

Probabilistic Graphical Models

10-708

Learning Partially Observed Graphical Models

Eric Xing

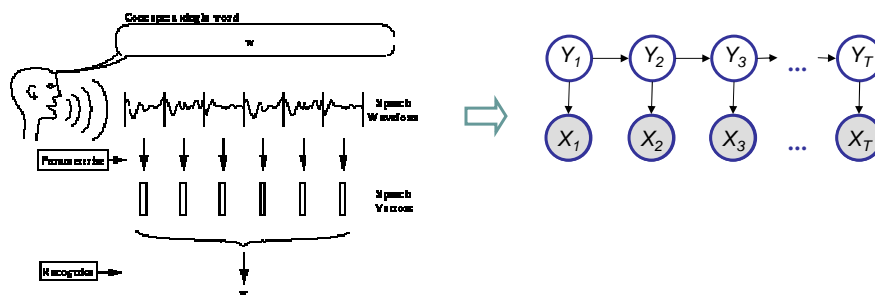
Lecture 13, Oct 26, 2005

Reading: MJ-Chap. 5,10,11



Partially observed GMs

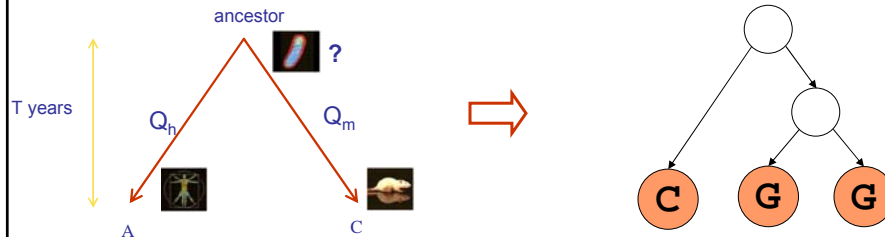
- Speech recognition



Partially observed GM



- Biological Evolution



Unobserved Variables

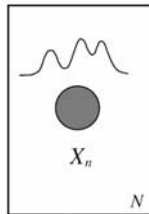


- A variable can be unobserved (latent) because:
 - it is an imaginary quantity meant to provide some simplified and abstractive view of the data generation process
 - e.g., speech recognition models, mixture models ...
 - it is a real-world object and/or phenomena, but difficult or impossible to measure
 - e.g., the temperature of a star, causes of a disease, evolutionary ancestors ...
 - it is a real-world object and/or phenomena, but sometimes wasn't measured, because of faulty sensors, etc.
- Discrete latent variables can be used to partition/cluster data into sub-groups.
- Continuous latent variables (factors) can be used for dimensionality reduction (factor analysis, etc).

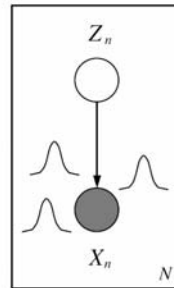
Mixture models



- A density model $p(x)$ may be multi-modal.
- We may be able to model it as a mixture of uni-modal distributions (e.g., Gaussians).
- Each mode may correspond to a different sub-population (e.g., male and female).



(a)



(b)

Gaussian Mixture Models (GMMs)



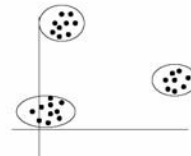
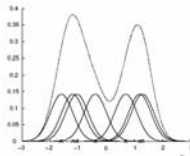
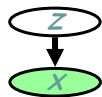
- Consider a mixture of K Gaussian components:
 - Z is a latent class indicator vector: $p(z_n) = \text{multi}(z_n : \pi) = \prod (\pi_k)^{z_n^k}$
 - X is a conditional Gaussian variable with a class-specific mean/covariance

$$p(x_n | z_n^k = 1, \mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2}(x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)\right\}$$

- The likelihood of a sample:

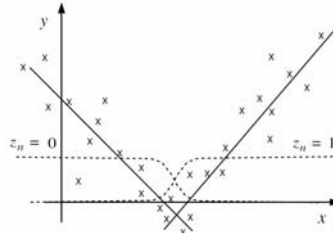
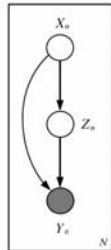
$$p(x_n | \mu, \Sigma) = \sum_k p(z^k = 1 | \pi) p(x_n | z^k = 1, \mu, \Sigma) \quad \begin{matrix} \text{mixture proportion} & \text{mixture component} \\ \downarrow & \downarrow \end{matrix}$$

$$= \sum_{z_n} \prod_k (\pi_k)^{z_n^k} \mathcal{N}(x_n; \mu_k, \Sigma_k)^{z_n^k} = \sum_k \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k)$$



- This model can be used for unsupervised clustering.
 - This model (fit by AutoClass) has been used to discover new kinds of stars in astronomical data, etc.

Conditional mixture model: Mixture of experts



- We will model $p(Y|X)$ using different experts, each responsible for different regions of the input space.

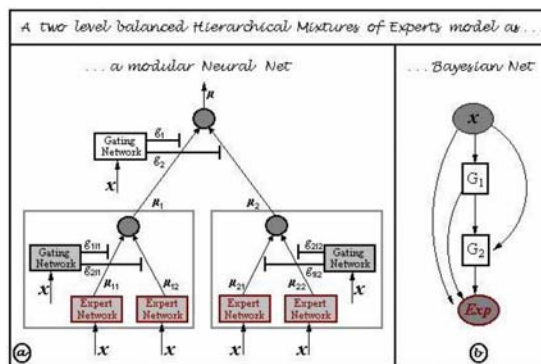
- Latent variable Z chooses expert using softmax gating function:

$$P(z^k = 1|x) = \text{Softmax}(\xi^T x)$$

- Each expert can be a linear regression model: $P(y|x, z^k = 1) = \mathcal{N}(y; \theta_k^T x, \sigma_k^2)$
- The posterior expert responsibilities are

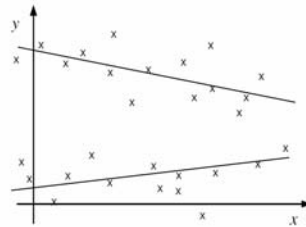
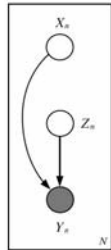
$$P(z^k = 1|x, y, \theta) = \frac{p(z^k = 1|x) p_k(y|x, \theta_k, \sigma_k^2)}{\sum_j p(z^j = 1|x) p_j(y|x, \theta_j, \sigma_j^2)}$$

Hierarchical mixture of experts



- This is like a soft version of a depth-2 classification/regression tree.
- $P(Y|X, G_1, G_2)$ can be modeled as a GLIM, with parameters dependent on the values of G_1 and G_2 (which specify a "conditional path" to a given leaf in the tree).

Mixture of overlapping experts



- By removing the $X \rightarrow Z$ arc, we can make the partitions independent of the input, thus allowing overlap.
- This is a mixture of linear regressors; each subpopulation has a different conditional mean.

$$p(z^k = 1 | x, y, \theta) = \frac{p(z^k = 1) p_k(y | x, \theta_k, \sigma_k^2)}{\sum_j p(z^j = 1) p_j(y | x, \theta_j, \sigma_j^2)}$$

Why is Learning Harder?

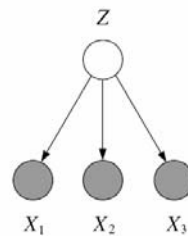
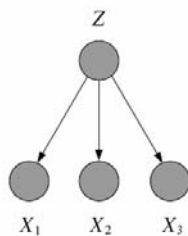


- In fully observed iid settings, the log likelihood decomposes into a sum of local terms (at least for directed models).

$$\ell_c(\theta; \mathcal{D}) = \log p(\mathbf{x}, \mathbf{z} | \theta) = \log p(\mathbf{z} | \theta_z) + \log p(\mathbf{x} | \mathbf{z}, \theta_x)$$

- With latent variables, all the parameters become coupled together via marginalization

$$\ell_c(\theta; \mathcal{D}) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z} | \theta) = \log \sum_{\mathbf{z}} p(\mathbf{z} | \theta_z) p(\mathbf{x} | \mathbf{z}, \theta_x)$$



Gradient Learning for mixture models



- We can learn mixture densities using gradient descent on the log likelihood. The gradients are quite interesting:

$$\begin{aligned}\mathcal{L}(\theta) &= \log p(x|\theta) = \log \sum_k \pi_k p_k(x|\theta_k) \\ \frac{\partial \mathcal{L}}{\partial \theta} &= \frac{1}{p(x|\theta)} \sum_k \pi_k \frac{\partial p_k(x|\theta_k)}{\partial \theta} \\ &= \sum_k \frac{\pi_k}{p(x|\theta)} p_k(x|\theta_k) \frac{\partial \log p_k(x|\theta_k)}{\partial \theta} \\ &= \sum_k \pi_k \frac{p_k(x|\theta_k)}{p(x|\theta)} \frac{\partial \log p_k(x|\theta_k)}{\partial \theta_k} = \sum_k r_k \frac{\partial \ell_k}{\partial \theta_k}\end{aligned}$$

- In other words, the gradient is the responsibility weighted sum of the individual log likelihood gradients.
- Can pass this to a conjugate gradient routine.

Parameter Constraints



- Often we have constraints on the parameters, e.g. $\sum_k \pi_k = 1$, Σ being symmetric positive definite (hence $\Sigma_{ii} > 0$).
- We can use constrained optimization, or we can reparameterize in terms of unconstrained values.

- For normalized weights, use the softmax transform: $\pi_k = \frac{\exp(\gamma_k)}{\sum_j \exp(\gamma_j)}$
- For covariance matrices, use the Cholesky decomposition:

$$\Sigma^{-1} = \mathbf{A}^T \mathbf{A}$$

where \mathbf{A} is upper diagonal with positive diagonal:

$$\mathbf{A}_{ii} = \exp(\lambda_i) > 0 \quad \mathbf{A}_{ij} = \eta_{ij} \quad (j > i) \quad \mathbf{A}_{ij} = 0 \quad (j < i)$$

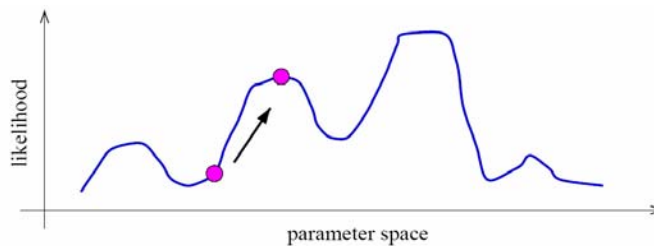
the parameters $\gamma_i, \lambda_i, \eta_{ij} \in \mathbb{R}$ are unconstrained.

- Use chain rule to compute $\frac{\partial \mathcal{L}}{\partial \pi}, \frac{\partial \mathcal{L}}{\partial \mathbf{A}}$.

Identifiability



- A mixture model induces a multi-modal likelihood.
- Hence gradient ascent can only find a local maximum.
- Mixture models are unidentifiable, since we can always switch the hidden labels without affecting the likelihood.
- Hence we should be careful in trying to interpret the “meaning” of latent variables.



Expectation-Maximization (EM) Algorithm



- EM is an optimization strategy for objective functions that can be interpreted as likelihoods in the presence of missing data.
- It is much simpler than gradient methods:
 - No need to choose step size.
 - Enforces constraints automatically.
 - Calls inference and fully observed learning as subroutines.
- EM is an Iterative algorithm with two linked steps:
 - E-step: fill-in hidden values using inference, $p(z|x, \theta)$.
 - M-step: update parameters $t+1$ using standard MLE/MAP method applied to completed data
- We will prove that this procedure monotonically improves (or leaves it unchanged). Thus it always converges to a local optimum of the likelihood.

Complete & Incomplete Log Likelihoods



- Complete log likelihood

Let X denote the observable variable(s), and Z denote the latent variable(s).
If Z could be observed, then

$$\ell_c(\theta; \mathbf{x}, \mathbf{z}) \stackrel{\text{def}}{=} \log p(\mathbf{x}, \mathbf{z} | \theta)$$

- Usually, optimizing $\ell_c()$ given both \mathbf{z} and \mathbf{x} is straightforward (c.f. MLE for fully observed models).
- Recalled that in this case the objective for, e.g., MLE, decomposes into a sum of factors, the parameter for each factor can be estimated separately.
- **But given that Z is not observed, $\ell_c()$ is a random quantity, cannot be maximized directly.**

- Incomplete log likelihood

With \mathbf{z} unobserved, our objective becomes the log of a marginal probability:

$$\ell(\theta; \mathbf{x}) = \log p(\mathbf{x} | \theta) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z} | \theta)$$

- **This objective won't decouple**

Expected Complete Log Likelihood



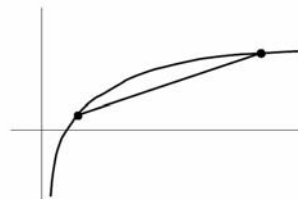
- For **any** distribution $q(\mathbf{z})$, define *expected complete log likelihood*:

$$\langle \ell_c(\theta; \mathbf{x}, \mathbf{z}) \rangle_q \stackrel{\text{def}}{=} \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}, \theta) \log p(\mathbf{x}, \mathbf{z} | \theta)$$

- A deterministic function of θ
- Linear in $\ell_c()$ --- inherit its factorizability
- Does maximizing this surrogate yield a maximizer of the likelihood?

- Jensen's inequality

$$\begin{aligned} \ell(\theta; \mathbf{x}) &= \log p(\mathbf{x} | \theta) \\ &= \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z} | \theta) \\ &= \log \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{q(\mathbf{z} | \mathbf{x})} \\ &= \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \log \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{q(\mathbf{z} | \mathbf{x})} \end{aligned}$$



$$\Rightarrow \ell(\theta; \mathbf{x}) \geq \langle \ell_c(\theta; \mathbf{x}, \mathbf{z}) \rangle_q + H_q$$

Lower Bounds and Free Energy

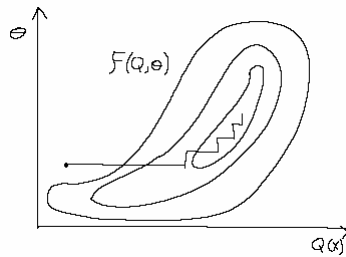


- For fixed data \mathbf{x} , define a functional called the free energy:

$$F(q, \theta) \stackrel{\text{def}}{=} \sum_z q(z | \mathbf{x}) \log \frac{p(\mathbf{x}, z | \theta)}{q(z | \mathbf{x})} \leq \ell(\theta; \mathbf{x})$$

- The EM algorithm is coordinate-ascent on F :

- E-step:** $q^{t+1} = \arg \max_q F(q, \theta^t)$
- M-step:** $\theta^{t+1} = \arg \max_{\theta} F(q^{t+1}, \theta)$



E-step: maximization of expected ℓ_c w.r.t. q



- Claim: $q^{t+1} = \arg \max_q F(q, \theta^t) = p(z | \mathbf{x}, \theta^t)$

- This is the posterior distribution over the latent variables given the data and the parameters. Often we need this at test time anyway (e.g. to perform classification).

- Proof (easy): this setting attains the bound $\ell(\theta; \mathbf{x}) \geq F(q, \theta)$

$$\begin{aligned} F(p(z | \mathbf{x}, \theta^t), \theta^t) &= \sum_z p(z | \mathbf{x}, \theta^t) \log \frac{p(\mathbf{x}, z | \theta^t)}{p(z | \mathbf{x}, \theta^t)} \\ &= \sum_z q(z | \mathbf{x}) \log p(\mathbf{x} | \theta^t) \\ &= \log p(\mathbf{x} | \theta^t) = \ell(\theta^t; \mathbf{x}) \end{aligned}$$

- Can also show this result using variational calculus or the fact that $\ell(\theta; \mathbf{x}) - F(q, \theta) = \text{KL}(q \| p(z | \mathbf{x}, \theta))$

E-step \equiv plug in posterior expectation of latent variables



- Without loss of generality: assume that $p(\mathbf{x}, \mathbf{z} | \theta)$ is a generalized exponential family distribution:

$$p(\mathbf{x}, \mathbf{z} | \theta) = \frac{1}{Z(\theta)} h(\mathbf{x}, \mathbf{z}) \exp\left\{ \sum_i \theta_i f_i(\mathbf{x}, \mathbf{z}) \right\}$$

- Special cases: if $p(\mathbf{x} | \mathbf{Z})$ are GLIMs, then $f_i(\mathbf{x}, \mathbf{z}) = \eta_i^T(\mathbf{z}) \xi_i(\mathbf{x})$
- The expected complete log likelihood under $q^{t+1} = p(\mathbf{z} | \mathbf{x}, \theta^t)$ is

$$\begin{aligned} \langle \ell_c(\theta^t; \mathbf{x}, \mathbf{z}) \rangle_{q^{t+1}} &= \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}, \theta^t) \log p(\mathbf{x}, \mathbf{z} | \theta^t) - A(\theta) \\ &= \sum_i \theta_i^t \langle f_i(\mathbf{x}, \mathbf{z}) \rangle_{q(\mathbf{z} | \mathbf{x}, \theta^t)} - A(\theta) \\ &= \sum_i \theta_i^t \langle \eta_i(\mathbf{z}) \rangle_{q(\mathbf{z} | \mathbf{x}, \theta^t)} \xi_i(\mathbf{x}) - A(\theta) \end{aligned}$$

M-step: maximization of expected ℓ_c w.r.t. θ



- Note that the free energy breaks into two terms:

$$\begin{aligned} F(q, \theta) &= \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \log \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{q(\mathbf{z} | \mathbf{x})} \\ &= \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \log p(\mathbf{x}, \mathbf{z} | \theta) - \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \log q(\mathbf{z} | \mathbf{x}) \\ &= \langle \ell_c(\theta; \mathbf{x}, \mathbf{z}) \rangle_q + H_q \end{aligned}$$

- The first term is the expected complete log likelihood (energy) and the second term, which does not depend on θ , is the entropy.
- Thus, in the M-step, maximizing with respect to θ for fixed q we only need to consider the first term:

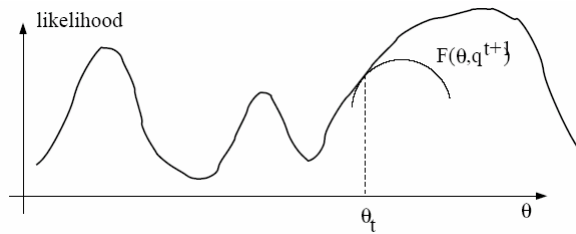
$$\theta^{t+1} = \arg \max_{\theta} \langle \ell_c(\theta; \mathbf{x}, \mathbf{z}) \rangle_{q^{t+1}} = \arg \max_{\theta} \sum_{\mathbf{z}} q(\mathbf{z} | \mathbf{x}) \log p(\mathbf{x}, \mathbf{z} | \theta)$$

- Under optimal q^{t+1} , this is equivalent to solving a standard MLE of fully observed model $p(\mathbf{x}, \mathbf{z} | \theta)$, with the **sufficient statistics** involving \mathbf{z} replaced by their expectations w.r.t. $p(\mathbf{z} | \mathbf{x}, \theta)$.

EM Constructs Sequential Convex Lower Bounds



- Consider the likelihood function and the function $F(q^{t+1}, \cdot)$.



- A hill-climbing algorithm

Summary: EM Algorithm



- A way of maximizing likelihood function for latent 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:
 - E-step: $q^{t+1} = \arg \max_q F(q, \theta^t)$
 - M-step: $\theta^{t+1} = \arg \max_{\theta} F(q^{t+1}, \theta)$
- In the M-step we optimize a lower bound on the likelihood. In the E-step we close the gap, making bound=likelihood.

Example: Gaussian mixture model



- A mixture of K Gaussians:

- Z is a latent class indicator vector

$$p(\mathbf{z}_n) = \text{multi}(\mathbf{z}_n; \boldsymbol{\pi}) = \prod_k (\pi_k)^{z_n^k}$$

- X is a conditional Gaussian variable with a class-specific mean/covariance

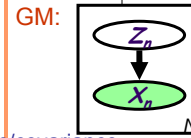
$$p(x_n | z_n^k = 1, \mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2}(x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)\right\}$$

- The likelihood of a sample:

$$\begin{aligned} p(x_n | \mu, \Sigma) &= \sum_k p(z^k = 1 | \boldsymbol{\pi}) p(x_n | z^k = 1, \mu, \Sigma) \\ &= \sum_{z_n} \prod_k (\pi_k)^{z_n^k} \mathcal{N}(x_n; \mu_k, \Sigma_k)^{z_n^k} = \sum_k \pi_k \mathcal{N}(x_n; \mu_k, \Sigma_k) \end{aligned}$$

- The expected complete log likelihood

$$\begin{aligned} \langle \ell_c(\boldsymbol{\theta}; \mathbf{x}, \mathbf{z}) \rangle &= \sum_n \langle \log p(\mathbf{z}_n | \boldsymbol{\pi}) \rangle_{p(\mathbf{z}|\mathbf{x})} + \sum_n \langle \log p(x_n | \mathbf{z}_n, \mu, \Sigma) \rangle_{p(\mathbf{z}|\mathbf{x})} \\ &= \sum_n \sum_k \langle z_n^k \rangle \log \pi_k - \frac{1}{2} \sum_n \sum_k \langle z_n^k \rangle (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) + \log |\Sigma_k| + \mathcal{C} \end{aligned}$$



E-step

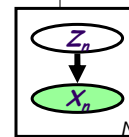


- We maximize $\langle \ell_c(\boldsymbol{\theta}) \rangle$ iteratively using the following iterative procedure:

- **Expectation step:** computing the expected value of the hidden variables (i.e., \mathbf{z}) given current est. of the parameters (i.e., $\boldsymbol{\pi}$ and μ).

$$\tau_n^{k(t)} = \langle z_n^k \rangle_{q^{(t)}} = p(z_n^k = 1 | \mathbf{x}, \mu^{(t)}, \Sigma^{(t)}) = \frac{\pi_k^{(t)} \mathcal{N}(x_n; \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_i \pi_i^{(t)} \mathcal{N}(x_n; \mu_i^{(t)}, \Sigma_i^{(t)})}$$

- Here we are essentially doing **inference**

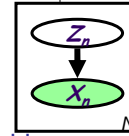


M-step



- We maximize $\langle l_c(\theta) \rangle$ iteratively using the following iterative procedure:

- Maximization step:** compute the parameters under current results of the expected value of the hidden variables



$$\pi_k^* = \arg \max \langle l_c(\theta) \rangle \Rightarrow \frac{\partial}{\partial \pi_k} \langle l_c(\theta) \rangle = 0, \forall k, \quad \text{s.t. } \sum_k \pi_k = 1$$

$$\Rightarrow \pi_k^* = \frac{\sum_n \langle z_n^k \rangle_{q^{(t)}}}{N} = \frac{\sum_n \tau_n^{k(t)}}{N} = n_k / N$$

$$\mu_k^* = \arg \max \langle l(\theta) \rangle \Rightarrow \mu_k^{(t+1)} = \frac{\sum_n \tau_n^{k(t)} x_n}{\sum_n \tau_n^{k(t)}}$$

$$\Sigma_k^* = \arg \max \langle l(\theta) \rangle \Rightarrow \Sigma_k^{(t+1)} = \frac{\sum_n \tau_n^{k(t)} (x_n - \mu_k^{(t+1)})(x_n - \mu_k^{(t+1)})^T}{\sum_n \tau_n^{k(t)}}$$

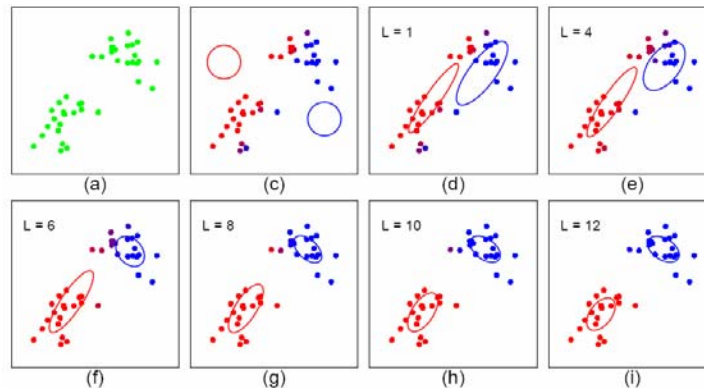
Fact:

$$\frac{\partial \log |A^{-1}|}{\partial A^{-1}} = A^T$$

$$\frac{\partial x^T A x}{\partial A} = x x^T$$

- This is isomorphic to **MLE** except that the variables that are hidden are replaced by their expectations (in general they will be replaced by their corresponding "**sufficient statistics**")

EM for MOG



Compare: K-means

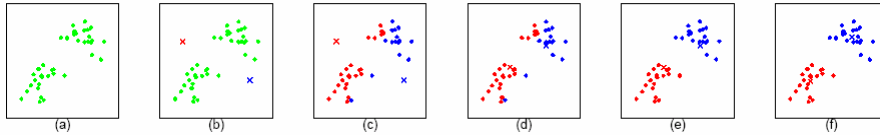


- The EM algorithm for mixtures of Gaussians is like a "soft version" of the K-means algorithm.
- In the K-means "E-step" we do hard assignment:

$$z_n^{(t)} = \arg \max_k (x_n - \mu_k^{(t)})^T \Sigma_k^{-1(t)} (x_n - \mu_k^{(t)})$$

- In the K-means "M-step" we update the means as the weighted sum of the data, but now the weights are 0 or 1:

$$\mu_k^{(t+1)} = \frac{\sum_n \delta(z_n^{(t)}, k) x_n}{\sum_n \delta(z_n^{(t)}, k)}$$



EM for conditional mixture model



- Model:

$$p(y|x) = \sum_k p(z^k = 1 | x, \xi) p(y | z^k = 1, x, \theta_k, \sigma)$$

- The objective function

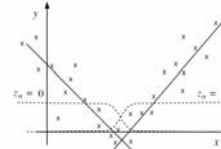
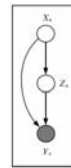
$$\begin{aligned} \langle \ell_c(\theta; x, y, z) \rangle &= \sum_n \langle \log p(z_n | x_n, \xi) \rangle_{p(z|x,y)} + \sum_n \langle \log p(y_n | x_n, z_n, \theta, \sigma) \rangle_{p(z|x,y)} \\ &= \sum_n \sum_k \langle z_n^k \rangle \log(\text{softmax}(\xi_k^T x_n)) - \frac{1}{2} \sum_n \sum_k \langle z_n^k \rangle \left(\frac{(y_n - \theta_k^T x_n)^2}{\sigma_k^2} + \log \sigma_k^2 + C \right) \end{aligned}$$

- EM:

- E-step: $\tau_n^{k(t)} = P(z_n^k = 1 | x_n, y_n, \theta) = \frac{p(z_n^k = 1 | x_n) p_k(y_n | x_n, \theta_k, \sigma_k^2)}{\sum_j p(z_n^j = 1 | x_n) p_j(y_n | x_n, \theta_j, \sigma_j^2)}$

- M-step:

- using the normal equation for standard LR $\theta = (X^T X)^{-1} X^T Y$, but with the data re-weighted by τ (homework)
- IRLS and/or weighted IRLS algorithm to update $\{\xi_k, \theta_k, \sigma_k\}$ based on data pair (x_n, y_n) , with weights $\tau_n^{k(t)}$ (homework)



EM for general BNs



while not converged

% E-step

for each node i

$ESS_i = 0$ **% reset expected sufficient statistics**

for each data sample n

do inference with $X_{n,H}$

for each node i

$$ESS_i += \left\langle SS_i(x_{n,i}, x_{n,\pi_i}) \right\rangle_{p(x_{n,H} | x_{n,-H})}$$

% M-step

for each node i

$\theta_i := \text{MLE}(ESS_i)$

Partially Hidden Data



- Of course, we can learn when there are missing (hidden) variables on some cases and not on others.

- In this case the cost function is:

$$\ell_c(\theta; D) = \sum_{n \in \text{Complete}} \log p(x_n, y_n | \theta) + \sum_{m \in \text{Missing}} \log \sum_{y_m} p(x_m, y_m | \theta)$$

- Note that Y_m do not have to be the same in each case --- the data can have different missing values in each different sample
- Now you can think of this in a new way: in the E-step we estimate the hidden variables on the incomplete cases only.
- The M-step optimizes the log likelihood on the complete data plus the expected likelihood on the incomplete data using the E-step.

EM Variants



- Sparse EM:
Do not re-compute exactly the posterior probability on each data point under all models, because it is almost zero. Instead keep an “active list” which you update every once in a while.
- Generalized (Incomplete) EM:
It might be hard to find the ML parameters in the M-step, even given the completed data. We can still make progress by doing an M-step that improves the likelihood a bit (e.g. gradient step). Recall the IRLS step in the mixture of experts model

A Report Card for EM



- Some good things about EM:
 - no learning rate (step-size) parameter
 - automatically enforces parameter constraints
 - very fast for low dimensions
 - each iteration guaranteed to improve likelihood
- Some bad things about EM:
 - can get stuck in local minima
 - can be slower than conjugate gradient (especially near convergence)
 - requires expensive inference step
 - is a maximum likelihood/MAP method