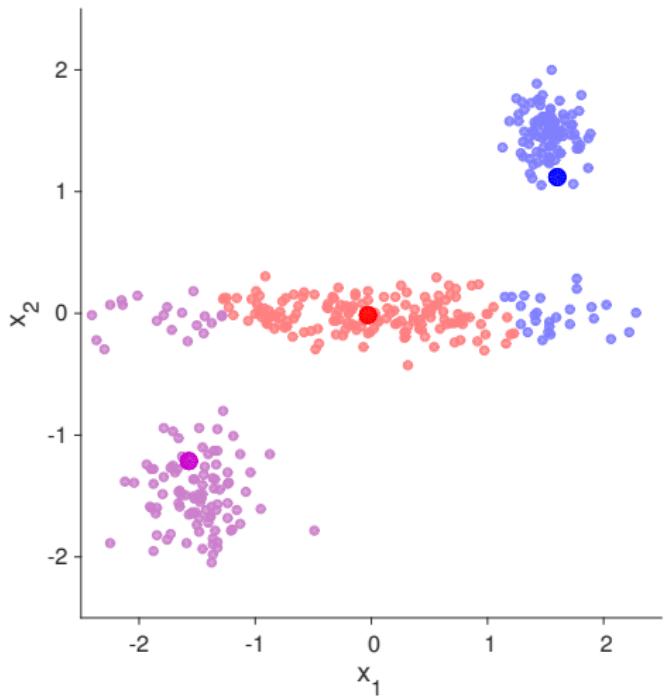
```
input:  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ ,  $\mathbf{x}_n \in \mathbb{R}^D$ 
initialise:  $\mathbf{m}_k \in \mathbb{R}^D$  for  $k = 1 \dots K$ 
repeat
    for  $n = 1 \dots N$ 
         $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$ 
    endfor
    for  $k = 1 \dots K$ 
         $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$ 
    endfor
until convergence ( $s_n$  fixed)
```

Where will K-means converge to when run on these data?



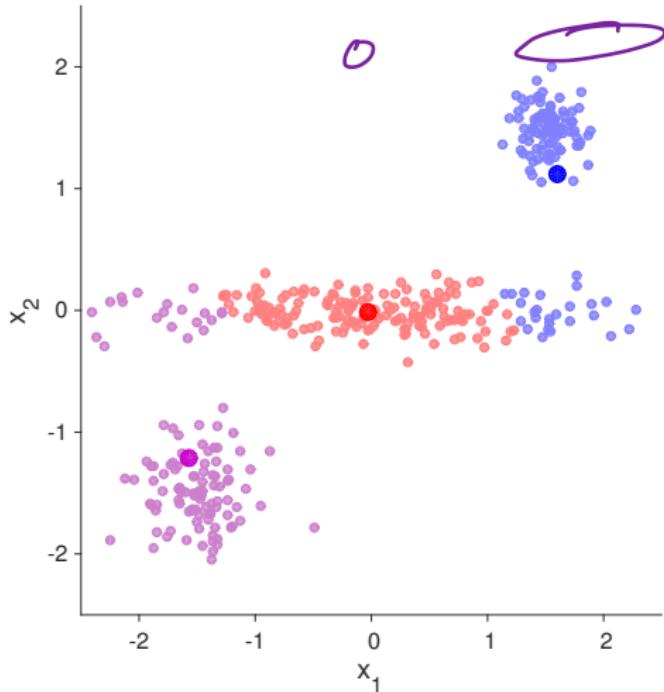
input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$
initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$
repeat
 for $n = 1 \dots N$
 $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$
 endfor
 for $k = 1 \dots K$
 $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$
 endfor
until convergence (s_n fixed)

Where will K-means converge to when run on these data?



input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$
initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$
repeat
 for $n = 1 \dots N$
 $s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$
 endfor
 for $k = 1 \dots K$
 $\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$
 endfor
until convergence (s_n fixed)

Where will K-means converge to when run on these data?



input: $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_n \in \mathbb{R}^D$

initialise: $\mathbf{m}_k \in \mathbb{R}^D$ for $k = 1 \dots K$

repeat

for $n = 1 \dots N$

$s_n = \arg \min_k \|\mathbf{x}_n - \mathbf{m}_k\|$

endfor

for $k = 1 \dots K$

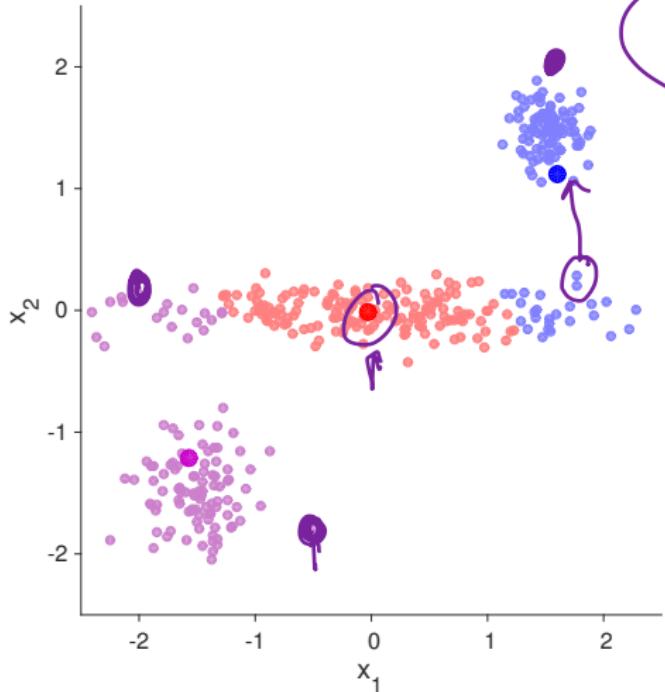
$\mathbf{m}_k = \text{mean}(\mathbf{x}_n : s_n = k)$

endfor

until convergence (s_n fixed)

$$\|\mathbf{x}_{d,n} - \mathbf{m}_d\|$$

Pathologies of K-means

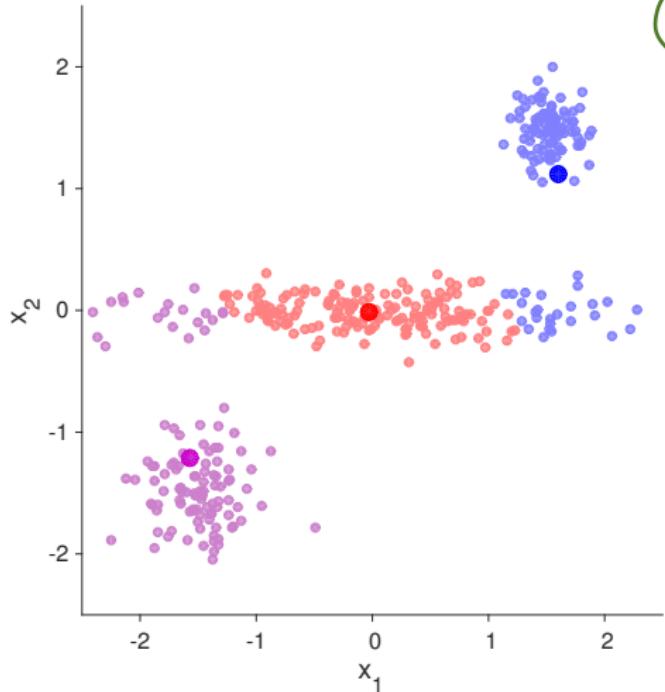


- ▶ strange results when clusters have very different shapes

- ▶ returns hard assignments
 - ▶ soft assignments more sensible

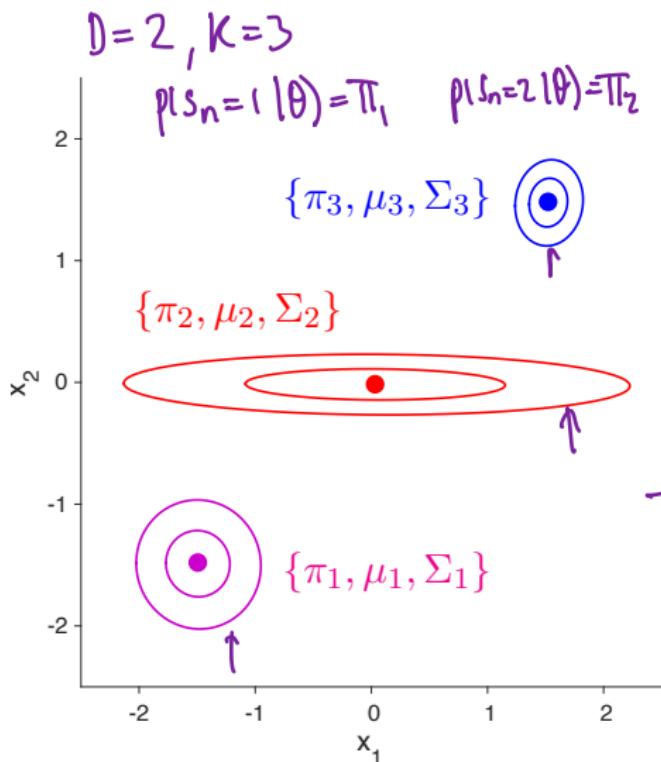
- ▶ initialisation delicate
 - ▶ k++ means algorithm: spreads out initial centres

Pathologies of K-means



- ▶ strange results when clusters have very different shapes
- ▶ returns hard assignments
 - ▶ soft assignments more sensible
- ▶ initialisation delicate
 - ▶ k++ means algorithm: spreads out initial centres

Mixture of Gaussians: Generative Model



For each data point $1 \dots N$

sample cluster membership:

$$p(s_n = k|\theta) = \pi_k \quad \text{note} \quad \sum_{k=1}^K \pi_k = 1$$

sample data-value given cluster mem.:

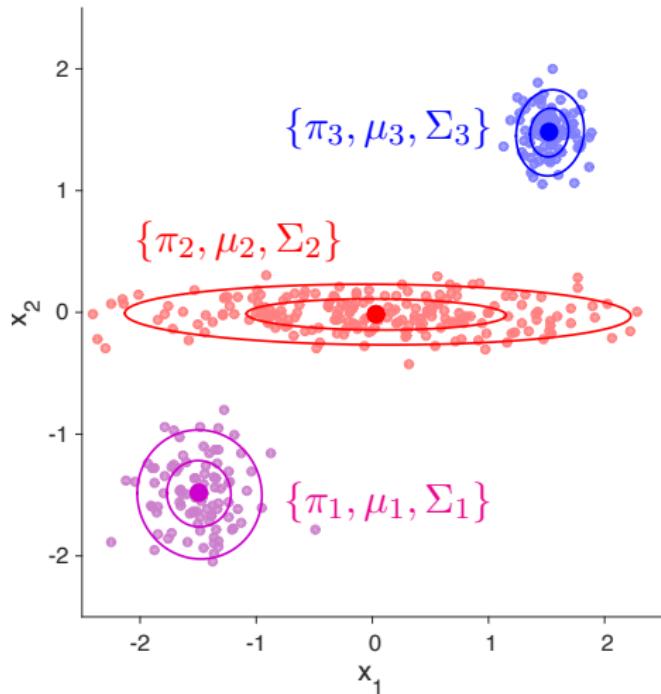
$$p(\mathbf{x}_n | s_n = k, \theta) = \mathcal{N}(\mathbf{x}_n; \mathbf{m}_k, \Sigma_k)$$

$$\prod_{n=1}^N p(\mathbf{x}_n, s_n | \theta) = \prod_{n=1}^N p(s_n | \theta) p(\mathbf{x}_n | s_n, \theta)$$

$$\prod_{n=1}^N p(\mathbf{x}_n | \theta)$$

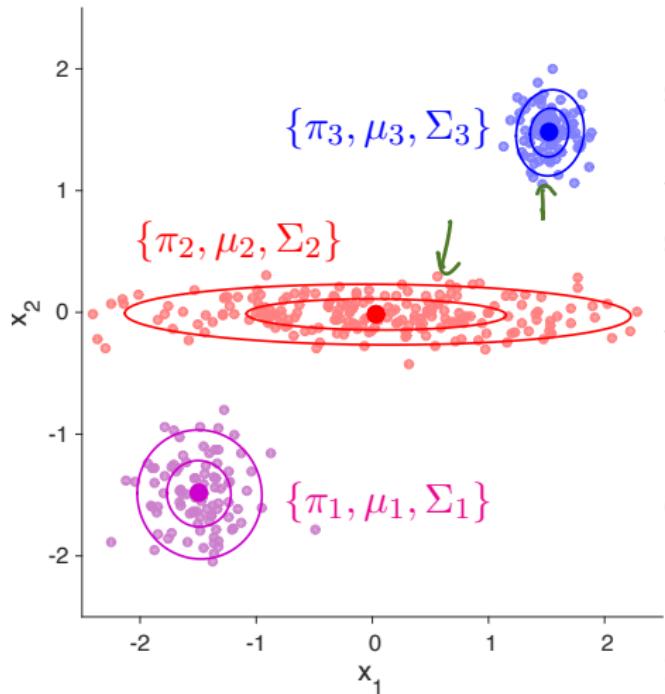
$$\Theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$$

Mixture of Gaussians: Generative Model



For each data point $1 \dots N$
sample cluster membership:
 $p(s_n = k|\theta) = \pi_k$ note $\sum_{k=1}^K \pi_k = 1$
sample data-value given cluster mem.:
 $p(\mathbf{x}_n | s_n = k, \theta) = \mathcal{N}(\mathbf{x}_n; \mathbf{m}_k, \Sigma_k)$

Mixture of Gaussians: Generative Model



- For each data point $1 \dots N$
- sample cluster membership:
- $$p(s_n = k | \theta) = \pi_k \quad \text{note} \quad \sum_{k=1}^K \pi_k = 1$$
- sample data-value given cluster mem.:
- $$p(\mathbf{x}_n | s_n = k, \theta) = \mathcal{N}(\mathbf{x}_n; \mathbf{m}_k, \Sigma_k)$$
- Perform clustering by:**
1. Learning the parameters from data using maximum-likelihood
$$\theta_{\text{ML}} = \arg \max_{\theta} \log p(\{\mathbf{x}_n\}_{n=1}^N | \theta)$$
 2. Inferring the cluster membership from the observed data
$$p(s_n = k | \mathbf{x}_n, \theta_{\text{ML}})$$

Which plot shows the density $p(x|\theta)$ of this 1D Mixture of Gaussians?

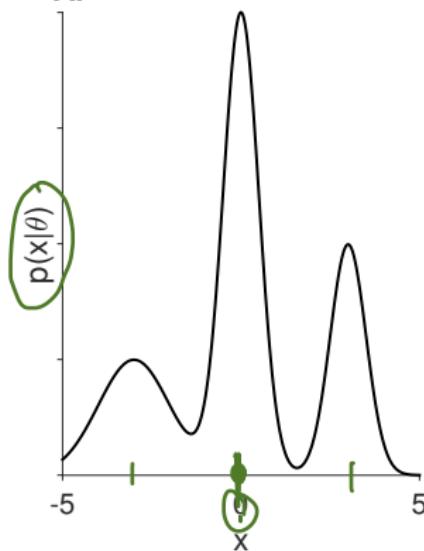
mixing proportions: $\pi_1 = \frac{1}{4}$, $\pi_2 = \frac{1}{2}$ and $\pi_3 = \frac{1}{4}$

means: $\mu_1 = 0$, $\mu_2 = 3$ and $\mu_3 = -3$

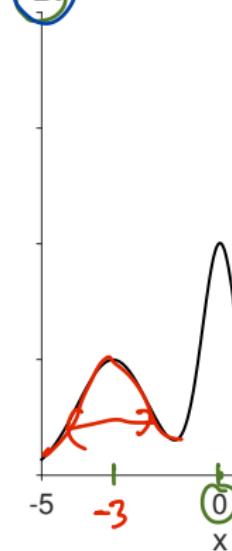
standard deviations: $\sigma_1 = \frac{1}{2}$, $\sigma_2 = \frac{1}{2}$ and $\sigma_3 = 1$

$\parallel \theta$

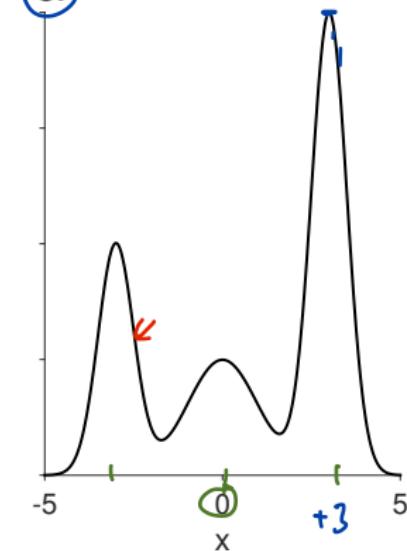
A.



B.



C.



Mixture of Gaussians: $p(s = k|\theta) = \pi_k$ $p(x|s = k, \theta) = \mathcal{N}(x; \mu_k, \sigma_k^2)$

$$\begin{aligned}
 p(\underline{x} | \theta) &= \sum_{k=1}^K p(s=k, \underline{x} | \theta) // \\
 &= \prod_{k=1}^K p(s=k | \theta) p(\underline{x} | s=k, \theta) \\
 &= \sum_{k=1}^K \pi_k N(\underline{x}; \underline{\mu}_k, \underline{\Sigma}_k)
 \end{aligned}$$

↑ ↑
Mixing proportion

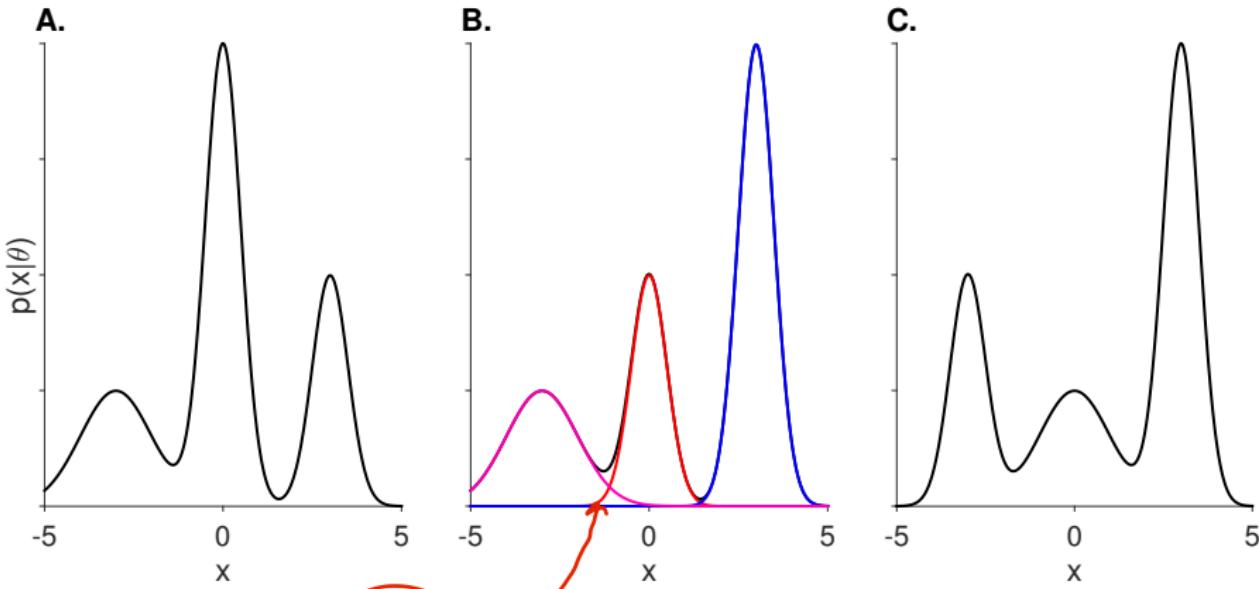
Mixture of Gaussians Model

Which plot shows the density $p(x|\theta)$ of this 1D Mixture of Gaussians?

mixing proportions: $\pi_1 = \frac{1}{4}$, $\pi_2 = \frac{1}{2}$ and $\pi_3 = \frac{1}{4}$

means: $\mu_1 = 0$, $\mu_2 = 3$ and $\mu_3 = -3$

standard deviations: $\sigma_1 = \frac{1}{2}$, $\sigma_2 = \frac{1}{2}$ and $\sigma_3 = 1$



$$\underline{\underline{p(x|\theta)}} = \frac{1}{4} \mathcal{N}(x; 0, \frac{1}{2}^2) + \frac{1}{2} \mathcal{N}(x; 3, \frac{1}{2}^2) + \frac{1}{4} \mathcal{N}(x; -3, 1)$$

Different ways of maximising the likelihood

$$\log p(\underbrace{\{\mathbf{x}_n\}_{n=1}^N}_{\text{---}} | \theta)$$

Different ways of maximising the likelihood

$$\log p(\{\mathbf{x}_n\}_{n=1}^N | \theta)$$

$$= \log \prod_{n=1}^N p(\mathbf{x}_n | \theta) = \sum_{n=1}^N \log p(\mathbf{x}_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K p(\mathbf{x}_n, s_n = k | \theta)$$

Different ways of maximising the likelihood

$$\log p(\{\mathbf{x}_n\}_{n=1}^N | \theta)$$

$$= \log \prod_{n=1}^N p(\mathbf{x}_n | \theta) = \sum_{n=1}^N \log p(\mathbf{x}_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K p(\mathbf{x}_n, s_n = k | \theta)$$

$$= \sum_{n=1}^N \log \sum_{k=1}^K p(s_n = k | \theta) \underbrace{p(\mathbf{x}_n | s_n = k, \theta)}_{\text{product rule}} = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \underbrace{\mathcal{N}(\mathbf{x}_n; \mu_k, \sigma_k^2)}$$

Different ways of maximising the likelihood

$$\begin{aligned} & \log p(\{\mathbf{x}_n\}_{n=1}^N | \theta) \\ &= \log \prod_{n=1}^N p(\mathbf{x}_n | \theta) = \sum_{n=1}^N \log p(\mathbf{x}_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K p(\mathbf{x}_n, s_n = k | \theta) \\ &= \sum_{n=1}^N \log \sum_{k=1}^K p(s_n = k | \theta) p(\mathbf{x}_n | s_n = k, \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n; \mu_k, \sigma_k^2) \end{aligned}$$

- ▶ Method 1: direct gradient based optimisation of the log-likelihood

Different ways of maximising the likelihood

$$\log p(\{\mathbf{x}_n\}_{n=1}^N | \theta)$$

$$= \log \prod_{n=1}^N p(\mathbf{x}_n | \theta) = \sum_{n=1}^N \log p(\mathbf{x}_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K p(\mathbf{x}_n, s_n = k | \theta)$$

$$= \sum_{n=1}^N \log \sum_{k=1}^K p(s_n = k | \theta) p(\mathbf{x}_n | s_n = k, \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n; \mu_k, \sigma_k^2)$$

- ▶ Method 1: direct gradient based optimisation of the log-likelihood
- ▶ Method 2: the **expectation-maximisation (EM) algorithm**

Different ways of maximising the likelihood

$$\log p(\{\mathbf{x}_n\}_{n=1}^N | \theta)$$

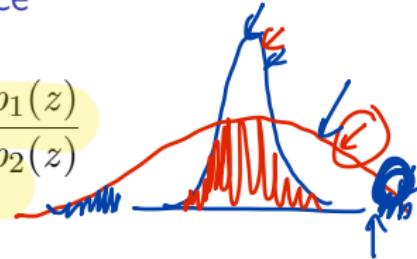
$$\begin{aligned} &= \log \prod_{n=1}^N p(\mathbf{x}_n | \theta) = \sum_{n=1}^N \log p(\mathbf{x}_n | \theta) = \sum_{n=1}^N \log \sum_{k=1}^K p(\mathbf{x}_n, s_n = k | \theta) \\ &= \sum_{n=1}^N \log \sum_{k=1}^K p(s_n = k | \theta) p(\mathbf{x}_n | s_n = k, \theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n; \mu_k, \sigma_k^2) \end{aligned}$$

- ▶ Method 1: direct gradient based optimisation of the log-likelihood
- ▶ Method 2: the **expectation-maximisation (EM) algorithm**

Direct optimisation is generally faster, but the EM algorithm (i) is sometimes **simpler to implement**, (ii) is more **widely used**, and (iii) leads to **important generalisations**.

A brief introduction to the Kullback-Leibler divergence

$$KL(p_1(z) || p_2(z)) = \sum_{z \in \mathcal{Z}} p_1(z) \log \frac{p_1(z)}{p_2(z)}$$



$$K = 2$$

$$p(s=2) = \pi = \pi_1, \quad p(s=1) = 1 - \pi = \pi_2 \quad \Leftarrow \quad \pi_1 + \pi_2 = 1$$

$$q(s=2) = \underline{\rho} \quad q(s=1) = \underline{1-\rho}$$

$$KL(q(s) || p(s)) = (1-\rho) \log \frac{(1-\rho)}{1-\pi} + \rho \log \frac{\rho}{\pi}$$

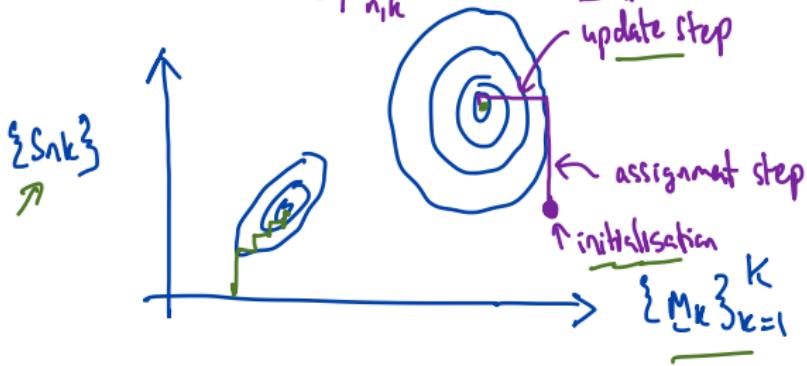
$$KL(p(s) || q(s)) = (1-\pi) \log \frac{(1-\pi)}{1-\rho} + \pi \log \frac{\pi}{\rho}$$

Summary K-means Algorithm

Algorithm pseudo-code

- ① initialise cluster centres $\{\underline{m}_k\}_{k=1}^K \leftarrow \underline{x}_n$
 - ② assign each point to "nearest" cluster centre
 - ③ update centres to mean of assigned points
- iterate

$$E = \frac{1}{2} p \sum_{n,k} s_{nk} \|\underline{x}_n - \underline{m}_k\|^2$$



Mixture of Gaussians Model for Clustering : Summary

Mixture of Gaussians Model

$$p(s_n=k | \theta) = \pi_k$$

$$p(\underline{x}_n | s_n=k, \theta) = N(\underline{x}_n; \underline{\mu}_k, \underline{\Sigma}_k)$$

Clustering approach

$$p(x_1) p(x_2|x_1) p(x_3|x_2, x_1) \dots$$

[latent variable]

$$\begin{matrix} \underline{\mu}_k \\ \underline{\Sigma}_k \end{matrix}$$

$$\theta = \{\pi_k, \underline{\mu}_k, \underline{\Sigma}_k\}_{k=1}^K$$

$$\begin{aligned} \textcircled{1} \quad \theta_{ML} &= \arg \max_{\theta} \log p(\{\underline{x}_n\}_{n=1}^N | \theta) \\ &= \arg \max_{\theta} \sum_{n=1}^N \log p(\underline{x}_n | \theta) \\ &= \sum_{n=1}^N \log \sum_{k=1}^K p(s_n=k | \theta) p(\underline{x}_n | s_n=k) \end{aligned}$$

2

$$p(s_n=k | \underline{x}_n, \theta_{ML}) = p(s_n=k | \theta_{ML}) p(\underline{x}_n | s_n=k, \theta_{ML}) / p(\underline{x}_n | \theta_{ML})$$

A Question about likelihood functions

① From last lecture $p(\underline{x}_n | \theta)$ = likelihood of parameters

$$s_{nk} = [0/1]$$

From notebooks

$$\begin{aligned} p(\underline{x}_n | s_n, \theta) &= \underbrace{\prod_k \pi_k N(x_n; \mu_k, \Sigma_k)}_{\text{likelihood of latent variables } s_n} \\ &\quad \downarrow s_{nk} \end{aligned}$$

What's going on ???!

A brief introduction to the Kullback-Leibler divergence

$$\mathcal{KL}(\underbrace{p_1(z)}_{\rho_2(z)} || \underbrace{p_2(z)}_{\rho_1(z)}) = \sum_{z \in \mathcal{Z}} p_1(z) \log \frac{p_1(z)}{p_2(z)}$$

Important properties:

A brief introduction to the Kullback-Leibler divergence

$$\mathcal{KL}(p_1(z) \parallel p_2(z)) = \sum_{z \in \mathcal{Z}} p_1(z) \log \frac{p_1(z)}{p_2(z)}$$

Important properties:

- ▶ Gibb's inequality: $\mathcal{KL}(p_1(z) \parallel p_2(z)) \geq 0$, equality at $p_1(z) = p_2(z)$
 - ▶ proof via Jensen's inequality or differentiation (see MacKay pg. 35)

A brief introduction to the Kullback-Leibler divergence

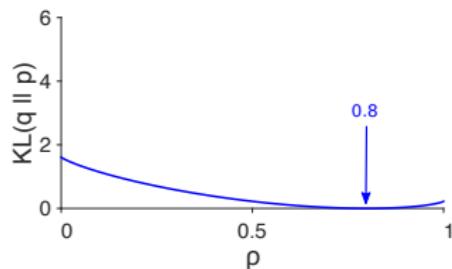
$$\mathcal{KL}(p_1(z) \parallel p_2(z)) = \sum_{z \in \mathcal{Z}} p_1(z) \log \frac{p_1(z)}{p_2(z)}$$

Important properties:

- ▶ Gibb's inequality: $\mathcal{KL}(p_1(z) \parallel p_2(z)) \geq 0$, equality at $p_1(z) = p_2(z)$
 - ▶ proof via Jensen's inequality or differentiation (see MacKay pg. 35)

Example:

- ▶ binary variables $z \in \{0, 1\}$
- ▶ $p(z = 1) = 0.8$ and $q(z = 1) = \rho$



A brief introduction to the Kullback-Leibler divergence

$$\mathcal{KL}(p_1(z) \parallel p_2(z)) = \sum_{z \in \mathcal{Z}} p_1(z) \log \frac{p_1(z)}{p_2(z)}$$

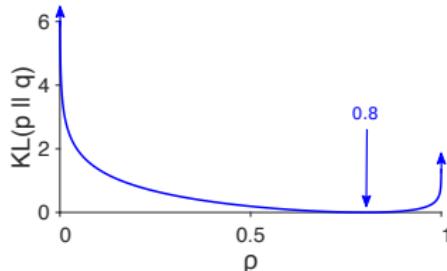
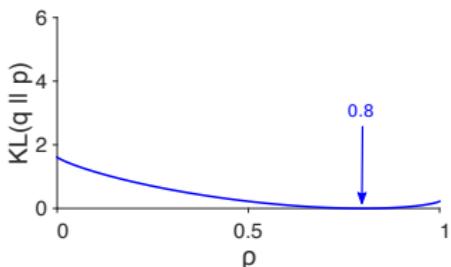
~~1~~

Important properties:

- ▶ Gibb's inequality: $\mathcal{KL}(p_1(z) \parallel p_2(z)) \geq 0$, equality at $p_1(z) = p_2(z)$
 - ▶ proof via Jensen's inequality or differentiation (see MacKay pg. 35)

Example:

- ▶ binary variables $z \in \{0, 1\}$
- ▶ $p(z = 1) = 0.8$ and $q(z = 1) = \rho$



A brief introduction to the Kullback-Leibler divergence

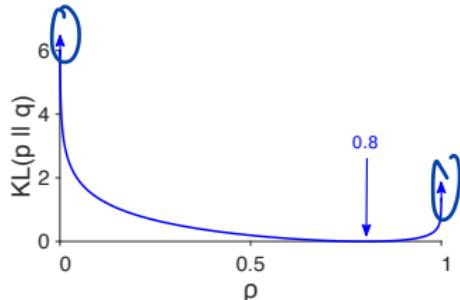
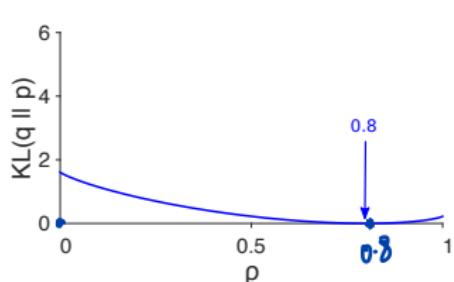
$$\mathcal{KL}(p_1(z)||p_2(z)) = \sum_{z \in \mathcal{Z}} p_1(z) \log \frac{p_1(z)}{p_2(z)}$$

Important properties:

- ▶ Gibb's inequality: $\mathcal{KL}(p_1(z)||p_2(z)) \geq 0$, equality at $p_1(z) = p_2(z)$
 - ▶ proof via Jensen's inequality or differentiation (see MacKay pg. 35)
- ▶ Non-symmetric: $\mathcal{KL}(p_1(z)||p_2(z)) \neq \mathcal{KL}(p_2(z)||p_1(z))$
 - ▶ hence named *divergence* and not *distance*

Example:

- ▶ binary variables $z \in \{0, 1\}$
- ▶ $p(z=1) = 0.8$ and $q(z=1) = \rho$



The Expectation Maximisation Algorithm and the Variational Free-energy

- ▶ Build toward **Expectation Maximisation (EM) Algorithm**

$$KL(q(s) \parallel p(s)) = \sum_{k=1}^K \frac{q(s=k)}{p(s=k)} \log \frac{q(s=k)}{p(s=k)}$$

$s \in \{1 \dots K\}$

$$\frac{d}{dq(s=l)} \left[\underbrace{KL(q(s) \parallel p(s))}_{\text{evaluates to 0 at solution}} - \lambda \left(\sum_{k=1}^K q(s=k) - 1 \right) \right] = 0$$

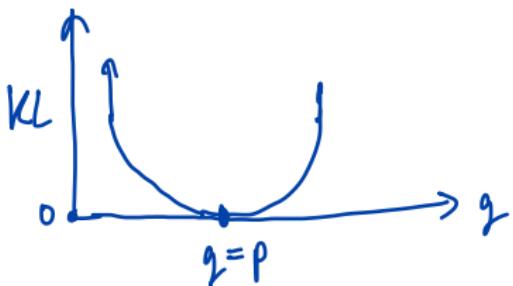
$$\underbrace{\log \frac{q(s=l)}{p(s=l)}} + \frac{q(s=l)}{\cancel{q(s=l)}} - \lambda = 0$$

$$\left. \begin{aligned} q(s=l) &= \underbrace{p(s=l)} e^{\lambda-1} \\ \sum_{l=1}^K q(s=l) &= 1 \end{aligned} \right\} \Rightarrow e^{\lambda-1} = 1 \Rightarrow q(s=l) \parallel p(s=l)$$

$$\frac{d^2}{dq(s=l)^2} \text{KL}(q(s) \parallel p(s)) = \frac{d}{dq(s=l)} \left[\log \frac{q(s=l)}{p(s=l)} + 1 \right]$$

$$= \frac{1}{q(s=l)} > 0$$

+ve



The Expectation Maximisation Algorithm and the Variational Free-energy

- ▶ Build toward **Expectation Maximisation (EM) Algorithm**
- ▶ EM Algorithm can be used for maximum-likelihood parameter learning of **many latent variable models**

The Expectation Maximisation Algorithm and the Variational Free-energy

- ▶ Build toward **Expectation Maximisation (EM) Algorithm**
- ▶ EM Algorithm can be used for maximum-likelihood parameter learning of **many latent variable models**
- ▶ Leads to a variety of extensions such as **approximate variational inference**

The Expectation Maximisation Algorithm and the Variational Free-energy

- ▶ Build toward **Expectation Maximisation (EM) Algorithm**
- ▶ EM Algorithm can be used for maximum-likelihood parameter learning of **many latent variable models**
- ▶ Leads to a variety of extensions such as **approximate variational inference**

For generality, we simplify notation: let $\mathbf{x} = \{\mathbf{x}_n\}_{n=1}^N$ and $\mathbf{s} = \{s_n\}_{n=1}^N$ so



The Expectation Maximisation Algorithm and the Variational Free-energy

- ▶ Build toward **Expectation Maximisation (EM) Algorithm**
- ▶ EM Algorithm can be used for maximum-likelihood parameter learning of **many latent variable models**
- ▶ Leads to a variety of extensions such as **approximate variational inference**

For generality, we simplify notation: let $\mathbf{x} = \{\mathbf{x}_n\}_{n=1}^N$ and $\mathbf{s} = \{s_n\}_{n=1}^N$ so

$$\log p(\underbrace{\{\mathbf{x}_n\}}_{\mathbf{x}}_{n=1}^N | \theta) = \log \underbrace{p(\mathbf{x} | \theta)}_{\mathbf{s}} = \log \sum_{\mathbf{s}} p(\mathbf{x}, \mathbf{s} | \theta)$$

A lower bound on the log-likelihood $\log p(\mathbf{x}|\theta)$

$$\overbrace{\log p(\mathbf{x}|\theta)}^{\text{log-likelihood}}$$

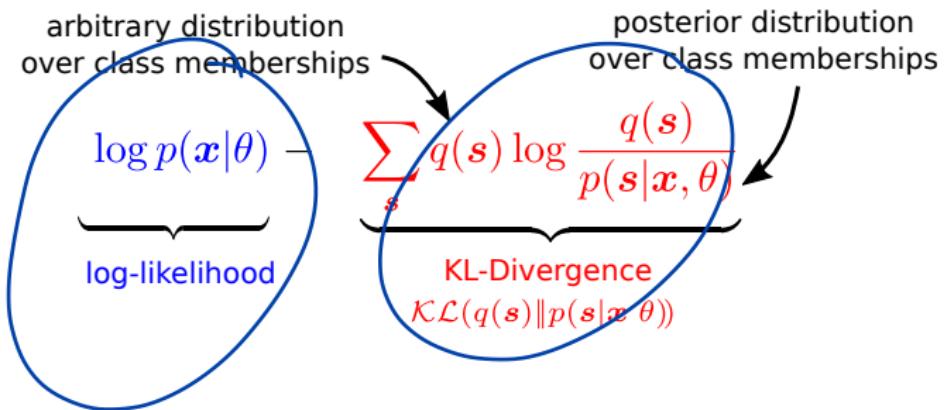
A lower bound on the log-likelihood $\log p(\mathbf{x}|\theta)$

$$\underbrace{\log p(\mathbf{x}|\theta)}_{\text{log-likelihood}} - \sum_s q(s) \log \frac{q(s)}{\underbrace{p(s|\mathbf{x}, \theta)}_{\text{posterior}}}$$

arbitrary distribution over class memberships

posterior distribution over class memberships

A lower bound on the log-likelihood $\log p(\mathbf{x}|\theta)$



A lower bound on the log-likelihood $\log p(\mathbf{x}|\theta)$

$$\mathcal{F}(q(\mathbf{s}), \theta) = \underbrace{\log p(\mathbf{x}|\theta)}_{\text{log-likelihood}} - \underbrace{\sum_{\mathbf{s}} q(\mathbf{s}) \log \frac{q(\mathbf{s})}{p(\mathbf{s}|\mathbf{x}, \theta)}}_{\text{KL-Divergence } \mathcal{KL}(q(\mathbf{s}) \| p(\mathbf{s}|\mathbf{x}, \theta))}$$

arbitrary distribution over class memberships posterior distribution over class memberships

free-energy log-likelihood

ELBO

A lower bound on the log-likelihood $\log p(\mathbf{x}|\theta)$

$$\mathcal{F}(q(\mathbf{s}), \theta) = \underbrace{\log p(\mathbf{x}|\theta)}_{\text{log-likelihood}} - \underbrace{\sum_{\mathbf{s}} q(\mathbf{s}) \log \frac{q(\mathbf{s})}{p(\mathbf{s}|\mathbf{x}, \theta)}}_{\text{KL-Divergence } \mathcal{KL}(q(\mathbf{s}) \| p(\mathbf{s}|\mathbf{x}, \theta))}$$

arbitrary distribution over class memberships

posterior distribution over class memberships

free-energy log-likelihood

A lower bound on the log-likelihood $\log p(\mathbf{x}|\theta)$

$$\mathcal{F}(q(\mathbf{s}), \theta) = \underbrace{\log p(\mathbf{x}|\theta)}_{\text{log-likelihood}} - \underbrace{\sum_{\mathbf{s}} q(\mathbf{s}) \log \frac{q(\mathbf{s})}{p(\mathbf{s}|\mathbf{x}, \theta)}}_{\text{KL-Divergence } \mathcal{KL}(q(\mathbf{s}) \| p(\mathbf{s}|\mathbf{x}, \theta))}$$

arbitrary distribution over class memberships posterior distribution over class memberships

Non-negativity
of KL-Divergence

$$\mathcal{KL}(q(\mathbf{s}) \| p(\mathbf{s}|\mathbf{x}, \theta)) \geq 0 \implies \mathcal{F}(q(\mathbf{s}), \theta) \leq \log p(\mathbf{x}|\theta)$$

Free-energy is lower bound on
log-likelihood

A lower bound on the log-likelihood $\log p(\mathbf{x}|\theta)$

$$\mathcal{F}(q(\mathbf{s}), \theta) = \underbrace{\log p(\mathbf{x}|\theta)}_{\text{log-likelihood}} - \underbrace{\sum_{\mathbf{s}} q(\mathbf{s}) \log \frac{q(\mathbf{s})}{p(\mathbf{s}|\mathbf{x}, \theta)}}_{\text{KL-Divergence } \mathcal{KL}(q(\mathbf{s}) \| p(\mathbf{s}|\mathbf{x}, \theta))}$$

arbitrary distribution over class memberships posterior distribution over class memberships

The diagram illustrates the decomposition of the free-energy function. It starts with the formula $\mathcal{F}(q(\mathbf{s}), \theta) = \log p(\mathbf{x}|\theta) - \sum_{\mathbf{s}} q(\mathbf{s}) \log \frac{q(\mathbf{s})}{p(\mathbf{s}|\mathbf{x}, \theta)}$. The term $\log p(\mathbf{x}|\theta)$ is labeled "log-likelihood". The term $\sum_{\mathbf{s}} q(\mathbf{s}) \log \frac{q(\mathbf{s})}{p(\mathbf{s}|\mathbf{x}, \theta)}$ is labeled "KL-Divergence" and is further broken down into "arbitrary distribution over class memberships" and "posterior distribution over class memberships". Arrows point from the labels to their corresponding parts in the formula.

Non-negativity
of KL-Divergence

$$\mathcal{KL}(q(\mathbf{s}) \| p(\mathbf{s}|\mathbf{x}, \theta)) \geq 0 \implies \mathcal{F}(q(\mathbf{s}), \theta) \leq \log p(\mathbf{x}|\theta)$$

Free-energy is lower bound on
log-likelihood

KL-Divergence equal to 0
when $q(\mathbf{s}) = p(\mathbf{s}|\mathbf{x}, \theta)$

Free-energy equal to log-likelihood
when $q(\mathbf{s}) = p(\mathbf{s}|\mathbf{x}, \theta)$

$$\mathcal{KL}(q(\mathbf{s}) \| p(\mathbf{s}|\mathbf{x}, \theta)) = 0 \implies \mathcal{F}(q(\mathbf{s}), \theta) = \log p(\mathbf{x}|\theta)$$

A lower bound on the log-likelihood $\log p(\mathbf{x}|\theta)$

$$q(\underline{s}) = \prod_{n=1}^N q_n(s_n) \stackrel{\text{if}}{=} q_n(s_n=1) = 1 - \rho$$

arbitrary distribution
over class memberships

$$\mathcal{F}(q(\mathbf{s}), \theta) = \underbrace{\log p(\mathbf{x}|\theta)}_{\text{free-energy}} - \underbrace{\sum_{\mathbf{s}} q(\mathbf{s}) \log \frac{q(\mathbf{s})}{p(\mathbf{s}|\mathbf{x}, \theta)}}_{\text{log-likelihood}}$$

posterior distribution
over class memberships

KL-Divergence
 $\mathcal{KL}(q(\mathbf{s}) \| p(\mathbf{s}|\mathbf{x}, \theta))$

$$\Rightarrow \mathcal{F}(q(\mathbf{s}), \theta) = \sum_{\mathbf{s}} q(\mathbf{s}) \log \frac{p(\mathbf{x}|\mathbf{s}, \theta)p(\mathbf{s}|\theta)}{q(\mathbf{s})} \xrightarrow{\text{simple to compute}}$$

Non-negativity
of KL-Divergence

Free-energy is lower bound on
log-likelihood

$$\mathcal{KL}(q(\mathbf{s}) \| p(\mathbf{s}|\mathbf{x}, \theta)) \geq 0 \implies \mathcal{F}(q(\mathbf{s}), \theta) \leq \log p(\mathbf{x}|\theta)$$

KL-Divergence equal to 0
when $q(\mathbf{s}) = p(\mathbf{s}|\mathbf{x}, \theta)$

Free-energy equal to log-likelihood
when $q(\mathbf{s}) = p(\mathbf{s}|\mathbf{x}, \theta)$

$$\mathcal{KL}(q(\mathbf{s}) \| p(\mathbf{s}|\mathbf{x}, \theta)) = 0 \implies \mathcal{F}(q(\mathbf{s}), \theta) = \log p(\mathbf{x}|\theta)$$

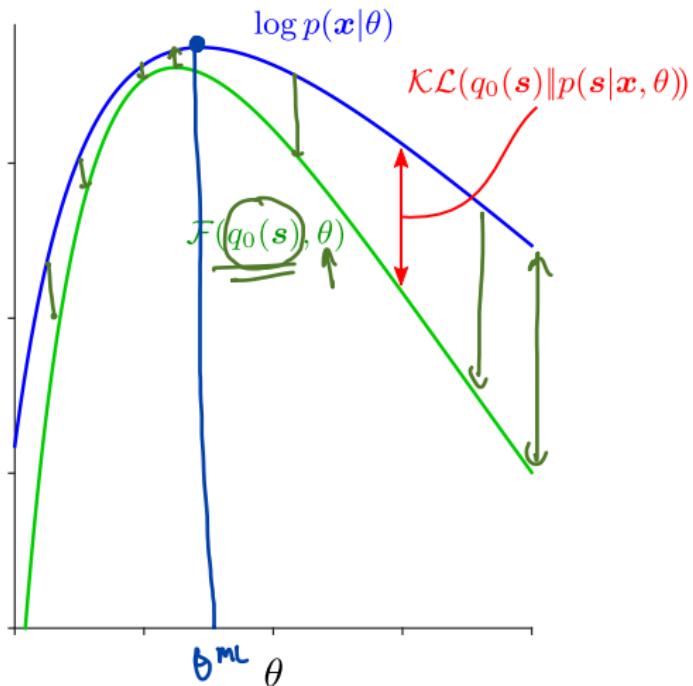
Visualising the free-energy lower bound

Visualising the free-energy lower bound

$$\mathcal{F}(q(s), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(s)\|p(s|\mathbf{x}, \theta))$$

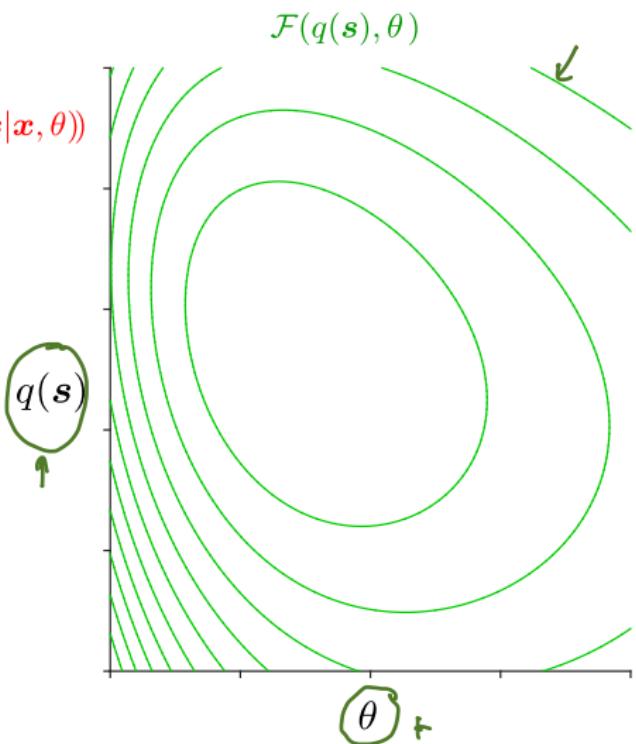
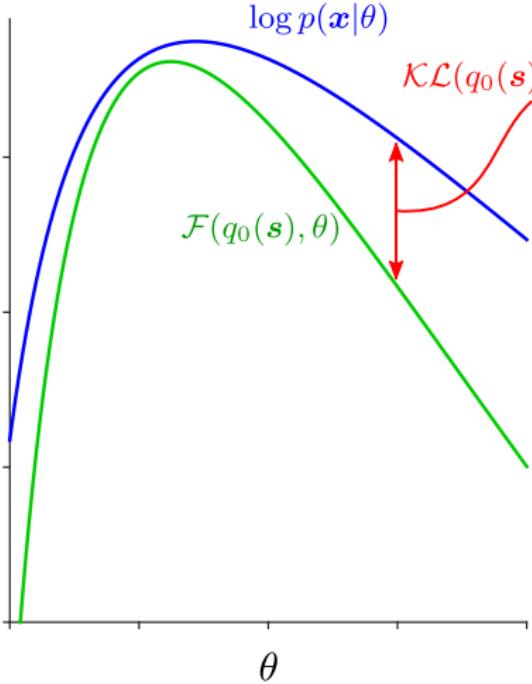
Visualising the free-energy lower bound

$$\mathcal{F}(q(s), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(s)\|p(s|\mathbf{x}, \theta))$$



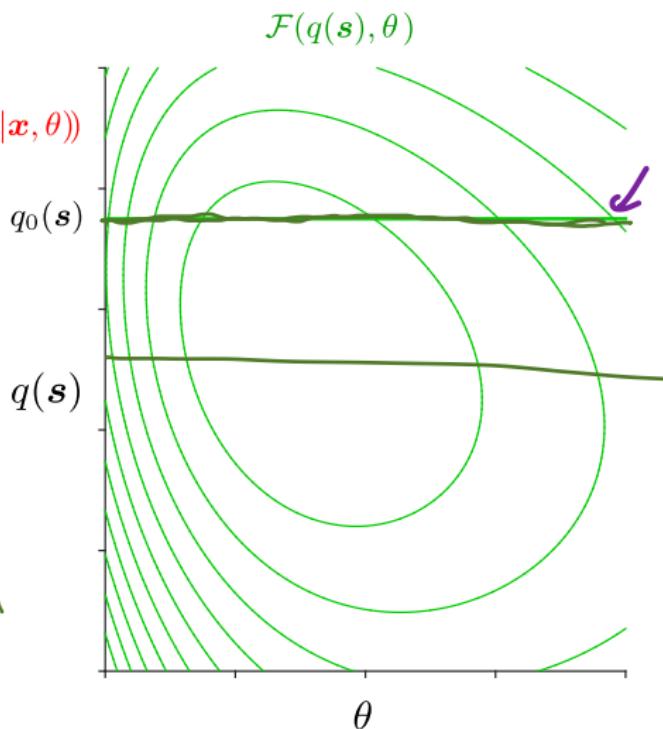
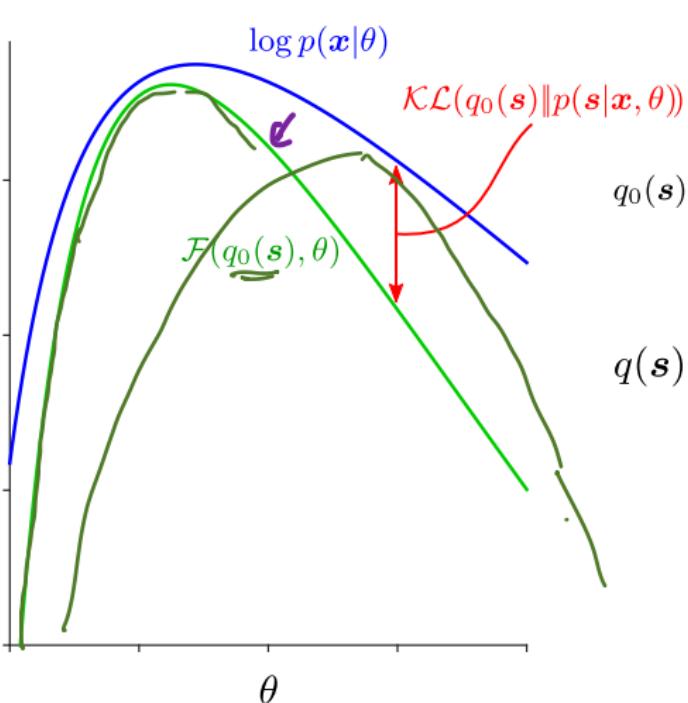
Visualising the free-energy lower bound

$$\mathcal{F}(q(s), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(s)\|p(s|\mathbf{x}, \theta))$$



Visualising the free-energy lower bound

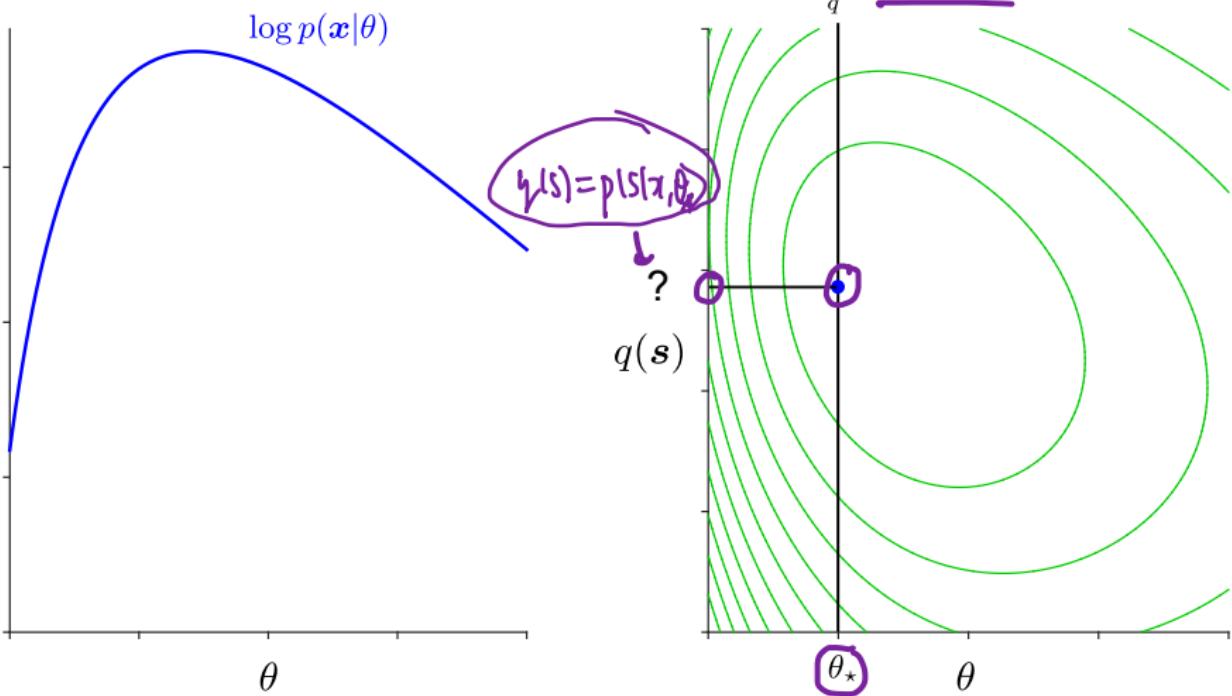
$$\mathcal{F}(q(s), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(s)\|p(s|\mathbf{x}, \theta))$$



What is the maximal value of the free-energy along this vertical slice?

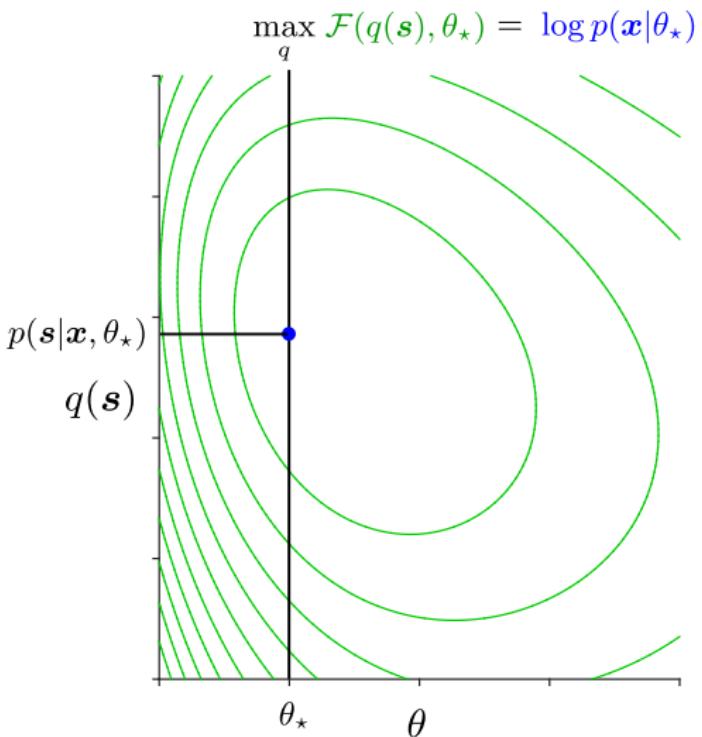
$$\Rightarrow \mathcal{F}(q(s), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(s) \| p(s|\mathbf{x}, \theta)) = 0$$

$$\max_q \mathcal{F}(q(s), \theta_*) = ?$$



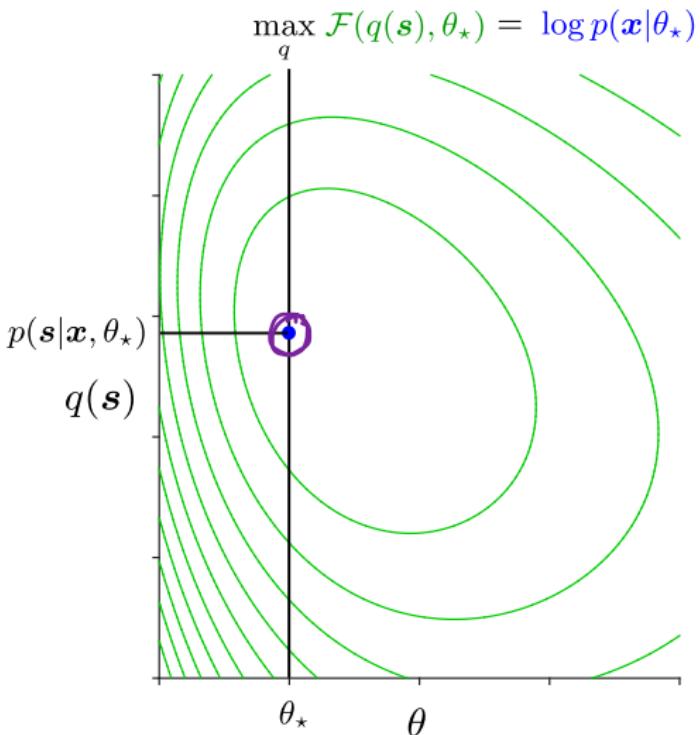
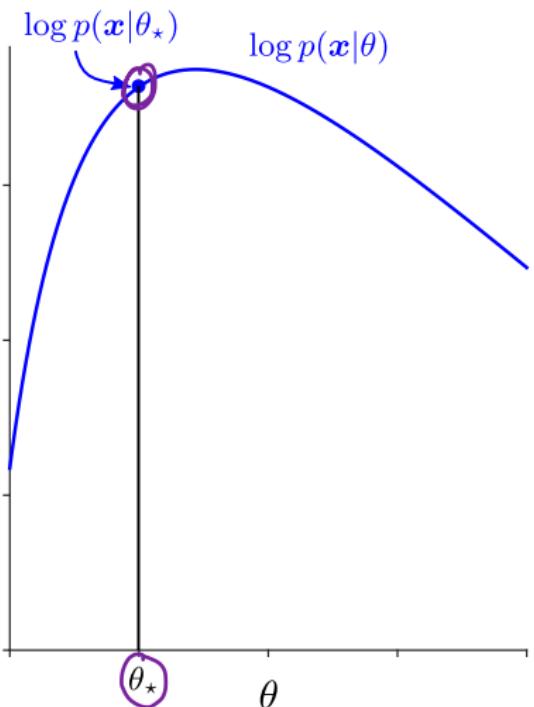
What is the maximal value of the free-energy along this vertical slice?

$$\mathcal{F}(q(s), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(s)\|p(s|\mathbf{x}, \theta))$$



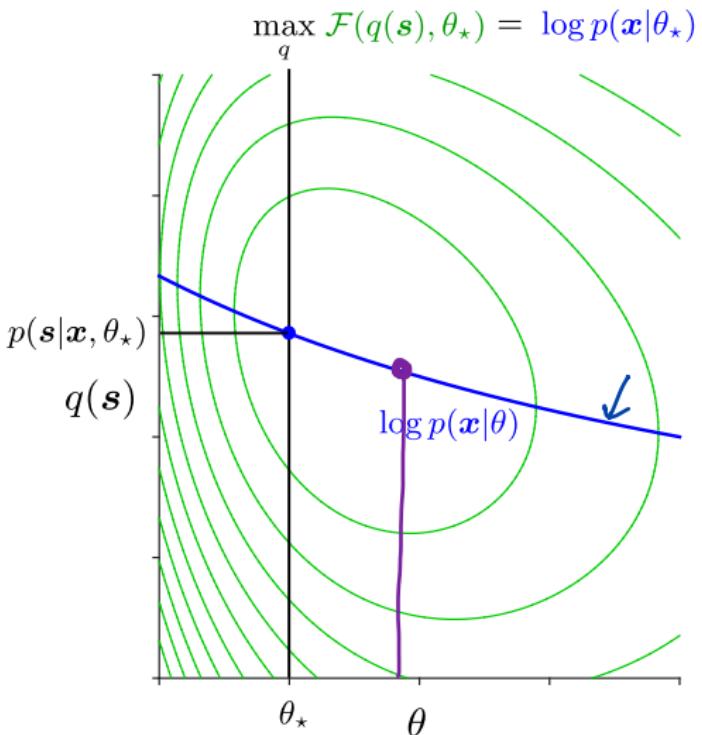
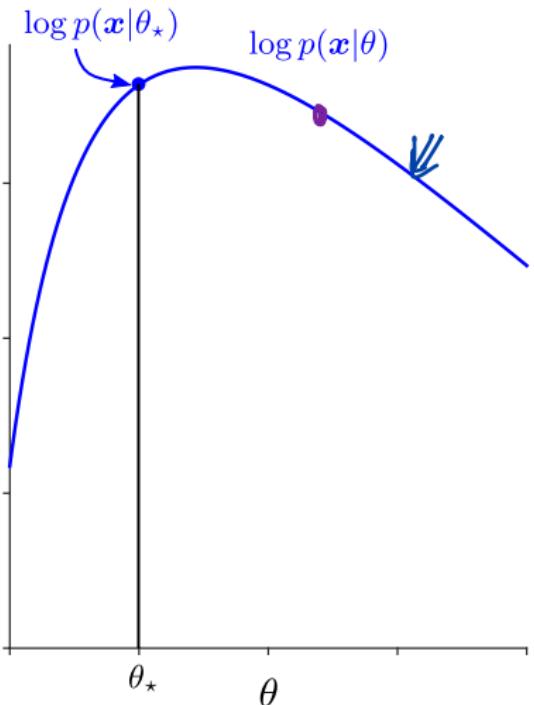
What is the maximal value of the free-energy along this vertical slice?

$$\mathcal{F}(q(s), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(s)\|p(s|\mathbf{x}, \theta))$$



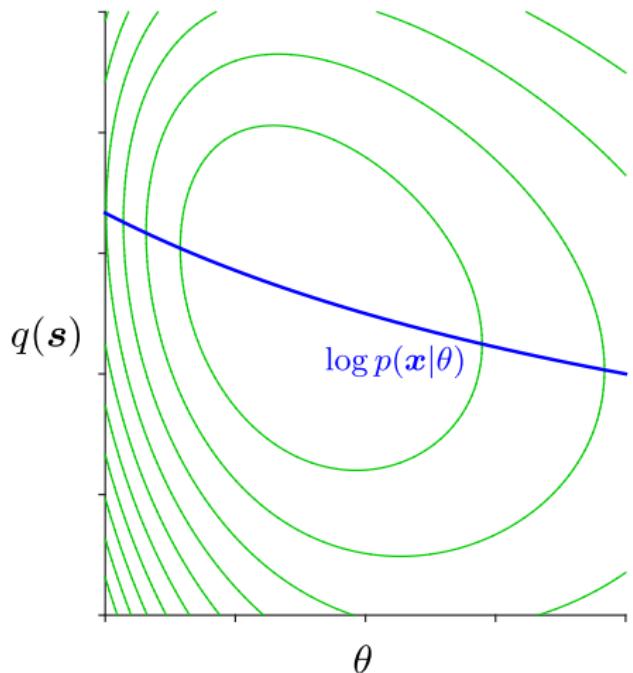
What is the maximal value of the free-energy along this vertical slice?

$$\mathcal{F}(q(s), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(s)\|p(s|\mathbf{x}, \theta))$$



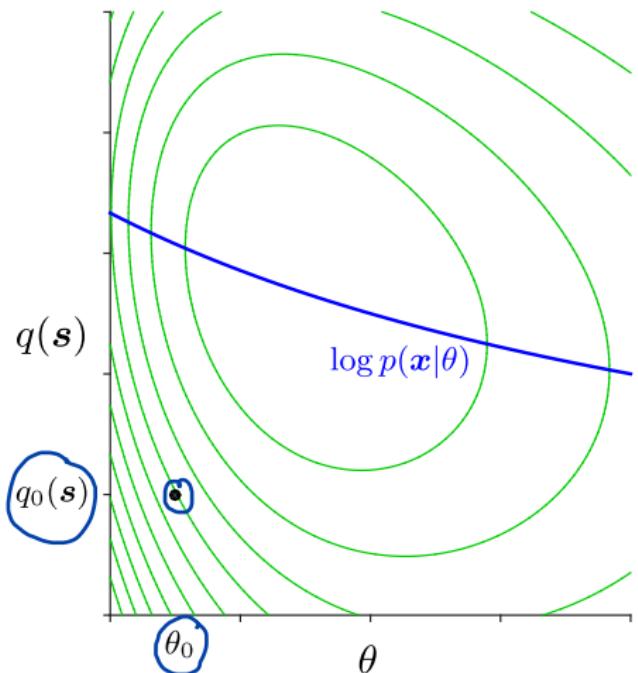
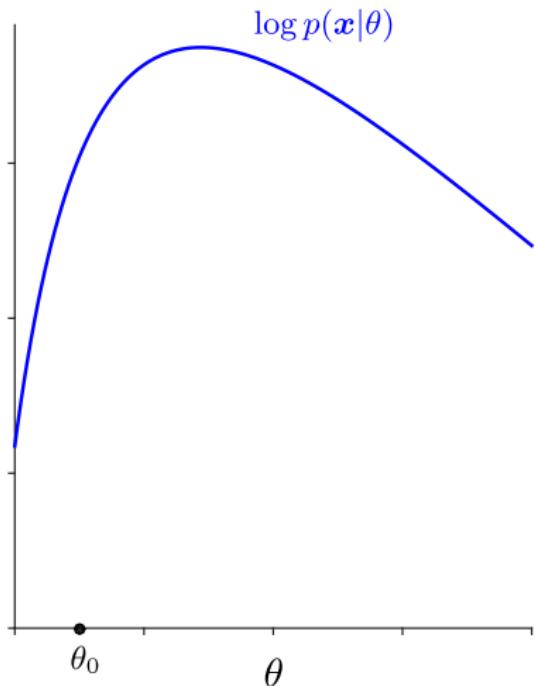
The Expectation Maximisation (EM) algorithm

$$\mathcal{F}(q(s), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(s)\|p(s|\mathbf{x}, \theta))$$



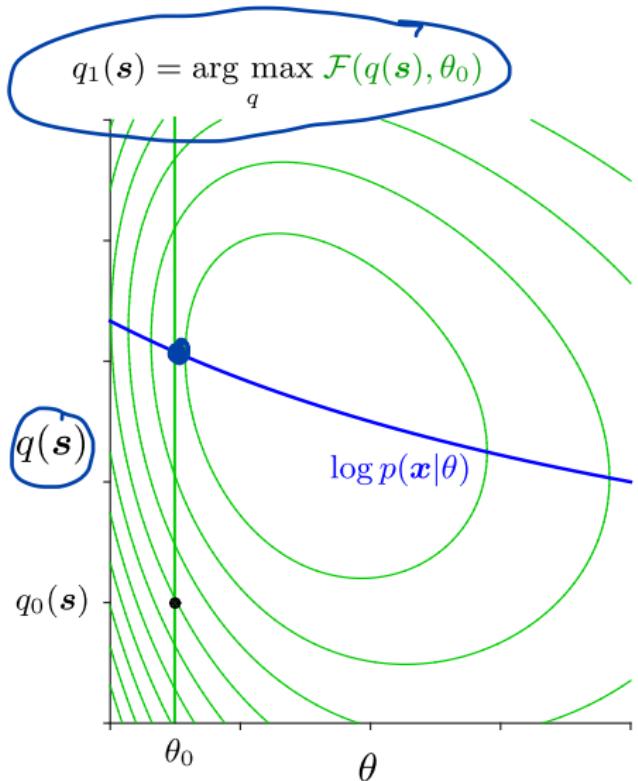
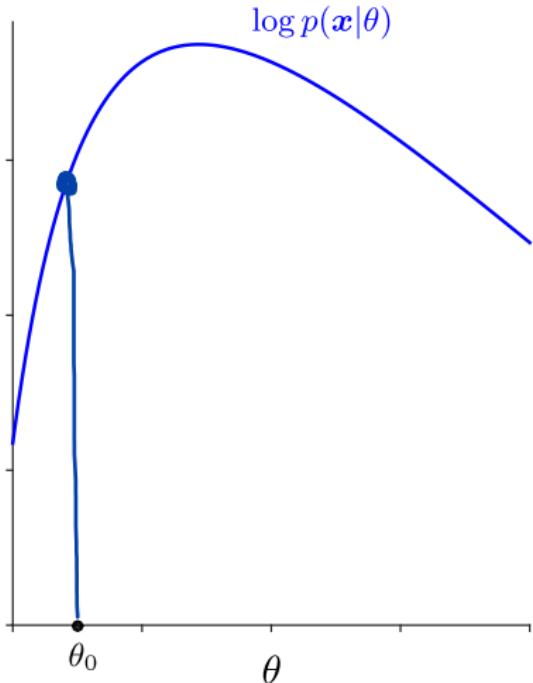
The Expectation Maximisation (EM) algorithm

$$\mathcal{F}(q(s), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(s)\|p(s|\mathbf{x}, \theta))$$



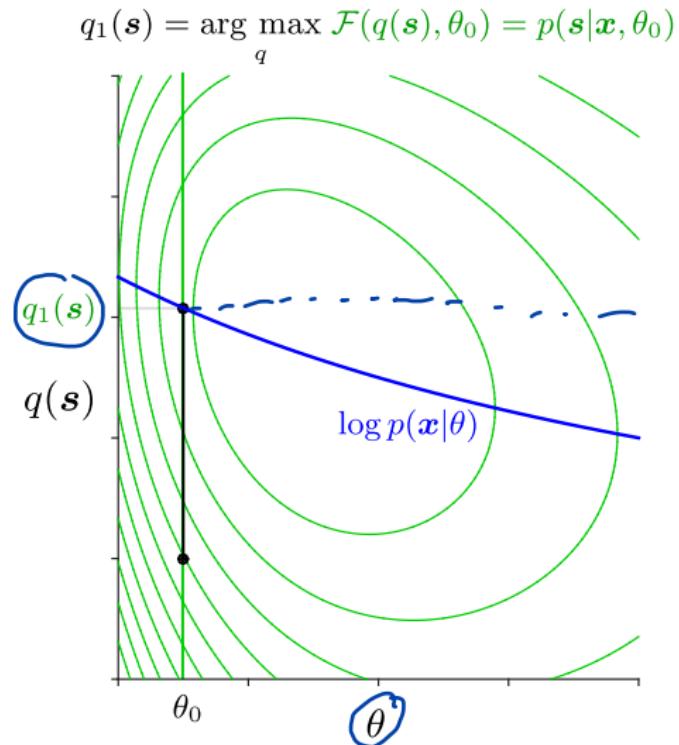
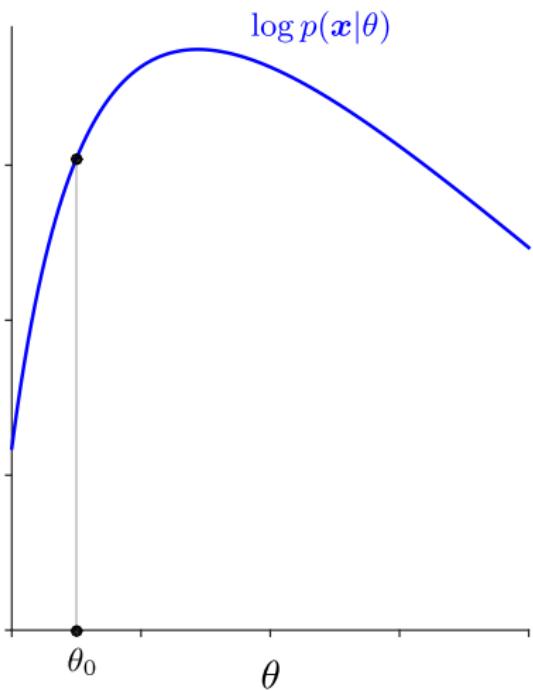
The Expectation Maximisation (EM) algorithm

$$\mathcal{F}(q(\mathbf{s}), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(\mathbf{s})\|p(\mathbf{s}|\mathbf{x}, \theta))$$



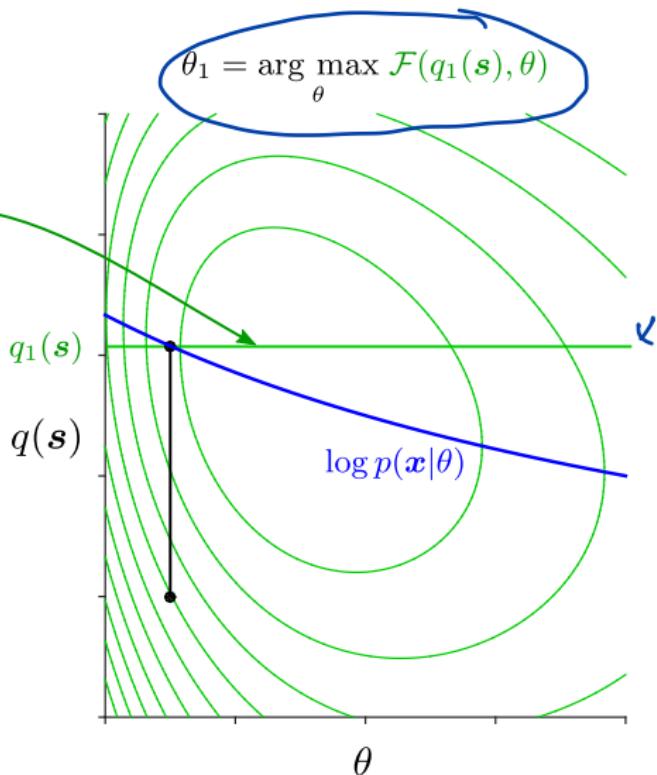
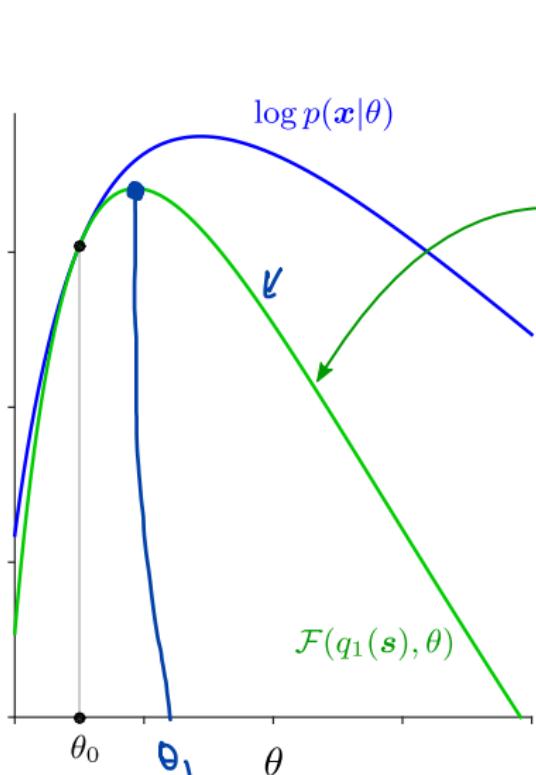
The Expectation Maximisation (EM) algorithm

$$\mathcal{F}(q(\mathbf{s}), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(\mathbf{s})\|p(\mathbf{s}|\mathbf{x}, \theta))$$



The Expectation Maximisation (EM) algorithm

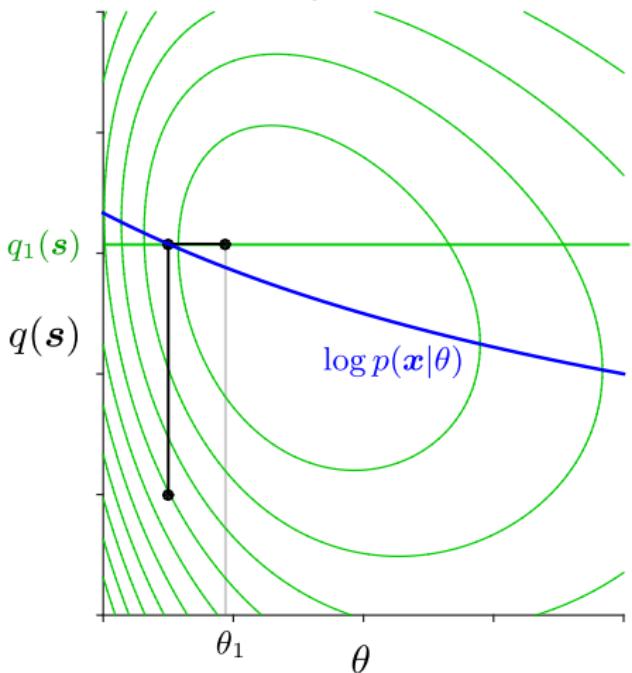
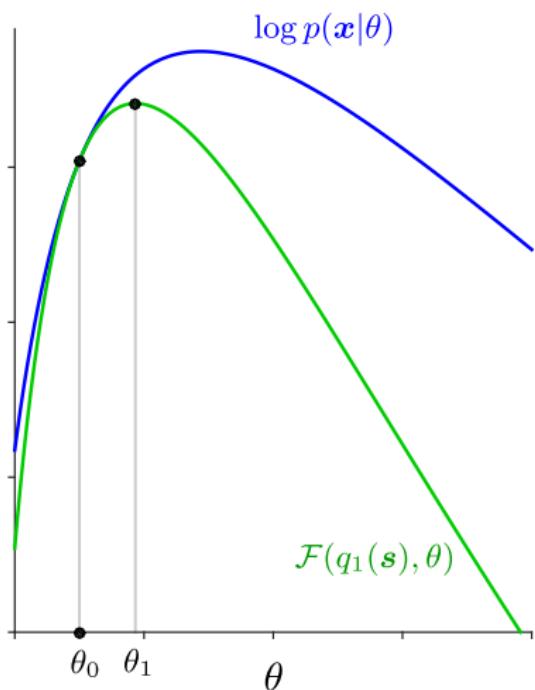
$$\mathcal{F}(q(s), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(s)\|p(s|\mathbf{x}, \theta))$$



The Expectation Maximisation (EM) algorithm

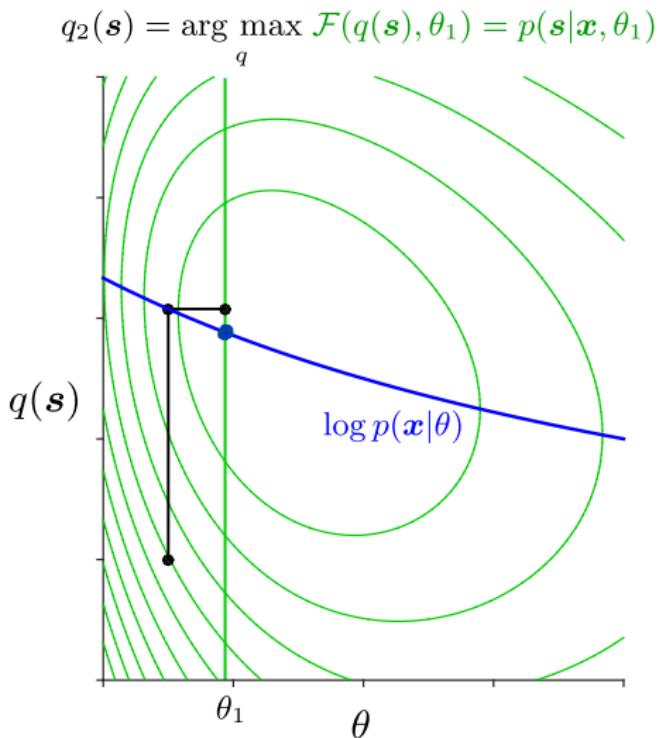
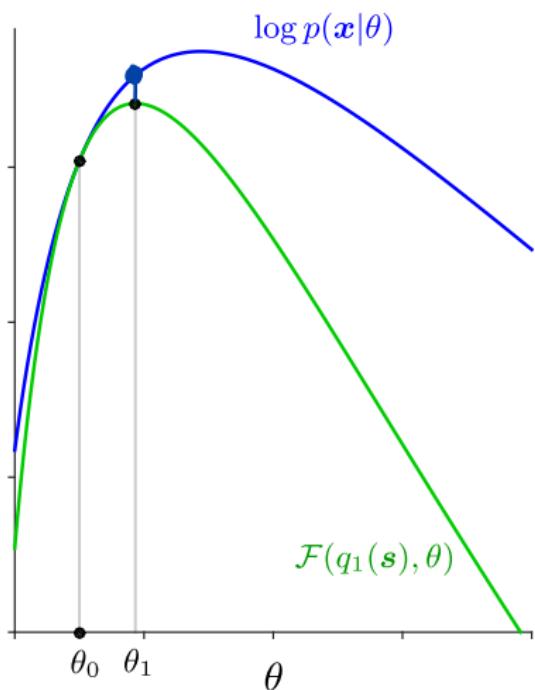
$$\mathcal{F}(q(s), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(s)\|p(s|\mathbf{x}, \theta))$$

$$\theta_1 = \arg \max_{\theta} \mathcal{F}(q_1(s), \theta)$$



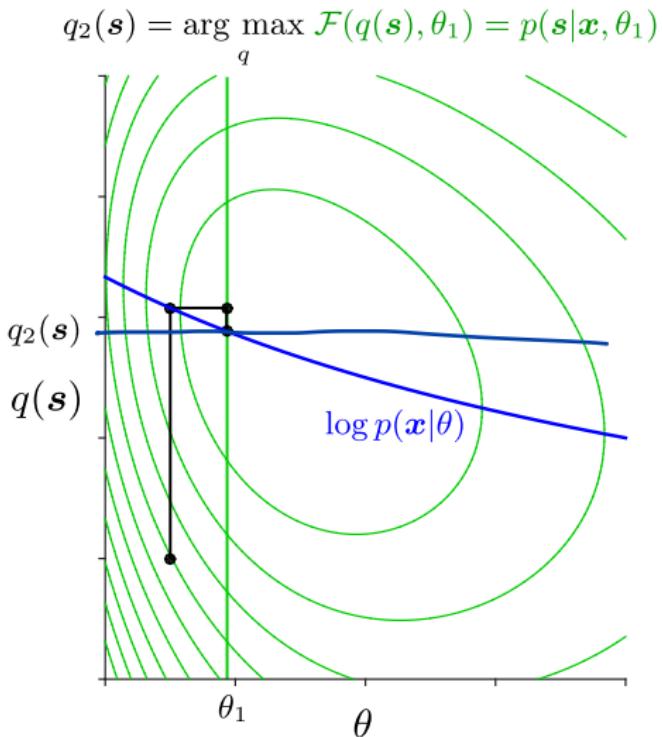
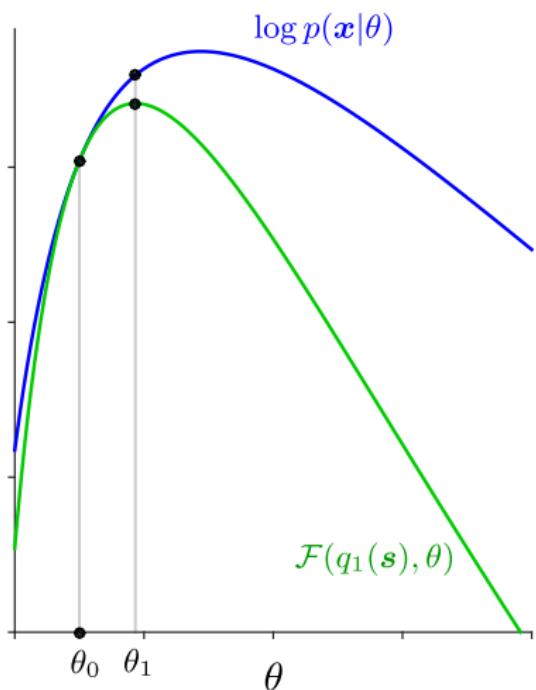
The Expectation Maximisation (EM) algorithm

$$\mathcal{F}(q(s), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(s)\|p(s|\mathbf{x}, \theta))$$



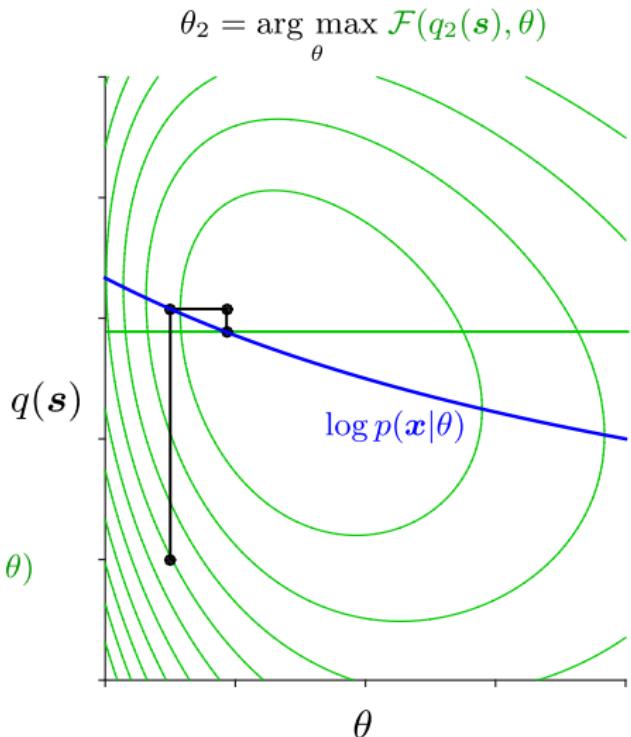
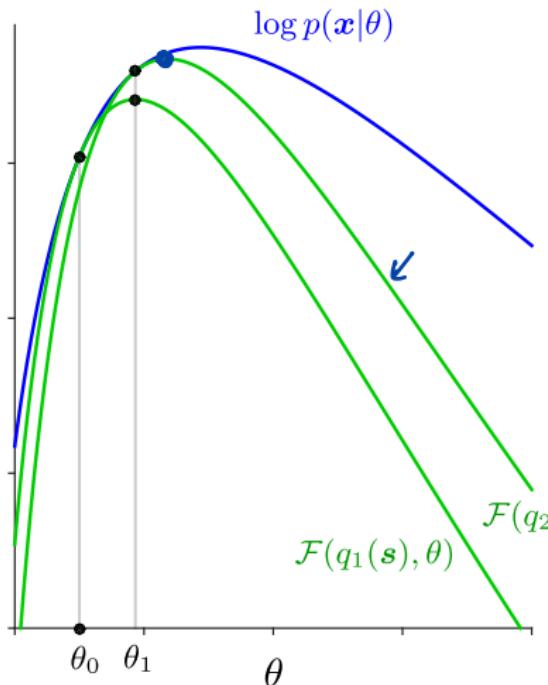
The Expectation Maximisation (EM) algorithm

$$\mathcal{F}(q(s), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(s)\|p(s|\mathbf{x}, \theta))$$



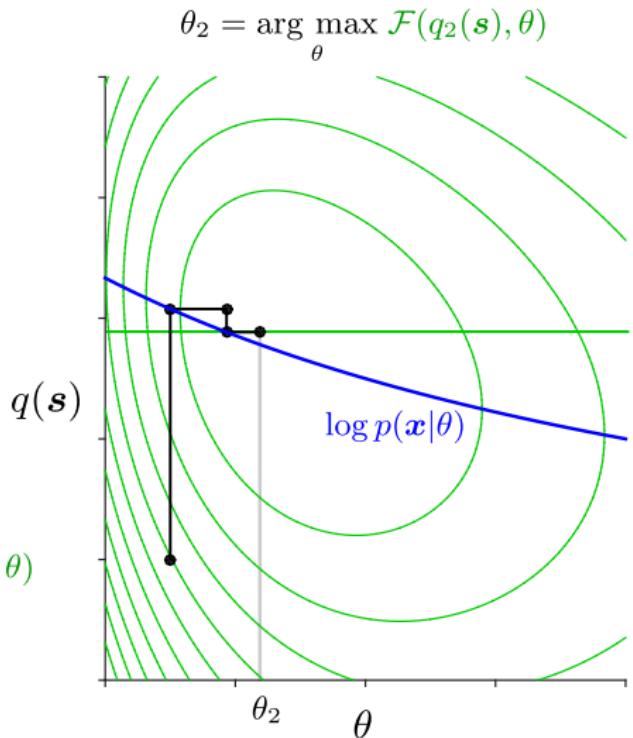
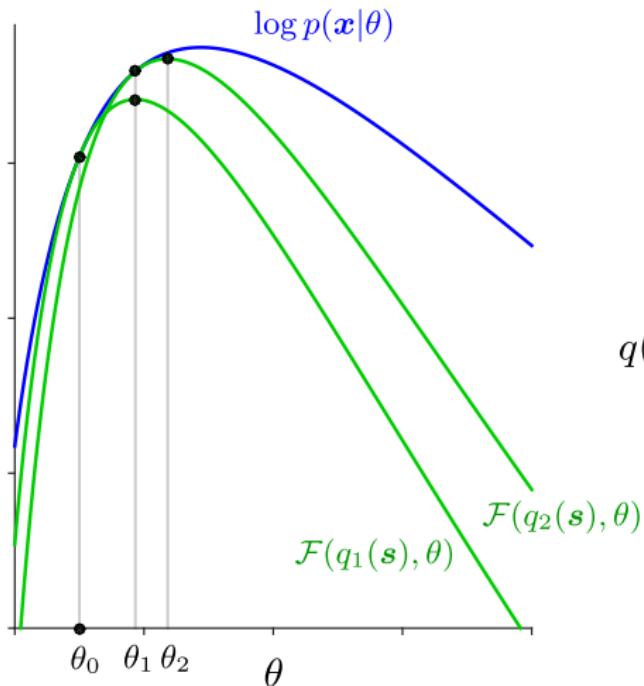
The Expectation Maximisation (EM) algorithm

$$\mathcal{F}(q(\mathbf{s}), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(\mathbf{s})\|p(\mathbf{s}|\mathbf{x}, \theta))$$



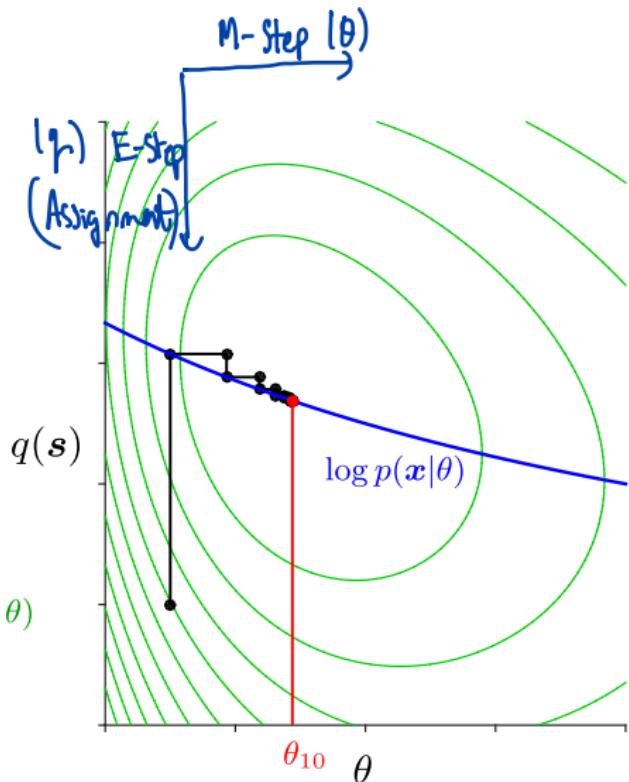
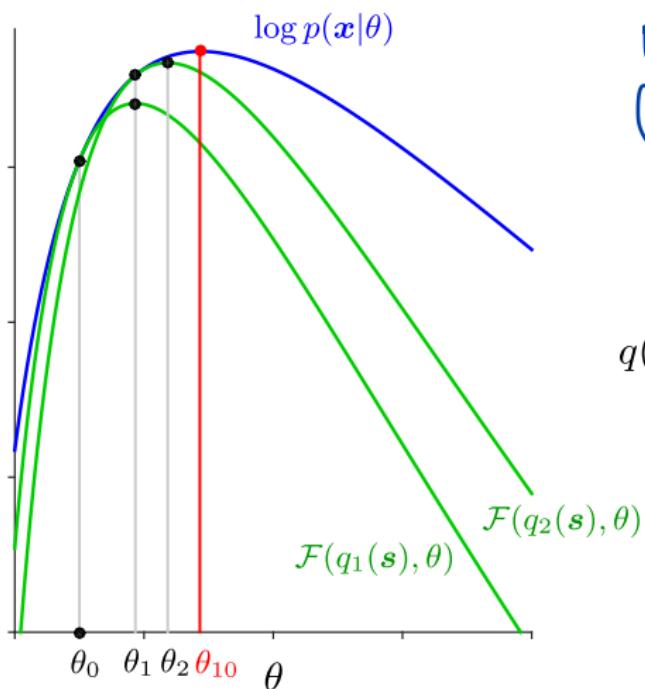
The Expectation Maximisation (EM) algorithm

$$\mathcal{F}(q(\mathbf{s}), \theta) = \log p(\mathbf{x}|\theta) - \mathcal{KL}(q(\mathbf{s})\|p(\mathbf{s}|\mathbf{x}, \theta))$$



The Expectation Maximisation (EM) algorithm

$$\mathcal{F}(q(s), \theta) = \underbrace{\log p(\mathbf{x}|\theta)}_{\text{log } p(\mathbf{x}|\theta)} - \underbrace{\mathcal{KL}(q(s)\|p(s|\mathbf{x}, \theta))}_{\text{M-Step } (\theta)}$$



The Expectation Maximisation (EM) algorithm

From initial (random) parameters θ_0 iterate $t = 1, \dots, T$ the two steps:

E step: for fixed θ_{t-1} , maximize lower bound $\mathcal{F}(q(s), \theta_{t-1})$ wrt $q(s)$.

As log likelihood $\log p(\mathbf{x}|\theta)$ is independent of $q(s)$ this is equivalent to minimizing $\mathcal{KL}(q(s)||p(s|\mathbf{x}, \theta_{t-1}))$, so $q_t(s) = p(s|\mathbf{x}, \underline{\theta_{t-1}})$.

M step: for fixed $q_t(s)$ maximize the lower bound $\mathcal{F}(q_t(s), \theta)$ wrt θ .

$$\Rightarrow \mathcal{F}(q(s), \theta) = \sum_s q(s) \log (p(\mathbf{x}|s, \theta)p(s|\theta)) - \sum_s q(s) \log q(s)$$
$$\begin{aligned} \mathcal{F}(q(s), \theta) &= \log p(\mathbf{x}|\theta) - \text{KL}(q(s) || p(s|\mathbf{x}, \theta)) \\ &= \log p(\mathbf{x}|\theta) - \sum_s q(s) \log \frac{q(s)}{p(s|\mathbf{x}, \theta)} \\ &= \sum_s q(s) \log p(\mathbf{x}|\theta) + \sum_s q(s) \log p(s|\mathbf{x}, \theta) - \sum_s q(s) \log q(s) \\ &= \sum_s q(s) \log [p(\mathbf{x}|\theta) p(s|\mathbf{x}, \theta)] - \sum_s q(s) \log q(s) \end{aligned}$$

Entropy
product rule
 $p(\mathbf{x}, s|\theta) = p(s|\mathbf{x}).$
 $p(\mathbf{x}|s, \theta)$
Model
specifies
this!

The Expectation Maximisation (EM) algorithm

From initial (random) parameters θ_0 iterate $t = 1, \dots, T$ the two steps:

E step: for fixed θ_{t-1} , maximize lower bound $\mathcal{F}(q(s), \theta_{t-1})$ wrt $q(s)$.

As log likelihood $\log p(\mathbf{x}|\theta)$ is independent of $q(s)$ this is equivalent to minimizing $\mathcal{KL}(q(s)||p(s|\mathbf{x}, \theta_{t-1}))$, so $q_t(s) = p(s|\mathbf{x}, \theta_{t-1})$.

M step: for fixed $q_t(s)$ maximize the lower bound $\mathcal{F}(q_t(s), \theta)$ wrt θ .

$$\mathcal{F}(q(s), \theta) = \sum_s q(s) \log (p(\mathbf{x}|s, \theta)p(s|\theta)) - \sum_s q(s) \log q(s),$$

the second term is the entropy of $q(s)$, independent of θ , so the M step is

$$\theta_t = \operatorname{argmax}_{\theta} \sum_s q_t(s) \log (p(\mathbf{x}|s, \theta)p(s|\theta)).$$

The Expectation Maximisation (EM) algorithm

From initial (random) parameters θ_0 iterate $t = 1, \dots, T$ the two steps:

E step: for fixed θ_{t-1} , maximize lower bound $\mathcal{F}(q(s), \theta_{t-1})$ wrt $q(s)$.

As log likelihood $\log p(\mathbf{x}|\theta)$ is independent of $q(s)$ this is equivalent to minimizing $\mathcal{KL}(q(s)||p(s|\mathbf{x}, \theta_{t-1}))$, so $q_t(s) = p(s|\mathbf{x}, \theta_{t-1})$.

M step: for fixed $q_t(s)$ maximize the lower bound $\mathcal{F}(q_t(s), \theta)$ wrt θ .

$$\mathcal{F}(q(s), \theta) = \underbrace{\left(\sum_s q(s) \log (p(\mathbf{x}|s, \theta)p(s|\theta)) - \sum_s q(s) \log q(s) \right)}$$

the second term is the entropy of $q(s)$, independent of θ , so the M step is

$$\theta_t = \operatorname{argmax}_{\theta} \sum_s q_t(s) \log (p(\mathbf{x}|s, \theta)p(s|\theta)).$$

Although the steps work with the lower bound (Lyupanov* function), each iteration cannot decrease the log likelihood as

$$\underbrace{\log p(\mathbf{x}|\theta_{t-1})}_{\text{E step}} = \underbrace{\mathcal{F}(q_t(s), \theta_{t-1})}_{\text{M step}} \leq \underbrace{\mathcal{F}(q_t(s), \theta_t)}_{\text{lower bound}} \leq \underbrace{\log p(\mathbf{x}|\theta_t)}$$

Application of EM to Mixture of Gaussians (E Step)

- ▶ Assume $D = 1$ dimensional data for x simplicity
- ▶ Gaussian mixture model parameters: $\theta = \{\mu_k, \sigma_k^2, \pi_k\}_{k=1\dots K}$
- ▶ One latent variable per datapoint $s_n, n = 1 \dots N$ takes values $1 \dots K$.

Probability of the observations given the latent variables and the parameters, and the prior on latent variables are:

$$\left| \begin{array}{l} p(x_n | s_n = k, \theta) = \frac{1}{\sqrt{2\pi\sigma_k^2}} e^{-\frac{1}{2\sigma_k^2}(x_n - \mu_k)^2} \\ p(s_n = k | \theta) = \pi_k \end{array} \right.$$

Application of EM to Mixture of Gaussians (E Step)

- ▶ Assume $D = 1$ dimensional data for x simplicity
- ▶ Gaussian mixture model parameters: $\theta = \{\mu_k, \sigma_k^2, \pi_k\}_{k=1\dots K}$
- ▶ One latent variable per datapoint $s_n, n = 1 \dots N$ takes values $1 \dots K$.

Probability of the observations given the latent variables and the parameters, and the prior on latent variables are:

$$p(x_n | s_n = k, \theta) = \frac{1}{\sqrt{2\pi\sigma_k^2}} e^{-\frac{1}{2\sigma_k^2}(x_n - \mu_k)^2}$$

$$p(s_n = k | \theta) = \pi_k$$

(Bayes' Rule)

so the E step becomes:

$$q(s_n = k) = p(s_n = k | x_n, \theta) \propto p(x_n, s_n = k | \theta) = \frac{\pi_k}{\sqrt{2\pi\sigma_k^2}} e^{-\frac{1}{2\sigma_k^2}(x_n - \mu_k)^2} = u_{nk}$$

That is: $q(s_n = k) = r_{nk} = \frac{u_{nk}}{u_n}$ where $u_n = \sum_{k=1}^K u_{nk}$ $p(x_n | \theta)$

Responsibilities

Application of EM to Mixture of Gaussians (E Step)

- ▶ Assume $D = 1$ dimensional data for x simplicity
- ▶ Gaussian mixture model parameters: $\theta = \{\mu_k, \sigma_k^2, \pi_k\}_{k=1\dots K}$
- ▶ One latent variable per datapoint $s_n, n = 1 \dots N$ takes values $1 \dots K$.

Probability of the observations given the latent variables and the parameters, and the prior on latent variables are:

$$p(x_n | s_n = k, \theta) = \frac{1}{\sqrt{2\pi\sigma_k^2}} e^{-\frac{1}{2\sigma_k^2}(x_n - \mu_k)^2} \quad p(s_n = k | \theta) = \pi_k$$

so the E step becomes:

$$q(s_n = k) = p(s_n = k | x_n, \theta) \propto p(x_n, s_n = k | \theta) = \frac{\pi_k}{\sqrt{2\pi\sigma_k^2}} e^{-\frac{1}{2\sigma_k^2}(x_n - \mu_k)^2} = u_{nk}$$

That is: $q(s_n = k) = r_{nk} = \frac{u_{nk}}{u_n}$ where $u_n = \sum_{k=1}^K u_{nk}$

Posterior for each latent variable, s_n follows a categorical distribution with probability given by the product of the prior and likelihood, renormalised. r_{nk} is called the **responsibility** that component k takes for data point n .

Application of EM to Mixture of Gaussians (M Step)

The lower bound is $\textcircled{1} \sum_s q(s) \log [p(s|\theta) p(x|s,\theta)] - \sum_s q(s) \log q(s)$

$$\textcircled{2} \mathcal{F}(q(s), \theta) = \sum_{n=1}^N \sum_{k=1}^K q(s_n = k) \left[\log(\pi_k) - \frac{1}{2\sigma_k^2} (x_n - \mu_k)^2 - \frac{1}{2} \log(\sigma_k^2) \right] + \text{const.}$$

$$\textcircled{A} \text{ iid } q(\underline{s}) = q(s_1, s_2, s_3, \dots, s_N) = p(\underline{s} | \underline{x}, \theta)$$

$$\prod_n q(s_n)$$

$$\prod_{n=1}^N p(s_n | \underline{x}_n, \theta)$$

Application of EM to Mixture of Gaussians (M Step)

The lower bound is

$$\mathcal{F}(q(s), \theta) = \sum_{n=1}^N \sum_{k=1}^K q(s_n = k) [\log(\pi_k) - \frac{1}{2\sigma_k^2} (x_n - \underline{\mu}_k)^2 - \frac{1}{2} \log(\sigma_k^2)] + \text{const.}$$

The M step, optimizing $\mathcal{F}(q(s), \theta)$ wrt the parameters, θ

$$\frac{\partial \mathcal{F}}{\partial \mu_j} = \sum_{n=1}^N q(s_n = j) \frac{x_n - \mu_j}{\sigma_j^2} = 0 \Rightarrow \underline{\mu}_j = \frac{\sum_{n=1}^N q(s_n = j) x_n}{\sum_{n=1}^N q(s_n = j)}$$

$$q(s_n = j) \quad \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \quad \underline{\mu}_1 = \frac{\underline{x}_1 + \underline{x}_4}{2}$$

= mean of the datapoints
that were assigned
to this cluster

Application of EM to Mixture of Gaussians (M Step)

The lower bound is

$$\mathcal{F}(q(s), \theta) = \sum_{n=1}^N \sum_{k=1}^K q(s_n = k) \left[\log(\pi_k) - \frac{1}{2\sigma_k^2} (x_n - \mu_k)^2 - \frac{1}{2} \log(\sigma_k^2) \right] + \text{const.}$$

The M step, optimizing $\mathcal{F}(q(s), \theta)$ wrt the parameters, θ

$$\frac{\partial \mathcal{F}}{\partial \mu_j} = \sum_{n=1}^N q(s_n = j) \frac{x_n - \mu_j}{\sigma_j^2} = 0 \Rightarrow \mu_j = \frac{\sum_{n=1}^N q(s_n = j) x_n}{\sum_{n=1}^N q(s_n = j)},$$

E step fills in the values of the hidden variables: M step just like performing **supervised learning** with known (soft) cluster assignments.

Application of EM to Mixture of Gaussians (M Step)

The lower bound is

$$\mathcal{F}(q(s), \theta) = \sum_{n=1}^N \sum_{k=1}^K q(s_n = k) \left[\log(\pi_k) - \frac{1}{2\sigma_k^2} (x_n - \mu_k)^2 - \frac{1}{2} \log(\sigma_k^2) \right] + \text{const.}$$

The M step, optimizing $\mathcal{F}(q(s), \theta)$ wrt the parameters, θ

$$\frac{\partial \mathcal{F}}{\partial \mu_j} = \sum_{n=1}^N q(s_n = j) \frac{x_n - \mu_j}{\sigma_j^2} = 0 \Rightarrow \mu_j = \frac{\sum_{n=1}^N q(s_n = j) x_n}{\sum_{n=1}^N q(s_n = j)},$$

$$\frac{\partial \mathcal{F}}{\partial \sigma_j^2} = \sum_{n=1}^N q(s_n = j) \left[\frac{(x_n - \mu_j)^2}{2\sigma_j^4} - \frac{1}{2\sigma_j^2} \right] = 0 \Rightarrow \sigma_j^2 = \frac{\sum_{n=1}^N q(s_n = j) (x_n - \mu_j)^2}{\sum_{n=1}^N q(s_n = j)}$$

E step fills in the values of the hidden variables: M step just like performing **supervised learning** with known (soft) cluster assignments.

Application of EM to Mixture of Gaussians (M Step)

The lower bound is

$$\mathcal{F}(q(s), \theta) = \sum_{n=1}^N \sum_{k=1}^K q(s_n = k) [\log(\pi_k) - \frac{1}{2\sigma_k^2} (x_n - \mu_k)^2 - \frac{1}{2} \log(\sigma_k^2)] + \text{const.}$$

The M step, optimizing $\mathcal{F}(q(s), \theta)$ wrt the parameters, θ

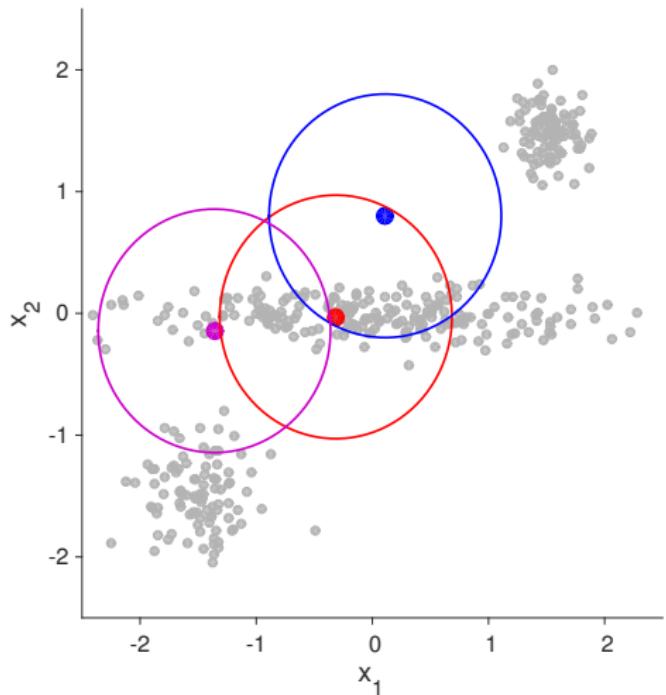
$$\frac{\partial \mathcal{F}}{\partial \mu_j} = \sum_{n=1}^N q(s_n = j) \frac{x_n - \mu_j}{\sigma_j^2} = 0 \Rightarrow \mu_j = \frac{\sum_{n=1}^N q(s_n = j) x_n}{\sum_{n=1}^N q(s_n = j)},$$

$$\frac{\partial \mathcal{F}}{\partial \sigma_j^2} = \sum_{n=1}^N q(s_n = j) \left[\frac{(x_n - \mu_j)^2}{2\sigma_j^4} - \frac{1}{2\sigma_j^2} \right] = 0 \Rightarrow \sigma_j^2 = \frac{\sum_{n=1}^N q(s_n = j) (x_n - \mu_j)^2}{\sum_{n=1}^N q(s_n = j)},$$

$$\frac{\partial [\mathcal{F} + \lambda(1 - \sum_k \pi_k)]}{\partial \pi_j} = \sum_{n=1}^N \frac{q(s_n = j)}{\pi_j} - \lambda = 0 \Rightarrow \pi_j = \frac{1}{N} \sum_{n=1}^N q(s_n = j)$$

E step fills in the values of the hidden variables: M step just like performing **supervised learning** with known (soft) cluster assignments.

EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

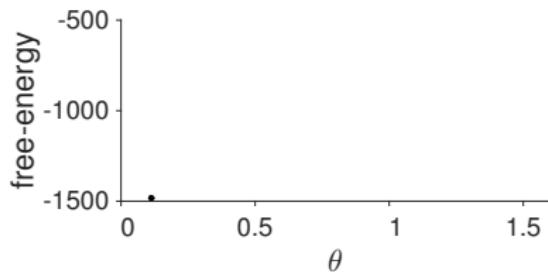
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

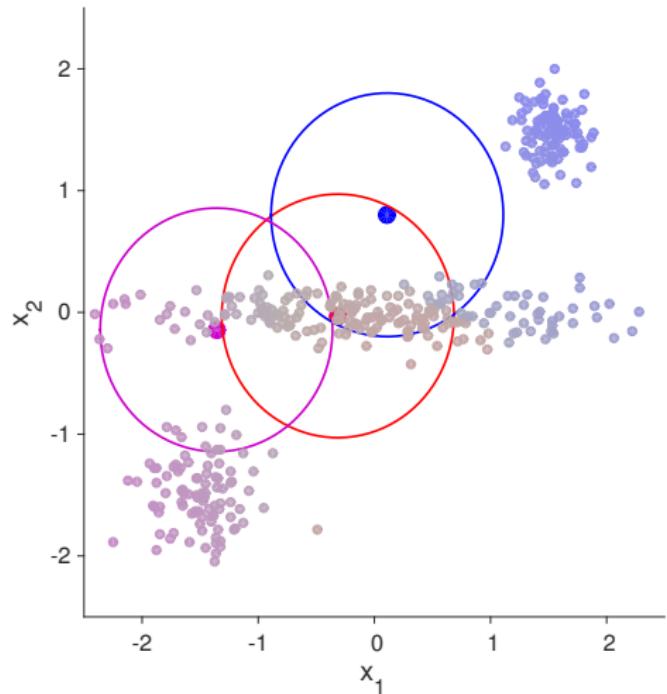
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

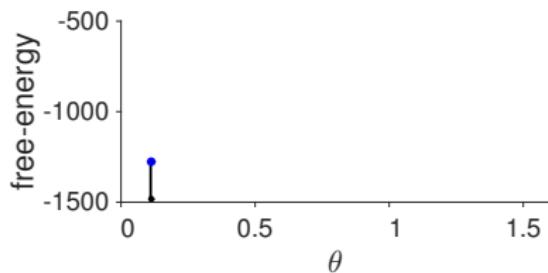
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

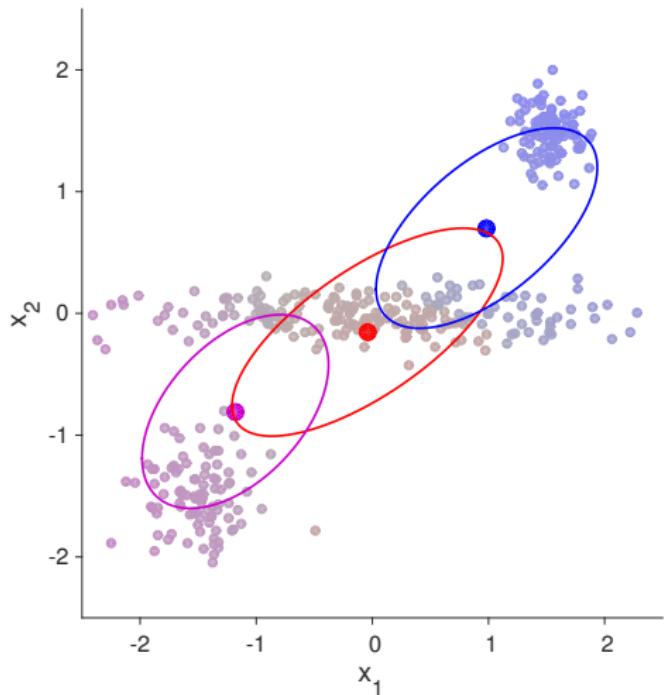
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

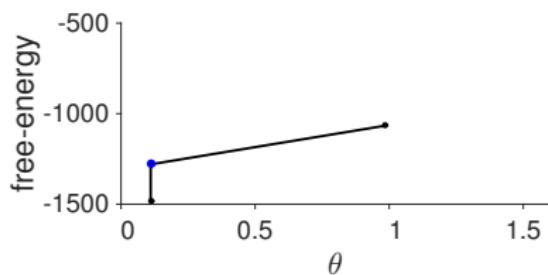
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

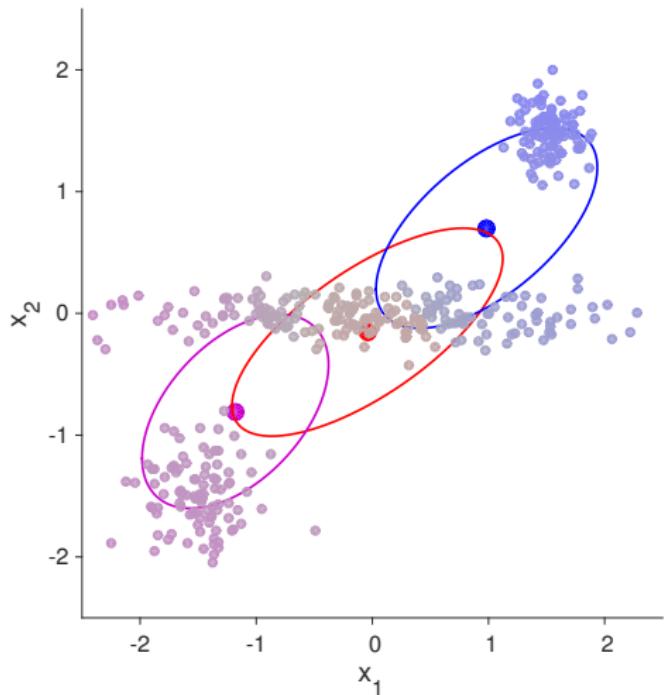
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

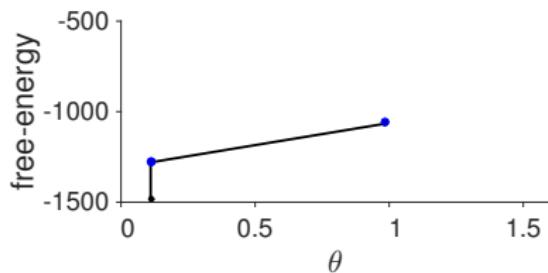
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

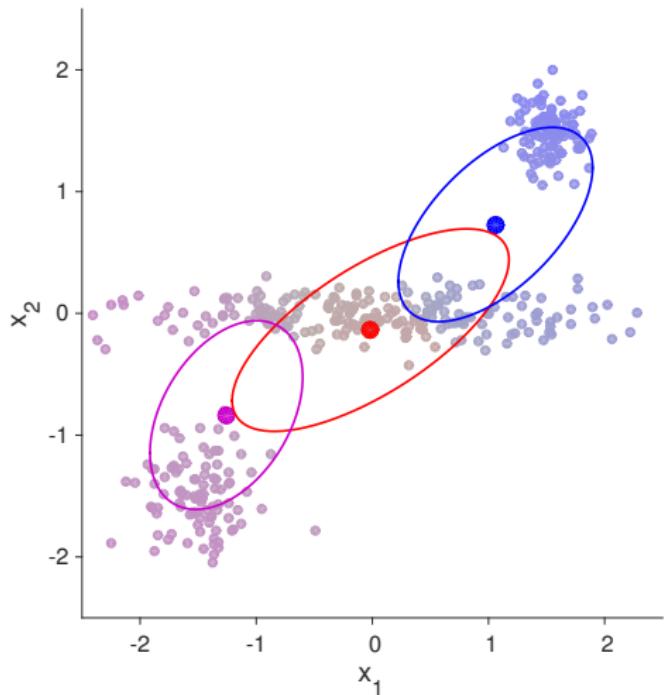
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

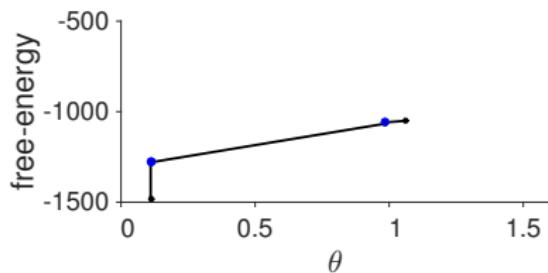
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

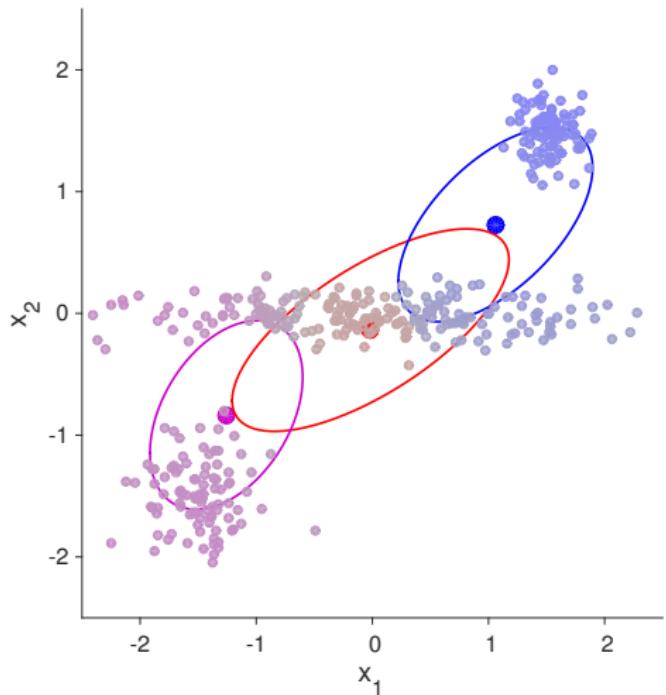
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

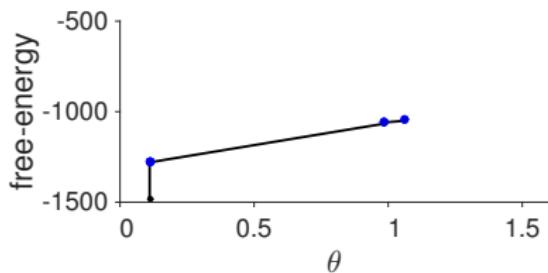
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

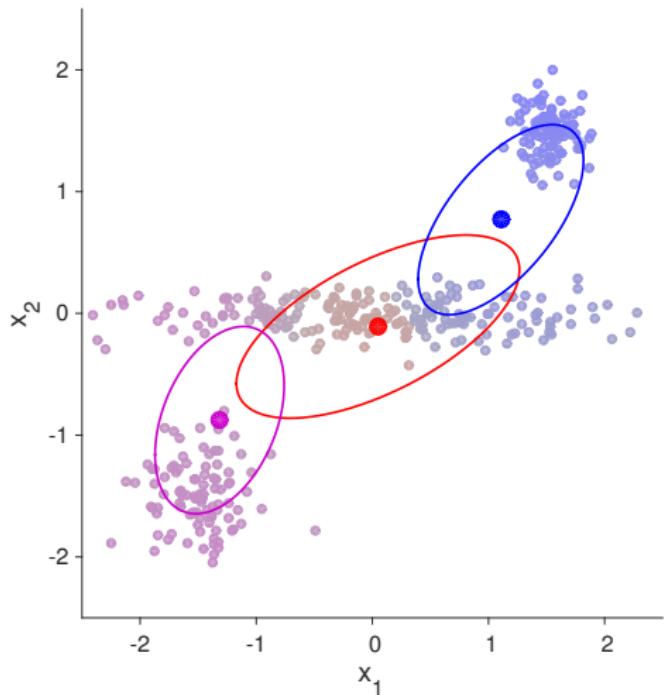
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

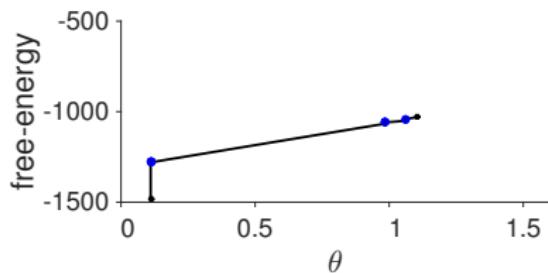
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

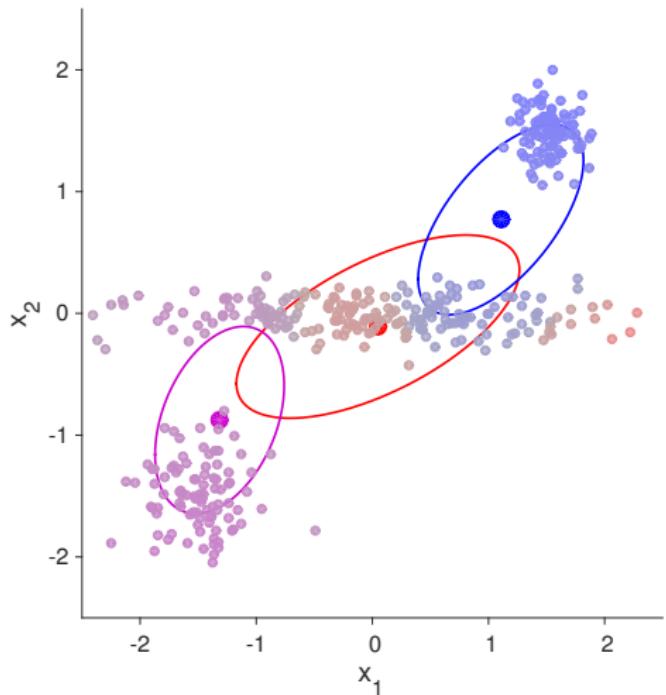
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

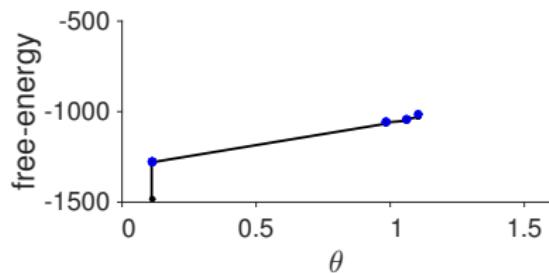
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

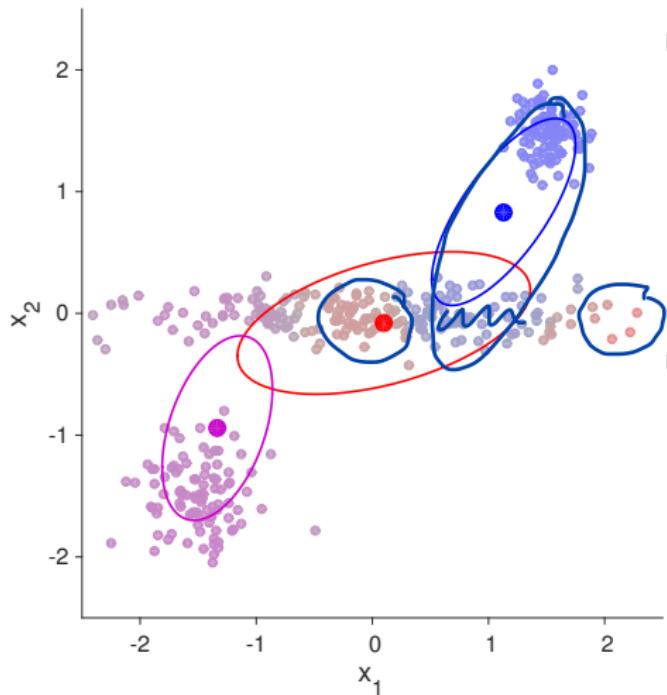
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

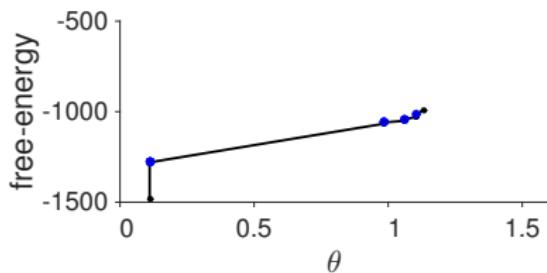
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

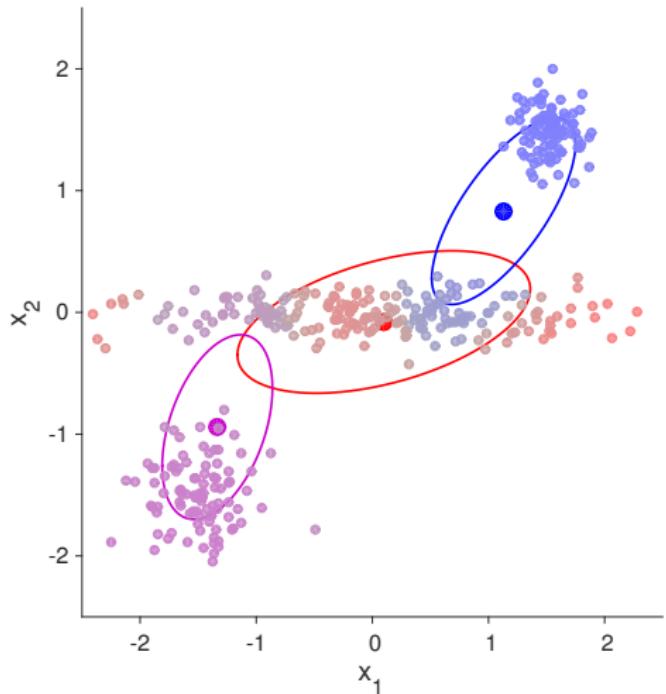
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

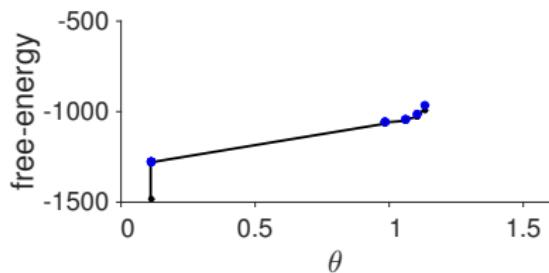
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

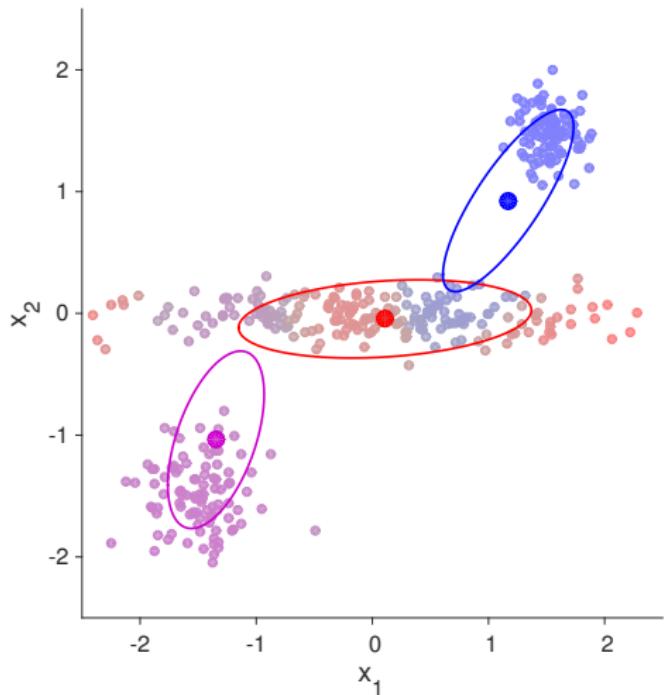
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

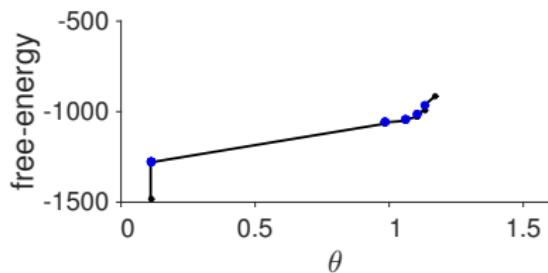
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

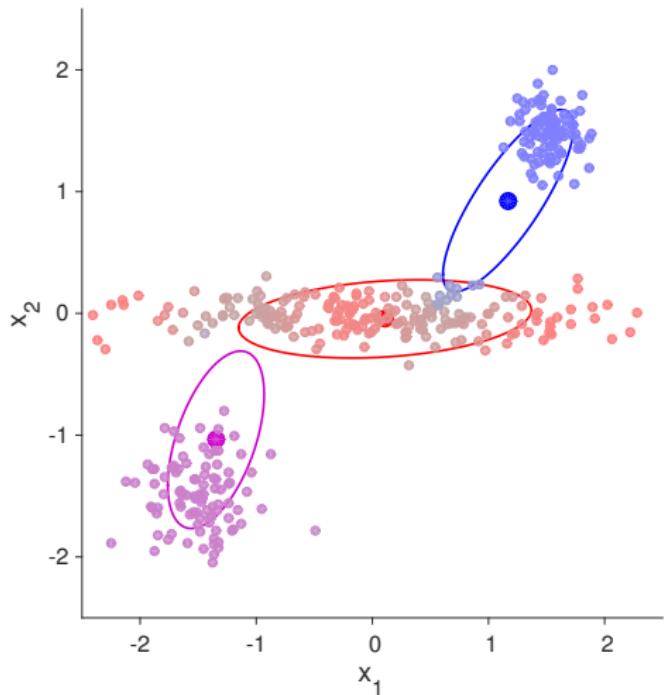
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

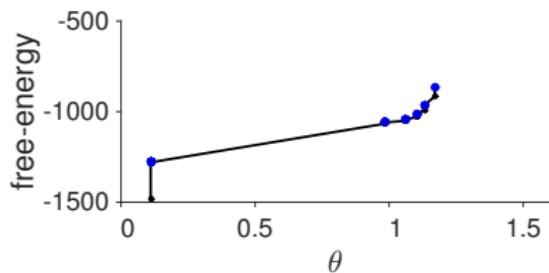
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

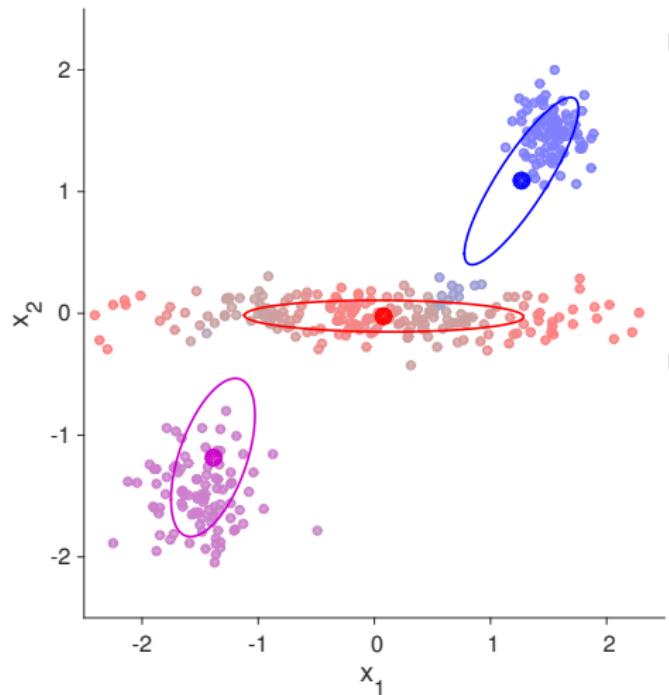
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

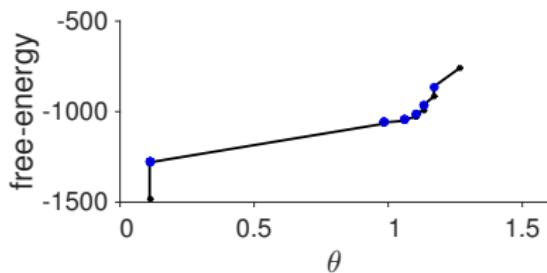
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

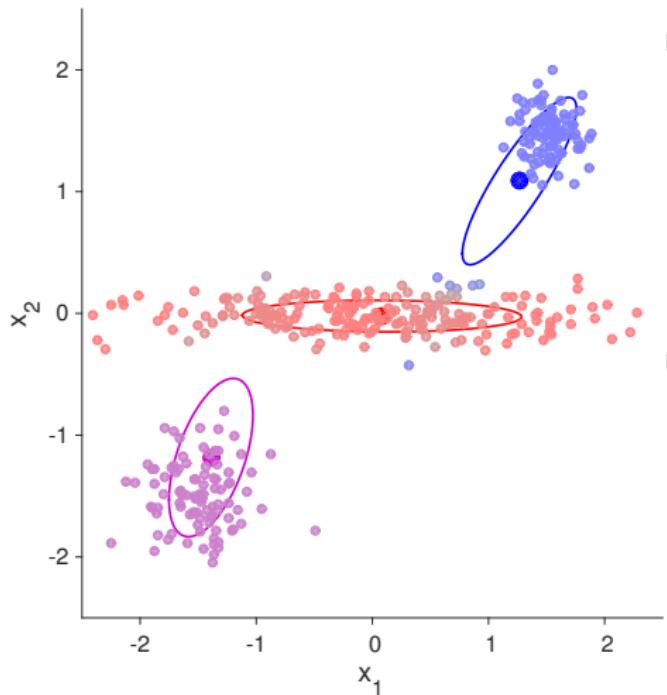
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

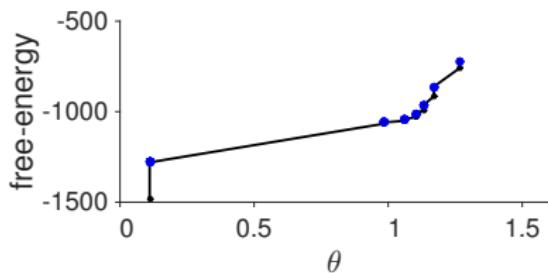
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

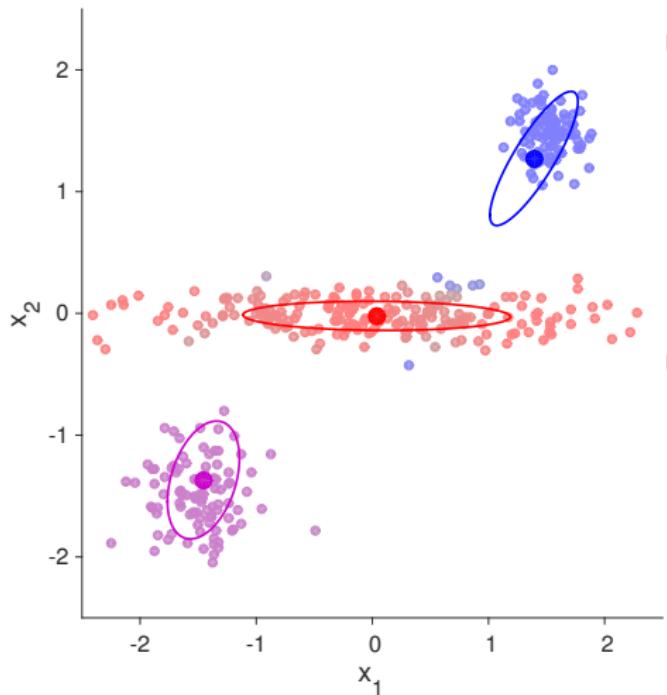
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

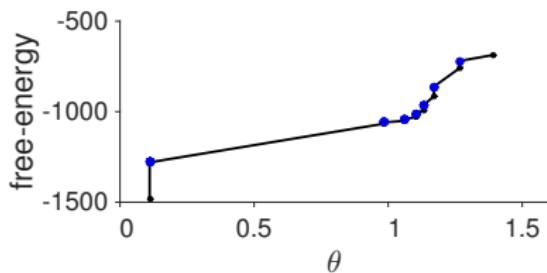
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

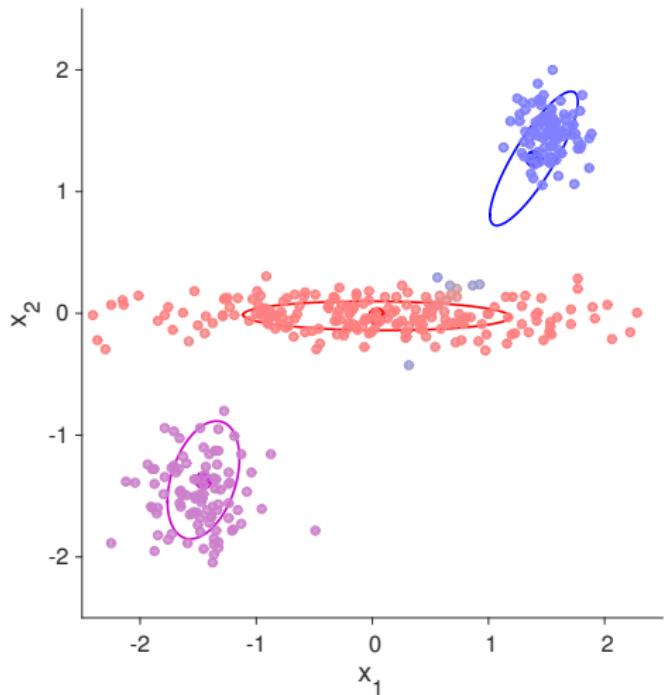
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

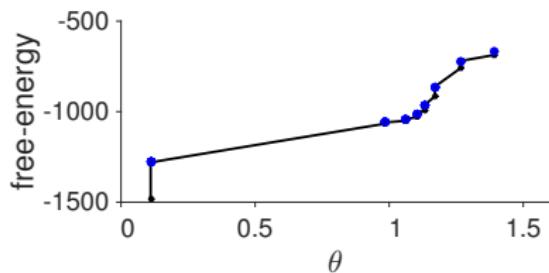
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

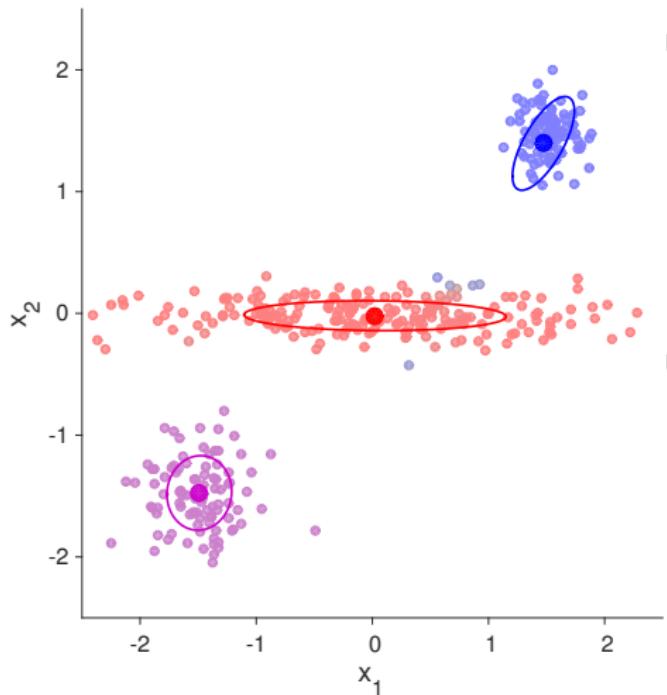
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

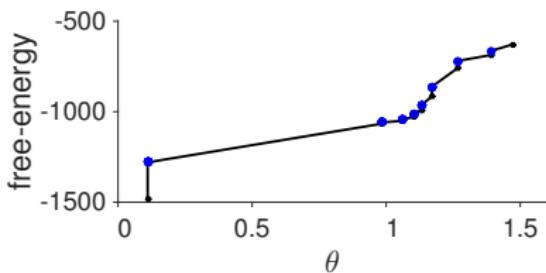
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

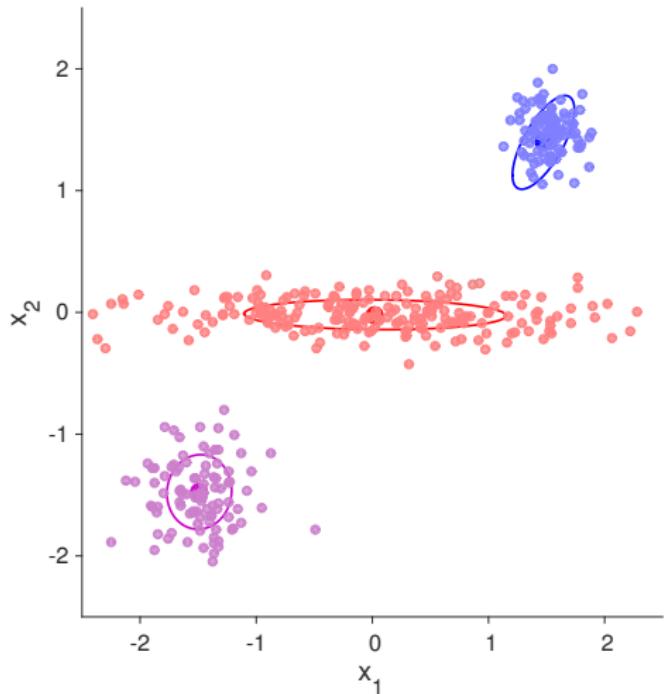
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

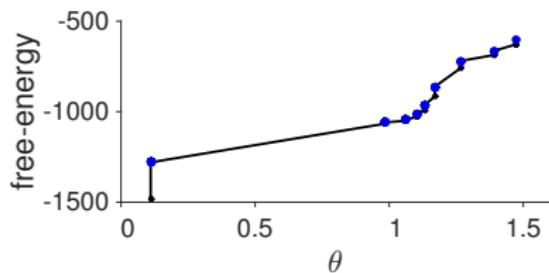
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

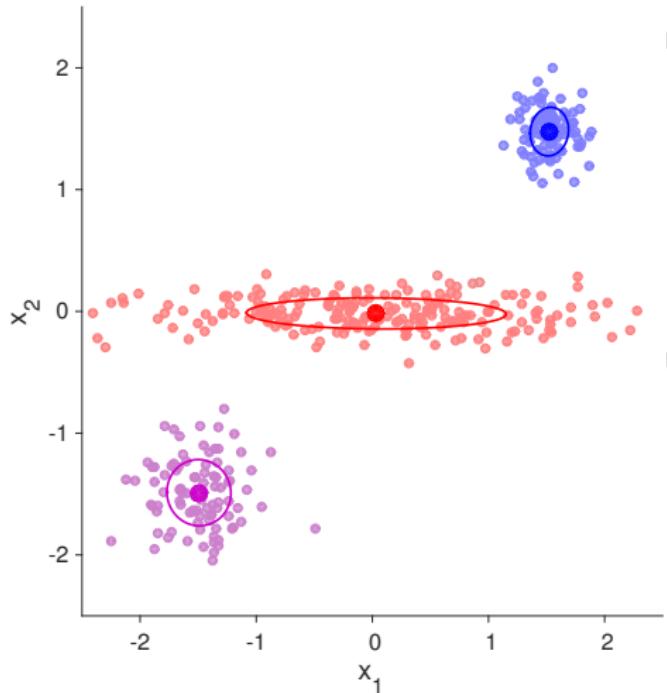
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

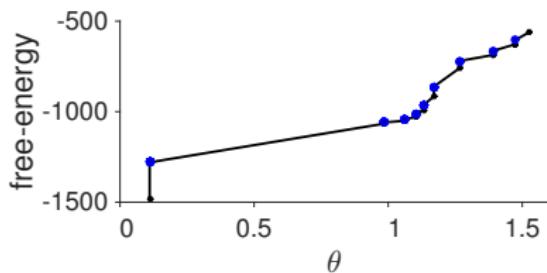
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

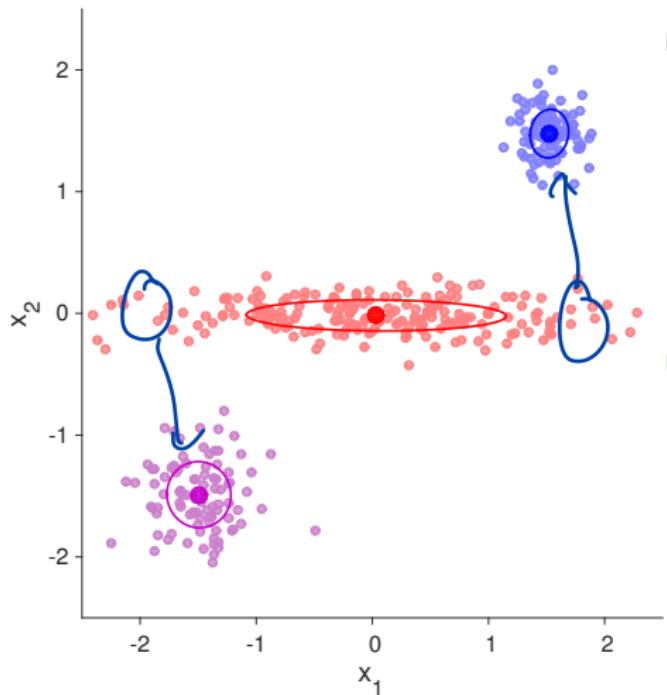
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



EM for MoGs: soft, non-axis aligned K-means



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

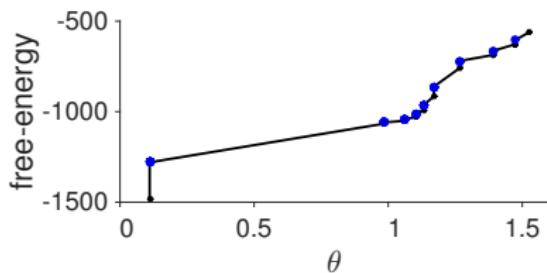
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

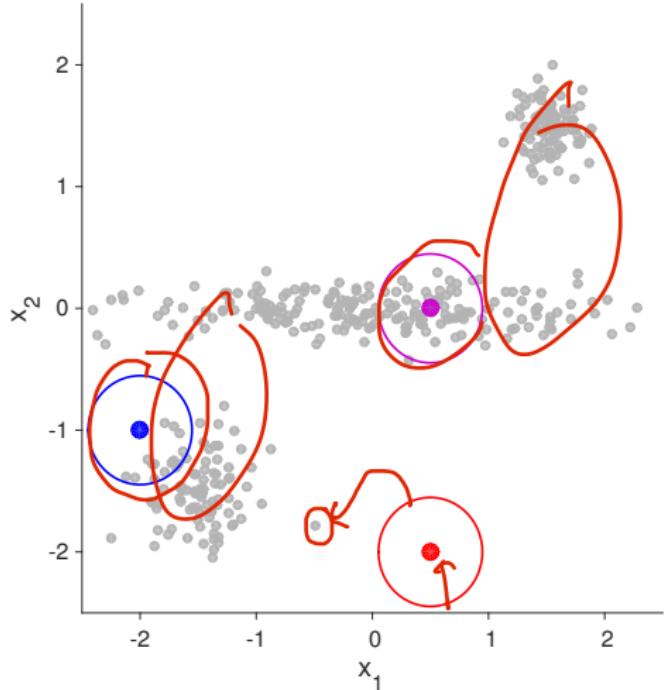
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

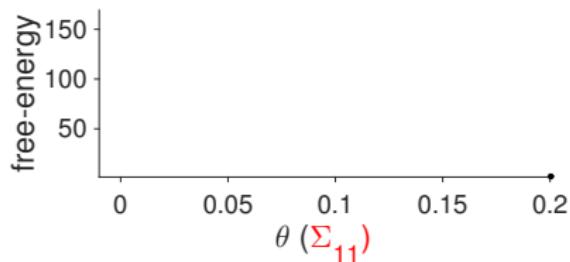
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

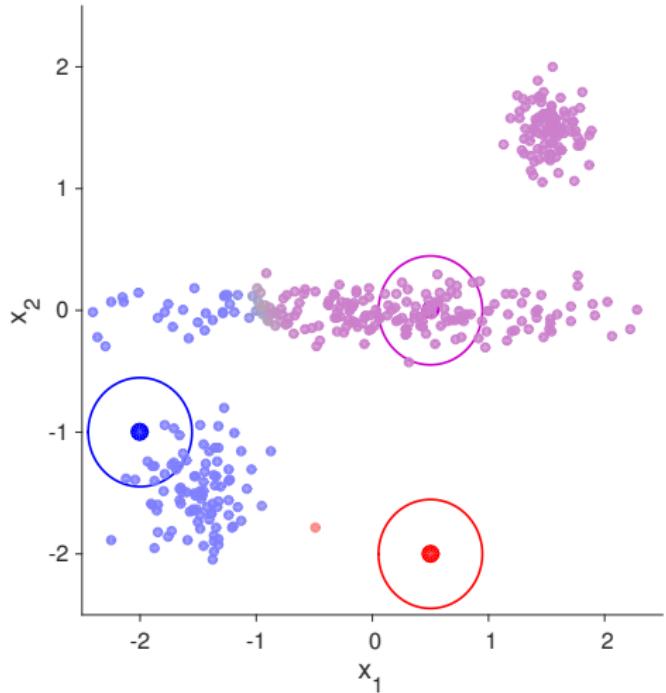
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

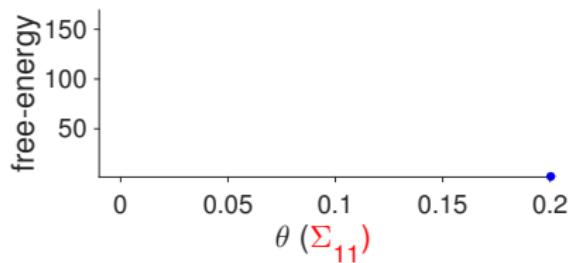
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

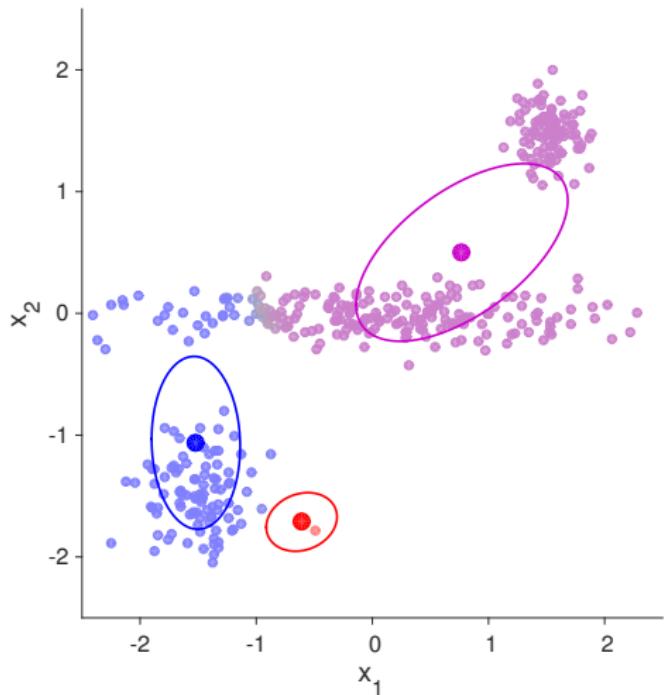
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

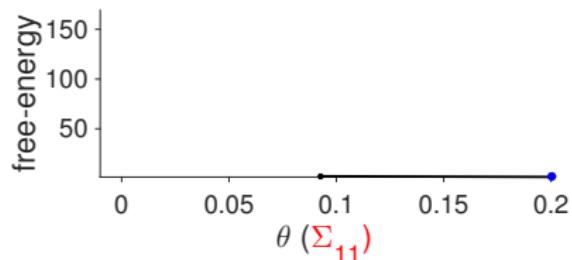
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

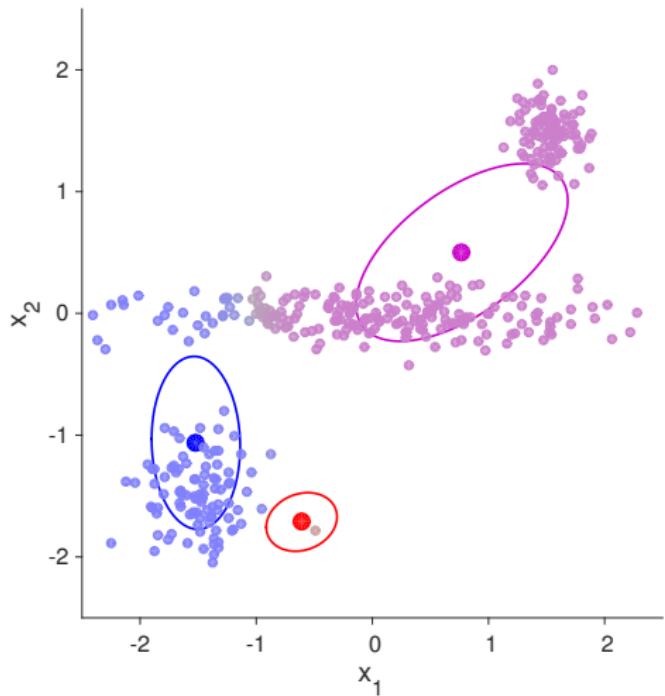
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

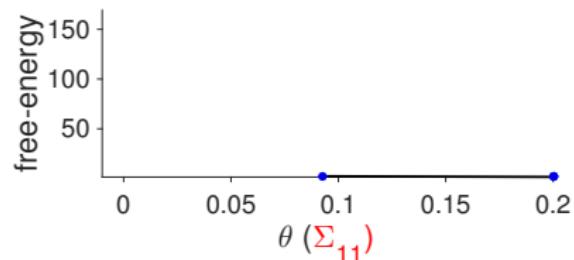
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

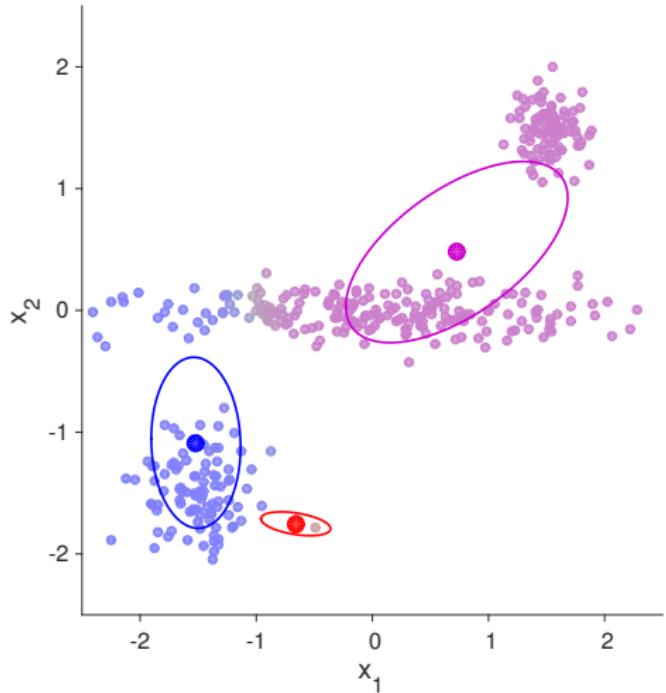
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

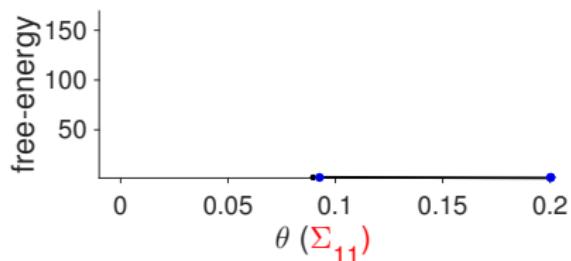
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

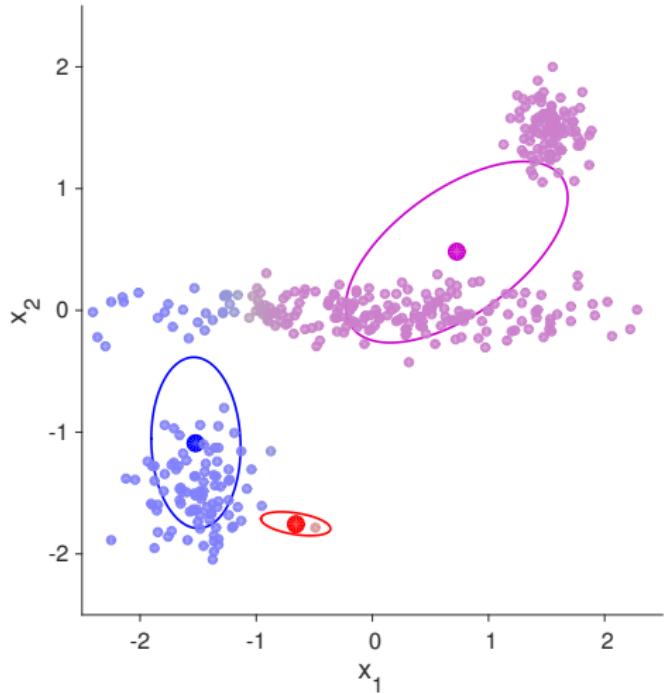
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

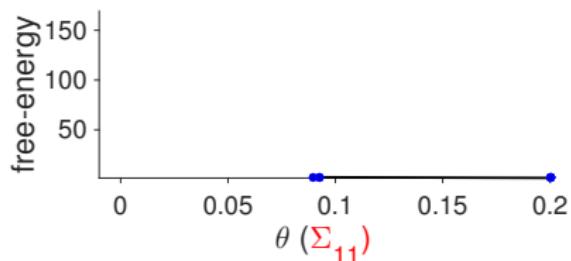
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

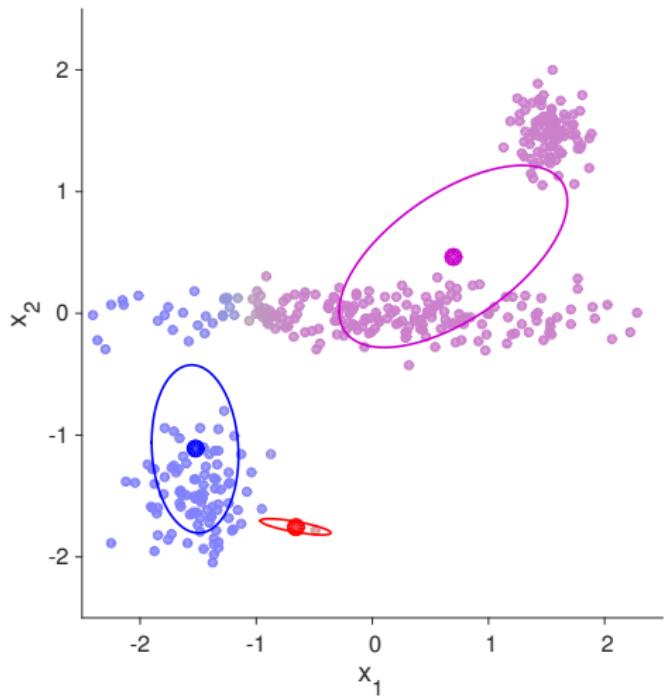
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

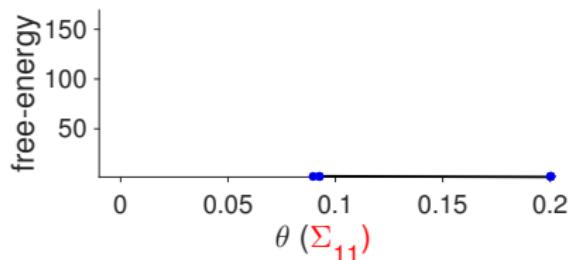
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

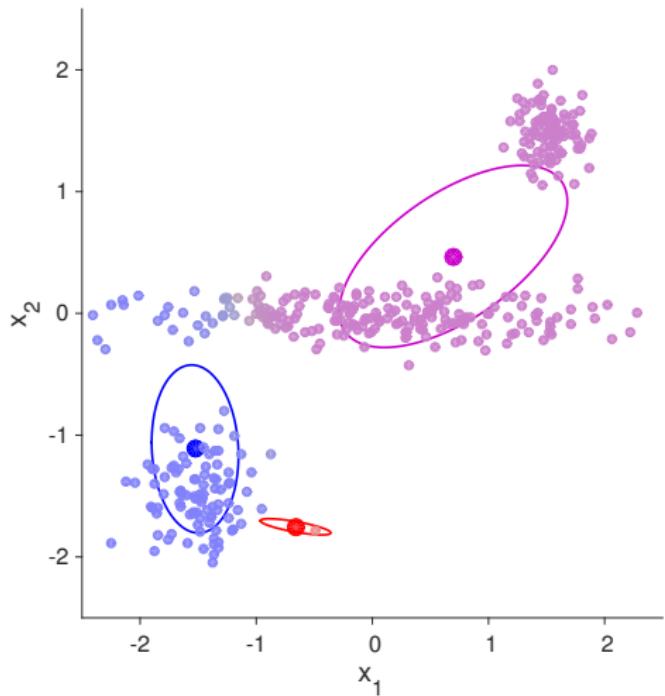
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

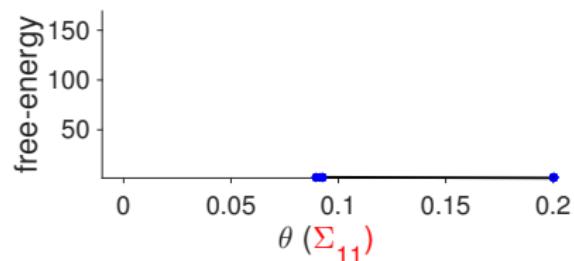
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

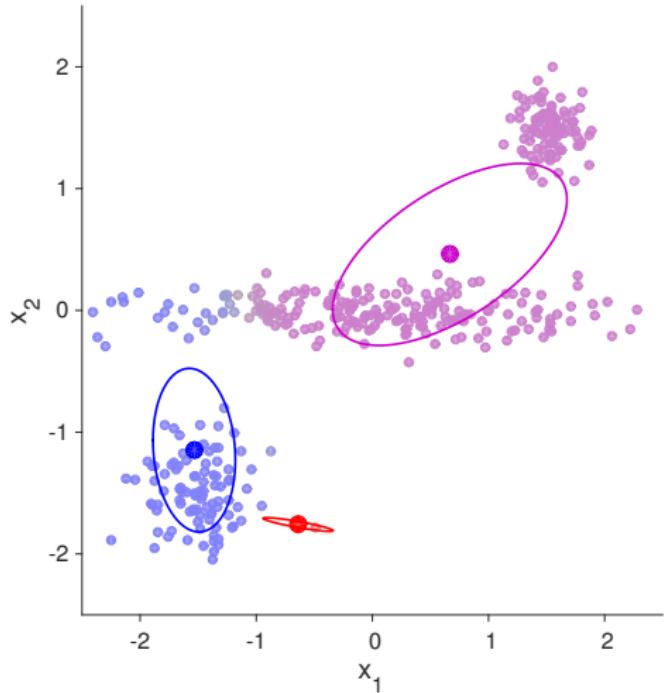
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

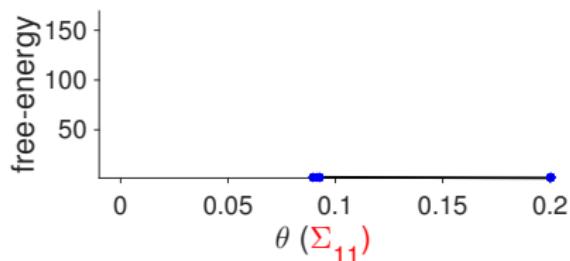
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

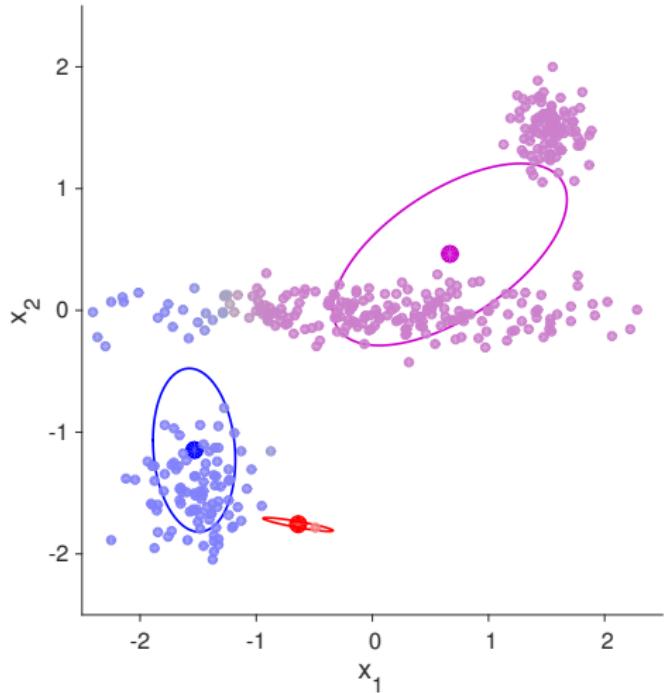
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

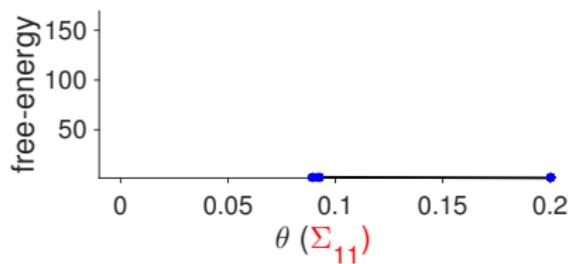
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

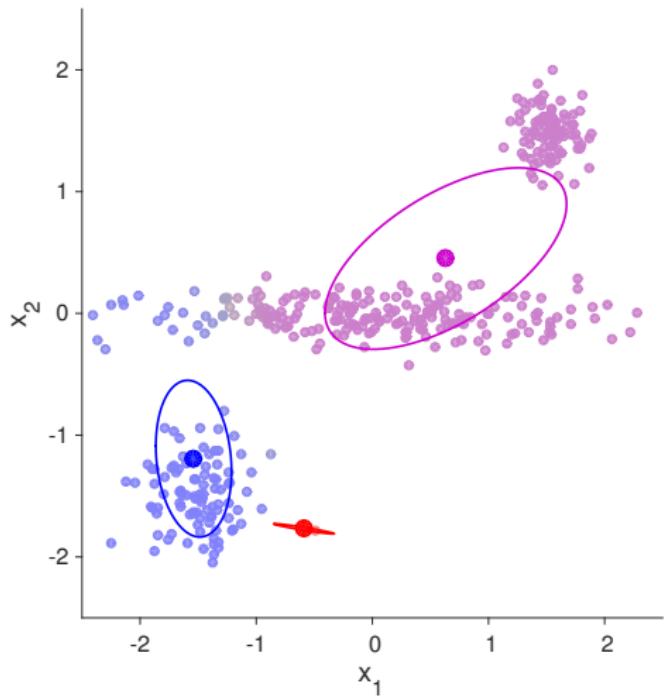
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

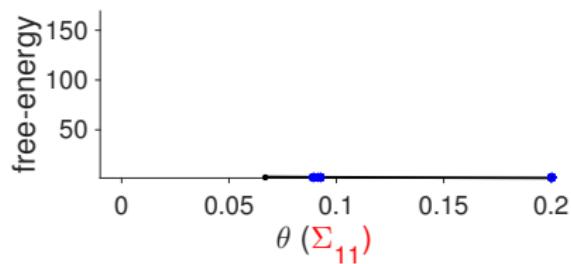
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

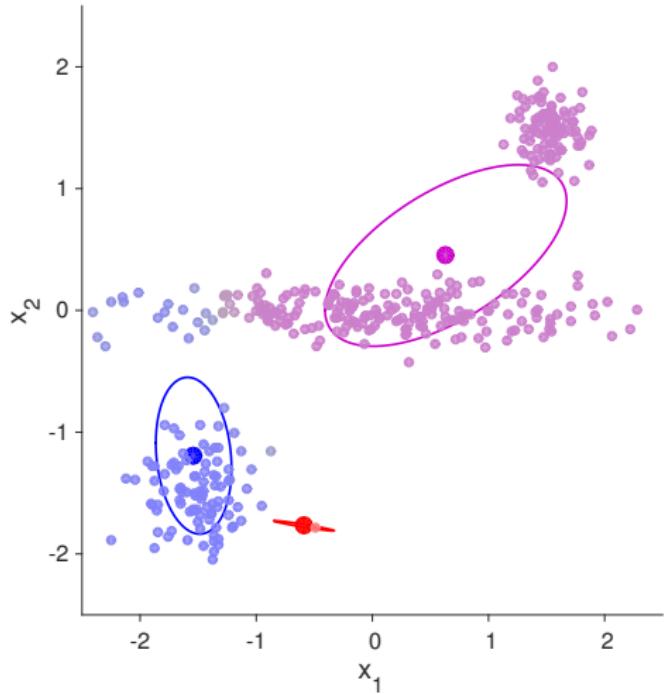
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

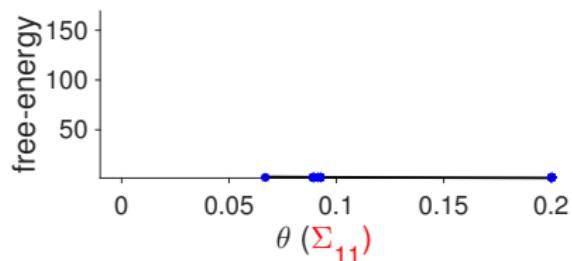
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

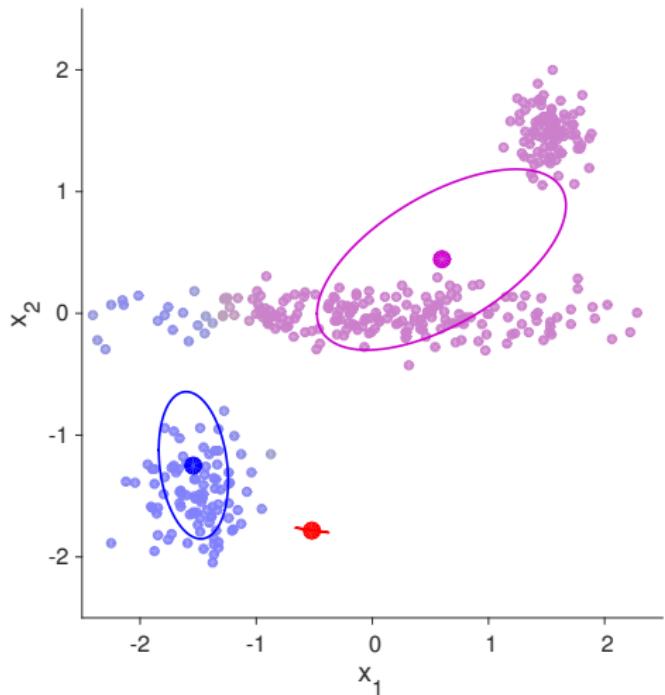
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

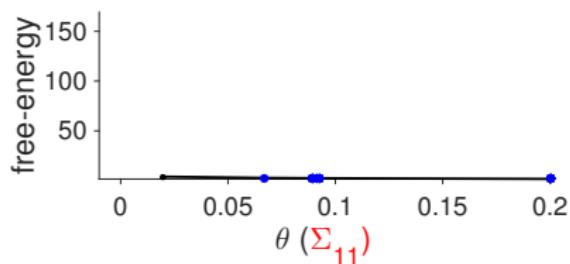
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

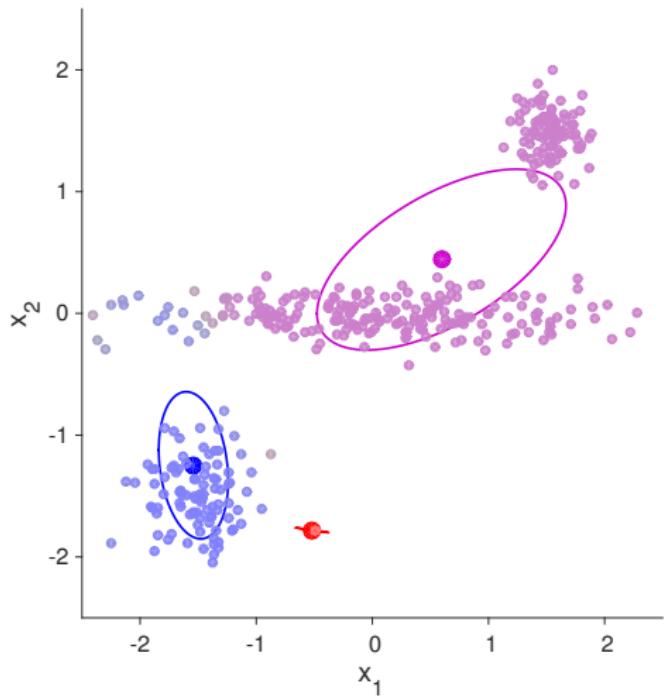
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

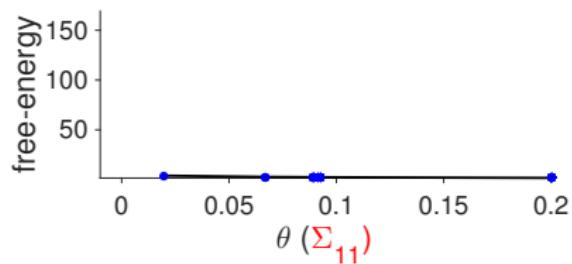
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

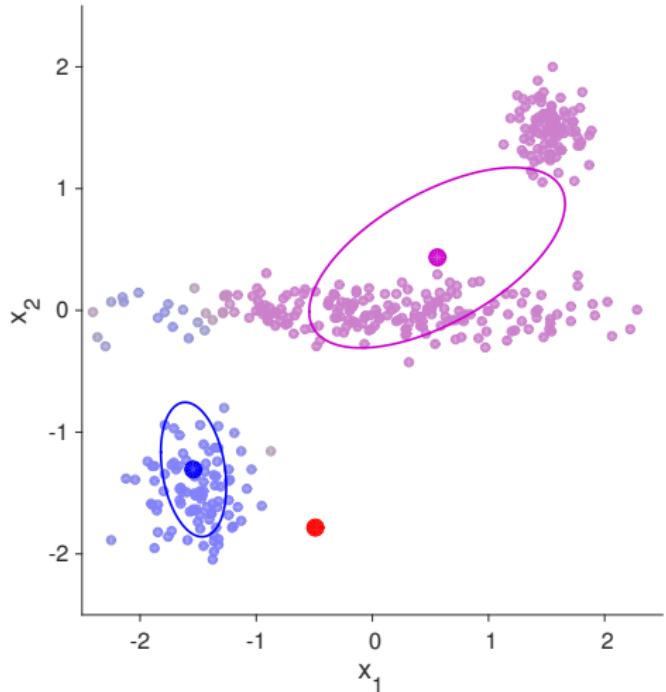
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

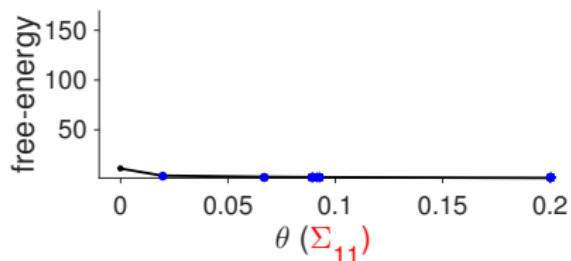
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

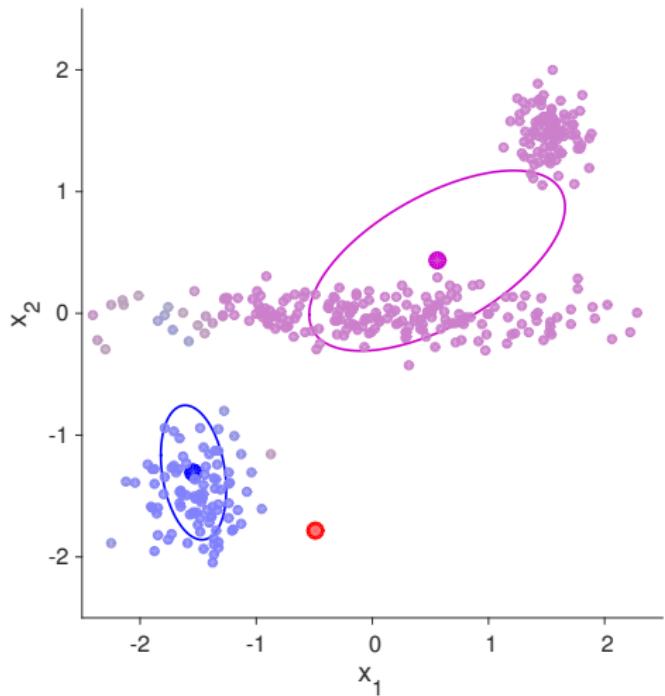
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

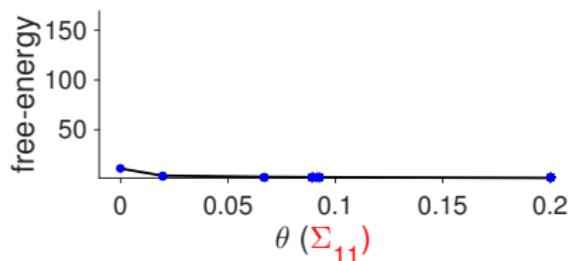
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

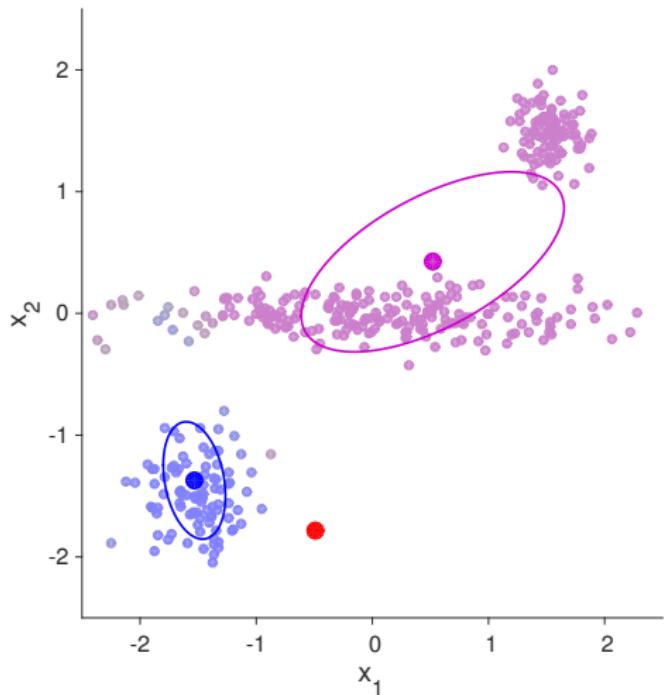
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

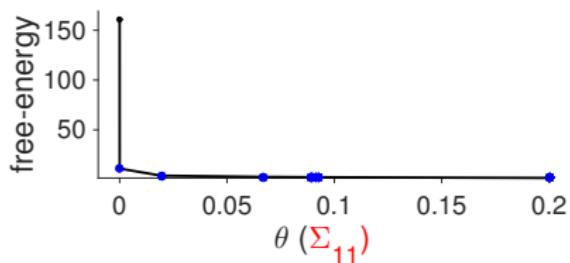
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

M-Step:

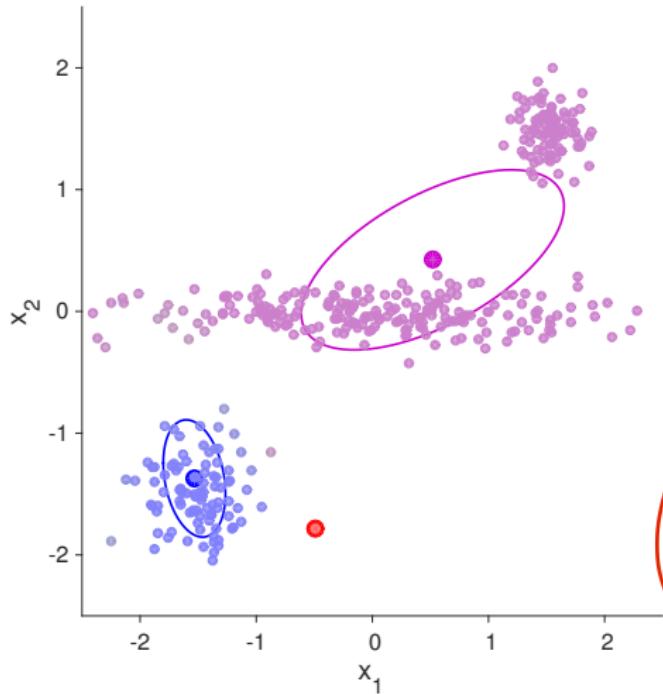
$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?

KABOOM!



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$

repeat

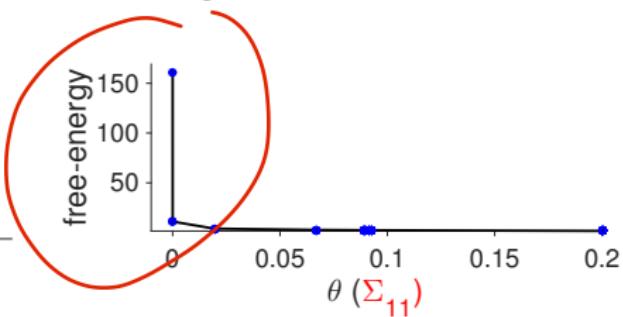
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

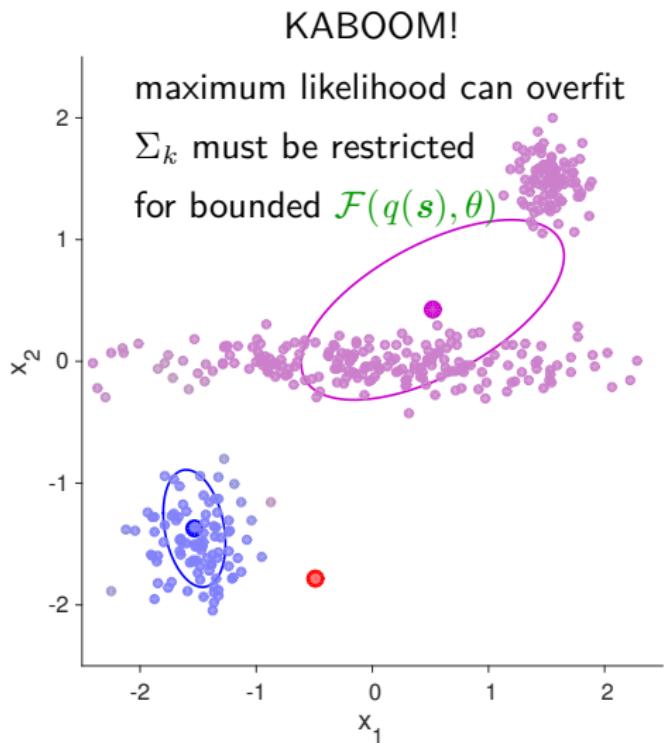
M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



What will happen with this initialisation?



initialise: $\theta = \{\pi_k, \mathbf{m}_k, \Sigma_k\}_{k=1}^K$
repeat

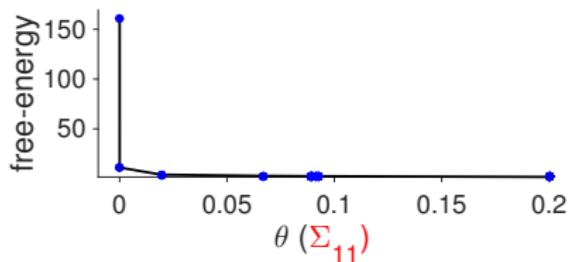
E-Step:

$$r_{nk} = p(s_n = k | \mathbf{x}_n, \theta) \text{ for } n = 1 \dots N$$

M-Step:

$$\arg \max_{\theta} \sum_{n,k} r_{nk} \log p(s_n = k, \mathbf{x} | \theta)$$

until convergence



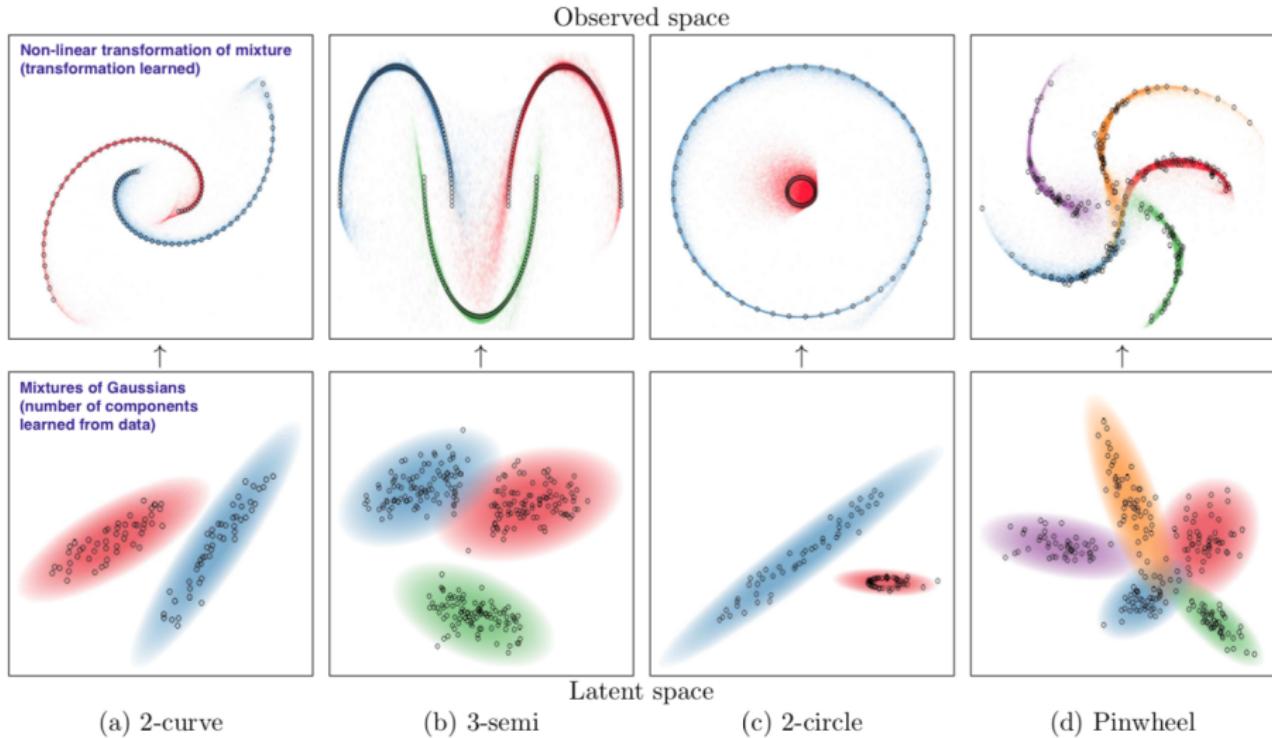
Summary

- ▶ MoG + EM algorithm = soft k-means clustering with non-axis aligned, non-equally weighted clusters
- ▶ EM can be used to fit **latent variable models** e.g. PCA, Factor analysis, MoGs, HMMs, ...
 - ▶ requires tractable posterior $p(s|x, \theta)$ entropy and average log-joint $\mathbb{E}_{q(s|\theta)} [\log p(s, x|\theta)]$

Limitations

- ▶ MoG clusters still have simple shapes (ellipses)
 - ▶ a single real cluster might be described by many components
 - ▶ more complex cluster models have been developed
- ▶ maximum-likelihood can overfit
 - ▶ Bayesian approaches avoid overfitting $p(\theta, s|x)$
- ▶ co-ordinate ascent is often slow to converge (lots of iterations required)
 - ▶ joint optimisation of $\mathcal{F}(q(s), \theta)$ faster (variational inference methods)
 - ▶ direct optimisation of log-likelihood $\log p(x|\theta)$

Research frontiers: variable numbers of non-Gaussian clusters



Warped Mixtures for Nonparametric Cluster Shapes, Iwata et al., UAI, 2013

Jupyter Notebooks

<https://github.com/cambridge-mlg/mphil-intro-module>

Appendix: proof of KL divergence properties

Minimise Kullback Leibler divergence (relative entropy) $\mathcal{KL}(q(x)||p(x))$: add Lagrange multiplier (enforce $q(x)$ normalises), take variational derivatives:

$$\frac{\delta}{\delta q(x)} \left[\int q(x) \log \frac{q(x)}{p(x)} dx + \lambda \left(1 - \int q(x) dx \right) \right] = \log \frac{q(x)}{p(x)} + 1 - \lambda.$$

Find stationary point by setting the derivative to zero:

$$q(x) = \exp(\lambda-1)p(x), \text{ normalization condition } \lambda = 1, \text{ so } q(x) = p(x),$$

which corresponds to a minimum, since the second derivative is positive:

$$\frac{\delta^2}{\delta q(x) \delta q(x)} \mathcal{KL}(q(x)||p(x)) = \frac{1}{q(x)} > 0.$$

The minimum value attained at $q(x) = p(x)$ is $\mathcal{KL}(p(x)||p(x)) = 0$, showing that $\mathcal{KL}(q(x)||p(x))$

- ▶ is non-negative
- ▶ attains its minimum 0 when $p(x)$ and $q(x)$ are equal