

A Gentle Introduction to the EM Algorithm

Ted Pedersen
Department of Computer Science
University of Minnesota Duluth
tpederse@d.umn.edu

EMNLP, June 2001

Ted Pedersen - EM Panel

1

A unifying methodology

- Dempster, Laird & Rubin (1977) unified many strands of apparently unrelated work under the banner of **The EM Algorithm**
- EM had gone incognito for many years
 - Newcomb (1887)
 - McKendrick (1926)
 - Hartley (1958)
 - Baum et. al. (1970)

EMNLP, June 2001

Ted Pedersen - EM Panel

2

A general framework for solving many kinds of problems

- Filling in missing data in a sample
- Discovering the value of latent variables
- Estimating parameters of HMMs
- Estimating parameters of finite mixtures
- Unsupervised learning of clusters
- ...

EMNLP, June 2001

Ted Pedersen - EM Panel

3

EM allows us to make MLE under adverse circumstances

- What are Maximum Likelihood Estimates?
- What are these adverse circumstances?
- How does EM triumph over adversity?
- PANEL: When does it really work?

EMNLP, June 2001

Ted Pedersen - EM Panel

4

Maximum Likelihood Estimates

- Parameters describe the characteristics of a population. Their values are estimated from samples collected from that population.
- A MLE is a parameter estimate that is most consistent with the sampled data. It maximizes the likelihood function.

EMNLP, June 2001

Ted Pedersen - EM Panel

5

Coin Tossing!

- How likely am I to toss a head? A series of 10 trials/tosses yields (h,t,t,t,h,t,t,h,t,t)
– ($x_1=3, x_2=7$), $n=10$
- Probability of tossing a head = $3/10$
- That's a MLE! This estimate is absolutely consistent with the observed data.
- A few underlying details are masked...

EMNLP, June 2001

Ted Pedersen - EM Panel

6

Coin tossing unmasked

- Coin tossing is well described by the binomial distribution since there are n independent trials with two outcomes.
- Given 10 tosses, how likely is 3 heads?

$$L(\theta) = \binom{10}{3} \theta^3 (1-\theta)^7$$

EMNLP, June 2001

Ted Pedersen - EM Panel

7

Maximum Likelihood Estimates

- We seek to estimate the parameter such that it maximizes the likelihood function.
- Take the first derivative of the likelihood function with respect to the parameter theta and solve for 0. This value maximizes the likelihood function and is the MLE.

EMNLP, June 2001

Ted Pedersen - EM Panel

8

Maximizing the likelihood

$$L(\theta) = \binom{10}{3} \theta^3 (1-\theta)^7$$

$$\log L(\theta) = \log \binom{10}{3} + 3 \log \theta + 7 \log(1-\theta)$$

$$\frac{d \log L(\theta)}{d \theta} = \frac{3}{\theta} - \frac{7}{1-\theta} = 0$$

$$\frac{3}{\theta} = \frac{7}{1-\theta} \Rightarrow \theta = \frac{3}{10}$$

EMNLP, June 2001

Ted Pedersen - EM Panel

9

Multinomial MLE example

- There are n animals classified into one of four possible categories (Rao 1973).
 - Category counts are the sufficient statistics to estimate multinomial parameters
- Technique for finding MLEs is the same
 - Take derivative of likelihood function
 - Solve for zero

EMNLP, June 2001

Ted Pedersen - EM Panel

10

Multinomial MLE example

There are $n = 197$ animals classified into one of 4 categories:
 $Y = (y_1, y_2, y_3, y_4) = (125, 18, 20, 34)$

The probability associated with each category is given as:

$$\Theta = \left(\frac{1}{2} + \frac{1}{4}\pi, \frac{1}{4}(1-\pi), \frac{1}{4}(1-\pi), \frac{1}{4}\pi \right)$$

The resulting likelihood function for this multinomial is:

$$L(\pi) = \frac{n!}{y_1! y_2! y_3! y_4!} \left(\frac{1}{2} + \frac{1}{4}\pi \right)^{y_1} \left(\frac{1}{4}(1-\pi) \right)^{y_2} \left(\frac{1}{4}(1-\pi) \right)^{y_3} \left(\frac{1}{4}\pi \right)^{y_4}$$

EMNLP, June 2001

Ted Pedersen - EM Panel

11

Multinomial MLE example

$$\log L(\pi) = y_1 \log \left(\frac{1}{2} + \frac{1}{4}\pi \right) + y_2 \log \left(\frac{1}{4}(1-\pi) \right) + y_3 \log \left(\frac{1}{4}(1-\pi) \right) + y_4 \log \left(\frac{1}{4}\pi \right)$$

$$\frac{d \log L(\pi)}{d \pi} = \frac{y_1}{2+\pi} - \frac{y_2+y_3}{1-\pi} + \frac{y_4}{\pi} = 0$$

$$\frac{d \log L(\pi)}{d \pi} = \frac{125}{2+\pi} - \frac{38}{1-\pi} + \frac{34}{\pi} = 0 \Rightarrow \pi = 0.627$$

EMNLP, June 2001

Ted Pedersen - EM Panel

12

Multinomial MLE runs aground?

- Adversity strikes! The observed data is incomplete. There are really 5 categories.
- y_1 is the composite of 2 categories (x_1+x_2)
 - $p(y_1) = \frac{1}{2} + \frac{1}{4} * \pi$, $p(x_1) = \frac{1}{2}$, $p(x_2) = \frac{1}{4} * \pi$
- How can we make a MLE, since we can't observe category counts x_1 and x_2 ?!
 - Unobserved sufficient statistics!?

EM triumphs over adversity!

- E-STEP: Find the expected values of the sufficient statistics for the complete data X , given the incomplete data Y and the current parameter estimates
- M-STEP: Use those sufficient statistics to make a MLE as usual!

MLE for complete data

$X = (x_1, x_2, x_3, x_4, x_5) = (x_1, x_2, 18, 20, 34)$ where $x_1 + x_2 = 125$

$$\Theta = \left(\frac{1}{2}, \frac{1}{4}\pi, \frac{1}{4}(1-\pi), \frac{1}{4}(1-\pi), \frac{1}{4}\pi \right)$$

$$L(\pi) = \frac{n!}{x_1!x_2!x_3!x_4!x_5!} * \left(\frac{1}{2}\right)^{x_1} * \left(\frac{1}{4}\pi\right)^{x_2} * \left(\frac{1}{4}(1-\pi)\right)^{x_3} * \left(\frac{1}{4}(1-\pi)\right)^{x_4} * \left(\frac{1}{4}\pi\right)^{x_5}$$

MLE for complete data

$$\log L(\pi) = x_2 * \log\left(\frac{1}{4}\pi\right) + x_3 * \log\left(\frac{1}{4}(1-\pi)\right) + x_4 * \log\left(\frac{1}{4}(1-\pi)\right) + x_5 * \log\left(\frac{1}{4}\pi\right)$$

$$\frac{d \log L(\pi)}{d\pi} = \frac{x_2 + x_5}{\pi} - \frac{x_3 + x_4}{1-\pi} = 0$$

$$\frac{d \log L(\pi)}{d\pi} = \frac{\mathbf{x_2} + 34}{\pi} - \frac{38}{1-\pi} = 0$$

E-step

- What are the sufficient statistics?
 - $X_1 \Rightarrow X_2 = 125 - x_1$
- How can their expected value be computed?
 - $E[x_1 | y_1] = n * p(x_1)$
- The unobserved counts x_1 and x_2 are the categories of a binomial distribution with a sample size of 125.
 - $p(x_1) + p(x_2) = p(y_1) = \frac{1}{2} + \frac{1}{4} * \pi$

E-Step

- $E[x_1 | y_1] = n * p(x_1)$
 - $p(x_1) = \frac{1}{2} / (\frac{1}{2} + \frac{1}{4} * \pi)$
- $E[x_2 | y_1] = n * p(x_2) = 125 - E[x_1 | y_1]$
 - $p(x_2) = \frac{1}{4} * \pi / (\frac{1}{2} + \frac{1}{4} * \pi)$
- Iteration 1? Start with $\pi = 0.5$ (this is just a random guess...)

E-Step Iteration 1

- $E[x_1|y_1] = 125 * (\frac{1}{2} / (\frac{1}{2} + \frac{1}{4} * 0.5)) = 100$
- $E[x_2|y_1] = 125 - 100 = 25$
- These are the expected values of the sufficient statistics, given the observed data and current parameter estimate (which was just a guess)

M-Step iteration 1

- Given sufficient statistics, make MLEs as usual

$$\frac{d \log L(\pi)}{d\pi} = \frac{x_2}{\pi} + \frac{34}{\pi} - \frac{38}{1-\pi} = 0$$

$$\frac{25}{\pi} + \frac{34}{\pi} - \frac{38}{1-\pi} = 0$$

$$\pi = .608$$

E-Step Iteration 2

- $E[x_1|y_1] = 125 * (\frac{1}{2} / (\frac{1}{2} + \frac{1}{4} * 0.608)) = 95.86$
- $E[x_2|y_1] = 125 - 95.86 = 29.14$
- These are the expected values of the sufficient statistics, given the observed data and current parameter estimate (from iteration 1)

M-Step iteration 2

- Given sufficient statistics, make MLEs as usual

$$\frac{d \log L(\pi)}{d\pi} = \frac{x_2}{\pi} + \frac{34}{\pi} - \frac{38}{1-\pi} = 0$$

$$\frac{29.14}{\pi} + \frac{34}{\pi} - \frac{38}{1-\pi} = 0$$

$$\pi = .624$$

Result?

- Converge in 4 iterations to $\pi = .627$
 - $E[x_1|y_1] = 95.2$
 - $E[x_2|y_1] = 29.8$

Conclusion

- Distribution must be appropriate to problem
- Sufficient statistics should be identifiable and have computable expected values
- Maximization operation should be possible
- Initialization should be good or lucky to avoid saddle points and local maxima
- Then...it might be safe to proceed...