Announcements

Assignments

- HW9 (online)
 - Due Thu 4/16, 11:59 pm

TA Applications:

Please apply to TA with us! Link for MLD applications will be in piazza.

Participation points

 Starting now, we're capping the denominator (57 polls) in the participation points calculation

Introduction to Machine Learning

Clustering
GMM and EM

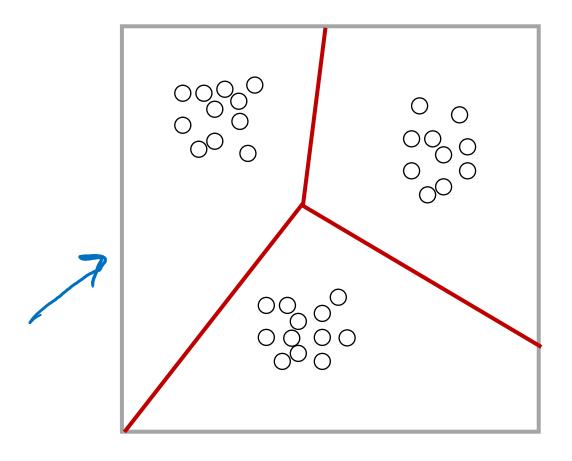
Instructor: Pat Virtue

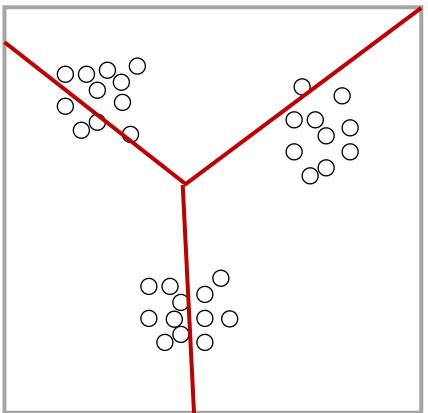
Slide credits: Aarti Singh, Eric Xing, Carlos Guestrin

K-means Optimization

Question: Which of these partitions is "better"?

min





K-means Optimization

Input:
$$\vec{x}^{(1)}$$
 ... $\vec{x}^{(N)}$ $\vec{x}^$

Slide credit: CMU MLD Matt Gormley

K-means Optimization

Alternating minimization

a)
$$\vec{z} = \arg\min_{z} ||x' - c_z||_2^2$$

b) $C = \arg\min_{z} ||x' - c_z||_2^2$

a) $z^{(i)} = \arg\min_{z} ||x'' - c_z||_2^2$

b) $\vec{c}_1 = \arg\min_{z \in Z^{(i)}} ||x^{(i)} - c_1||_2^2$

$$\vec{c}_k = \arg\min_{z \in Z^{(i)}} ||x^{(i)} - c_1||_2^2$$

Slide crefft: CMU MLD Matt Gormley

Piazza Poll 1

[True/False] The alternating minimization algorithm will find the global minimum of the k-means objective.

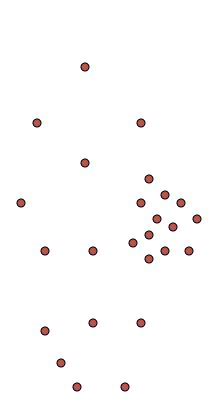
$$C, z = \underset{C, z}{\operatorname{argmin}} \sum_{i=1}^{N} ||x^{(i)} - c_{z^{(i)}}||_{2}^{2}$$

A.

B

C.

(One) bad case for K-means



- Clusters may overlap
- Some clusters may be "wider" than others
- Clusters may not be linearly separable

(One) bad case for K-means

- Clusters may overlap
- Some clusters may be "wider" than others
- Clusters may not be linearly separable

Partitioning Algorithms

- K-means
 - hard assignment: each object belongs to only one cluster

- Mixture modeling
 - soft assignment: probability that an object belongs to a cluster

Generative approach

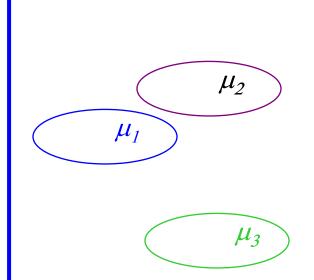
Gaussian Mixture Model

Mixture of K Gaussian distributions: (Multi-modal distribution)

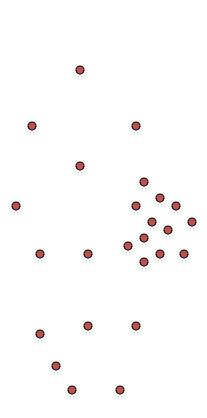
$$p(x|y=i) \sim N(\mu_i, \sigma^2 I)$$

$$p(x) = \sum_{i} p(x|y=i) P(y=i)$$

$$\downarrow \qquad \qquad \downarrow$$
Mixture
$$component \qquad proportion$$



(One) bad case for K-means



- Clusters may overlap
- Some clusters may be "wider" than others
- Clusters may not be linearly separable

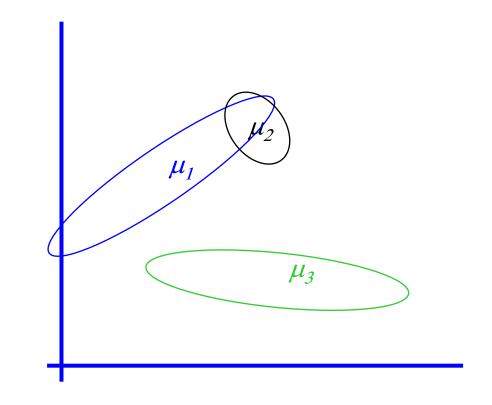
General GMM

GMM – Gaussian Mixture Model (Multi-modal distribution)

$$p(x|y=i) \sim N(\mu_i, \Sigma_i)$$

$$p(x) = \sum_i p(x|y=i) P(y=i)$$

$$\downarrow \qquad \qquad \downarrow$$
Mixture
$$component \qquad proportion$$



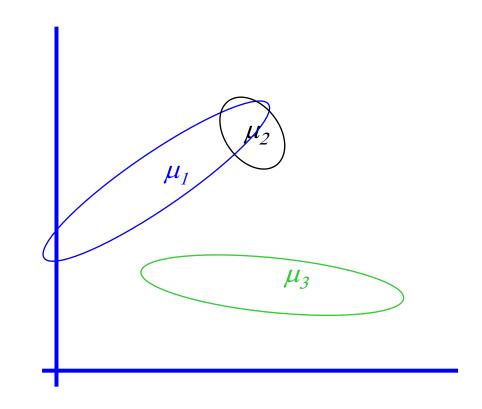
General GMM

GMM – Gaussian Mixture Model (Multi-modal distribution)

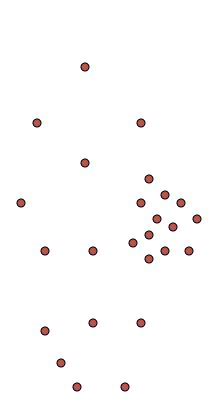
- There are k components
- Component i has an associated mean vector μ_i
- Each component generates data from a Gaussian with mean μ_i and covariance matrix Σ_i

Each data point is generated according to the following recipe:

- Pick a component at random: Choose component i with probability P(y=i)
- 2) Datapoint $x \sim N(\mu_i, \Sigma_i)$



(One) bad case for K-means



- Clusters may overlap
- Some clusters may be "wider" than others
- Clusters may not be linearly separable

General GMM

GMM - Gaussian Mixture Model (Multi-modal distribution)

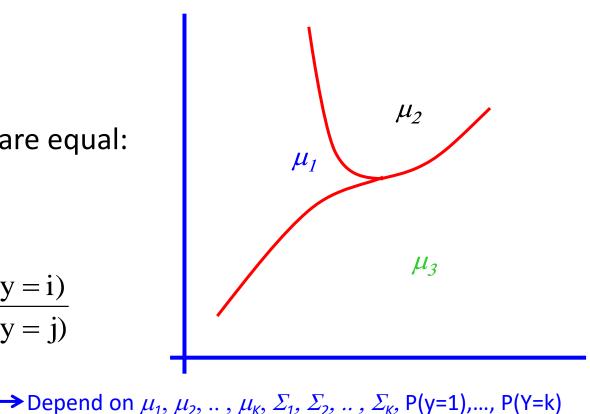
$$p(x|y=i) \sim N(\mu_i, \Sigma_i)$$

Decision boundary when probabilities are equal:

$$\log \frac{P(y=i \mid x)}{P(y=j \mid x)}$$

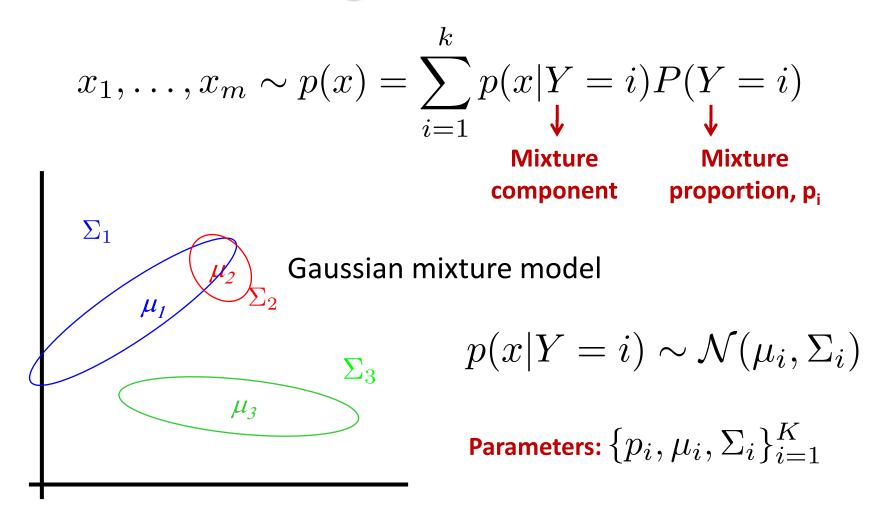
$$= \log \frac{p(x \mid y=i)P(y=i)}{p(x \mid y=j)P(y=j)}$$

$$= x^{T} W x + W^{T} x$$



"Quadratic Decision boundary" – second-order terms don't cancel out

Learning General GMM



 How to estimate parameters? Max Likelihood But don't know labels Y

Learning General GMM

Maximize marginal likelihood:

$$argmax \prod_{j} P(x_{j}) = argmax \prod_{j} \sum_{i=1}^{K} P(y_{j}=i,x_{j})$$
$$= argmax \prod_{j} \sum_{i=1}^{K} P(y_{j}=i)p(x_{j}|y_{j}=i)$$

 $P(y_i=i) = P(y=i)$ Mixture component i is chosen with prob P(y=i)

$$= \arg \max \prod_{j=1}^{m} \sum_{i=1}^{k} P(y=i) \frac{1}{\sqrt{\det(\sum_{i})}} \exp \left[-\frac{1}{2} (x_{j} - \mu_{i})^{T} \sum_{i} (x_{j} - \mu_{i}) \right]$$

How do we find the μ_i , Σ_i s and P(y=i)s which give max. marginal likelihood?

* Set $\frac{\partial}{\partial \mu_i}$ log Prob (....) = 0 and solve for μ_i 's. Non-linear not-analytically solvable

* Use gradient descent: Doable, but often slow

MLE Optimization

MLE Optimization

GMM vs. k-means

Maximize marginal likelihood:

$$\underset{\text{argmax }\prod_{j}\sum_{i=1}^{K}P(y_{j}=i,x_{j})}{\text{argmax }\prod_{j}\sum_{i=1}^{K}P(y_{j}=i,x_{j})}$$

$$=\underset{\text{argmax }\prod_{j}\sum_{i=1}^{K}P(y_{j}=i)p(x_{j}|y_{j}=i)}{\text{product}}$$

What happens if we assume Hard assignment?

$$P(y_j = i) = 1 \text{ if } i = C(j)$$

= 0 otherwise

$$\arg\max\prod_{j} P(x_j) = \arg\max\prod_{j} p(x_j|y_j = C(j))$$
 k-means!
$$= \arg\max\prod_{j=1}^n \exp(\frac{-1}{2\sigma^2} ||x_j - \mu_{C(j)}||^2)$$

$$= \arg\min\sum_{j=1}^n ||x_j - \mu_{C(j)}||^2) = \arg\min_{\mu, C} F(\mu, C)$$

Expectation-Maximization (EM)

A general algorithm to deal with hidden data, but we will study it in the context of unsupervised learning (hidden labels) first

- No need to choose step size as in Gradient methods.
- EM is an Iterative algorithm with two linked steps:

E-step: fill-in hidden data (Y) using inference
M-step: apply standard MLE/MAP method to estimate parameters

 $\{p_i, \mu_i, \Sigma_i\}_{i=1}^k$

• We will see that this procedure monotonically improves the likelihood (or leaves it unchanged). Thus it always converges to a local optimum of the likelihood.

EM for spherical, same variance GMMs

E-step

Compute "expected" classes of all datapoints for each class

M-step

Compute Max. like μ given our data's class membership distributions (weights)

Iterate.

EM for spherical, same variance GMMs

E-step

Compute "expected" classes of all datapoints for each class

$$P(y=i|x_j,\mu_1...\mu_k) \propto exp\left(-\frac{1}{2\sigma^2} \|x_j - \mu_i\|^2\right) P(y=i)$$
In K-means "E-step" we do hard assignment

In K-means "E-step"

EM does soft assignment

M-step

Compute Max. like μ given our data's class membership distributions (weights)

$$\mu_{i} = \frac{\sum_{j=1}^{m} P(y=i|x_{j})x_{j}}{\sum_{j=1}^{m} P(y=i|x_{j})}$$

Exactly same as MLE with weighted data

Iterate.

EM for general GMMs

Iterate. On iteration t let our estimates be

$$\lambda_t = \{ \, \mu_1^{(t)}, \, \mu_2^{(t)} \, ... \, \, \mu_k^{(t)}, \, \Sigma_1^{(t)}, \, \Sigma_2^{(t)} \, ... \, \, \Sigma_k^{(t)}, \, \rho_1^{(t)}, \, \rho_2^{(t)} \, ... \, \, \rho_k^{(t)} \, \}$$

 $p_i^{(t)}$ is shorthand for estimate of P(y=i) on t'th iteration

E-step

Compute "expected" classes of all datapoints for each class

$$P(y = i | x_j, \lambda_t) \propto p_i^{(t)} p(x_j | \mu_i^{(t)}, \Sigma_i^{(t)})$$

Just evaluate a Gaussian at x_i

M-step

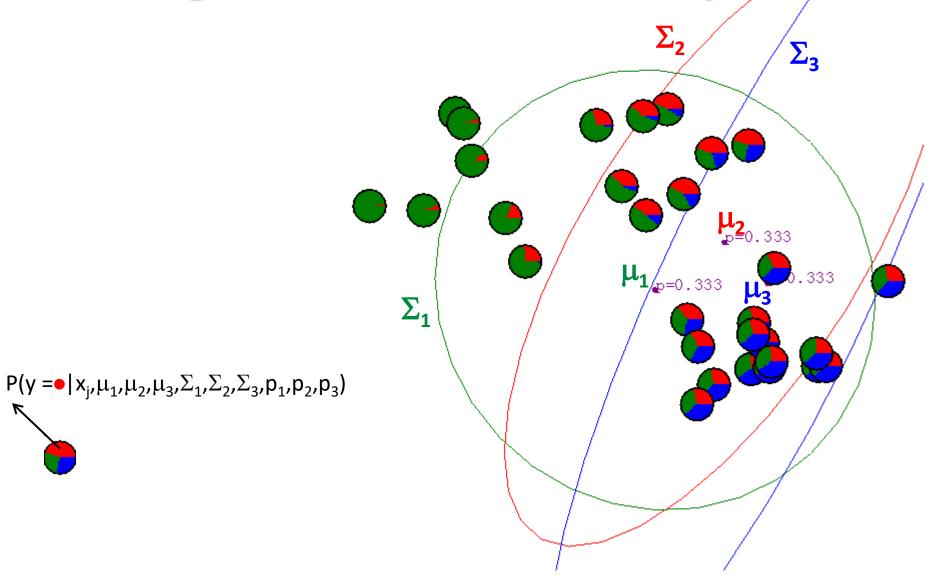
Compute MLEs given our data's class membership distributions (weights)

$$\mu_{i}^{(t+1)} = \frac{\sum_{j} P(y = i | x_{j}, \lambda_{t}) x_{j}}{\sum_{j} P(y = i | x_{j}, \lambda_{t})} \qquad \sum_{i} \frac{\sum_{j} P(y = i | x_{j}, \lambda_{t}) (x_{j} - \mu_{i}^{(t+1)}) (x_{j} - \mu_{i}^{(t+1)})^{T}}{\sum_{j} P(y = i | x_{j}, \lambda_{t})}$$

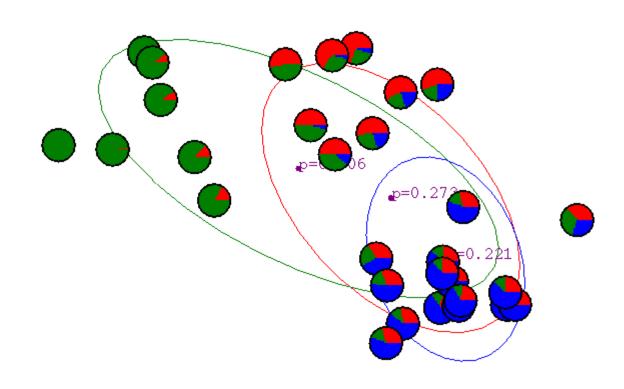
$$p_{i}^{(t+1)} = \frac{\sum_{j} P(y = i | x_{j}, \lambda_{t})}{m} \qquad m = \#data points$$

EM Convergence

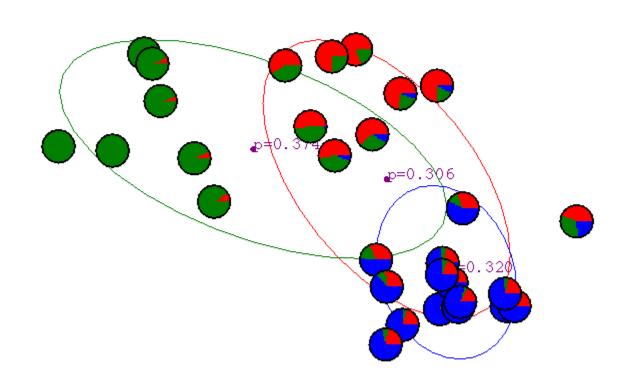
EM for general GMMs: Example



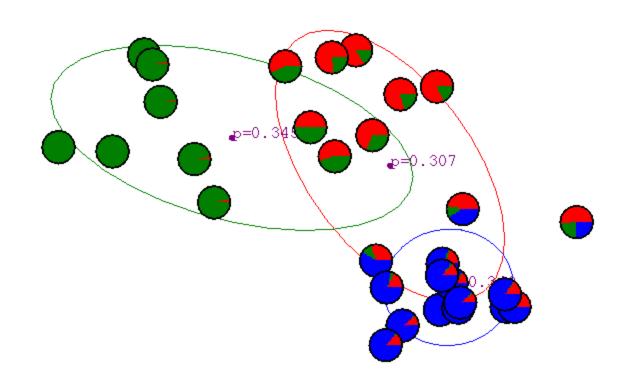
After 1st iteration



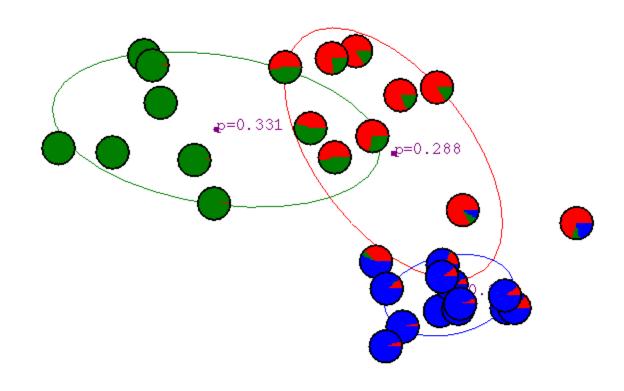
After 2nd iteration



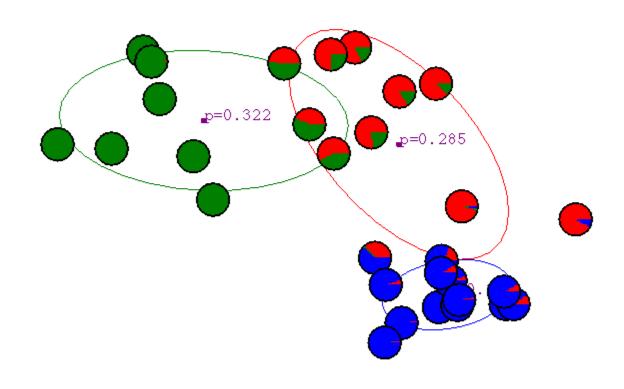
After 3rd iteration



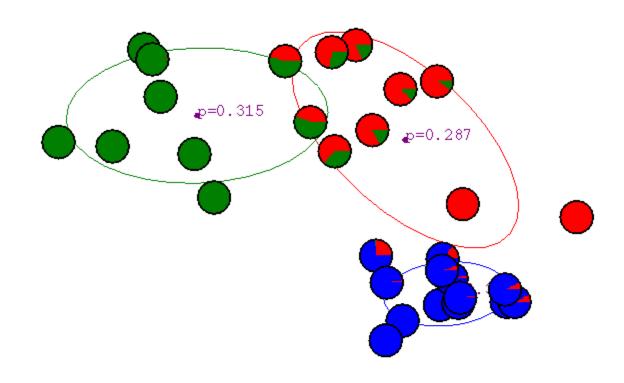
After 4th iteration



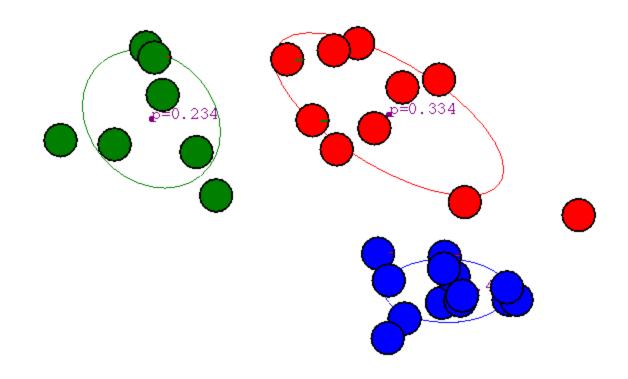
After 5th iteration



After 6th iteration

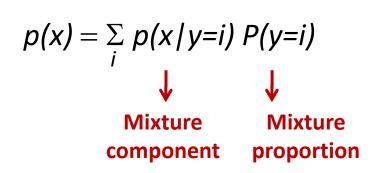


After 20th iteration

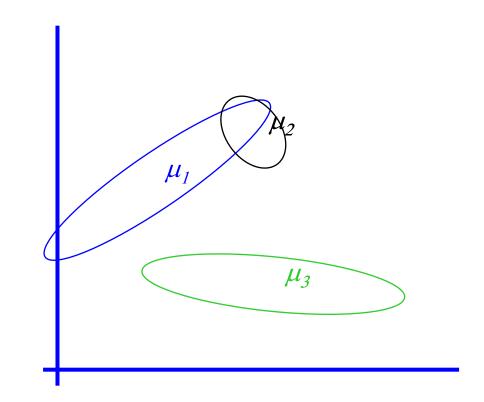


General GMM

GMM – Gaussian Mixture Model (Multi-modal distribution)



$$p(x|y=i) \sim N(\mu_i, \Sigma_i)$$



What you need to know...

- Hierarchical clustering algorithms
 - Single-linkage
 - Complete-linkage
 - Centroid-linkage
 - Average-linkage
- Partition based clustering algorithms
 - K-means
 - Coordinate descent
 - Seeding
 - Choosing K
 - Mixture modelsEM algorithm

General EM algorithm

Marginal likelihood – \mathbf{x} is observed, \mathbf{z} is missing:

$$\log P(D; \theta) = \log \prod_{j=1}^{m} P(\mathbf{x}_{j} | \theta)$$

$$= \sum_{j=1}^{m} \log P(\mathbf{x}_{j} | \theta)$$

$$= \sum_{j=1}^{m} \log \sum_{\mathbf{z}} P(\mathbf{x}_{j}, \mathbf{z} | \theta)$$

$$= \sum_{j=1}^{m} \log \sum_{\mathbf{z}} P(\mathbf{x}_{j}, \mathbf{z} | \theta)$$

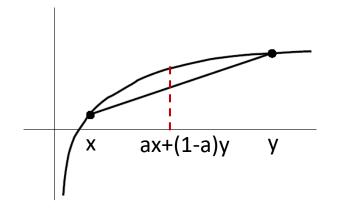
How to maximize marginal likelihood using EM?

Lower-bound on marginal likelihood

$$\log P(D; \theta) = \sum_{j=1}^{m} \log \sum_{\mathbf{z}} P(\mathbf{x}_{j}, \mathbf{z} \mid \theta)$$

$$= \sum_{j=1}^{m} \log \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_{j}) \frac{P(\mathbf{z}, \mathbf{x}_{j} \mid \theta)}{Q(\mathbf{z} \mid \mathbf{x}_{j})}$$
 Variational approach

Jensen's inequality: $\log \sum_{z} P(z) f(z) \ge \sum_{z} P(z) \log f(z)$



log: concave function

$$log(ax+(1-a)y) \ge a log(x) + (1-a) log(y)$$

Lower-bound on marginal likelihood

$$\log P(D; \theta) = \sum_{j=1}^{m} \log \sum_{\mathbf{z}} P(\mathbf{x}_{j}, \mathbf{z} \mid \theta)$$

$$= \sum_{j=1}^{m} \log \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_{j}) \frac{P(\mathbf{z}, \mathbf{x}_{j} \mid \theta)}{Q(\mathbf{z} \mid \mathbf{x}_{j})}$$

$$= \sum_{j=1}^{m} \log \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_{j}) \frac{P(\mathbf{z}, \mathbf{x}_{j} \mid \theta)}{Q(\mathbf{z} \mid \mathbf{x}_{j})}$$

Jensen's inequality: $\log \sum_{z} P(z) f(z) \ge \sum_{z} P(z) \log f(z)$

$$\geq \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j \mid \theta)}{Q(\mathbf{z} \mid \mathbf{x}_j)} =: F(\theta, Q)$$

EM as Coordinate Ascent

$$\log P(D; \theta) \geq F(\theta, Q)$$

E-step: Fix θ , maximize F over Q

$$Q^{t+1} = \arg\max_{Q} F(\theta^t, Q)$$

M-step: Fix Q, maximize F over θ

$$\theta^{t+1} = \arg\max_{\theta} F(\theta, Q^{t+1})$$

Convergence of EM

$$\log P(D; \theta) \geq F(\theta, Q)$$

E-step: Fix θ , maximize F over Q

$$Q^{t+1} = \arg\max_{Q} F(\theta^t, Q)$$

E-step maximizes lower bound F on marginal likelihood => doesn't decrease the marginal likelihood

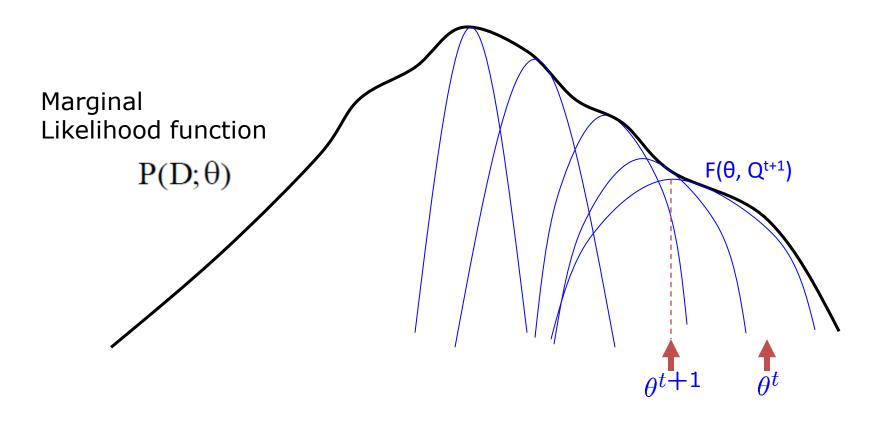
M-step: Fix Q, maximize F over θ

$$\theta^{t+1} = \arg\max_{\theta} F(\theta, Q^{t+1})$$

M-step maximizes lower bound F on marginal likelihood => doesn't decrease the marginal likelihood

Since marginal likelihood is bounded, convergence follows!

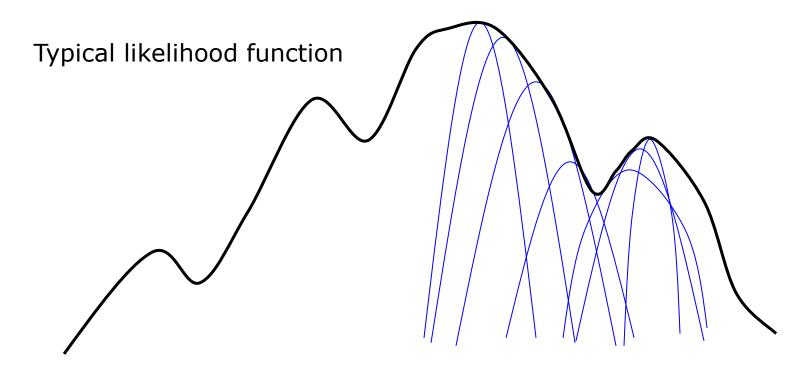
Convergence of EM



Sequence of EM lower bound F-functions

EM monotonically converges to a local maximum of likelihood!

EM & Local Maxima



Different sequence of EM lower bound F-functions depending on initialization

Use multiple, randomized initializations in practice

EM as Coordinate Ascent

$$\log P(D; \theta) \geq F(\theta, Q)$$

E-step: Fix θ , maximize F over Q

$$Q^{t+1} = \arg\max_{Q} F(\theta^t, Q)$$

M-step: Fix Q, maximize F over θ

$$\theta^{t+1} = \arg\max_{\theta} F(\theta, Q^{t+1})$$

E step

$$\log P(D; \theta) \geq F(\theta, Q)$$

E-step: Fix θ , maximize F over Q

$$\log P(D; \theta^{(t)}) \geq F(\theta^{(t)}, Q) = \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_{j}) \log \frac{P(\mathbf{z}, \mathbf{x}_{j} \mid \theta^{(t)})}{Q(\mathbf{z} \mid \mathbf{x}_{j})}$$

$$= \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_{j}) \log \frac{P(\mathbf{z} \mid \mathbf{x}_{j}, \theta^{(t)}) P(\mathbf{x}_{j} \mid \theta^{(t)})}{Q(\mathbf{z} \mid \mathbf{x}_{j})}$$

$$= \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_{j}) \log \frac{P(\mathbf{z} \mid \mathbf{x}_{j}, \theta^{(t)})}{Q(\mathbf{z} \mid \mathbf{x}_{j})} + \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_{j}) \log P(\mathbf{x}_{j} \mid \theta^{(t)})$$

$$-KL(Q(\mathbf{z} \mid \mathbf{x}_{j}), P(\mathbf{z} \mid \mathbf{x}_{j}, \theta^{(t)})) \qquad \log P(D; \theta^{(t)})$$

KL divergence between two distributions

E step

$$\log P(D; \theta) \geq F(\theta, Q)$$

E-step: Fix θ , maximize F over Q

$$\log \mathbf{P}(\mathbf{D}; \boldsymbol{\theta}^{(t)}) \geq F(\boldsymbol{\theta}^{(t)}, Q) = \sum_{j=1}^{m} \sum_{\mathbf{z}} Q(\mathbf{z} \mid \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j \mid \boldsymbol{\theta}^{(t)})}{Q(\mathbf{z} \mid \mathbf{x}_j)}$$
$$= \sum_{j=1}^{m} -KL(Q(\mathbf{z} | \mathbf{x}_j), P(\mathbf{z} | \mathbf{x}_j, \boldsymbol{\theta}^{(t)})) + \log \mathbf{P}(\mathbf{D}; \boldsymbol{\theta}^{(t)})$$

KL>=0, above expression is maximized if KL divergence = 0

$$KL(Q,P) = 0 \text{ iff } Q = P$$

Therefore,

Estep:
$$Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) = P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)})$$

E step

$$\log P(D; \theta) \geq F(\theta, Q)$$

E-step: Fix θ , maximize F over Q

$$\log \mathbf{P}(\mathbf{D}; \boldsymbol{\theta}^{(t)}) \ge F(\boldsymbol{\theta}^{(t)}, Q) = \sum_{j=1}^{m} -KL(Q(\mathbf{z}|\mathbf{x}_j), P(\mathbf{z}|\mathbf{x}_j, \boldsymbol{\theta}^{(t)})) + \log \mathbf{P}(\mathbf{D}; \boldsymbol{\theta}^{(t)})$$

Compute probability of missing data z given current choice of θ

Re-aligns F with marginal likelihood!!

$$F(\theta^{(t)}, Q^{(t+1)}) = \log P(D; \theta^{(t)})$$

M step

$$\log P(D; \theta) \geq F(\theta, Q)$$

M-step: Fix Q, maximize F over θ

$$\log \mathbf{P}(\mathbf{D}; \boldsymbol{\theta}) \ge F(\boldsymbol{\theta}, Q^{(t+1)}) = \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \log \frac{P(\mathbf{z}, \mathbf{x}_j \mid \boldsymbol{\theta})}{Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j)}$$

$$= \sum_{j=1}^{m} \sum_{\mathbf{z}} Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_{j}) \log P(\mathbf{z}, \mathbf{x}_{j} \mid \theta) + \sum_{j=1}^{m} H(Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_{j}))$$
Fixed (Independent of θ)
$$\sum_{\mathbf{z}} \sum_{j=1}^{m} \log P(\mathbf{z}, \mathbf{x}_{j} \mid \theta) Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_{j}) \qquad \text{Log likelihood if } \mathbf{z} \text{ was known}$$
Expected log likelihood wrt Q

M step

$$\log P(D; \theta) \geq F(\theta, Q)$$

M-step: Fix Q, maximize F over θ

$$\log \mathbf{P}(\mathbf{D}; \boldsymbol{\theta}) \geq F(\boldsymbol{\theta}, Q^{(t+1)}) = \sum_{\mathbf{z}} \sum_{j=1}^{m} \log P(\mathbf{z}, \mathbf{x}_j \mid \boldsymbol{\theta}) Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) \\ + \sum_{j=1}^{m} H(Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j))$$
Fixed (Independent of $\boldsymbol{\theta}$)

$$\theta^{(t+1)} \leftarrow \arg\max_{\theta} \sum_{\mathbf{z}} \sum_{j=1}^{m} \log P(\mathbf{z}, \mathbf{x}_{j} \mid \theta) Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_{j})$$
Expected log likelihood wrt Q^(t+1)

Use expected counts instead of counts when computing MLE: If learning requires Count(\mathbf{x} , \mathbf{z}), Use $E_{Q(t+1)}[Count(\mathbf{x}$, \mathbf{z})]

EM as Coordinate Ascent

$$\log P(D; \theta) \geq F(\theta, Q)$$

E-step: Fix θ , maximize F over Q

$$Q^{t+1} = \arg\max_{Q} F(\theta^t, Q)$$

$$Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j) = P(\mathbf{z} \mid \mathbf{x}_j, \theta^{(t)}) \qquad \text{E.g., } P(y = i \mid x_j, m_t)$$

Compute probability of missing data given current choice of θ

M-step: Fix Q, maximize F over θ

$$\theta^{t+1} = \arg\max_{\theta} F(\theta, Q^{t+1})$$

$$\theta^{(t+1)} \leftarrow \arg\max_{\theta} \sum_{\mathbf{z}} \sum_{j=1}^{m} \log P(\mathbf{z}, \mathbf{x}_j \mid \theta) Q^{(t+1)}(\mathbf{z} \mid \mathbf{x}_j)$$

Compute estimate of θ by maximizing marginal likelihood using $Q^{(t+1)}(z|x_i)$

Summary: EM Algorithm

- A way of maximizing likelihood function for hidden variable models. Finds MLE of parameters when the original (hard) problem can be broken up into two (easy) pieces:
 - 1. Estimate some "missing" or "unobserved" data from observed data and current parameters.
 - 2. Using this "complete" data, find the maximum likelihood parameter estimates.
- Alternate between filling in the latent variables using the best guess (posterior) and updating the parameters based on this guess:
 - 1. E-step: $Q^{t+1} = arg \max_{Q} F(\theta^{t}, Q)$ 2. M-step: $\theta^{t+1} = arg \max_{Q} F(\theta, Q^{t+1})$
- In the M-step we optimize a lower bound on the likelihood. In the E-step we close the gap, making bound=likelihood.
- EM performs coordinate ascent on F, can get stuck in local minima.
- BUT Extremely popular in practice.