

Notes on Neal and Hinton's Generalized Expectation Maximization (GEM) Algorithm

Mark Johnson

Brown University

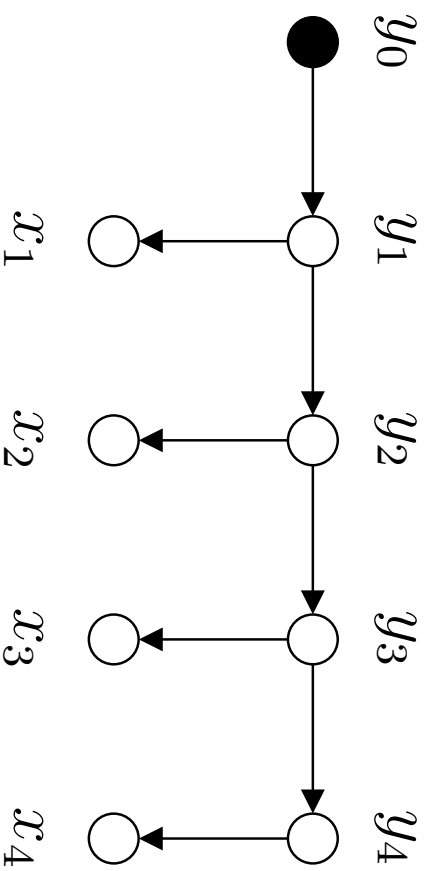
February 2005, updated November 2008

Talk overview

- What kinds of problems does expectation maximization solve?
- An example of EM
- Relaxation, and proving that EM converges
- Sufficient statistics and EM
- The generalized EM algorithm

Hidden Markov Models

States (e.g., parts of speech)



Observations (e.g., words)

$$P(\mathbf{Y}, \mathbf{X} | \boldsymbol{\theta}) = \prod_{i=1}^n P(Y_i | Y_{i-1}, \boldsymbol{\theta}) P(X_i | Y_i, \boldsymbol{\theta})$$

$$P(y_i | y_{i-1}, \boldsymbol{\theta}) = \theta_{y_i, y_{i-1}}$$

$$P(x_i | y_i, \boldsymbol{\theta}) = \theta_{x_i, y_i}$$

3

Maximum likelihood estimation

- Given *visible data* (\mathbf{y}, \mathbf{x}) , how can we estimate θ ?
- Maximum likelihood principle:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} L_{(\mathbf{y}, \mathbf{x})}(\theta), \text{ where:}$$

$$L_{(\mathbf{y}, \mathbf{x})}(\theta) = \log P_{\theta}(\mathbf{y}, \mathbf{x}) = \log P(\mathbf{y}, \mathbf{x} | \theta)$$

- For a HMM, these are simple to calculate:

$$\begin{aligned} \hat{\theta}_{y_i, y_j} &= \frac{n_{y_i, y_j}(\mathbf{y}, \mathbf{x})}{\sum_{y'_i} n_{y'_i, y_j}(\mathbf{y}, \mathbf{x})} \\ \hat{\theta}_{x_i, y_i} &= \frac{n_{x_i, y_i}(\mathbf{y}, \mathbf{x})}{\sum_{x'_i} n_{x'_i, y_i}(\mathbf{y}, \mathbf{x})} \end{aligned}$$

ML estimation from hidden data

- Our model defines $P(\mathbf{Y}, \mathbf{X})$, but our data only contains values for \mathbf{X} , i.e., the variable \mathbf{Y} is *hidden*
 - HMM example: D only contains words \mathbf{x} but not their labels \mathbf{y}
- Maximum likelihood principle still applies:

$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} L_{\mathbf{x}}(\boldsymbol{\theta})$, where:

$$L_{\mathbf{x}}(\boldsymbol{\theta}) = \log P(\mathbf{x}|\boldsymbol{\theta}) = \log \sum_{\mathbf{y} \in \mathcal{Y}} P(\mathbf{y}, \mathbf{x}|\boldsymbol{\theta})$$

- But maximizing $L_{\mathbf{x}}(\boldsymbol{\theta})$ may now be a non-trivial problem!

What does Expectation Maximization do?

- Expectation Maximization (EM) is a *maximum likelihood estimation procedure* for problems with hidden variables
- EM is good for problems where:
 - our model $P(Y, X|\theta)$ involves variables Y and X
 - our training data contains x but not y
 - maximizing $P(x|\theta)$ is hard
 - maximizing $P(y, x|\theta)$ is easy
- In HMM example: if training data consists of words x alone, and does not contain their labels

The EM algorithm

- The EM algorithm:
 - Guess an initial model $\theta^{(0)}$
 - For $t = 1, 2, \dots$, compute $Q^{(t)}(y)$ and $\theta^{(t)}$, where

$$Q^{(t)}(y) = P(y|x, \theta^{(t-1)}) \quad (\text{E-step})$$

$$\begin{aligned} \theta^{(t)} &= \operatorname{argmax}_{\theta} E_{Y \sim Q^{(t)}} [\log P(Y, x|\theta)] \quad (\text{M-step}) \\ &= \operatorname{argmax}_{\theta} \sum_{y \in \mathcal{Y}} Q^{(t)}(y) \log P(y, x|\theta) \\ &= \operatorname{argmax}_{\theta} \prod_{y \in \mathcal{Y}} P(y, x|\theta)^{Q^{(t)}(y)} \end{aligned}$$

- $Q^{(t)}(y)$ is probability of “pseudo-data” y using model $\theta^{(t-1)}$
- $\theta^{(t)}$ is the MLE based on pseudo-data (y, x) , where each (y, x) is weighted according to $Q^{(t)}(y)$

HMM example

- For a HMM, the EM formulae are:

$$\begin{aligned} Q^{(t)}(\mathbf{y}) &= P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}^{(t-1)}) \\ &= \frac{P(\mathbf{y}, \mathbf{x}|\boldsymbol{\theta}^{(t-1)})}{\sum_{\mathbf{y} \in \mathcal{Y}} P(\mathbf{y}, \mathbf{x}|\boldsymbol{\theta}^{(t-1)})} \\ \theta_{y_i, y_j}^{(t)} &= \frac{\sum_{\mathbf{y} \in \mathcal{Y}} Q^{(t)}(\mathbf{y}) n_{y_i, y_j}(\mathbf{y}, \mathbf{x})}{\sum_{y_i'} \sum_{\mathbf{y} \in \mathcal{Y}} Q^{(t)}(\mathbf{y}) n_{y_i', y_j}(\mathbf{y}, \mathbf{x})} \\ \theta_{x_i, y_i}^{(t)} &= \frac{\sum_{\mathbf{y} \in \mathcal{Y}} Q^{(t)}(\mathbf{y}) n_{x_i, y_i}(\mathbf{y}, \mathbf{x})}{\sum_{x_i'} \sum_{\mathbf{y} \in \mathcal{Y}} Q^{(t)}(\mathbf{y}) n_{x_i', y_i}(\mathbf{y}, \mathbf{x})} \end{aligned}$$

EM converges — overview

- Neal and Hinton define a function $F(Q, \theta)$ where:
 - $Q(Y)$ is a probability distribution over the hidden variables
 - θ are the model parameters

$$\operatorname{argmax}_{\theta} \max_Q F(Q, \theta) = \hat{\theta}, \text{ the MLE of } \theta$$

$$\max_Q F(Q, \theta) = L_x(\theta), \text{ the log likelihood of } \theta$$

$$\operatorname{argmax}_Q F(Q, \theta) = P(Y|x, \theta) \text{ for all } \theta$$

- The EM algorithm is an *alternating maximization* of F

$$Q^{(t)} = \operatorname{argmax}_Q F(Q, \theta^{(t-1)}) \quad (\text{E-step})$$

$$\theta^{(t)} = \operatorname{argmax}_{\theta} F(Q^{(t)}, \theta) \quad (\text{M-step})$$

The EM algorithm converges

$$\begin{aligned} F(Q, \theta) &= E_{Y \sim Q} [\log P(Y, x|\theta)] + H(Q) \\ &= L_x(\theta) - \text{KL}(Q(Y) \| P(Y|x, \theta)) \end{aligned}$$

$H(Q)$ = entropy of Q

$L_x(\theta)$ = $\log P(x|\theta)$ = log likelihood of θ

$\text{KL}(Q \| P)$ = KL divergence between Q and P

$$Q^{(t)}(Y) = P(Y|x, \theta^{(t-1)}) \quad = \underset{Q}{\text{argmax}} F(Q, \theta^{(t-1)}) \quad (\text{E-st})$$

$$\theta^{(t)} = \underset{\theta}{\text{argmax}} E_{Y \sim Q^{(t)}} [\log P(Y, x|\theta)] = \underset{\theta}{\text{argmax}} F(Q^{(t)}, \theta) \quad (\text{M-s})$$

- The maximum value of F is achieved at $\theta = \hat{\theta}$ and $Q(Y) = P(Y|x, \hat{\theta})$.
- The sequence of F values produced by the EM algorithm is *non-decreasing* and *bounded above* by $L(\hat{\theta})$.

Generalized EM

- Idea: anything that increases F gets you closer to $\hat{\theta}$
- Idea: insert any additional operations you want into the EM algorithm so long as they don't decrease F
 - Update θ after each data item has been processed
 - Visit some data items more often than others
 - Only update some components of θ on some iterations

Incremental EM for factored models

- Data and model both factor: $Y = (Y_1, \dots, Y_n), X = (X_1, \dots, X_n)$

$$P(Y, X|\theta) = \prod_{i=1}^n P(Y_i, X_i|\theta)$$

- Incremental EM algorithm:
 - Initialize $\theta^{(0)}$ and $Q_i^{(0)}(Y_i)$ for $i = 1, \dots, n$
 - E-step: Choose some data item i to be updated

$$Q_j^{(t)} = Q_j^{(t-1)} \text{ for all } j \neq i$$

$$Q_i^{(t)}(Y_i) = P(Y_i|x_i, \theta^{(t-1)})$$

- M-step:

$$\theta^{(t)} = \operatorname{argmax}_{\theta} E_{Y \sim Q^{(t)}} [\log P(Y, x|\theta)]$$

EM using sufficient statistics

- Model parameters θ estimated from *sufficient statistics* S :

$$(Y, X) \rightarrow S \rightarrow \theta$$

- In HMM example, pseudo-counts are sufficient statistics
- EM algorithm with sufficient statistics:

$$\tilde{g}^{(t)} = E_{Y \sim P(Y|x, \theta^{(t-1)})} [S] \quad (\text{E-step})$$

$$\theta^{(t)} = \text{maximum likelihood value for } \theta \text{ based on } \tilde{g}^{(t)} \quad (\text{M-step})$$

Incremental EM using sufficient statistics

- Incremental EM algorithm with sufficient statistics:

$$\boxed{(Y_i, X_i)} \rightarrow S_i \rightarrow S \rightarrow \theta \qquad S = \sum_i S_i$$

- Initialize $\theta^{(0)}$ and $\tilde{\mathbf{s}}_i^{(0)}$ for $i = 1, \dots, n$
- E-step: Choose some data item i to be updated

$$\tilde{\mathbf{s}}_j^{(t)} = \tilde{\mathbf{s}}_j^{(t-1)} \text{ for all } j \neq i$$

$$\tilde{\mathbf{s}}_i^{(t)} = E_{Y_i \sim P(Y_i | x_i, \theta^{(t-1)})} [S_i]$$

$$\tilde{\mathbf{s}}^{(t)} = \tilde{\mathbf{s}}^{(t-1)} + (\tilde{\mathbf{s}}_i^{(t)} - \tilde{\mathbf{s}}_i^{(t-1)})$$

- M-step:

$$\theta^{(t)} = \text{maximum likelihood value for } \theta \text{ based on } \tilde{\mathbf{s}}^{(t)}$$

Conclusion

- The Expectation-Maximization algorithm is a general technique for using supervised maximum likelihood estimators to solve unsupervised estimation problems
- The E-step and the M-step can be viewed as steps of an *alternating maximization procedure*
 - The functional F is bounded above by the log likelihood
 - Each E and M step increases F

⇒ Eventually the EM algorithm converges to a *local optimum* (not necessarily a global optimum)
- We can insert any steps we like into the EM algorithm so long as they do not decrease F

⇒ Incremental versions of the EM algorithm