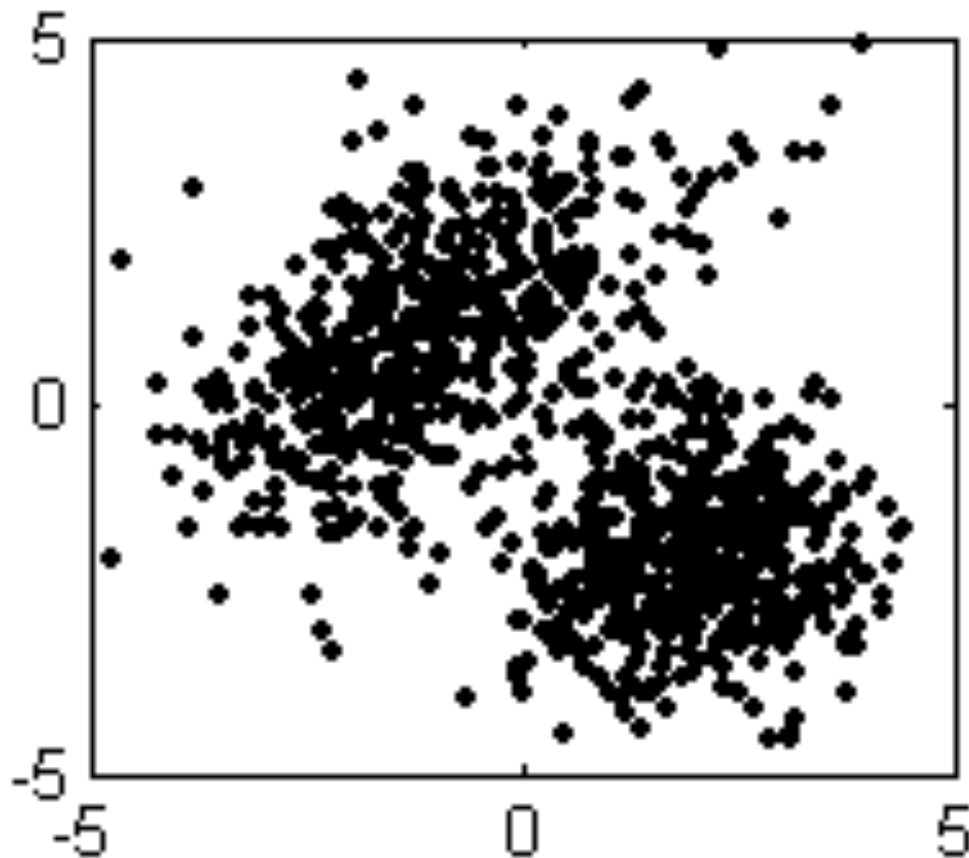


Machine Learning

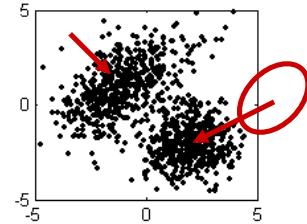
4771

Instructor: Itsik Pe'er

Reminder: Mixture of Gaussians



Kmeans: Notation Reminder



$\vec{z} \in \mathbf{B}^K$, $\sum_{i=1}^K \vec{z}(i) = 1$ or

$\vec{z} \in \{\vec{\delta}_1, \dots, \vec{\delta}_K\}$ where $\vec{\delta}_i(i) = 1, \vec{\delta}_i(j) = 0$ for $i \neq j$

mixing proportions (prior) $= \pi = p(\vec{z} = \vec{\delta}_i | \pi)$

mixture components $= p(\vec{x} | \vec{z} = \vec{\delta}_i, \theta)$

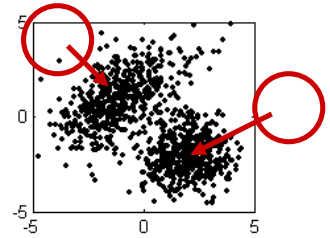
posteriors (responsibilities) $= \tau_{n,i} = p(\vec{z} = \vec{\delta}_i | \vec{x}_n, \theta) = \frac{p(\vec{x}_n | \vec{z} = \vec{\delta}_i, \theta) p(\vec{z} = \vec{\delta}_i | \theta)}{p(\vec{x}_n | \theta)}$

K-Means Clustering

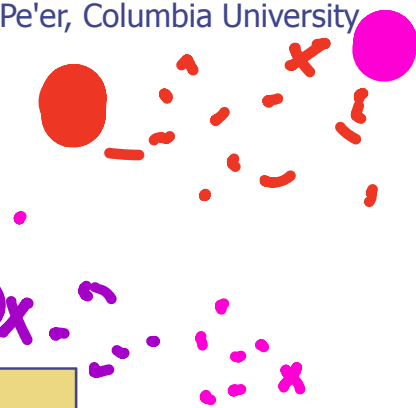
- An old “heuristic” clustering algorithm
- Gobble up data with a divide & conquer scheme
- Assume each point x has a discrete multinomial vector z
- Chicken and Egg problem:

If know classes, we can get model (max likelihood!)

If know the model, we can predict the classes (classifier!)



The K-Means Algorithm

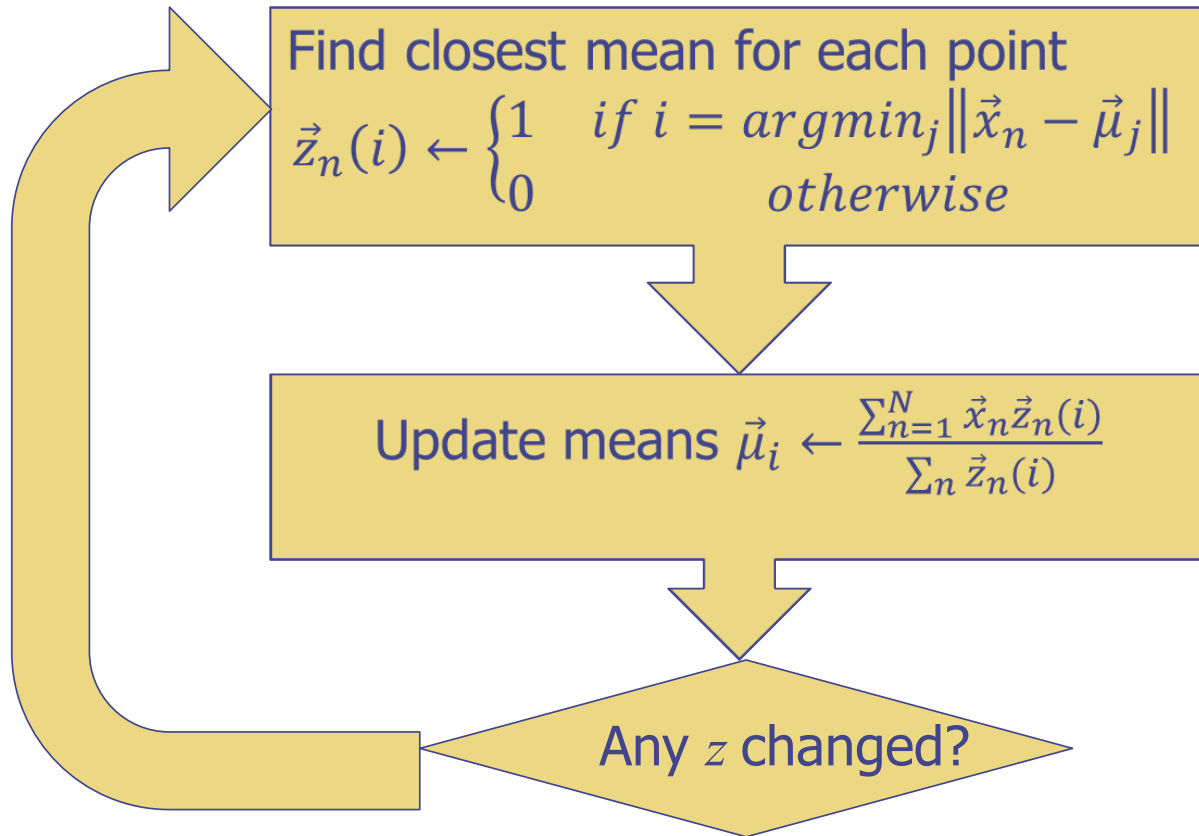


Find closest mean for each point

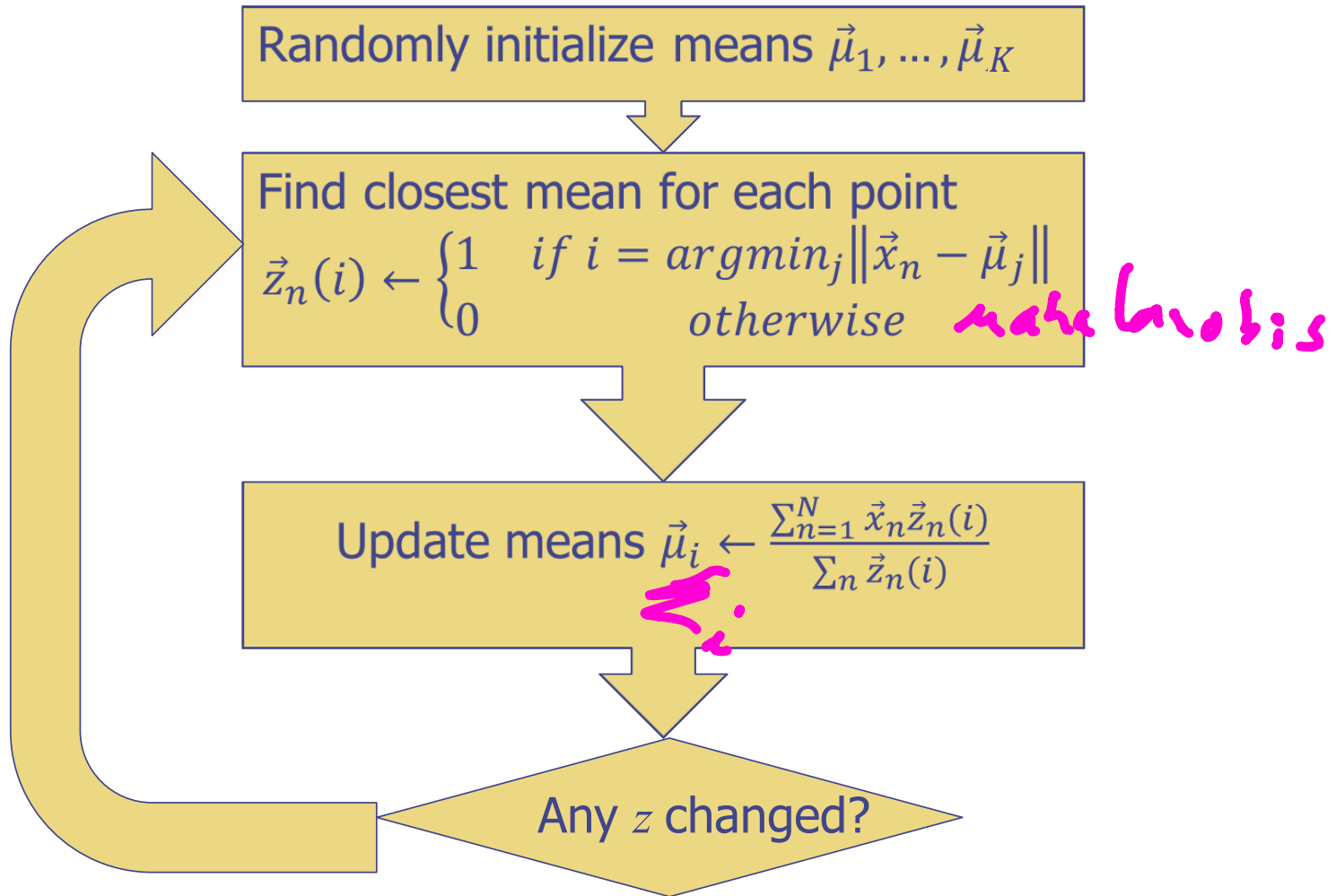
$$\vec{z}_n(i) \leftarrow \begin{cases} 1 & \text{if } i = \operatorname{argmin}_j \|\vec{x}_n - \vec{\mu}_j\| \\ 0 & \text{otherwise} \end{cases}$$

Update means $\vec{\mu}_i \leftarrow \frac{\sum_{n=1}^N \vec{x}_n \vec{z}_n(i)}{\sum_n \vec{z}_n(i)}$

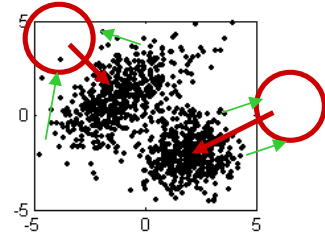
The K-Means Algorithm



The K-Means Algorithm



K-Means Clustering

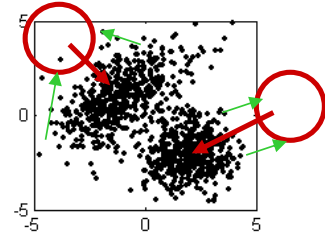


- Geometric, each point goes to closest Gaussian
- Recompute the means by their assigned points
- Minimizing $\min_{\mu} \min_{\vec{z}} J(\vec{\mu}_1, \dots, \vec{\mu}_K, \vec{z}_1, \dots, \vec{z}_N)$ cost function:

$$J(\vec{\mu}_1, \dots, \vec{\mu}_K, \vec{z}_1, \dots, \vec{z}_N) = \sum_{n=1}^N \sum_{i=1}^K \vec{z}_n(i) \|\vec{x}_n - \vec{\mu}_i\|^2$$

$$\vec{z}_n(i) = \begin{cases} 1 & \text{if } i = \operatorname{argmin}_j \|\vec{x}_n - \vec{\mu}_j\| \\ 0 & \text{otherwise} \end{cases} \quad \text{and} \quad \vec{\mu}_i = \frac{\sum_{n=1}^N \vec{x}_n \vec{z}_n(i)}{\sum_n \vec{z}_n(i)}$$

K-Means Clustering

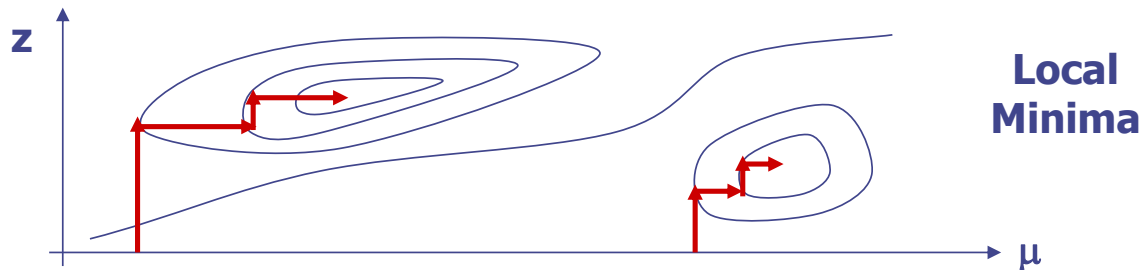


- Geometric, each point goes to closest Gaussian
- Recompute the means by their assigned points
- Minimizing $\min_{\mu} \min_z J(\vec{\mu}_1, \dots, \vec{\mu}_K, \vec{z}_1, \dots, \vec{z}_N)$ cost function:

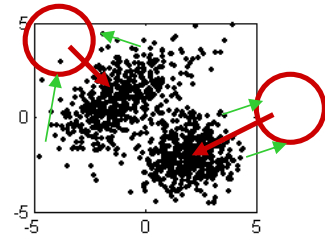
$$J(\vec{\mu}_1, \dots, \vec{\mu}_K, \vec{z}_1, \dots, \vec{z}_N) = \sum_{n=1}^N \sum_{i=1}^K \vec{z}_n(i) \|\vec{x}_n - \vec{\mu}_i\|^2$$

$$\vec{z}_n(i) = \begin{cases} 1 & \text{if } i = \operatorname{argmin}_j \|\vec{x}_n - \vec{\mu}_j\| \\ 0 & \text{otherwise} \end{cases} \quad \text{and} \quad \vec{\mu}_i = \frac{\sum_{n=1}^N \vec{x}_n \vec{z}_n(i)}{\sum_n \vec{z}_n(i)}$$

- Guaranteed to improve per iteration and converge
- Like **Coordinate Descent** (lock one var, maximize the other)
- A.k.a. **Axis-Parallel Optimization** or **Alternating Minimization**



K-Means Clustering

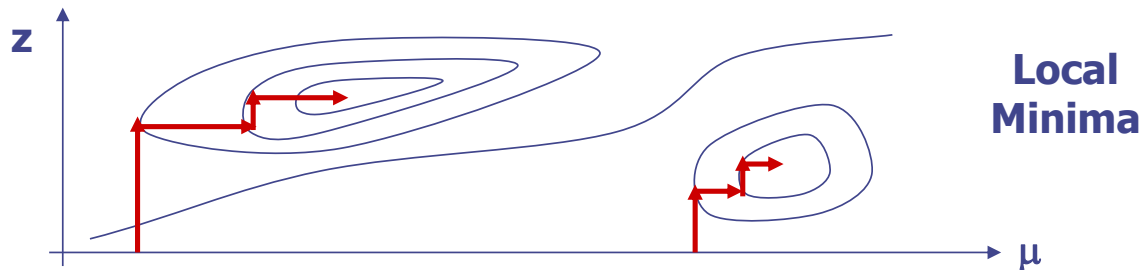


- Geometric, each point goes to closest Gaussian
- Recompute the means by their assigned points
- Minimizing $\min_{\mu, \Sigma} \min_z J(\vec{\mu}_1, \dots, \vec{\mu}_K, \Sigma_1, \dots, \Sigma_K, \vec{z}_1, \dots, \vec{z}_N)$ cost:

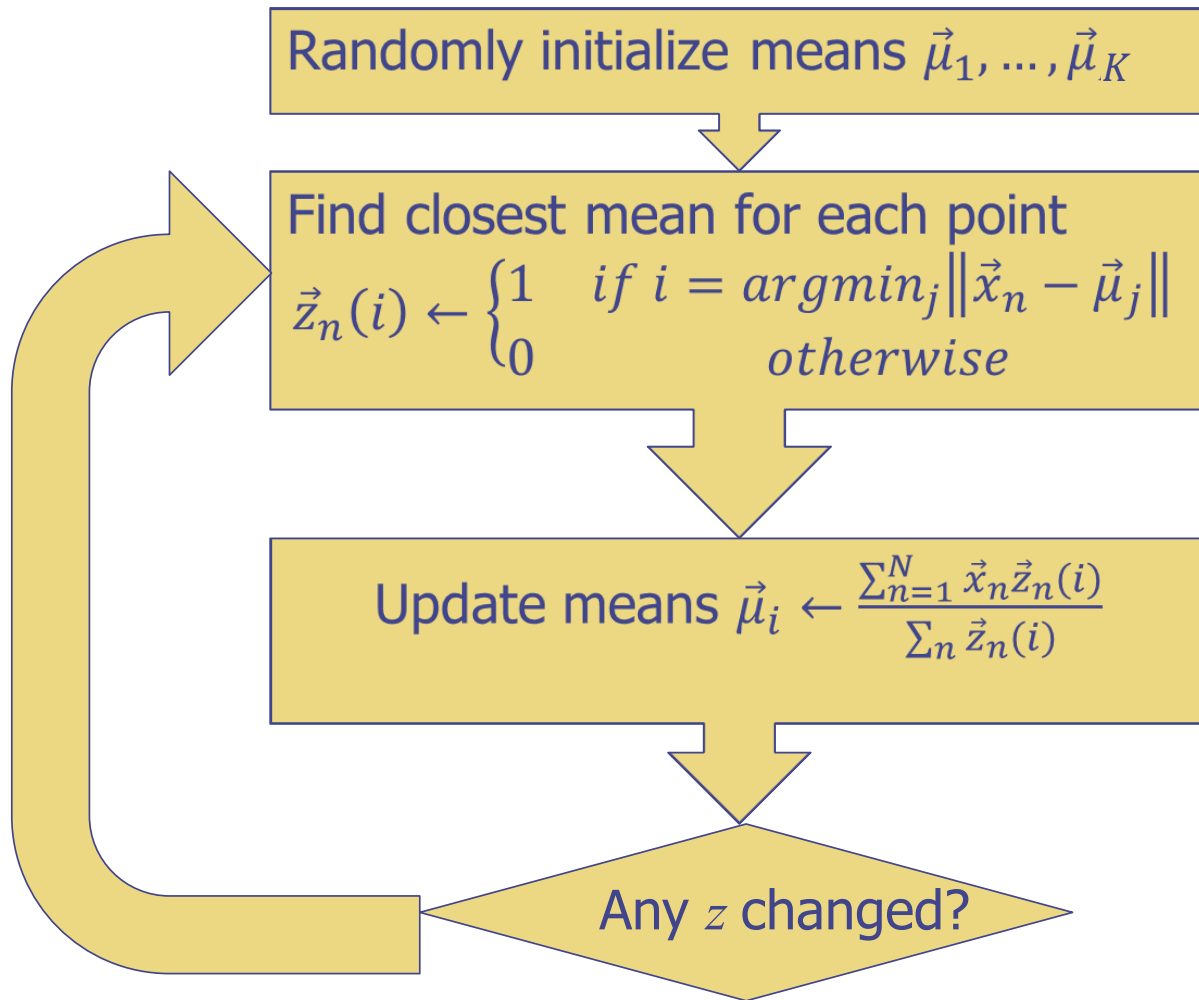
$$J(\vec{\mu}_1, \dots, \vec{\mu}_K, \vec{z}_1, \dots, \vec{z}_N) = \sum_{n=1}^N \sum_{i=1}^K \vec{z}_n(i) (\vec{x}_n - \vec{\mu}_i)^T \Sigma^{-1} (\vec{x}_n - \vec{\mu}_i)$$

$$\vec{z}_n(i) = \begin{cases} 1 & \text{if } i = \operatorname{argmin}_j (\vec{x}_n - \vec{\mu}_j)^T \Sigma^{-1} (\vec{x}_n - \vec{\mu}_j) \\ 0 & \text{otherwise} \end{cases}$$

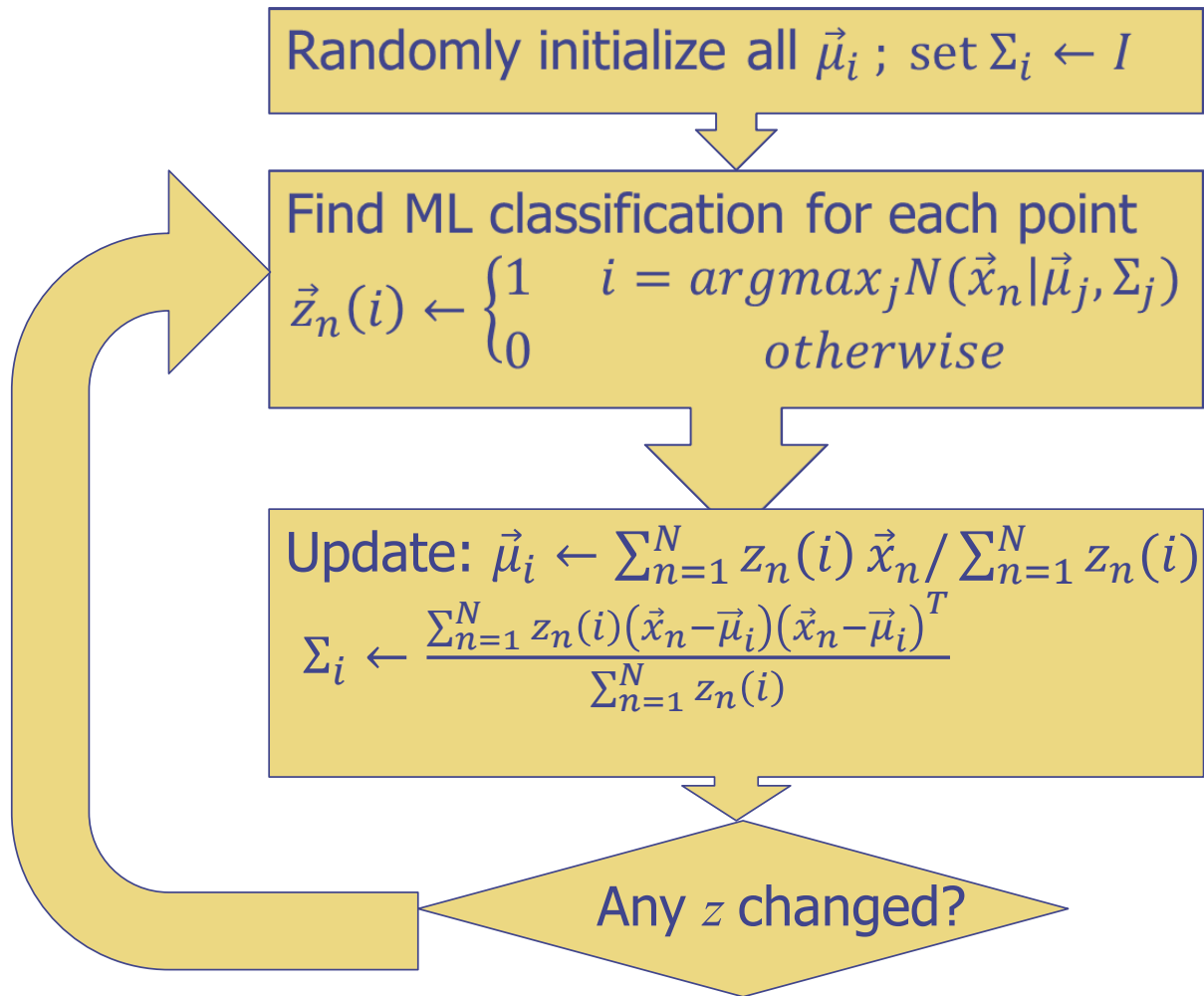
- Guaranteed to improve per iteration and converge
- Like **Coordinate Descent** (lock one var, maximize the other)
- A.k.a. **Axis-Parallel Optimization** or **Alternating Minimization**



The K-Means Algorithm

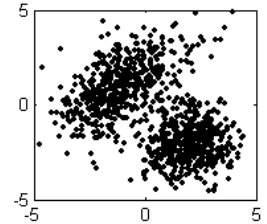


The Gaussian K-Means Algorithm



Kmeans→Expectation Maximization

At each stage of K-means, an element chooses the most likely (responsible) cluster



Hedging allows smoother optimization surface

Take expectation over responsible cluster

mixing proportions (prior) $= \pi = p(\vec{z} = \vec{\delta}_i | \pi)$

mixture components $= p(\vec{x} | \vec{z} = \vec{\delta}_i, \theta)$

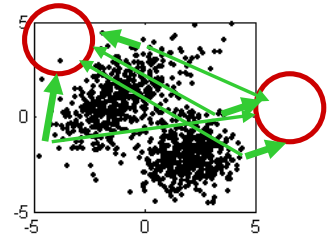
posteriors (responsibilities) $= \tau_{n,i} = p(\vec{z} = \vec{\delta}_i | \vec{x}_n, \theta) = \frac{p(\vec{x}_n | \vec{z} = \vec{\delta}_i, \theta) p(\vec{z} = \vec{\delta}_i | \theta)}{p(\vec{x}_n | \theta)}$

Expectation-Maximization (EM)

- EM is a soft/fuzzy version of K-Means (which does winner-takes-all, closest Gaussian Mean completely wins datapoint)

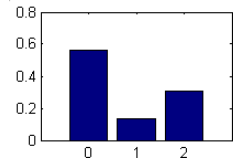
$$\vec{z}_n(i) = \begin{cases} 1 & i = \operatorname{argmax}_j N(\vec{x}_n | \vec{\mu}_j, \Sigma_j) \\ 0 & \text{otherwise} \end{cases}$$

- Instead, consider soft percentage assignment of datapoint (responsibility)



$$\tau_{n,i} = p(\vec{z}_n = \delta(j) | \vec{x}_n, \theta) \propto \pi_j \frac{\exp\left(-\frac{1}{2}(\vec{x}_n - \vec{\mu}_j)^T \Sigma^{-1}(\vec{x}_n - \vec{\mu}_j)\right)}{(2\pi)^{\frac{D}{2}} \sqrt{|\Sigma|}}$$

$$\tau_{n,1}, \dots, \tau_{n,K} =$$



- Update for the means are then

$$\vec{\mu}_i = \frac{\sum_n \tau_{n,i} \cdot \vec{x}_n}{\sum_n \tau_{n,i}}$$

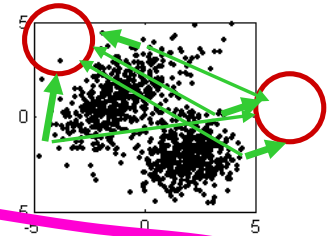
Expectation-Maximization (EM)

- EM is a soft/fuzzy version of K-Means (which does winner-takes-all, closest Gaussian Mean completely wins datapoint)

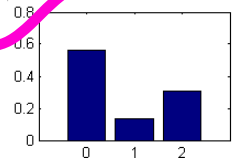
$$\vec{z}_n(i) = \begin{cases} 1 & i = \operatorname{argmax}_j N(\vec{x}_n | \vec{\mu}_j, \Sigma_j) \\ 0 & \text{otherwise} \end{cases}$$

- Instead, consider soft percentage assignment of datapoint (responsibility)

$$\tau_{n,i} = p(\vec{z}_n = \delta(j) | \vec{x}_n, \theta) \propto \pi_j \frac{\exp\left(-\frac{1}{2}(\vec{x}_n - \vec{\mu}_j)^T \Sigma^{-1}(\vec{x}_n - \vec{\mu}_j)\right)}{(2\pi)^{\frac{D}{2}} \sqrt{|\Sigma|}}$$



$$\tau_{n,1}, \dots, \tau_{n,K} =$$



- Update for the means are then

'weighted' by responsibilities:
$$\vec{\mu}_i = \frac{\sum_{n=1}^N \tau_{n,i} \vec{x}_n}{\sum_n \tau_{n,i}}$$

- Same for other parameters

Expectation-Maximization

- EM uses expected value of $\vec{z}_n(i)$ rather than max

$$\tau_{n,i} = E\{\vec{z}_n(i) | \vec{x}_n\} = p(\vec{z}_n = \vec{\delta}(j) | \vec{x}_n, \theta)$$

- EM updates covariances, mixing proportions AND means...
- The algorithm for Gaussian mixtures:

EXPECTATION: $\tau_{n,i}^{(t)} \leftarrow$

MAXIMIZATION: $\vec{\mu}_i^{(t+1)} \leftarrow$
 $\Sigma_i^{(t+1)} \leftarrow$

$\pi_i^{(t+1)} \leftarrow$

Expectation-Maximization

- EM uses expected value of $\vec{z}_n(i)$ rather than max

$$\tau_{n,i} = E\{\vec{z}_n(i) | \vec{x}_n\} = p(\vec{z}_n = \vec{\delta}(j) | \vec{x}_n, \theta)$$

- EM updates covariances, mixing proportions AND means...
- The algorithm for Gaussian mixtures:

→ EXPECTATION: $\tau_{n,i}^{(t)} \leftarrow \frac{\pi_i^{(t)} N(\vec{x}_n | \vec{\mu}_i^{(t)}, \Sigma_i^{(t)})}{\sum_j \pi_j^{(t)} N(\vec{x}_n | \vec{\mu}_j^{(t)}, \Sigma_j^{(t)})}$

→ MAXIMIZATION: $\vec{\mu}_i^{(t+1)} \leftarrow \frac{\sum_{n=1}^N \tau_{n,i}^{(t)} \vec{x}_n}{\sum_n \tau_{n,i}^{(t)}}$

$$\Sigma_i^{(t+1)} \leftarrow$$

Expectation-Maximization

- EM uses expected value of $\vec{z}_n(i)$ rather than max

$$\tau_{n,i} = E\{\vec{z}_n(i) | \vec{x}_n\} = p(\vec{z}_n = \vec{\delta}(j) | \vec{x}_n, \theta)$$

- EM updates covariances, mixing proportions AND means...
- The algorithm for Gaussian mixtures:

→ EXPECTATION: $\tau_{n,i}^{(t)} \leftarrow \frac{\pi_i^{(t)} N(\vec{x}_n | \vec{\mu}_i^{(t)}, \Sigma_i^{(t)})}{\sum_j \pi_j^{(t)} N(\vec{x}_n | \vec{\mu}_j^{(t)}, \Sigma_j^{(t)})}$

→ MAXIMIZATION: $\vec{\mu}_i^{(t+1)} \leftarrow \frac{\sum_{n=1}^N \tau_{n,i}^{(t)} \vec{x}_n}{\sum_n \tau_{n,i}^{(t)}} \quad \pi_i^{(t+1)} \leftarrow \frac{\sum_n \tau_{n,i}^{(t)}}{N}$

$$\Sigma_i^{(t+1)} \leftarrow \frac{\sum_{n=1}^N \tau_{n,i}^{(t)} (\vec{x}_n - \vec{\mu}_i^{(t+1)})^T (\vec{x}_n - \vec{\mu}_i^{(t+1)})}{\sum_n \tau_{n,i}^{(t)}}$$

The EM Clustering Algorithm

Initialize: random $\vec{\mu}_i^{(1)}, \Sigma_i^{(1)} = I, \pi_i^{(1)} = \frac{1}{K}$

Find expected assignment to each \vec{x}_n

$$\tau_{n,i}^{(t)} \leftarrow \frac{\pi_i^{(t)} N(\vec{x}_n | \vec{\mu}_i^{(t)}, \Sigma_i^{(t)})}{\sum_j \pi_j^{(t)} N(\vec{x}_n | \vec{\mu}_j^{(t)}, \Sigma_j^{(t)})}$$

Expectation

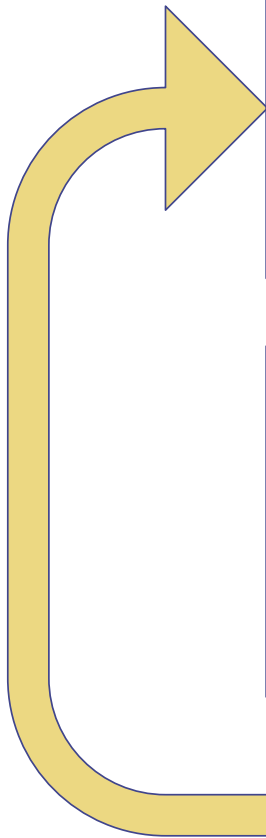
Update $\vec{\mu}_i^{(t+1)} \leftarrow \sum_{n=1}^N \tau_{n,i}^{(t)} \vec{x}_n / \sum_{n=1}^N \tau_{n,i}^{(t)}$

$$\Sigma_i^{(t+1)} \leftarrow \frac{\sum_{n=1}^N \tau_{n,i}^{(t)} (\vec{x}_n - \vec{\mu}_i^{(t+1)}) (\vec{x}_n - \vec{\mu}_i^{(t+1)})^T}{\sum_{n=1}^N \tau_{n,i}^{(t)}}$$

$$\pi_i^{(t+1)} \leftarrow \sum_{n=1}^N \tau_{n,i}^{(t)} / N$$

Maximization

Little changed?



The EM Clustering Algorithm

Initialize: random $\vec{\mu}_i^{(1)}, \Sigma_i^{(1)} = I, \pi_i^{(1)} = \frac{1}{K}$

Find expected assignment to each \vec{x}_n

$$\tau_{n,i}^{(t)} \leftarrow E\{\vec{z}_n(i)\}$$

Expectation

Update

$$\vec{\mu}_i^{(t+1)} \leftarrow \operatorname{argmax}_{\mu^*} E_{\vec{z}}\{l(\theta, z \text{ s.t. } \vec{\mu}_i = \mu^*)\}$$

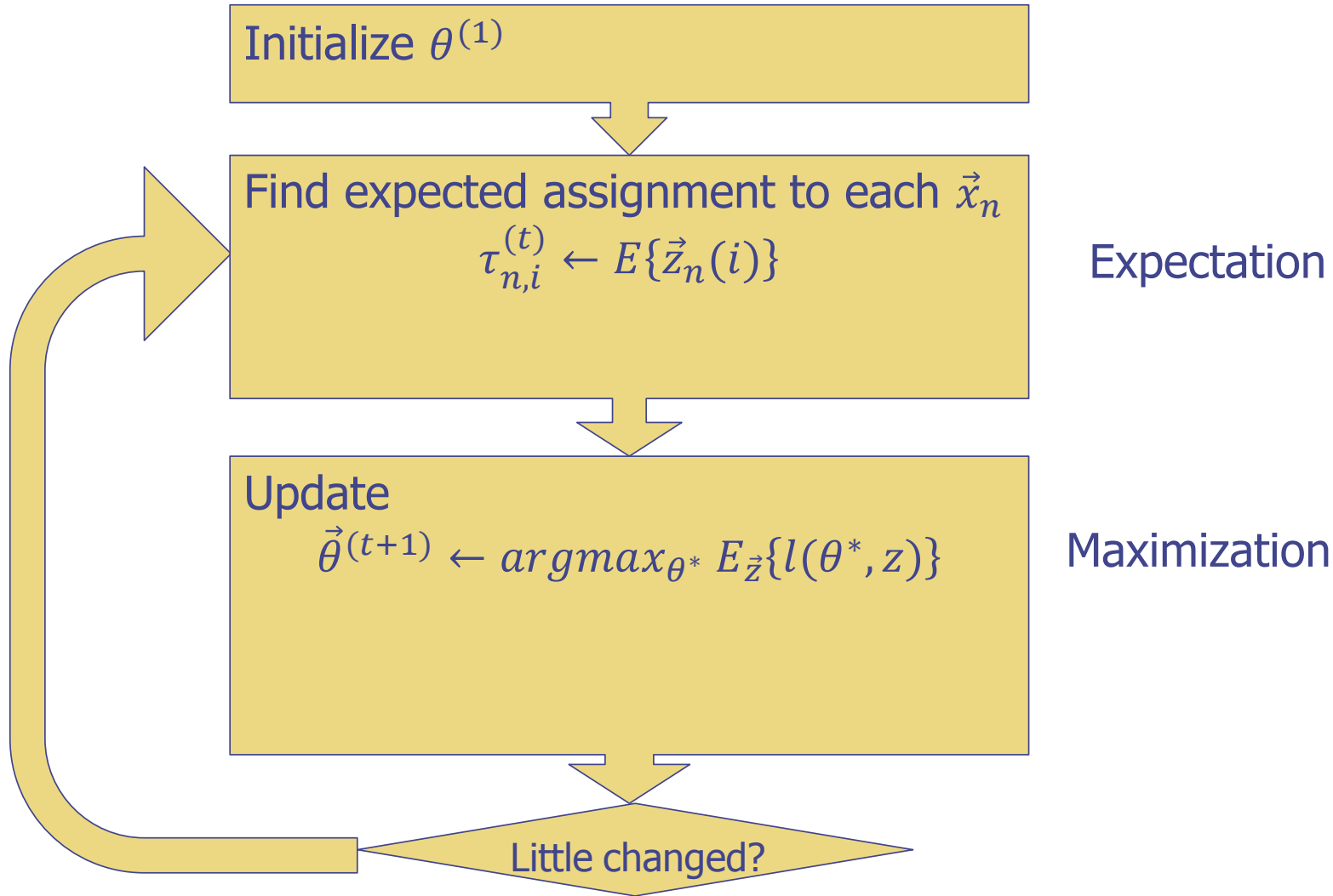
$$\Sigma_i^{(t+1)} \leftarrow \operatorname{argmax}_{\Sigma^*} E_{\vec{z}}\{l(\theta, z \text{ s.t. } \Sigma_i = \Sigma^*)\}$$

$$\pi_i^{(t+1)} \leftarrow \operatorname{argmax}_{\pi^*} E_{\vec{z}}\{l(\theta, z \text{ s.t. } \pi_i = \pi^*)\}$$

Maximization

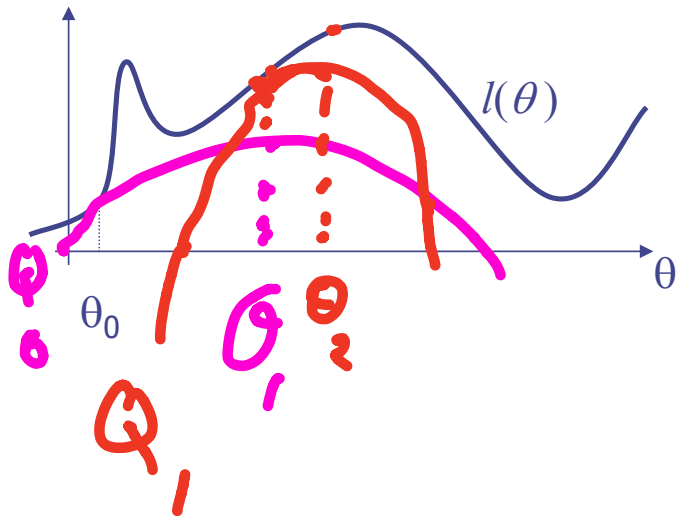
Little changed?

The General EM Algorithm



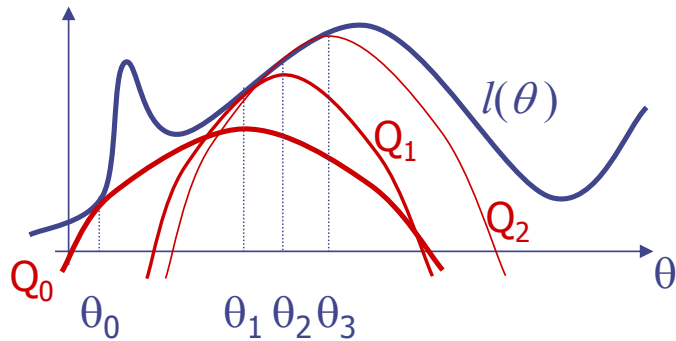
EM as Bound Maximization

- Let's now show that EM indeed maximizes likelihood
- **Bound Maximization:** optimize a lower bound on $l(\theta)$
- Since log-likelihood $l(\theta)$ not concave, can't max it directly
- Consider an auxiliary function $Q(\theta)$ which is concave
- $Q(\theta)$ kisses $l(\theta)$ at a point and is less than it elsewhere



EM as Bound Maximization

- Let's now show that EM indeed maximizes likelihood
- **Bound Maximization:** optimize a lower bound on $l(\theta)$
- Since log-likelihood $l(\theta)$ not concave, can't max it directly
- Consider an auxiliary function $Q(\theta)$ which is concave
- $Q(\theta)$ kisses $l(\theta)$ at a point and is less than it elsewhere



EM as Bound Maximization

- Let's now show that EM indeed maximizes likelihood
- **Bound Maximization:** optimize a lower bound on $l(\theta)$
- Since log-likelihood $l(\theta)$ not concave, can't max it directly
- Consider an auxiliary function $Q(\theta)$ which is concave
- $Q(\theta)$ kisses $l(\theta)$ at a point and is less than it elsewhere

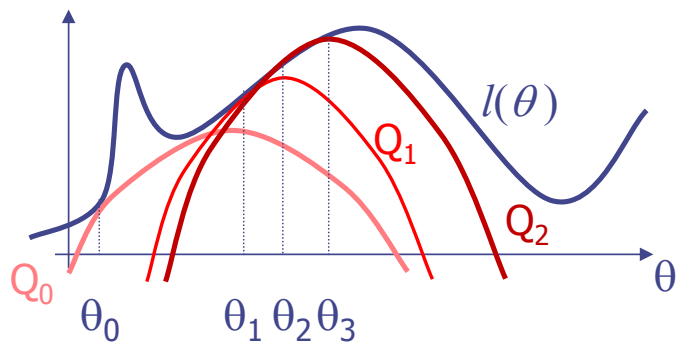
$$\forall \theta \forall t : l(\theta) \geq Q_t(\theta)$$

$$l(\theta_t) = Q_t(\theta_t)$$

$$Q_t(\theta_{t+1}) > Q_t(\theta_t)$$

$$\text{because } \theta_{t+1} = \arg \max_{\theta} Q_t(\theta)$$

$$l(\theta_{t+1}) \geq Q_t(\theta_{t+1}) > Q_t(\theta_t) = l(\theta_t)$$

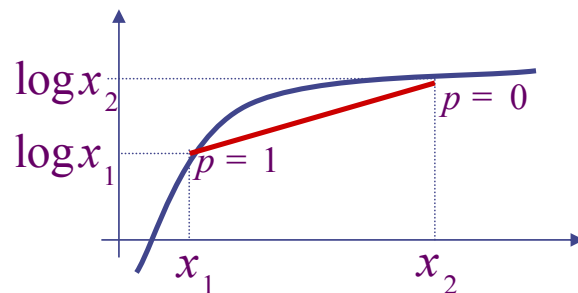


- Monotonically increases log-likelihood
- But how to find a bound and guarantee we max it?

Jensen's Inequality



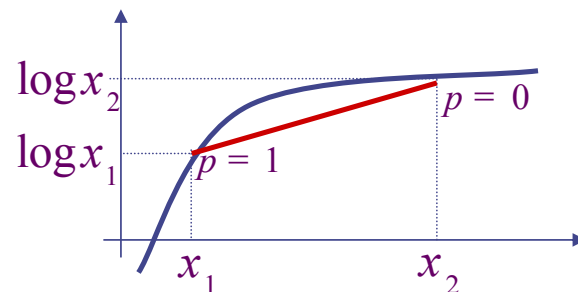
- Example: $f(x) = \log(x)$ = concave and $M=2$
$$\log(px_1 + (1-p)x_2) \geq p \log x_1 + (1-p) \log x_2$$
- Bound $\log(\text{sum})$ with $\text{sum}(\log)$



Jensen's Inequality



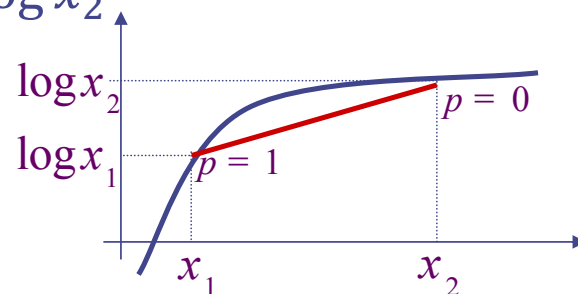
- Expectation in discrete case is sum weight by probability
- For convex f : $f\left(\sum_{i=1}^M p_i x_i\right) \leq \sum_{i=1}^M p_i f(x_i)$ when $\sum_{i=1}^M p_i = 1, p_i \geq 0$
- For concave f : $f\left(\sum_{i=1}^M p_i x_i\right) \geq \sum_{i=1}^M p_i f(x_i)$ when $\sum_{i=1}^M p_i = 1, p_i \geq 0$
- Example: $f(x) = \log(x)$ = concave and $M=2$
 $\log(px_1 + (1-p)x_2) \geq p \log x_1 + (1-p) \log x_2$
- Bound $\log(\text{sum})$ with $\text{sum}(\log)$





Jensen's Inequality

- An important general bound from Jensen (1906)
- For convex f : $f(E\{x\}) \leq E\{f(x)\}$
- For concave f : $f(E\{x\}) \geq E\{f(x)\}$
- Expectation in discrete case is sum weight by probability
- For convex f : $f(\sum_{i=1}^M p_i x_i) \leq \sum_{i=1}^M p_i f(x_i)$ when $\sum_{i=1}^M p_i = 1, p_i \geq 0$
- For concave f : $f(\sum_{i=1}^M p_i x_i) \geq \sum_{i=1}^M p_i f(x_i)$ when $\sum_{i=1}^M p_i = 1, p_i \geq 0$
- Example: $f(x) = \log(x)$ = concave and $M=2$
 $\log(px_1 + (1-p)x_2) \geq p \log x_1 + (1-p) \log x_2$
- Bound $\log(\text{sum})$ with $\text{sum}(\log)$
- How to apply this to mixture models?



Expectation-Maximization

$$\begin{aligned}
 l(\theta) &= \sum_n \log P(z_n | \theta) = \\
 &= \sum_n \log \sum_z P(x_n, z | \theta) \frac{P(z | x_n, \theta)}{P(z | x_n, \theta)} \\
 &= \sum_n \log \sum_z P(x_n, z | \theta) \cdot \frac{P(x_n, z | \theta)}{P(z | x_n, \theta)} \\
 &= \sum_n \log E \left[P(x_n, z | \theta) \right]
 \end{aligned}$$

Expectation-Maximization

$$\begin{aligned}
 l(\theta) &= \sum_{n=1}^N \log p(x_n | \theta) && \leftarrow \text{Original Log-Likelihood} \\
 &= \sum_{n=1}^N \log \sum_z p(x_n, z | \theta) && \leftarrow \text{Has Hidden Variables (messy)} \\
 &= \sum_{n=1}^N \log \sum_z p(x_n, z | \theta) \frac{p(z | x_n, \theta_t)}{p(z | x_n, \theta_t)} && \leftarrow \begin{array}{l} \text{Multiply by 1} \\ \text{Ratio of hidden} \\ \text{posterior density} \end{array} \\
 &= \sum_{n=1}^N \log \sum_z p(z | x_n, \theta_t) \frac{p(x_n, z | \theta)}{p(z | x_n, \theta_t)} && \leftarrow \text{Rearrange}
 \end{aligned}$$

Expectation-Maximization

$$\begin{aligned}
 l(\theta) &= \sum_{n=1}^N \log p(x_n | \theta) && \text{Original Log-Likelihood} \\
 &= \sum_{n=1}^N \log \sum_z p(x_n, z | \theta) && \text{Has Hidden Variables (messy)} \\
 &= \sum_{n=1}^N \log \sum_z p(x_n, z | \theta) \frac{p(z | x_n, \theta_t)}{p(z | x_n, \theta_t)} && \text{Multiply by 1 Ratio of hidden posterior density} \\
 &= \sum_{n=1}^N \log \sum_z p(z | x_n, \theta_t) \frac{p(x_n, z | \theta)}{p(z | x_n, \theta_t)} && \text{Rearrange} \\
 &\geq \sum_{n=1}^N \sum_z p(z | x_n, \theta_t) \log \frac{p(x_n, z | \theta)}{p(z | x_n, \theta_t)} && \text{Jensen log}(\sum_i p_i x_i) \\
 &= \sum_{n=1}^N \sum_z p(z | x_n, \theta_t) \log p(x_n, z | \theta) - \sum_{n=1}^N \sum_z p(z | x_n, \theta_t) \log p(z | x_n, \theta_t) \\
 &= Q(\theta | \theta_t) - \text{const} && \text{New auxiliary function called Q (not messy)}
 \end{aligned}$$

EM as Bound Maximization

- Now have the following bound and maximize it:

$$l(\theta) \geq$$

EM as Bound Maximization

- Now have the following bound and maximize it:

$$l(\theta) \geq Q(\theta|\theta_t) - \sum_{n=1}^N \sum_z p(z|x_n, \theta_t) \log p(z|x_n, \theta_t)$$

$$\theta_{t+1} = \arg \max Q(\theta|\theta_t) = \arg \max \sum_{n=1}^N \sum_z p(z|x_n, \theta_t) \log p(x_n, z|\theta_t)$$

$$= \arg \max \sum_{n=1}^N \sum_z \tau_{n,z} \log p(x_n, z|\theta_t)$$

- $Q(\theta|\theta_t)$ is called **Auxiliary Function**... take derivatives of it
- This is easy for many functions... just weighted max likelihood!
- For example, Gaussian mixture:

$$\frac{\partial Q(\theta)}{\partial \vec{\mu}_k} =$$

EM as Bound Maximization

- Now have the following bound and maximize it:

$$l(\theta) \geq Q(\theta|\theta_t) - \sum_{n=1}^N \sum_z p(z|x_n, \theta_t) \log p(z|x_n, \theta_t)$$

$$\theta_{t+1} = \arg \max Q(\theta|\theta_t) = \arg \max \sum_{n=1}^N \sum_z p(z|x_n, \theta_t) \log p(x_n, z|\theta_t)$$

$$= \arg \max \sum_{n=1}^N \sum_z \tau_{n,z} \log p(x_n, z|\theta_t)$$

- $Q(\theta|\theta_t)$ is called **Auxiliary Function**... take derivatives of it
- This is easy for many functions... just weighted max likelihood!
- For example, Gaussian mixture:

$$\frac{\partial Q(\theta)}{\partial \vec{\mu}_k} = \frac{\partial}{\partial \vec{\mu}_k} \sum_{n=1}^N \sum_{j=1}^K \tau_{n,k} \log \pi_k N(\vec{x}_n | \vec{\mu}_k, \Sigma_k)$$

$$0 = \sum_{n=1}^N \tau_{n,k} \frac{\partial}{\partial \vec{\mu}_k} \left(-\frac{1}{2} (\vec{x}_n - \vec{\mu}_k)^T \Sigma^{-1} (\vec{x}_n - \vec{\mu}_k) \right)$$

$$\vec{\mu}_k = \frac{\sum_{n=1}^N \tau_{n,k} \vec{x}_n}{\sum_{n=1}^N \tau_{n,k}}$$

... similarly get π_k and Σ_k

EM as Expected Log-Likelihood

- Incomplete Log-Likelihood
- Complete Log-Likelihood

EM as Expected Log-Likelihood

- Incomplete Log-Likelihood

$$l(\theta) = \log p(\text{observed}|\theta) = \sum_{n=1}^N \log \sum_z p(x_n, z_n|\theta)$$

- Complete Log-Likelihood

$$l(\theta) = \log p(\text{observed}, \text{hidden}|\theta) = \sum_{n=1}^N \log p(x_n, z_n|\theta)$$

- We don't know the hidden variables z
- EM computes expected values of hidden z under current θ_t
- EM chooses Q to be the Expected Complete Log-Likelihood

EM as Expected Log-Likelihood

- Incomplete Log-Likelihood

$$l(\theta) = \log p(\text{observed}|\theta) = \sum_{n=1}^N \log \sum_z p(x_n, z_n|\theta)$$

- Complete Log-Likelihood

$$l(\theta) = \log p(\text{observed}, \text{hidden}|\theta) = \sum_{n=1}^N \log p(x_n, z_n|\theta)$$

- We don't know the hidden variables z
- EM computes expected values of hidden z under current θ_t
- EM chooses Q to be the Expected Complete Log-Likelihood

$$E\{l^c(\theta)\} =$$

EM as Expected Log-Likelihood

• Incomplete Log-Likelihood

$$l(\theta) = \log p(\text{observed}|\theta) = \sum_{n=1}^N \log \sum_z p(x_n, z_n|\theta)$$

• Complete Log-Likelihood

$$l(\theta) = \log p(\text{observed}, \text{hidden}|\theta) = \sum_{n=1}^N \log p(x_n, z_n|\theta)$$

- We don't know the hidden variables z
- EM computes expected values of hidden z under current θ_t
- EM chooses Q to be the **Expected Complete Log-Likelihood**

$$\begin{aligned} E\{l^c(\theta)\} &= \sum_{\text{hidden}} p(\text{hidden}|\text{observed}, \theta_t) l^c(\theta) \\ &= \sum_{z_1} \sum_{z_2} \dots \sum_{z_N} p(z_1, \dots, z_N | x_1, \dots, x_N, \theta_t) l^c(\theta) \\ &= \sum_{z_1} \sum_{z_2} \dots \sum_{z_N} \prod_n p(z_n | x_n, \theta_t) l^c(\theta) \\ &= \sum_{z_1} \sum_{z_2} \dots \sum_{z_N} \prod_n p(z_n | x_n, \theta_t) \sum_n \log p(x_n, z_n | \theta) \\ &= \sum_n \sum_{z_n} p(z_n | x_n, \theta_t) \log p(x_n, z_n | \theta) \sum_{z_1} \dots \sum_{z_{i \neq n}} \dots \sum_{z_N} \prod_{i \neq n} p(z_i | x_i, \theta_t) \\ &= \sum_n \sum_n p(z_n | x_n, \theta_t) \log p(x_n, z_n | \theta) = Q(\theta | \theta_t) \end{aligned}$$

Summary

- Mixture Models and Hidden Variables
 - Clustering
 - K-Means
 - Expectation Maximization
-
- Happy Spring Break!