

# The EM Algorithm

Michael Collins

In fulfillment of the Written Preliminary Exam II requirement, September 1997

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Preliminaries</b>	<b>4</b>
2.1	Notation . . . . .	4
2.2	Maximum-likelihood Estimation . . . . .	4
2.2.1	An example . . . . .	4
2.3	Sufficient Statistics . . . . .	5
2.3.1	An example . . . . .	5
2.4	Exponential Families . . . . .	5
2.4.1	An example: the normal distribution . . . . .	6
2.4.2	Other important properties . . . . .	6
<b>3</b>	<b>The EM algorithm</b>	<b>7</b>
3.1	An example . . . . .	8
3.2	Proof that $L(\Theta)$ is non-decreasing at each iteration . . . . .	9
3.2.1	Proof of equation 25 . . . . .	11
3.2.2	Proof of equation 26 . . . . .	11
3.3	Proof that $L(\Theta)$ is <b>increasing</b> if $\Theta$ is not a stationary point of $L$ . . . . .	11
3.4	Generalised EM (GEM) algorithms . . . . .	12
3.5	Special Cases of the EM Algorithm . . . . .	12
3.5.1	Exponential Families . . . . .	13
3.5.2	Algebraic Models . . . . .	14
3.6	Summary of the 4 Theorems in <b>DLR</b> . . . . .	16
3.6.1	Theorem 1 . . . . .	16
3.6.2	Theorems 2 and 3 . . . . .	16
3.6.3	Theorem 4 . . . . .	17
<b>4</b>	<b>(Wu 83)'s Commentary on the EM algorithm</b>	<b>18</b>
4.1	Is $L^*$ a global maximum, local maximum or stationary value? . . . . .	19
4.1.1	Theorem 1 . . . . .	19
4.1.2	Theorem 2 . . . . .	19
4.1.3	Theorem 3 . . . . .	20
4.1.4	Summary of Theorems 1, 2 and 3 . . . . .	20
4.1.5	Example of Convergence to a Saddle Point . . . . .	20
4.1.6	Proof of Theorem 1 . . . . .	20
4.1.7	Corollary 1 . . . . .	21
4.2	Does $\Theta$ Converge to a point $\Theta^*$ ? . . . . .	21
4.2.1	Theorem 4 . . . . .	21
4.2.2	Theorem 5 . . . . .	21
4.3	The Non-convergent GEM Algorithm given in (Boyles 83) . . . . .	22

<b>5</b>	<b>(Jamshidian and Jennrich 93)</b>	<b>22</b>
5.1	Optimisation of Quadratic Functions . . . . .	22
5.1.1	Conjugate Gradient Methods . . . . .	23
5.1.2	Generalised Conjugate Gradient Methods . . . . .	24
5.2	Accelerating EM using Generalised Conjugate Gradients . . . . .	24
5.3	Discussion . . . . .	24
<b>6</b>	<b>Conclusions</b>	<b>25</b>

# 1 Introduction

The Expectation Maximization (EM) algorithm is a parameter estimation method which falls into the general framework of maximum-likelihood estimation, and is applied in cases where part of the data can be considered to be incomplete, or “hidden”. It is essentially an iterative optimisation algorithm which, at least under certain conditions, will converge to parameter values at a local maximum of the likelihood function. There are many statistical models which turn out to be special cases of EM, for example: Hidden Markov Models (HMMs) (Baum 71); the generalisation of HMMs to Stochastic Context-Free Grammars (Baker 79); mixture models; and estimation in cases of missing data.

(Dempster, Laird and Rubin) (from here on referred to as **DLR**) defined the EM algorithm, and proved certain properties, in particular that at each iteration the log-likelihood of the observed data is guaranteed to be non-decreasing. That is, if  $L(\Theta)$  is the likelihood of the observed data given parameter values  $\Theta$ , and  $\Theta_t, \Theta_{t+1}$  are the parameter values at the  $t$ 'th and  $t+1$ 'th iterations respectively, then  $L(\Theta_{t+1}) \geq L(\Theta_t)$ . They also defined Generalised EM (GEM) algorithms, which include EM as a special case, and can be more computationally efficient, while still guaranteeing that  $L(\Theta_{t+1}) \geq L(\Theta_t)$ .

(Wu 1983) addressed two issues:

1. Given that  $L$  converges to some value  $L^*$ , then is  $L^*$  a global maximum, local maximum, saddle point or some other point? It is well known that  $L^*$  can not, in general, be guaranteed to be a global maximum.  $L(\Theta_{t+1}) \geq L(\Theta_t)$  is one condition for convergence to a stationary point of  $L$ , (Wu 83) defines additional conditions for convergence of an EM or GEM algorithm to a stationary point. At least for EM algorithms, these conditions are quite mild. He also gave a condition for convergence to a local maximum as opposed to a saddle point, but this condition is difficult to verify in practice (and does not hold in many practical applications).
2. Under what conditions do the parameter estimates  $\Theta$  also converge to some point  $\Theta^*$ ? Convergence of  $L$  to a point  $L^*$  does not guarantee convergence of the parameter estimates to some  $\Theta^*$ , particularly if there is more than one point  $\Theta$  satisfying  $L(\Theta) = L^*$ .

(JJ 93) emphasise that EM is an optimisation algorithm for  $L$ , and show that it is approximately a steepest descent algorithm, an optimisation method which often converges slowly. They show that with a relatively minor increase in complexity the EM algorithm can be modified to a conjugate-gradient descent method, which is known to be an improved optimisation algorithm. They give

experimental results showing that their algorithm typically converges around 3-10 times faster than standard EM, and can in some cases be 25-100 times faster.

The remainder of this paper gives some background about maximum-likelihood estimation in section 2; considers the major results of **DLR**, (Wu 83) and (JJ 77) in sections 3, 4 and 5; and concludes in section 6. For a summary of the major points of this paper the reader should refer at this point to the bullet points in section 6.

## 2 Preliminaries

Most of the results in this section are taken from [BD 77].

### 2.1 Notation

We use bold-face throughout to denote matrices, normal typeface to denote scalars. Given a vector  $\mathbf{X}$ , we write its  $i$ 'th component as  $X_i$ . We use the  $\mathbf{D}$  operator to denote differentiation. Where there is ambiguity regarding which variable differentiation is with respect to, we use superscripts on the  $\mathbf{D}$  operator. For example,  $\mathbf{D}^{10}Q(\boldsymbol{\Theta}_1, \boldsymbol{\Theta}_2)$  is the first derivative of  $Q$  w.r.t.  $\boldsymbol{\Theta}_1$ ,  $\mathbf{D}^{01}Q(\boldsymbol{\Theta}_1, \boldsymbol{\Theta}_2)$  is the first derivative w.r.t.  $\boldsymbol{\Theta}_2$ .

### 2.2 Maximum-likelihood Estimation

In general we have

- a *sample*  $\mathbf{X} = \{X_1, X_2, \dots, X_n\}$  where each  $X_i$  is a random variable (a single value, or vector of values).
- A vector of *parameters*  $\boldsymbol{\Theta}$  such that we can define the *likelihood* of the data  $P(\mathbf{X}|\boldsymbol{\Theta})$ . We can also define the *log-likelihood*  $L(\mathbf{X}|\boldsymbol{\Theta}) = \log P(\mathbf{X}|\boldsymbol{\Theta})$ . Often the  $X_i$ s are independently identically distributed (i.i.d.) so that  $L(\mathbf{X}|\boldsymbol{\Theta}) = \sum_{i=1 \dots n} \log P(X_i|\boldsymbol{\Theta})$ .

If  $\Omega$  is the parameter space, maximum-likelihood (ML) estimation involves setting the ML estimate  $\boldsymbol{\Theta}_{ML}$  such that

$$\boldsymbol{\Theta}_{ML} = \arg \max_{\boldsymbol{\Theta} \in \Omega} L(\mathbf{X}|\boldsymbol{\Theta}) \quad (1)$$

#### 2.2.1 An example

Suppose we toss a coin 6 times, and  $X_i = 1$  if the  $i$ 'th toss is heads, 0 if it is tails. Say our sample  $\mathbf{x} = \{1, 0, 0, 0, 1, 0\}$ . Assume the coin has a probability  $p$  of being heads,  $1 - p$  of being tails, so that  $\boldsymbol{\Theta} = p$ . Then

$$\begin{aligned} L(\mathbf{X} = \mathbf{x}|\boldsymbol{\Theta}) &= \sum_{i=1}^n \log(P(X_i = x_i|p)) \\ &= 2 \log p + 4 \log(1 - p) \end{aligned} \quad (2)$$

We can maximize  $\mathbf{L}$  by setting the derivative w.r.t.  $p$  equal to 0:

$$\frac{d L(\mathbf{X} = \mathbf{x}|\Theta)}{d p} = \frac{2}{p} - \frac{4}{1-p} = 0 \quad (3)$$

Solving this gives  $p = \frac{2}{6}$ , which is the “intuitive” estimate for  $p$ , the proportion of heads which have been seen in the sample.

Another common example of maximum-likelihood estimation is when the components of  $\mathbf{X}$  are drawn i.i.d. from a normal distribution with unknown mean  $\mu$  and known variance  $\sigma^2$ . It’s simple enough to prove that the ML estimate for  $\mu$  is  $\frac{\sum X_i}{n}$ , i.e., the sample mean.

### 2.3 Sufficient Statistics

A statistic  $\mathbf{T}(\mathbf{X})$  is any real or vector-valued function of the data  $\mathbf{X}$ . Note that if  $\mathbf{T}(\mathbf{X}_1) = \mathbf{T}(\mathbf{X}_2)$  for two samples  $\mathbf{X}_1$  and  $\mathbf{X}_2$  such that  $\mathbf{X}_1 \neq \mathbf{X}_2$  then  $\mathbf{T}$  reduces the data, by mapping different samples to the same value.  $\mathbf{T}$  is *sufficient* if there are functions  $g(\mathbf{T}(\mathbf{X}), \Theta)$  and  $h(\mathbf{X})$  s.t.

$$P(\mathbf{X}|\Theta) = g(\mathbf{T}(\mathbf{X}), \Theta)h(\mathbf{X}) \quad (4)$$

Typically,  $g(\mathbf{T}(\mathbf{X}), \Theta) = P(\mathbf{T}(\mathbf{X})|\Theta)$  and  $h(\mathbf{X}) = P(\mathbf{X}|\mathbf{T}(\mathbf{X}))$ . The crucial point is that when maximizing  $P(\mathbf{X}|\Theta)$  w.r.t.  $\Theta$  we can simply maximize  $g(\mathbf{T}(\mathbf{X}), \Theta)$ , so the sufficient statistics *summarise* the data – for ML estimation, once we know  $\mathbf{T}$  we don’t need to know anything else about the data.

#### 2.3.1 An example

For the coin-tossing example, if the sample size is  $n$  and the number of heads in the sample is  $N_h$ , then

$$P(\mathbf{X}|\Theta) = p^{N_h}(1-p)^{(n-N_h)} \quad (5)$$

So  $\mathbf{T} = (N_h, n)$  is sufficient.

### 2.4 Exponential Families

An important class of distributions is the exponential family, where the likelihood can be written

$$P(\mathbf{X}|\Theta) = \{\exp[\sum C_i(\Theta)T_i(\mathbf{X}) + d(\Theta) + S(\mathbf{X})]\}I_A(\mathbf{X}) \quad (6)$$

$I_A$  is the indicator function over the set  $A$ , and  $A$  cannot depend on  $\Theta$ . Note that  $\mathbf{T}(\mathbf{X}) = \{T_1(\mathbf{X}), T_2(\mathbf{X}) \dots T_n(\mathbf{X})\}$  is sufficient.

If we define the parameters  $\Theta = \{\Theta_1, \Theta_2, \dots \Theta_n\}$  such that  $C_i(\Theta) = \Theta_i$  then these are called the *natural* parameters. This can be a useful simplification, for example if when maximizing  $\mathbf{L}$  we differentiate w.r.t.  $\Theta$ , where for the natural parameters the derivative is then a simple function involving  $\mathbf{T}$ .

### 2.4.1 An example: the normal distribution

$$\begin{aligned} P(X|\Theta) &= \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(X-\mu)^2}{2\sigma^2}\right] \\ &= \exp\left[-\frac{X^2}{2\sigma^2} + \frac{\mu}{\sigma^2}X - \frac{\mu^2}{2\sigma^2} - \log \sigma\sqrt{2\pi}\right] \end{aligned} \quad (7)$$

In this case  $\mathbf{C}(\Theta) = \{-\frac{1}{2\sigma^2}, \frac{\mu}{\sigma^2}\}$ ,  $\mathbf{T}(\mathbf{X}) = \{X^2, X\}$ ,  $d(\Theta) = \frac{\mu^2}{2\sigma^2} - \log \sigma\sqrt{2\pi}$ . The natural parameters are  $\{-\frac{1}{2\sigma^2}, \frac{\mu}{\sigma^2}\}$ , being functions of the conventional parameters  $\{\mu, \sigma\}$ .

### 2.4.2 Other important properties

By noting that (by the definition of probability)

$$\int P(\mathbf{X}|\Theta) d\mathbf{X} = 1 \quad (8)$$

it is easy to show

$$d(\Theta) = -\log \int \{\exp[\sum C_i(\Theta)T_i(\mathbf{X}) + S(\mathbf{X})]\} I_A(\mathbf{X}) d\mathbf{X} \quad (9)$$

Using<sup>1</sup>  $\nabla(-\log f(\Theta)) = -\frac{f'(\Theta)}{f(\Theta)}$ , and assuming that we're using natural parameters (hence  $\nabla \sum C_i(\Theta)T_i(\mathbf{X}) = \mathbf{T}(\mathbf{X})$ )

$$\begin{aligned} d'(\Theta) &= -\frac{\nabla \int \{\exp[\sum C_i(\Theta)T_i(\mathbf{X}) + S(\mathbf{X})]\} I_A(\mathbf{X}) d\mathbf{X}}{\int \{\exp[\sum C_i(\Theta)T_i(\mathbf{X}) + S(\mathbf{X})]\} I_A(\mathbf{X}) d\mathbf{X}} \\ &= -\frac{\nabla \int \{\exp[\sum C_i(\Theta)T_i(\mathbf{X}) + S(\mathbf{X})]\} I_A(\mathbf{X}) d\mathbf{X}}{\exp[-d(\Theta)]} \\ &= -\frac{\int \mathbf{T}(\mathbf{X}) \{\exp[\sum C_i(\Theta)T_i(\mathbf{X}) + S(\mathbf{X})]\} I_A(\mathbf{X}) d\mathbf{X}}{\exp[-d(\Theta)]} \\ &= -\int \mathbf{T}(\mathbf{X}) \{\exp[\sum C_i(\Theta)T_i(\mathbf{X}) + S(\mathbf{X}) + d(\Theta)]\} I_A(\mathbf{X}) d\mathbf{X} \\ &= -\int \mathbf{T}(\mathbf{X}) P(\mathbf{X}|\Theta) d\mathbf{X} \\ &= -\mathbf{E}[\mathbf{T}(\mathbf{X})|\Theta] \end{aligned} \quad (10)$$

Now note that the log-likelihood

$$L(\mathbf{X}|\Theta) = \sum C_i(\Theta)T_i(\mathbf{X}) + d(\Theta) + S(\mathbf{X}) \quad (11)$$

---

<sup>1</sup>  $\nabla$  refers to differentiation w.r.t.  $\Theta$

So to obtain the ML estimates by differentiating w.r.t.  $\Theta$  (again, assuming natural parameters)

$$\begin{aligned} L'(\mathbf{X}|\Theta) &= \mathbf{T}(\mathbf{X}) + d'(\Theta) \\ &= \mathbf{T}(\mathbf{X}) - \mathbf{E}[\mathbf{T}(\mathbf{X})|\Theta] \end{aligned} \quad (12)$$

So setting  $\mathbf{T}(\mathbf{X}) = \mathbf{E}[\mathbf{T}(\mathbf{X})|\Theta]$  will give  $L'(\mathbf{X}|\Theta) = 0$ , and maximize the log-likelihood. For example, for a binomial distribution, the sufficient statistic  $\mathbf{T}(\mathbf{X}) = \sum X_i$  and  $\mathbf{E}[\sum \mathbf{X}_i|\mathbf{p}] = np$  where  $n$  is the sample size and  $p$  is the binomial parameter. So solving  $\sum X_i = np$  gives the ML estimate of  $p$ .

If we assume non-natural parameters, then (12) is modified to give

$$L'(\mathbf{X}|\Theta) = \frac{d\mathbf{C}(\Theta)}{d\Theta} [\mathbf{T}(\mathbf{X}) - \mathbf{E}[\mathbf{T}(\mathbf{X})|\Theta]] \quad (13)$$

Solving  $\mathbf{T}(\mathbf{X}) = \mathbf{E}[\mathbf{T}(\mathbf{X})|\Theta]$  is also a solution to (13), but this solution may not always exist – it may be necessary to also solve (13) as it stands (see section 3.5.1 for an example where  $\mathbf{T}(\mathbf{X}) = \mathbf{E}[\mathbf{T}(\mathbf{X})|\Theta]$  has no solution, but  $\frac{d\mathbf{C}(\Theta)}{d\Theta} [\mathbf{T}(\mathbf{X}) - \mathbf{E}[\mathbf{T}(\mathbf{X})|\Theta]] = 0$  does have a solution.)

### 3 The EM algorithm

The EM algorithm assumes the following problem definition: we have two sample spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , such that there is a many-one mapping  $\mathbf{Y} = f(\mathbf{X})$  from an observation  $\mathbf{X}$  in  $\mathcal{X}$  to an observation  $\mathbf{Y}$  in  $\mathcal{Y}$ . We define

$$\mathcal{X}(\mathbf{Y}) = \{\mathbf{X} : f(\mathbf{X}) = \mathbf{Y}\} \quad (14)$$

$\mathbf{X}$  is the *complete* data, and  $\mathbf{Y}$  is the *observed* data. If the distribution  $f(\mathbf{X}|\Theta)$  is well defined then the probability of  $\mathbf{Y}$  given  $\Theta$  is

$$g(\mathbf{Y}|\Theta) = \int_{\mathcal{X}(\mathbf{Y})} f(\mathbf{X}|\Theta) d\mathbf{X} \quad (15)$$

EM attempts to solve the following problem: given a sample from  $\mathbf{Y}$  is observed, but the corresponding  $\mathbf{X}$  are unobserved, or *hidden*, find the maximum-likelihood estimate  $\Theta_{ML}$  which maximizes  $L(\Theta) = \log g(\mathbf{Y}|\Theta)$ . In general,  $\log f(\mathbf{X}|\Theta)$  will have an easily-defined, analytically solvable maximum, but maximization of  $L(\Theta)$  has no analytic solution. EM is an iterative optimisation algorithm which defines a sequence of parameter settings through a mapping  $\Theta_t \rightarrow \Theta_{t+1}$  such that  $L(\Theta_{t+1}) \geq L(\Theta_t)$  with equality holding only at stationary points of  $L(\Theta)$ . Thus EM is a hill-climbing algorithm which, at least under certain conditions, will converge to a stationary point of  $L(\Theta)$ .

The mapping  $\Theta_t \rightarrow \Theta_{t+1}$  is defined in two steps:

1. The **Estimation** step. Define  $\tilde{p}(\mathbf{X}) = p(\mathbf{X}|\mathbf{Y}, \Theta_t)$ . (Note that  $\tilde{p}(\mathbf{X}) = 0$  outside  $\mathcal{X}(\mathbf{Y})$ .) Calculate

$$Q(\Theta', \Theta_t) = E[\log f(\mathbf{X}|\Theta') | \tilde{p}(\mathbf{X})] = \int \tilde{p}(\mathbf{X}) \log f(\mathbf{X}|\Theta') d\mathbf{X} \quad (16)$$

2. The **Maximization** step. Set  $\Theta_{t+1} = \arg \max_{\Theta'} Q(\Theta', \Theta_t)$ .

The intuition is as follows: if we had the complete data, we would simply estimate  $\Theta'$  to maximize  $\log f(\mathbf{X}|\Theta')$ . But with some of the complete data missing we instead maximize the *expectation* of  $\log f(\mathbf{X}|\Theta')$  given the observed data and the current value of  $\Theta$ .

### 3.1 An example

Say we observe a series of coin-tosses which we assume have been generated in the following way: a person has two coins in her pocket. Coin 1 has probability of heads =  $p_1$ , coin 2 has probability  $p_2$ . At each point she chooses coin 1 with probability  $\lambda$ , coin 2 with probability  $1 - \lambda$ , and tosses it 3 times. Thus the observed data is a sequence of triples of coin tosses, e.g.  $\mathbf{Y} = \{\langle HHH \rangle, \langle TTT \rangle, \langle HHH \rangle, \langle TTT \rangle\}$ . The complete data  $\mathbf{X}$ , if we could observe it, would additionally show the coin chosen at each step, e.g.  $\mathbf{X} = \{\langle HHH, 1 \rangle, \langle TTT, 2 \rangle, \langle HHH, 1 \rangle, \langle TTT, 2 \rangle\}$ . The parameters, all of which are to be estimated, are  $\Theta = \{\lambda, p_1, p_2\}$ .

Assume that  $\mathbf{X}$  is unobserved. Then the EM steps are as follows.

1. The **estimation** step: define  $\tilde{p}_i = P(X_i = \langle Y_i, 1 \rangle \mid Y_i, \Theta)$ , i.e. the probability of the  $i$ 'th coin being coin 1, given the observed data and the current parameter settings. If  $P_c(Y_i|p)$  is the probability of seeing  $Y_i$  given a coin with prob of heads =  $p$ , then we have

$$f_i(X_i = \langle Y_i, 1 \rangle \mid \Theta) = \lambda P_c(Y_i|p_1) \quad (17)$$

$$g_i(Y_i \mid \Theta) = \lambda P_c(Y_i|p_1) + (1 - \lambda) P_c(Y_i|p_2) \quad (18)$$

$$\tilde{p}_i = \frac{f_i(X_i = \langle Y_i, 1 \rangle \mid \Theta)}{g_i(Y_i \mid \Theta)} \quad (19)$$

$$= \frac{\lambda P_c(Y_i|p_1)}{\lambda P_c(Y_i|p_1) + (1 - \lambda) P_c(Y_i|p_2)} \quad (20)$$

$\tilde{p}_i$  is the *posterior* probability of coin 1 having generated the  $i$ 'th observation. If we define  $H_i$  as the number of heads in  $Y_i$  then  $P_c(Y_i|p) = p^{H_i} (1 - p)^{3-H_i}$ . Say  $\Theta' = \{\lambda', p'_1, p'_2\}$ . As the samples are i.i.d., we can write

$$\begin{aligned} & E [\log f(\mathbf{X}|\Theta') \mid \tilde{p}(\mathbf{X})] \\ &= \sum E [\log f_i(\mathbf{X}|\Theta') \mid \tilde{p}_i] \\ &= \sum \tilde{p}_i \log(f_i(X_i = \langle Y_i, 1 \rangle \mid \Theta')) + (1 - \tilde{p}_i) \log(f_i(X_i = \langle Y_i, 2 \rangle \mid \Theta')) \\ &= \sum \tilde{p}_i \log \lambda' P_c(Y_i|p'_1) + (1 - \tilde{p}_i) \log(1 - \lambda') P_c(Y_i|p'_2) \\ &= \sum \tilde{p}_i \log \lambda' p_1'^{H_i} (1 - p_1')^{3-H_i} + (1 - \tilde{p}_i) \log(1 - \lambda') p_2'^{H_i} (1 - p_2')^{3-H_i} \\ &= \sum \tilde{p}_i \log \lambda' + (1 - \tilde{p}_i) \log(1 - \lambda') + \tilde{p}_i \log p_1'^{H_i} (1 - p_1')^{3-H_i} + (1 - \tilde{p}_i) \log p_2'^{H_i} (1 - p_2')^{3-H_i} \end{aligned} \quad (21)$$



2. The **Maximization** step: Maximizing this function by setting the differentials w.r.t.  $\lambda'$ ,  $p'_1$  and  $p'_2$  respectively to 0 gives the following update formulae:

$$\lambda' = \frac{\sum \tilde{p}_i}{n} \quad (22)$$

$$p'_1 = \frac{\sum \frac{H_i}{3} \tilde{p}_i}{\sum \tilde{p}_i} \quad (23)$$

$$p'_2 = \frac{\sum \frac{H_i}{3} (1 - \tilde{p}_i)}{\sum (1 - \tilde{p}_i)} \quad (24)$$

These formulae have a nicely intuitive interpretation.  $\lambda$  is the average posterior probability of coin 1 having generated the  $i$ 'th sample.  $p_1$  is a weighted average over the observations of the usual ML estimate,  $\frac{H_i}{3}$ , where the weight corresponds to  $\tilde{p}_i$ , the posterior probability of coin 1 for  $Y_i$ . Similarly,  $p_2$  is a weighted average over the observations, where the weight corresponds to  $1 - \tilde{p}_i$ , the posterior probability of coin 2 generating  $Y_i$ . See tables 1, 2 and 3 for examples of the EM algorithm for this problem.

### 3.2 Proof that $L(\Theta)$ is non-decreasing at each iteration

A crucial property of the EM algorithm is that the log-likelihood  $L(\Theta) = \log g(\mathbf{Y}|\Theta)$  is non-decreasing at each iteration. Formally, if we define the EM mapping as  $\Theta_t \rightarrow \Theta_{t+1}$  where  $\Theta_{t+1} = \arg \max_{\Theta'} Q(\Theta', \Theta_t)$  then  $L(\Theta_{t+1}) \geq L(\Theta_t)$ . The proof rests on two results:

1. Define  $k(\mathbf{X}|\mathbf{Y}, \Theta)$  to be the posterior likelihood of the complete data given the data  $\mathbf{Y}$  and the parameters  $\Theta$ , so that  $k(\mathbf{X}|\mathbf{Y}, \Theta) = \frac{f(\mathbf{X}|\Theta)}{g(\mathbf{Y}|\Theta)}$ . If we define  $H(\Theta', \Theta) = E[\log k(\mathbf{X}|\mathbf{Y}, \Theta') | \tilde{p}(\mathbf{X})]$ , (as before,  $\tilde{p}(\mathbf{X}) = p(\mathbf{X}|\mathbf{Y}, \Theta)$ ), then

$$L(\Theta') = Q(\Theta', \Theta) - H(\Theta', \Theta) \quad (25)$$

- 2.

$$\forall \Theta' \quad H(\Theta', \Theta) \leq H(\Theta, \Theta) \quad (26)$$

with equality iff  $\log k(\mathbf{X}|\mathbf{Y}, \Theta') = \log k(\mathbf{X}|\mathbf{Y}, \Theta)$  almost everywhere.

Given (25),

$$L(\Theta_{t+1}) - L(\Theta_t) = \{Q(\Theta_{t+1}, \Theta_t) - Q(\Theta_t, \Theta_t)\} - \{H(\Theta_{t+1}, \Theta_t) - H(\Theta_t, \Theta_t)\} \quad (27)$$

But  $\{Q(\Theta_{t+1}, \Theta_t) - Q(\Theta_t, \Theta_t)\} \geq 0$  (by the definition of the M step), and from (26)  $\{H(\Theta_{t+1}, \Theta_t) - H(\Theta_t, \Theta_t)\} \leq 0$ , so clearly  $L(\Theta_{t+1}) - L(\Theta_t) \geq 0$ .

Iteration	$\lambda$	$p_1$	$p_2$	$\tilde{p}_1$	$\tilde{p}_2$	$\tilde{p}_3$	$\tilde{p}_4$
0	0.3000	0.3000	0.6000	0.0508	0.6967	0.0508	0.6967
1	0.3738	0.0680	0.7578	0.0004	0.9714	0.0004	0.9714
2	0.4859	0.0004	0.9722	0.0000	1.0000	0.0000	1.0000
3	0.5000	0.0000	1.0000	0.0000	1.0000	0.0000	1.0000

Table 1: The coin example for  $\mathbf{Y} = \{\langle HHH \rangle, \langle TTT \rangle, \langle HHH \rangle, \langle TTT \rangle\}$ . The solution that EM reaches is intuitively correct: the coin-tosser has two coins, one which always shows heads, the other which always shows tails, and is picking between them with equal probability ( $\lambda = 0.5$ ). The posterior probabilities  $\tilde{p}_i$  show that we are certain that coin 1 (tail-biased) generated  $Y_2$  and  $Y_4$ , whereas coin 2 generated  $Y_1$  and  $Y_3$ .

Iteration	$\lambda$	$p_1$	$p_2$	$\tilde{p}_1$	$\tilde{p}_2$	$\tilde{p}_3$	$\tilde{p}_4$	$\tilde{p}_5$
0	0.3000	0.3000	0.6000	0.0508	0.6967	0.0508	0.6967	0.0508
1	0.3092	0.0987	0.8244	0.0008	0.9837	0.0008	0.9837	0.0008
2	0.3940	0.0012	0.9893	0.0000	1.0000	0.0000	1.0000	0.0000
3	0.4000	0.0000	1.0000	0.0000	1.0000	0.0000	1.0000	0.0000

Table 2: The coin example for  $\{\langle HHH \rangle, \langle TTT \rangle, \langle HHH \rangle, \langle TTT \rangle, \langle HHH \rangle\}$ .  $\lambda$  is now 0.4, indicating that the coin-tosser has probability 0.4 of selecting the tail-biased coin.

Iteration	$\lambda$	$p_1$	$p_2$	$\tilde{p}_1$	$\tilde{p}_2$	$\tilde{p}_3$	$\tilde{p}_4$
0	0.3000	0.3000	0.6000	0.1579	0.6967	0.0508	0.6967
1	0.4005	0.0974	0.6300	0.0375	0.9065	0.0025	0.9065
2	0.4632	0.0148	0.7635	0.0014	0.9842	0.0000	0.9842
3	0.4924	0.0005	0.8205	0.0000	0.9941	0.0000	0.9941
4	0.4970	0.0000	0.8284	0.0000	0.9949	0.0000	0.9949

Table 3: The coin example for  $\mathbf{Y} = \{\langle HHT \rangle, \langle TTT \rangle, \langle HHH \rangle, \langle TTT \rangle\}$ . EM selects a tails-only coin, and a coin which is heavily heads-biased ( $p_2 = 0.8284$ ). It's certain that  $Y_1$  and  $Y_3$  were generated by coin 2, as they contain heads.  $Y_2$  and  $Y_4$  could have been generated by either coin, but coin 1 is far more likely.

### 3.2.1 Proof of equation 25

By the rules of conditional probability,

$$\begin{aligned} k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta}') &= \frac{f(\mathbf{X}|\boldsymbol{\Theta}')}{g(\mathbf{Y}|\boldsymbol{\Theta}')} \\ \log k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta}') &= \log f(\mathbf{X}|\boldsymbol{\Theta}') - \log g(\mathbf{Y}|\boldsymbol{\Theta}') \end{aligned} \quad (28)$$

We can now take expectations w.r.t.  $\tilde{p}(\mathbf{X}) = p(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta})$ :

$$\begin{aligned} E[\log k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta}') | \tilde{p}(\mathbf{X})] &= E[\log f(\mathbf{X}|\boldsymbol{\Theta}') | \tilde{p}(\mathbf{X})] - E[\log g(\mathbf{Y}|\boldsymbol{\Theta}') | \tilde{p}(\mathbf{X})] \\ &= E[\log f(\mathbf{X}|\boldsymbol{\Theta}') | \tilde{p}(\mathbf{X})] - \log g(\mathbf{Y}|\boldsymbol{\Theta}') \end{aligned} \quad (29)$$

(Note that  $E[\log g(\mathbf{Y}|\boldsymbol{\Theta}') | \tilde{p}(\mathbf{X})] = \log g(\mathbf{Y}|\boldsymbol{\Theta}')$  as  $\log g(\mathbf{Y}|\boldsymbol{\Theta})$  does not depend on  $\mathbf{X}$ .) So by the definitions of  $H$ ,  $Q$  and  $L$ ,

$$H(\boldsymbol{\Theta}', \boldsymbol{\Theta}) = Q(\boldsymbol{\Theta}', \boldsymbol{\Theta}) - L(\boldsymbol{\Theta}') \quad (30)$$

### 3.2.2 Proof of equation 26

One thing to note is that  $H(\boldsymbol{\Theta}, \boldsymbol{\Theta}) - H(\boldsymbol{\Theta}', \boldsymbol{\Theta})$  is the Kullback-Liebler distance between  $k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta})$  and  $k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta}')$ , which is known to be  $\geq 0$  with equality only if the two distributions are equal.

A formal proof is through the following theorem stated in (Rao 1e.6.6): Let  $f(\mathbf{X})$  and  $g(\mathbf{X})$  be non-negative and integrable functions, and  $S$  be the region in which  $f(\mathbf{X}) > 0$ . The theorem states that if  $\int_S (f(\mathbf{X}) - g(\mathbf{X})) d\mathbf{X} \geq 0$ , then  $\int_S f(\mathbf{X}) \log \frac{f(\mathbf{X})}{g(\mathbf{X})} d\mathbf{X} \geq 0$ .

If we put  $f(\mathbf{X}) = k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta})$  and  $g(\mathbf{X}) = k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta}')$  then clearly  $\int_S (f(\mathbf{X}) - g(\mathbf{X})) d\mathbf{X} \geq 0$ , as  $\int_S f(\mathbf{X}) d\mathbf{X} = 1$  and by the laws of probability  $\int_S g(\mathbf{X}) d\mathbf{X} \leq 1$ . Hence

$$\int_S f(\mathbf{X}) \log \frac{f(\mathbf{X})}{g(\mathbf{X})} d\mathbf{X} = \int_S k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta}) \log \frac{k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta})}{k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta}')} d\mathbf{X} \geq 0$$

But

$$\begin{aligned} H(\boldsymbol{\Theta}, \boldsymbol{\Theta}) - H(\boldsymbol{\Theta}', \boldsymbol{\Theta}) &= E[\log k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta}) | \tilde{p}(\mathbf{X})] - E[\log k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta}') | \tilde{p}(\mathbf{X})] \\ &= \int_S k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta}) \log k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta}) - \int_S k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta}) \log k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta}') \\ &= \int_S k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta}) \log \frac{k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta})}{k(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta}')} \\ &\geq 0 \end{aligned} \quad (31)$$

## 3.3 Proof that $L(\boldsymbol{\Theta})$ is increasing if $\boldsymbol{\Theta}$ is not a stationary point of $L$

The result given in **DLR**, that  $L(\boldsymbol{\Theta}_{t+1}) \geq L(\boldsymbol{\Theta}_t)$ , is not all that useful, as the likelihood could remain at the same value at any iteration: for example the trivial mapping  $\boldsymbol{\Theta}_{t+1} = \boldsymbol{\Theta}_t$  would satisfy

it with equality. (Wu 1983) proves the more useful result that if  $\mathcal{L}$  is the set of stationary points of  $L$  in  $\Omega$ , where  $\Omega$  is the space of  $\Theta$ , then  $L(\Theta_{t+1}) > L(\Theta_t)$  for any  $\Theta_t \notin \mathcal{L}$ . So this says that  $L$  will *increase* unless the algorithm has already reached a stationary point of  $L$ . This is one necessary condition for convergence to a stationary point of  $L$ .

The proof is as follows: from (25) we can write

$$L(\Theta_t) = Q(\Theta_t, \Theta_t) - H(\Theta_t, \Theta_t) \quad (32)$$

Differentiating gives

$$\mathbf{D}L(\Theta_t) = \mathbf{D}^{10} Q(\Theta_t, \Theta_t) - \mathbf{D}^{10} H(\Theta_t, \Theta_t) \quad (33)$$

From (26),  $\Theta = \Theta_t$  maximizes  $H(\Theta, \Theta_t)$ , so  $\mathbf{D}^{10} H(\Theta_t, \Theta_t) = 0$ , therefore

$$\mathbf{D}L(\Theta_t) = \mathbf{D}^{10} Q(\Theta_t, \Theta_t) \quad (34)$$

If  $\Theta_t \notin \mathcal{L}$  then  $\mathbf{D}L(\Theta_t) \neq 0$ , so  $\mathbf{D}^{10} Q(\Theta_t, \Theta_t) \neq 0$ . Therefore we cannot be at a maximum of  $Q$ , hence given that  $\Theta_{t+1}$  maximizes  $Q(\Theta, \Theta_t)$  we have

$$\forall \Theta_t \notin \mathcal{L} \quad Q(\Theta_{t+1}, \Theta_t) > Q(\Theta_t, \Theta_t) \quad (35)$$

From (25), (26) and (35) it is clear that

$$\forall \Theta_t \notin \mathcal{L} \quad L(\Theta_{t+1}) > L(\Theta_t) \quad (36)$$

### 3.4 Generalised EM (GEM) algorithms

**DLR** defined a GEM algorithm to be any iterative scheme  $\Theta_t \rightarrow \Theta_{t+1}$  such that  $Q(\Theta_{t+1}, \Theta_t) \geq Q(\Theta_t, \Theta_t)$ . The point here is that it is not necessary to maximize  $Q$  at each step, instead it is sufficient for  $Q$  to simply increase at each step to ensure that  $L(\Theta_{t+1}) \geq L(\Theta_t)$ . In some situations it is less computationally demanding to increase  $Q$  at each step rather than to maximize it. Clearly EM algorithms are a special case of GEM algorithms.

As it stands, this definition is flawed. An additional criterion is required for  $L$  to converge to a stationary point, namely

$$\forall \Theta_t \notin \mathcal{L} \quad Q(\Theta_{t+1}, \Theta_t) > Q(\Theta_t, \Theta_t) \quad (37)$$

i.e.  $Q$  must be strictly increasing if we have not reached a stationary point of  $L$ . Note that EM algorithms automatically satisfy (37) by the proof in section 3.3. (Wu 83) identified this flaw in the definition of GEM algorithms, and states additional conditions (see Theorem 1 of (Wu 83), described in section 4.1) which guarantee convergence. Without the additional conditions non-convergent algorithms such as the trivial  $\Theta_{t+1} = \Theta_t$  satisfy the GEM definition.

### 3.5 Special Cases of the EM Algorithm

The generality of the EM formulation is extremely useful, but also means that  $Q$  has to be defined for each problem, and furthermore that a method for maximizing  $Q$  must be found for each case. This section describes a few special cases of EM problems where  $Q$  has been defined and the maximization step has a simple analytic solution.

### 3.5.1 Exponential Families

Here the complete data is generated from a distribution which is a member of the exponential family, that is  $f(\mathbf{X}|\boldsymbol{\Theta}) = \{\exp[\sum C_i(\boldsymbol{\Theta})T_i(\mathbf{X}) + d(\boldsymbol{\Theta}) + S(\mathbf{X})]\}I_A(\mathbf{X})$ . **DLR** show that in this case the following is an EM algorithm:

1. **Expectation** step. As before, define  $\tilde{p}(\mathbf{X}) = p(\mathbf{X}|\mathbf{Y}, \boldsymbol{\Theta}_t)$ . Calculate

$$\mathbf{T}^p = E[\mathbf{T}(\mathbf{X}) | \tilde{p}(\mathbf{X})] \quad (38)$$

2. **Maximization** step. Find  $\boldsymbol{\Theta}'$  such that

$$E[\mathbf{T}(\mathbf{X})|\boldsymbol{\Theta}'] = \mathbf{T}^p \quad (39)$$

Note that (39) is very similar to the usual ML solution for exponential families, setting  $\mathbf{T}(\mathbf{X}) = \mathbf{E}[\mathbf{T}(\mathbf{X})|\boldsymbol{\Theta}']$  (see section 2.4.2), except we set  $\mathbf{E}[\mathbf{T}(\mathbf{X})|\boldsymbol{\Theta}']$  to be the *expected value of the sufficient statistics* given the observed data and the current  $\boldsymbol{\Theta}$ , instead of calculating the sufficient statistics from the complete data which is unobserved.

The proof that this procedure maximizes  $Q$  is as follows:

$$\begin{aligned} Q(\boldsymbol{\Theta}', \boldsymbol{\Theta}) &= \int \tilde{p}(\mathbf{X}) \log f(\mathbf{X}|\boldsymbol{\Theta}') d\mathbf{X} \\ &= \int \tilde{p}(\mathbf{X}) \left[ \sum C_i(\boldsymbol{\Theta}')T_i(\mathbf{X}) + d(\boldsymbol{\Theta}') + S(\mathbf{X}) \right] d\mathbf{X} \\ \mathbf{D}^{10}Q(\boldsymbol{\Theta}', \boldsymbol{\Theta}) &= \int \tilde{p}(\mathbf{X}) \left[ \frac{d\mathbf{C}(\boldsymbol{\Theta}')}{d\boldsymbol{\Theta}} T(\mathbf{X}) + d'(\boldsymbol{\Theta}') \right] d\mathbf{X} \\ &= E \left[ \frac{d\mathbf{C}(\boldsymbol{\Theta}')}{d\boldsymbol{\Theta}} T(\mathbf{X}) | \tilde{p}(\mathbf{X}) \right] + E[d'(\boldsymbol{\Theta}') | \tilde{p}(\mathbf{X})] \\ &= \frac{d\mathbf{C}(\boldsymbol{\Theta}')}{d\boldsymbol{\Theta}} E[T(\mathbf{X}) | \tilde{p}(\mathbf{X})] + d'(\boldsymbol{\Theta}') \\ &= \frac{d\mathbf{C}(\boldsymbol{\Theta}')}{d\boldsymbol{\Theta}} (E[T(\mathbf{X}) | \tilde{p}(\mathbf{X})] - E[T(\mathbf{X}) | \boldsymbol{\Theta}']) \end{aligned} \quad (40)$$

So (38), (39) give  $\mathbf{D}^{10}Q(\boldsymbol{\Theta}', \boldsymbol{\Theta}) = 0$ , and therefore maximize  $Q$ , the required result.

A solution to (39) may not actually exist, and a more general formulation of the maximization step is to set

$$\frac{d\mathbf{C}(\boldsymbol{\Theta}')}{d\boldsymbol{\Theta}} [\mathbf{T}^p - E[\mathbf{T}(\mathbf{X})|\boldsymbol{\Theta}']] = 0 \quad (41)$$

**An example.** **DLR** present an initial motivating example for EM. Say the complete data  $\mathbf{X} = (x_1, x_2, x_3, x_4, x_5)$  is drawn from a multinomial distribution  $(\frac{1}{2}, \frac{\pi}{4}, \frac{1-\pi}{4}, \frac{1-\pi}{4}, \frac{\pi}{4})$ . The observed data  $\mathbf{Y} = (y_1, y_2, y_3, y_4) = (125, 18, 20, 34)$ , and  $y_1 = x_1 + x_2$ ,  $y_2 = x_3$ ,  $y_3 = x_4$ ,  $y_4 = x_5$ . Thus  $x_1$  and  $x_2$  are hidden, while  $x_1 + x_2$  is observed. The log-likelihood of the complete data can be written

$$\log f(\mathbf{X} | \pi) = x_1 \log \frac{1}{2} + x_2 \log \frac{\pi}{4} + x_3 \log \frac{1-\pi}{4} + x_4 \log \frac{1-\pi}{4} + x_5 \log \frac{\pi}{4} + S(\mathbf{X}) \quad (42)$$

where  $S(\mathbf{X})$  is a multinomial coefficient. This is then an exponential distribution, with sufficient statistics  $\mathbf{T}(\mathbf{X}) = (x_1, x_2, x_3, x_4, x_5)$  and  $\mathbf{C}(\boldsymbol{\Theta}) = (\log \frac{1}{2}, \log \frac{\pi}{4}, \log \frac{1-\pi}{4}, \log \frac{1-\pi}{4}, \frac{\pi}{4})$ . While  $x_3, x_4, x_5$  are observed, the expectation step involves estimating  $x_1^p$  and  $x_2^p$  given the current parameter settings and the observed data  $\mathbf{Y}$ . This gives

$$x_1^p = E[x_1 | \mathbf{Y}, \pi] = 125 \frac{\frac{1}{2}}{\frac{1}{2} + \frac{\pi}{4}} \quad (43)$$

$$x_2^p = E[x_2 | \mathbf{Y}, \pi] = 125 \frac{\frac{\pi}{4}}{\frac{1}{2} + \frac{\pi}{4}} \quad (44)$$

Then  $E[\mathbf{T}(\mathbf{X}) | \mathbf{Y}, \pi] = (x_1^p, x_2^p, x_3, x_4, x_5)$ . We can also calculate  $E[\mathbf{T}(\mathbf{X}) | \pi] = n(\frac{1}{2}, \frac{\pi}{4}, \frac{1-\pi}{4}, \frac{1-\pi}{4}, \frac{\pi}{4})$  where  $n = y_1 + y_2 + y_3 + y_4$ . **DLR**'s maximization step (39) gives

$$E[\mathbf{T}(\mathbf{X}) | \pi'] = E[\mathbf{T}(\mathbf{X}) | \mathbf{Y}, \pi] \\ n(\frac{1}{2}, \frac{\pi'}{4}, \frac{1-\pi'}{4}, \frac{1-\pi'}{4}, \frac{\pi'}{4}) = (x_1^p, x_2^p, x_3, x_4, x_5) \quad (45)$$

where  $\pi'$  is the new value for  $\pi$ . However this does not have a solution for  $\pi'$ , so it fails as a maximization step (for example, it requires  $x_3 = x_4 = n \frac{1-\pi'}{4}$ , but  $x_3 = 18, x_4 = 20$ ). If we instead use the more general formula in (41), noting that  $\frac{d\mathbf{C}(\boldsymbol{\Theta})}{d\boldsymbol{\Theta}} = (0, \frac{1}{\pi}, \frac{-1}{1-\pi}, \frac{-1}{1-\pi}, \frac{1}{\pi})$ , then the maximization step becomes

$$(0, \frac{1}{\pi'}, \frac{-1}{1-\pi'}, \frac{-1}{1-\pi'}, \frac{1}{\pi'}) \times \{n(\frac{1}{2}, \frac{\pi'}{4}, \frac{1-\pi'}{4}, \frac{1-\pi'}{4}, \frac{\pi'}{4}) - (x_1^p, x_2^p, x_3, x_4, x_5)\}^T = 0 \\ \frac{x_2^p + x_5}{\pi'} - \frac{x_3 + x_4}{1-\pi'} = 0 \\ \pi' = \frac{x_2^p + x_5}{x_2^p + x_3 + x_4 + x_5} \quad (46)$$

This is precisely the solution given in equation (1.5) of **DLR**.

### 3.5.2 Algebraic Models

(Lafferty) describes an important class of EM problem — EM applied to algebraic models. Say  $\boldsymbol{\Theta} = \{p_1, p_2, \dots, p_n\}$  is the combination of  $m$  multinomial distributions  $\Omega_1, \Omega_2, \dots, \Omega_m$  such that the  $\Omega$ s are disjoint subsets (forming a partition of  $\{1, 2, 3, \dots, n\}$ ) of the integers  $\{1, 2, \dots, n\}$  and  $\{p_i : i \in \Omega_j\}$  are the parameters of the  $j$ 'th multinomial, so that  $\sum_{i \in \Omega_j} p_i = 1$ . The probability of the complete data can be written

$$f(\mathbf{X} | \boldsymbol{\Theta}) = \prod_{i \in \Omega_1} p_i^{C(i, \mathbf{X})} \prod_{i \in \Omega_2} p_i^{C(i, \mathbf{X})} \dots \prod_{i \in \Omega_m} p_i^{C(i, \mathbf{X})} \quad (47)$$

where  $C(i, \mathbf{X})$  is the count of the event in  $\mathbf{X}$  which corresponds to  $p_i$  —  $C(i, \mathbf{X})$  are the sufficient statistics for  $\mathbf{X}$ . In this case, if  $\boldsymbol{\Theta}' = \{p_1, p_2, \dots, p_n\}$  then the  $Q$  function is:

$$Q(\boldsymbol{\Theta}', \boldsymbol{\Theta}) = \sum_{i \in \Omega_1} \tilde{p}(\mathbf{X}) \sum C(i, \mathbf{X}) \log p_i + \sum_{i \in \Omega_2} \tilde{p}(\mathbf{X}) \sum C(i, \mathbf{X}) \log p_i + \dots \sum_{i \in \Omega_m} \tilde{p}(\mathbf{X}) \sum C(i, \mathbf{X}) \log p_i \quad (48)$$

We can maximize  $Q$  by maximizing each of these double sums separately. This gives  $m$  constrained optimisation problems, for example maximize

$$Q_1(\Theta', \Theta) = \sum_{i \in \Omega_1} \tilde{p}(\mathbf{X}) \sum C(i, \mathbf{X}) \log p_i \quad (49)$$

subject to the constraint

$$\sum_{i \in \Omega_1} p_i = 1 \quad (50)$$

Using Lagrange multipliers, the unconstrained problem is to maximize

$$\hat{Q}_1(\Theta', \Theta) = \sum_{i \in \Omega_1} \tilde{p}(\mathbf{X}) \sum C(i, \mathbf{X}) \log p_i - \lambda \sum_{i \in \Omega_1} p_i \quad (51)$$

Taking partial derivatives w.r.t.  $p_j$  and setting them to 0 gives

$$\begin{aligned} \frac{d\hat{Q}_1(\Theta', \Theta)}{dp_j} &= \sum \frac{\tilde{p}(\mathbf{X}) C(j, \mathbf{X})}{p_j} - \lambda = 0 \\ \Rightarrow p_j &= \frac{\sum \tilde{p}(\mathbf{X}) C(j, \mathbf{X})}{\lambda} \end{aligned} \quad (52)$$

If we define  $\tilde{C}(j, \mathbf{X}) = \sum \tilde{p}(\mathbf{X}) C(j, \mathbf{X})$  and find  $\lambda$  s.t.  $\sum_{i \in \Omega_1} p_i = 1$ , then

$$p_j = \frac{\tilde{C}(j, \mathbf{X})}{\sum_{i \in \Omega_1} \tilde{C}(i, \mathbf{X})} \quad (53)$$

$\tilde{C}(j, \mathbf{X})$  can be interpreted as the *expected* count corresponding to parameter  $p_j$ , and  $p_j$  is then the normalised expected count.

**Examples of algebraic models** Hidden markov models are an important class of algebraic model. An EM algorithm for HMMs was first suggested by [Baum 71]. An HMM has  $n$  states  $\{s_1, s_2, \dots, s_n\}$ . There are  $n \times n$  transition probabilities,  $P(s_i \rightarrow s_j)$  such that  $\sum_{j=1..n} P(s_i \rightarrow s_j) = 1$ . Initial and final states are defined, without loss of generality we take the initial state to be  $s_1$ , the final state to be  $s_n$ . Given an output alphabet  $\Sigma$  s.t.  $|\Sigma| = m$ , there are also  $n \times m$  emission probabilities  $p(i \uparrow j)$ , where this is the probability of state  $s_i$  emitting symbol  $j$ . The observed data is a sequence of symbols from  $\Sigma$ , and the complete data is this sequence together with the underlying sequence of states which generated this data.

Say we observe an output sequence  $\{o_1, o_2, \dots, o_l\}$ , and the state sequence is  $\{q_1, q_2, \dots, q_l\}$  (excluding the initial and final states). Then the probability of the complete data is

$$\begin{aligned} &f(\{o_1, o_2, \dots, o_l\}, \{q_1, q_2, \dots, q_l\} \mid \Theta) = \\ &\prod p(s_1 \rightarrow q_1) p(q_1 \rightarrow q_2) p(q_2 \rightarrow q_3) \dots p(q_l \rightarrow s_n) \times \prod (q_1 \uparrow o_1)(q_2 \uparrow o_2) \dots (q_l \uparrow o_l) \end{aligned} \quad (54)$$

If  $C(s_i \rightarrow s_j)$  is the number of times we see  $p(s_i \rightarrow s_j)$  in the first product, and  $C(s_i \uparrow j)$  is the number of times we see  $p(s_i \uparrow j)$  in the second product, then we can re-write this as

$$\begin{aligned} & f(\{o_1, o_2, \dots, o_l\}, \{q_1, q_2, \dots, q_l\} \mid \Theta) = \\ & \prod_{i=1..n} \prod_{j=1..n} p(s_i \rightarrow s_j)^{C(s_i \rightarrow s_j)} \times \prod_{i=1..n} \prod_{j=1..m} p(s_i \uparrow j)^{C(s_i \uparrow j)} \end{aligned} \quad (55)$$

Thus this defines an algebraic model with  $2n$  multinomial distributions. We can define the EM algorithm for HMMs using the general method shown above. The result is a set of expected counts,  $\tilde{C}(s_i \rightarrow s_j)$  and  $\tilde{C}(s_i \uparrow j)$ , with

$$p(s_i \rightarrow s_j) = \frac{\tilde{C}(s_i \rightarrow s_j)}{\sum_{k=1..n} \tilde{C}(s_i \rightarrow s_k)} \quad p(s_i \uparrow j) = \frac{\tilde{C}(s_i \uparrow j)}{\sum_{k=1..m} \tilde{C}(s_i \uparrow k)} \quad (56)$$

Note that  $\sum_{k=1..n} \tilde{C}(s_i \rightarrow s_k) = \sum_{k=1..m} \tilde{C}(s_i \uparrow k)$  is the expected number of times the model was in state  $s_i$  given the observed data and the current parameter values. A naive algorithm for this problem would be the following: for every possible state sequence, calculate  $f(\{o_1, o_2, \dots, o_l\}, \{q_1, q_2, \dots, q_l\} \mid \Theta)$ , and the sufficient statistics for  $\mathbf{X} = \{\{o_1, o_2, \dots, o_l\}, \{q_1, q_2, \dots, q_l\}\}$ . From  $f$  calculate the marginal  $k(\{q_1, q_2, \dots, q_l\} \mid \{o_1, o_2, \dots, o_l\}, \Theta)$ . From the marginal probabilities and sufficient statistics for each state sequence, calculate the expected counts and from these the parameter values. This algorithm is unworkable though, given that the number of state sequences  $n^l$  is exponential. Fortunately the forward-backward algorithm given in [Baum 71] gives a dynamic programming algorithm for calculation of the expected counts, which runs in  $\mathcal{O}(n^2 l)$  time.

### 3.6 Summary of the 4 Theorems in DLR

In the following section we take  $M(\Theta)$  to define a GEM algorithm, that is  $Q(M(\Theta), \Theta) \geq Q(\Theta, \Theta)$  for all  $\Theta$ . **DLR** states 4 central theorems with respect to GEM algorithms.

#### 3.6.1 Theorem 1

For every GEM algorithm,  $L(M(\Theta)) \geq L(\Theta)$ , with equality if and only if both  $Q(M(\Theta), \Theta) = Q(\Theta, \Theta)$  and  $k(\mathbf{X} \mid \mathbf{Y}, M(\Theta)) = k(\mathbf{X} \mid \mathbf{Y}, \Theta)$  almost everywhere. The proof is in section 3.2.

**Corollary 1.** Suppose for some  $\Theta^*$ ,  $L(\Theta^*) \geq L(\Theta)$  for all  $\Theta$ , i.e.  $\Theta^*$  is a (possibly non-unique) global maximum of  $L$ . Then  $L(M(\Theta^*)) = L(\Theta^*)$ ,  $Q(M(\Theta^*), \Theta^*) = Q(\Theta^*, \Theta^*)$  and  $k(\mathbf{X} \mid \mathbf{Y}, M(\Theta^*)) = k(\mathbf{X} \mid \mathbf{Y}, \Theta^*)$  almost everywhere. Hence if EM reaches a global maximum, the likelihood remains fixed at this point.

**Corollary 2.** Suppose for some  $\Theta^*$ ,  $L(\Theta^*) > L(\Theta)$  for all  $\Theta \neq \Theta^*$ , i.e.  $\Theta^*$  is a unique global maximum of  $L$ . Then for every GEM algorithm  $M(\Theta^*) = \Theta^*$ . So if EM has reached a unique global maximum, the parameter values remain unchanged at each iteration.

#### 3.6.2 Theorems 2 and 3

Theorems 2 and 3 attempted to show that under certain conditions  $\Theta$  converges to some point  $\Theta^*$ . Note that  $L(\Theta)$  will usually converge to some value  $L^*$  ((Wu 83) defines the exact conditions



required for this to be the case), but this does not imply that  $\Theta$  also converges –  $\Theta$  could, for example, be oscillating between points nearby two local maxima with the same maximal value  $L^*$ . Unfortunately, as noted by (Wu 1983), a key step of the **DLR** proof is wrong, namely the application of the triangle inequality to go from step (3.13) to (3.14). Because of this both theorems 2 and 3 are invalid.

### 3.6.3 Theorem 4

It is useful to be able to calculate the *rate of convergence* of an EM algorithm. For the one parameter case where  $\Theta$  converges to a point  $\Theta^*$  we can define the rate of convergence  $R_{\Theta_p} = \frac{\Theta_{p+1} - \Theta^*}{\Theta_p - \Theta^*}$ . This can be interpreted as the reduction in the distance to  $\Theta^*$  when going from  $\Theta_p$  to  $\Theta_{p+1}$ , for example if  $\Theta^* = 0.5$ ,  $\Theta_p = 1.0$ ,  $\Theta_{p+1} = 0.75$  then  $R_{\Theta_p} = 0.5$ , i.e. the distance to  $\Theta^*$  is halved. It's possible to prove that  $\lim_{\Theta_p \rightarrow \Theta^*} R_{\Theta_p} = M'(\Theta^*)$  where  $M'(\Theta) = \frac{dM(\Theta)}{d\Theta}$ . So if we can calculate  $M'(\Theta^*)$ , then we have an estimate of the rate of convergence close to  $\Theta^*$ .

**Proof.** Say  $\Theta_p = \Theta^* + \delta$ . Then  $\Theta_{p+1} = M(\Theta^* + \delta)$ . Noting also that  $M(\Theta^*) = \Theta^*$ , as  $\Theta^*$  is a maximum of  $L$ ,

$$\begin{aligned} R_{\Theta_p} &= \frac{\Theta_{p+1} - \Theta^*}{\Theta_p - \Theta^*} \\ &= \frac{M(\Theta^* + \delta) - \Theta^*}{\Theta^* + \delta - \Theta^*} \\ &= \frac{M(\Theta^* + \delta) - M(\Theta^*)}{\delta} \end{aligned} \tag{57}$$

So,

$$\lim_{\Theta_p \rightarrow \Theta^*} R_{\Theta_p} = \lim_{\delta \rightarrow 0} \frac{M(\Theta^* + \delta) - M(\Theta^*)}{\delta} = \frac{dM(\Theta)}{d\Theta} \tag{58}$$

This can be generalised to the multiparameter case:  $\mathbf{R} = \frac{d\mathbf{M}(\Theta)}{d\Theta} = \mathbf{DM}(\Theta)$  is a vector of rates of convergence, where the  $i$ 'th component of  $\mathbf{R}$  is the rate of convergence of the  $i$ 'th parameter in  $\Theta$ .

**Theorem 4 of DLR** states that

$$\mathbf{DM}(\Theta^*) = \mathbf{D}^{20}\mathbf{H}(\Theta^*, \Theta^*) [\mathbf{D}^{20}\mathbf{Q}(\Theta^*, \Theta^*)]^{-1} \tag{59}$$

under the following conditions:

1.  $\Theta^{\mathbf{P}}$  converges to some  $\Theta^*$ .
2.  $\mathbf{D}^{10}\mathbf{Q}(\Theta_{\mathbf{p}+1}, \Theta_{\mathbf{p}}) = 0$ . This means that  $Q$  is maximized at each iteration, EM algorithms satisfy this condition, as do GEM algorithms which find a stationary point of  $Q$  at each iteration.
3.  $\mathbf{D}^{20}\mathbf{Q}(\Theta_{\mathbf{p}+1}, \Theta_{\mathbf{p}})$  is negative definite with eigenvalues bounded away from 0.

**An example.** If we return to the example in section 3.5.1, we can calculate  $\mathbf{D}^{20}Q(\pi', \pi)$  and  $\mathbf{D}^{20}H(\pi', \pi)$  as follows:

From (40)

$$\begin{aligned}\mathbf{D}^{10}Q(\pi', \pi) &= \left(0, \frac{1}{\pi'}, \frac{-1}{1-\pi'}, \frac{-1}{1-\pi'}, \frac{1}{\pi'}\right) \times \left\{n\left(\frac{1}{2}, \frac{\pi'}{4}, \frac{1-\pi'}{4}, \frac{1-\pi'}{4}, \frac{\pi'}{4}\right) - (x_1^p, x_2^p, x_3, x_4, x_5)\right\}^T \\ &= \frac{x_2^p + x_5}{\pi'} - \frac{x_3 + x_4}{1-\pi'}\end{aligned}\quad (60)$$

(This is the quantity we set to 0 in the maximization step (46)). Differentiating again gives

$$\mathbf{D}^{20}Q(\pi', \pi) = -\frac{x_2^p + x_5}{\pi'^2} - \frac{x_3 + x_4}{(1-\pi')^2}\quad (61)$$

This is the value given on page 10 of **DLR**. To calculate  $\mathbf{D}^{20}H(\pi', \pi)$  we note that  $L(\pi') = Q(\pi', \pi) - H(\pi', \pi)$ , so that  $\mathbf{D}^{20}H(\pi', \pi) = \mathbf{D}^{20}Q(\pi', \pi) - \mathbf{D}^2L(\pi')$ . We have

$$\begin{aligned}L(\pi') &= y_1 \log\left(\frac{2+\pi'}{4}\right) + x_3 \log\frac{1-\pi'}{4} + x_4 \log\frac{1-\pi'}{4} + x_5 \log\frac{\pi'}{4} + S(\mathbf{X}) \\ \mathbf{D}^1L(\pi') &= \frac{y_1}{2+\pi'} - \frac{x_3+x_4}{1-\pi'} + \frac{x_5}{\pi'} \\ \mathbf{D}^2L(\pi') &= -\frac{y_1}{(2+\pi')^2} - \frac{x_3+x_4}{(1-\pi')^2} - \frac{x_5}{\pi'^2}\end{aligned}\quad (62)$$

So

$$\begin{aligned}\mathbf{D}^{20}H(\pi', \pi) &= \mathbf{D}^{20}Q(\pi', \pi) - \mathbf{D}^2L(\pi') \\ &= -\frac{x_2^p + x_5}{\pi'^2} - \frac{x_3 + x_4}{(1-\pi')^2} - \left\{-\frac{y_1}{(2+\pi')^2} - \frac{x_3+x_4}{(1-\pi')^2} - \frac{x_5}{\pi'^2}\right\} \\ &= -\frac{x_2^p}{\pi'^2} + \frac{y_1}{(2+\pi')^2}\end{aligned}\quad (63)$$

Again, this is the value given on page 10 of **DLR**.

## 4 (Wu 83)'s Commentary on the EM algorithm

(Wu 83) addresses two points concerning the EM algorithm:

1. If  $L$  converges to some  $L^*$ , what is the nature of  $L^*$ ? (A global maximum, local maximum, stationary value or other point?) He shows that in general there can only be a guarantee that  $L^*$  is a stationary value (i.e. a local/global maximum or a saddle point), and specifies conditions under which  $L^*$  falls into these categories (without these conditions,  $L^*$  could potentially be *any* value).

2. Under what conditions does  $\Theta$  converge to some  $\Theta^*$ ? Note that even if  $L$  converges,  $\Theta$  may not converge, for example it could oscillate between points on two local maxima which have the same maximum.

(Wu 83) makes the following assumptions throughout, so they can be taken as preconditions of every theorem stated in this section:

- $\Omega$  is a subset in the  $r$ -dimensional Euclidean space  $R^r$ .  
( $\Omega$  is the parameter space so  $\Theta \in \Omega$ ).

(64)

- $\Omega_{\Theta_0} = \{\Theta \in \Omega : L(\Theta) \geq L(\Theta_0)\}$  is compact for any  $L(\Theta_0) > -\infty$ .

(65)

- $L$  is continuous in  $\Theta$  and differentiable in the interior of  $\Omega$ .

(66)

As a consequence of these conditions it follows that

- $\{L(\Theta_p)\}_{p \geq 0}$  is bounded above for any  $\Theta_0 \in \Omega$ .

(67)

#### 4.1 Is $L^*$ a global maximum, local maximum or stationary value?

We define the following subsets of  $\Omega$ :

- $\mathcal{M}$  is the set of local maxima in the interior of  $\Omega$ .
- $\mathcal{L}$  is the set of stationary points in the interior of  $\Omega$ .

From this it follows that  $\mathcal{L} \setminus \mathcal{M}$  is the set of saddle points in  $\Omega$  (the set of stationary points which are not local maxima).

##### 4.1.1 Theorem 1

Let  $\{\Theta_p\}$  be a GEM sequence generated by  $\Theta_{p+1} \in M(\Theta_p)$ . Then  $L$  converges monotonically to  $L^* = L(\Theta^*)$  for some  $\Theta^* \in \mathcal{L}$  under the following conditions:

- i)  $M$  is a closed point-to-set map over the complement of  $\mathcal{L}$  (we define closed point-to-set maps in section 4.1.6 below).
- ii)  $L(\Theta_{p+1}) > L(\Theta_p)$  for all  $\Theta_p \notin \mathcal{L}$

This theorem also holds if we replace every mention of  $\mathcal{L}$  with  $\mathcal{M}$  – this gives a similar theorem but for the conditions for convergence to a local maximum rather than just any stationary point.

##### 4.1.2 Theorem 2

Section 3.3 showed that for EM algorithms (as opposed to any GEM algorithm) condition (ii) of theorem 1 holds for the  $\mathcal{L}$  (stationary value) case.  $M$  can be shown to satisfy condition (i) if  $Q(\Theta', \Theta)$  is continuous in both  $\Theta'$  and  $\Theta$ . This leads to theorem 2:

**If  $Q(\Theta', \Theta)$  is continuous in both  $\Theta'$  and  $\Theta$ , then all the limit points of an EM algorithm are stationary points of  $L$  and  $L$  converges monotonically to  $L^* = L(\Theta^*)$  for some  $\Theta^* \in \mathcal{L}$ .**

#### 4.1.3 Theorem 3

Theorem 2 guarantees convergence to a stationary value, but this stationary value could be a saddle point (a member of  $\mathcal{L} \setminus \mathcal{M}$ ). The problem is that EM satisfies condition (ii) of Theorem 1 for stationary values  $\Theta_p \notin \mathcal{L}$ , but there may be saddle points  $\Theta_s \in \mathcal{L} \setminus \mathcal{M}$  such that  $L(M(\Theta_s)) = L(\Theta_s)$ . Theorem 3 states that convergence to a *local maximum* is guaranteed if every saddle point of  $L$  is *not* a global maximum of  $Q(\Theta', \Theta)$  w.r.t  $\Theta'$ . From (34)  $DQ = 0$  at any saddle point of  $L$ , so this must mean that any saddle point of  $L$  is a saddle point or local maximum of  $Q$ , but *not* the global maximum of  $Q$ . If this condition is satisfied then (given that EM maximizes  $Q$  at each step)  $Q$  will increase even at the saddle point, and  $L$  will also increase.

#### 4.1.4 Summary of Theorems 1, 2 and 3

To summarise these theorems, when designing a GEM or EM algorithm:

- For GEM algorithms check that conditions (i) and (ii) of theorem 1 hold, and the algorithm will then converge to some point in  $\mathcal{L}$  (or  $\mathcal{M}$  for the version of theorem 1 regarding  $\mathcal{M}$ ).
- For EM algorithms check that  $Q(\Theta', \Theta)$  is continuous in both  $\Theta'$  and  $\Theta$ , then by theorem 2 the algorithm will converge to some point in  $\mathcal{L}$ . In addition, if it can be shown that every saddle point of  $L$  is *not* a global maximum of  $Q$ , then  $L$  will converge to some point in  $\mathcal{M}$ .

#### 4.1.5 Example of Convergence to a Saddle Point

If we return to the example in section 3.1, but initialise  $p_1$  and  $p_2$  to the same value, we get the behaviour in table 4.

Theorem 3 is violated for this example, and in general it's hard to guarantee this theorem's requirement. (Wu 83) suggests that it's important when running EM to try several starting points, and to randomly select initial parameter values, in particular to avoid symmetries such as  $p_1 = p_2$  in the last example. Note that if we initialise  $p_1$  and  $p_2$  even slightly differently from each other we get convergence to the global maximum, see table 5.

#### 4.1.6 Proof of Theorem 1

The proof rests on the Global Convergence Theorem, which is stated and proved in (Zangwill 69):

Say  $\{x_k\}_{k=0}^\infty$  is generated by  $x_{k+1} \in M(x_k)$ , where  $M$  is a point-to-set map on  $X$ . (A point-to-set map on  $X$  is a function from points in  $X$  to subsets of  $X$ ). Let a solution set  $\Gamma$  be given. Suppose that:

1. All points  $x_k$  are contained in a compact set  $S \subset X$ .
2.  $M$  is closed over the complement of  $\Gamma$ . "Closedness" is a generalisation to point-to-set maps from continuity of a point-to-point map; for a point-to-point map continuity implies closedness.
3. There is a continuous function  $\alpha$  on  $X$  such that

- (a) if  $x \notin \Gamma$ ,  $\alpha(y) > \alpha(x)$  for all  $y \in M(x)$ .
- (b) if  $x \in \Gamma$ ,  $\alpha(y) \geq \alpha(x)$  for all  $y \in M(x)$ .

Then all limit points of  $\{x_k\}$  are in the solution set  $\Gamma$  and  $\alpha(x_k)$  converges monotonically to  $\alpha(x)$  for some  $x \in \Gamma$ .

If we take  $M$  as a GEM algorithm,  $\alpha = L$  and  $\Gamma = \mathcal{L}$  (or  $\mathcal{M}$  for Theorem 1 for  $\mathcal{M}$ ), GEM algorithms always satisfy the following conditions:

- **Condition (1).** If  $\Theta_0$  is such that  $L(\Theta_0) > -\infty$  then, given that  $L$  is non-decreasing for a GEM algorithm, all points  $x_k$  are contained in  $\Omega_{\Theta_0}$ , which by assumption (65) is compact.
- **Condition (3b).** This comes from  $L$  being non-decreasing at each iteration of a GEM algorithm.

The remaining conditions — (2) and (3a) — are conditions (i) and (ii) of Theorem 1, hence if Theorem 1 holds then the global convergence theorem holds. Theorems 2 and 3 follow from Theorem 1.

#### 4.1.7 Corollary 1

(Wu 83) states one additional result concerning the convergence of  $L$  to  $L^*$ , corollary 1 to Theorem 6:

**Suppose  $L(\Theta)$  is continuous in  $\Omega$  with  $\Theta^*$  being the only stationary point and that  $D^{10}Q(\Theta', \Theta)$  is continuous in both  $\Theta'$  and  $\Theta$ . Then for any EM sequence  $\{\Theta_p\}$ ,  $\Theta_p$  converges to the unique maximizer  $\Theta^*$  of  $L(\Theta)$ .**

## 4.2 Does $\Theta$ Converge to a point $\Theta^*$ ?

Say  $\{\Theta_p\}$  is an instance of a GEM algorithm which satisfies theorem 1. Define  $\mathcal{L}(\alpha) = \{\Theta : L(\Theta) = \alpha\}$ . Then by theorem 1  $L$  converges to  $L^*$  and all the limit points of  $\{\Theta_p\}$  are in  $\mathcal{L}(L^*)$ . Theorems 4 and 5 of (Wu 83) then give conditions where  $\Theta$  converges to a point  $\Theta^*$ .

### 4.2.1 Theorem 4

If  $\mathcal{L}(L^*) = \{\Theta^*\}$ , that is there is only one stationary point of  $L$  at which  $L$  is  $L^*$ , then  $\Theta \rightarrow \Theta^*$ .

### 4.2.2 Theorem 5

If  $\|\Theta_{p+1} - \Theta_p\| \rightarrow 0$  as  $p \rightarrow \infty$  then all the limit points of  $\Theta_p$  are in a connected and compact subset of  $\mathcal{L}(L^*)$ . In particular, if  $\mathcal{L}(L^*)$  is discrete, then  $\Theta \rightarrow \Theta^*$  where  $\Theta^*$  is some member of  $\mathcal{L}(L^*)$ .

### 4.3 The Non-convergent GEM Algorithm given in (Boyles 83)

(Boyles 83) gives an algorithm which satisfies the GEM definition in **DLR**, and is interesting in two respects:

1. It does not converge to a stationary point of the likelihood function.
2.  $\Theta$  does not converge to some point  $\Theta^*$ .

It is useful to examine this example in the context of (Wu 83)'s theorems. The basic idea of (Boyles 83) is to define a GEM algorithm for a two parameter problem, maximization of  $L(\theta_1, \theta_2)$ . If  $y_1, y_2$  are the ML-estimates for  $\theta_1, \theta_2$ , then in their example it can be shown that all points  $\theta_1, \theta_2$  which lie on a circle of radius  $r$  centered on  $y_1, y_2$  give the same value for  $L(\theta_1, \theta_2)$ . Moreover, the lower the value of the radius  $r$  the more  $L$  increases. The GEM algorithm given then defines the mapping  $\theta_1^k, \theta_2^k \rightarrow \theta_1^{k+1}, \theta_2^{k+1}$  such that the parameter values spiral in, with the radius decreasing at each point, such that the limit of the parameter values is to rotate on the circle at  $r = 1$ . Clearly the algorithm fails to reach the maximal point at  $r = 0$ , and  $\Theta$  does not converge, instead in the limit the parameters continually traverse the circle. So which convergence criteria in (Wu83) does the algorithm fail?

First, it fails Theorem 1 condition (ii). For any  $r \leq 1$  the radius remains constant at each iteration, hence the likelihood also remains constant. But  $r = 0$  is the only stationary point of  $L$  in the problem, so Theorem 1 condition (ii) is violated.

Second, it fails the requirements of both theorems 4 and 5.  $L$  converges to a value  $L^*$  which has many values of  $\theta_1, \theta_2$  such that  $L(\theta_1, \theta_2) = L^*$ , in fact any  $\theta_1, \theta_2$  which lie on the circle of radius 1. And it is clear that  $\|\Theta_{p+1} - \Theta_p\| \rightarrow 0$  is not satisfied.

## 5 (Jamshidian and Jennrich 93)

(JJ 93) further emphasise that the EM algorithm is an optimisation algorithm, and apply a standard optimisation algorithm, generalised conjugate gradient descent.

### 5.1 Optimisation of Quadratic Functions

(Zangwill 69 Chapter 6) describes Conjugate-Gradient methods for optimisation. Taking a Taylor's expansion (page 326 of Zangwill 69) about a local maximum of  $L$  at  $\Theta^*$ , gives

$$L(\mathbf{X}|\Theta^* + \delta\Theta) \approx L(\mathbf{X}|\Theta^*) + \mathbf{D}^1\mathbf{L}(\mathbf{X}|\Theta^*)\delta\Theta + \frac{1}{2}\delta\Theta^T\mathbf{D}^2\mathbf{L}(\mathbf{X}|\Theta^*)\delta\Theta \quad (68)$$

Where  $\mathbf{D}^1\dots$  and  $\mathbf{D}^2\dots$  are the first and second derivatives of  $L$  with respect to  $\Theta$ . Taylor's theorem shows that this approximation becomes increasingly accurate as  $\delta\Theta \rightarrow 0$ . Thus we have a quadratic (second-order) approximation to  $L$  close to a maximal point.

Gradient-based optimisation algorithms draw heavily on results for optimisation of a quadratic function such as

$$f(\mathbf{Z}) = c + \mathbf{b}^T\mathbf{Z} + \frac{1}{2}\mathbf{Z}^T\mathbf{A}\mathbf{Z} \quad (69)$$

The most obvious optimisation algorithm is steepest descent: at each point calculate the gradient  $\mathbf{D}f$ , and move in this direction to the point which maximizes  $f(\mathbf{Z} + \alpha \mathbf{D}f)$  where  $\alpha$  is the distance moved in the gradient's direction.

A crucial result is that steepest descent algorithms can be very poor as optimisation algorithms for quadratic functions, whereas conjugate gradient methods maximize a quadratic in at most  $p$  steps for a  $p$ -variate quadratic. If we assume that the function being optimised is well approximated by some quadratic, then we can assume that these results carry over to the function: steepest gradient will be considerably poorer than a conjugate gradient method based on the quadratic approximation of the function.

### 5.1.1 Conjugate Gradient Methods

Given an  $n \times n$  symmetric matrix  $\mathbf{A}$ , the directions  $d_1, d_2 \dots d_n$  are said to be  $\mathbf{A}$ -conjugate if they are linearly independent, and  $\mathbf{d}_i^T \mathbf{A} \mathbf{d}_j = 0$  for  $i \neq j$ .

It can then be shown that for a quadratic function  $f(\mathbf{Z}) = c + \mathbf{b}^T \mathbf{Z} + \frac{1}{2} \mathbf{Z}^T \mathbf{A} \mathbf{Z}$ , given  $n$  conjugate directions and a starting point  $\mathbf{Z}_0$ , then the following algorithm will maximize a quadratic in  $n$  steps (i.e.  $\mathbf{Z}_n$  maximizes  $f$ ):

- For  $k = 1 \dots n$ , find the  $\alpha_k$  which maximizes  $f(\mathbf{Z}_{k-1} + \alpha_k \mathbf{d}_k)$ . Set  $\mathbf{Z}_k = \mathbf{Z}_{k-1} + \alpha_k \mathbf{d}_k$ .

All that remains, then, is to find an algorithm which constructs  $d_1, d_2 \dots d_n$  for a given quadratic. Conjugate gradient methods are remarkable in that they find  $d_1, d_2 \dots d_n$  *without knowledge of  $\mathbf{A}$* . This is important when maximizing a function using a Taylor approximation such as (68), in that there is no need to calculate the Hessian of the function (in the example,  $\mathbf{D}^2 \mathbf{L}$ ), a step which can be computationally expensive. A conjugate algorithm proceeds as follows:

- For an arbitrary starting point  $\mathbf{Z}_0$ , let  $d_1 = f'(\mathbf{Z}_0)$ , where  $f'$  is the first derivative of  $f$ .
- For  $k = 1 \dots n$ ,
  - set  $\mathbf{Z}_k = \mathbf{Z}_{k-1} + \alpha_k d_k$ , where  $\alpha_k$  maximizes  $f(\mathbf{Z}_{k-1} + \alpha_k d_k)$ .
  - Calculate  $d_{k+1} = f'(\mathbf{Z}_k) + \beta_k d_k$ .

The most commonly used values for  $\beta_k$  are (where  $g_k = f'(\mathbf{Z}_k)$ ):

- The Fletcher-Reeves version  $\beta_k = \frac{\mathbf{g}_{k+1}^T \mathbf{g}_{k+1}}{\mathbf{g}_k^T \mathbf{g}_k}$
- The Polak-Ribiere algorithm  $\beta_k = \frac{(\mathbf{g}_{k+1} - \mathbf{g}_k)^T \mathbf{g}_{k+1}}{\mathbf{g}_k^T \mathbf{g}_k}$ . When maximizing a quadratic this is the same as the Fletcher-Reeves algorithm (as  $\mathbf{g}_k^T \mathbf{g}_{k+1} = 0$ ). The difference becomes important when maximizing an approximation to a quadratic, where Polak-Ribiere updates have been shown to be better in some cases.

### 5.1.2 Generalised Conjugate Gradient Methods

A generalised conjugate gradient method is identical to the algorithm in section 5.1.1, except the gradients are modified by some positive definite matrix  $\mathbf{W}$ . So every mention of  $f'$  in section 5.1.1 is replaced by  $\mathbf{W}^{-1}f'$ . (JJ 93) state that the use of an appropriate matrix  $\mathbf{W}$  can significantly improve the performance of any optimisation algorithm which uses gradients.

## 5.2 Accelerating EM using Generalised Conjugate Gradients

(JJ 93) first show that EM is approximately a steepest descent algorithm optimising  $L$  with the generalised gradient  $\tilde{\mathbf{g}} = \mathbf{W}^{-1}\mathbf{g}(\Theta)$  where  $\mathbf{W} = (-\mathbf{D}^{20}\mathbf{Q}(\Theta', \Theta'))$ . They do this by proving that an EM step  $\Theta \rightarrow \Theta'$  is such that

$$\Theta' - \Theta \approx (-\mathbf{D}^{20}\mathbf{Q}(\Theta', \Theta'))^{-1}\mathbf{g}(\Theta) \quad (70)$$

where  $\mathbf{g}(\Theta) = \mathbf{D}^1 L(\Theta)$ . (JJ 93) then define a conjugate gradient method based on the approximate generalised gradients  $\tilde{\mathbf{g}}(\Theta) = \Theta' - \Theta$ . The  $k$ 'th iteration of this algorithm is:

1. Perform the  $k$ 'th EM step to find  $\Theta'$  and calculate  $\tilde{\mathbf{g}}_{k+1} = \Theta' - \Theta_k$ .
2. At the first and every  $p^{th}$  step thereafter, where  $p$  = the number of parameters in  $\Theta$ , set  $d_{k+1} = \tilde{\mathbf{g}}(\Theta)$ . At other steps:
  - (a) Calculate  $\mathbf{g}_k = \mathbf{D}^1 L(\Theta_k)$ .
  - (b) Set  $\beta_k = \frac{\tilde{\mathbf{g}}_{k+1}^T(\mathbf{g}_{k+1} - \mathbf{g}_k)}{\mathbf{d}_k^T(\mathbf{g}_{k+1} - \mathbf{g}_k)}$ . This is an alternative to the Fletcher-Reeves and Polak-Ribiere updates.
  - (c) Set  $\mathbf{d}_{k+1} = \tilde{\mathbf{g}}_{k+1} - \beta_k \mathbf{d}_k$ .
3. Find  $\alpha_k$ , the value of  $\alpha$  which maximizes  $L(\Theta_k + \alpha \mathbf{d}_{k+1})$ , and set  $\Theta_{k+1} = \Theta_k + \alpha_k \mathbf{d}_{k+1}$ .

It can be seen that the algorithm involves a modest increase in complexity over the EM algorithm, mainly the calculation of  $\mathbf{g}_k = \mathbf{D}^1 L(\Theta_k)$ . (JJ 93) then pick a few example applications, and find that their algorithm is a around 3-10 times faster than EM on usual examples, and can be 25-100 times faster on cases where EM converges particularly slowly.

## 5.3 Discussion

While the results in (JJ 93) are impressive, there are a number of points which are unanswered in the paper:

- They justify the (unusual) choice of  $\beta_k$  because it gives conjugate gradients even when the line searches for the  $\alpha_k$ 's aren't exact. But they neither prove this, nor provide a citation.



- They stress that the choice of  $\mathbf{W}$  can be crucial when using generalised gradients. But they never given any justification of why  $\mathbf{W} = (-\mathbf{D}^{20}\mathbf{Q}(\boldsymbol{\Theta}', \boldsymbol{\Theta}'))$  might be a good choice for  $\mathbf{W}$ . In particular they do not compare their algorithm to a standard conjugate gradient algorithm, which would also require the calculation of  $g_k$  but would not require the EM step to calculate  $\tilde{g}_k$ .
- The method presumably becomes increasingly useful as the maximum of  $L$  is approached, as the quadratic approximation becomes increasingly accurate. The method runs a few EM iterations at the start of the algorithm until change in likelihood falls below a certain threshold. No real justification is given for this threshold, nor are tests done to see how robust the method is to the choice of this parameter (and no indication is given of how much this parameter was tuned to the particular problems used).

## 6 Conclusions

The following is a summary of the major points in this paper:

- To derive an EM or GEM algorithm first find the function  $Q(\boldsymbol{\Theta}', \boldsymbol{\Theta})$  as defined in section 3. An EM mapping  $f : \boldsymbol{\Theta}_k \rightarrow \boldsymbol{\Theta}_{k+1}$  is then  $f(\boldsymbol{\Theta}) = \arg \max_{\boldsymbol{\Theta}'} Q(\boldsymbol{\Theta}', \boldsymbol{\Theta})$ . A GEM mapping is any  $f : \boldsymbol{\Theta}_k \rightarrow \boldsymbol{\Theta}_{k+1}$  such that  $Q(f(\boldsymbol{\Theta}), \boldsymbol{\Theta}) \geq Q(\boldsymbol{\Theta}, \boldsymbol{\Theta})$ .
- Section 3.2 then contains the proof, given in **DLR**, that  $L(\boldsymbol{\Theta}_{k+1}) \geq L(\boldsymbol{\Theta}_k)$  for EM and GEM algorithms. (Wu 83) also proved that for any EM algorithm,  $L(\boldsymbol{\Theta}_{k+1}) > L(\boldsymbol{\Theta}_k)$  at any non-stationary point, (see section 3.3) though this is not necessarily true for any GEM algorithm.
- Section 3.5 gives a couple of special cases – namely the exponential family of distributions, and algebraic models – where it is easy to derive  $Q$  and the parameter values which maximize  $Q$ . Note that the M-step for exponential families given in **DLR**, (39), may not have a solution, and does not have a solution for the motivating multinomial example given at the start of **DLR**. When no solution is found, the more general (41) should be used.
- Section 3.6.3 describes a method for calculating the rate of convergence of an EM algorithm, as described in **DLR** Theorem 4. The theory requires calculation of the Hessians of the  $Q$  and  $H$  functions at the local maximum.
- (Wu 83) gives further conditions for the convergence of an EM or GEM algorithm to a stationary point of  $L$ . For an EM algorithm the most important criterion is that  $Q$  is continuous in both its parameters. See section 4.1 for the exact criteria for convergence. Major points are that convergence to a global maximum can not, in general, be guaranteed, and furthermore convergence to a saddle point is possible – theorem 3 of (Wu 83) gives criteria for convergence to a local maximum rather than a saddle point, but these criteria are unlikely to apply and are hard to prove. “Common wisdom” seems to be that parameter values should be randomly initialised to avoid symmetries which can lead to convergence to saddle points, and that EM should be run a few times with different starting points.

- Section 4.2 gives the theorems provided by (Wu 83) which state conditions for the estimate of  $\Theta$  to also converge. In most applications convergence of the estimate of  $\Theta$  is probably less important than convergence of  $L$  though.
- (JJ 93) show that EM is approximately a steepest descent algorithm, and describe a fairly simple modification that gives a conjugate-gradient optimisation algorithm. This is shown to improve performance on a number of problems. The major additional complexity of the method is the calculation of  $\mathbf{DL}$  at each iteration.

In conclusion, **DLR** provided a very general framework for defining iterative algorithms which find local maxima of the log-likelihood in incomplete data problems. However, as they defined EM and GEM there was not a guarantee of convergence to a stationary point of  $L$ , (Wu 83) considerably tightened the theory by giving strict conditions for convergence. (JJ 93) show that EM, while having the advantage of simplicity, may be a poor algorithm in some situations, and used existing algorithms in the optimisation literature to improve convergence rates.

A major weakness in EM-style algorithms is the guarantee concerning the nature of the limit point — global maximum, local maximum or saddle point? While examples like that in section 4.1.5 suggest that perturbing the parameters slightly will result in the parameters diverging from the saddle point, and **DLR** page 10 suggest that this will always be the case, as far as we know there is no general theory about the “stability” of a saddle point, or a general method for perturbing the parameters in such a way that the algorithm diverges from the saddle point. Nor, as far as we know, is there any general theory about the nature of the log-likelihood surface — conditions under which it has a single global maximum, or under which it has no saddle points. While it may be very difficult to formulate theories about the general case, it seems that it would be possible for special (but extremely common) cases like HMMs or mixture models.

## References

- [Baker 79] Baker, J. (1979). Trainable Grammars for Speech Recognition. In Jared J. Wolf and Dennis H. Klatt, editors, *Speech Communication Papers for the 97th Meeting of the Acoustical Society of America*, pages 547-550, MIT, Cambridge, Mass.
- [Baum 71] Baum, L.E. (1971). An Inequality and Associated Maximization Technique in Statistical Estimation for Probabilistic Functions of Markov Processes. In *Inequalities, III: Proceedings of a Symposium*. (Shish, Qved ed.). New York: Academic Press.
- [BD 77] Bickel and Docksum (1977). *Mathematical Statistics: Basic Ideas and Selected Topics*. Prentice Hall, Englewood Cliffs, New Jersey.
- [Boyles 83] Boyles, R.A. (1983). On the Convergence Properties of the EM Algorithm. *Journal of the Royal Statistical Society*, Ser B, 44, 47-50.
- [DLR] Dempster, A.P., Laird, N.M., and Rubin, D.B. (1977). Maximum Likelihood from Incomplete Data Via the EM Algorithm, *Journal of the Royal Statistical Society*, Ser B, 39, 1-38.

- [JJ 93] Jamshidian, M. and Jennrich, R.I. (1993). Conjugate Gradient Acceleration of the EM Algorithm , *Journal of the American Statistical Association*, Vol. 88, No 412, 221-228.
- [Lafferty] Lafferty, J. Notes on the EM Algorithm. *Unpublished Course Notes*.
- [Rao] Rao, C.R. (1965). *Linear Statistical Inference and its Applications*. New York, Wiley.
- [Wu 83] Wu, C.F.Jeff (1983). On the Convergence Properties of the EM Algorithm, *The Annals of Statistics*, 11, 95-103.
- [Zangwill 69] Zangwill, W.I. (1969). *Nonlinear Programming: A Unified Approach*. Prentice Hall, Englewood Cliffs, New Jersey.

Iteration	$\lambda$	$p_1$	$p_2$	$\tilde{p}_1$	$\tilde{p}_2$	$\tilde{p}_3$	$\tilde{p}_4$
0	0.3000	0.7000	0.7000	0.3000	0.3000	0.3000	0.3000
1	0.3000	0.5000	0.5000	0.3000	0.3000	0.3000	0.3000
2	0.3000	0.5000	0.5000	0.3000	0.3000	0.3000	0.3000
3	0.3000	0.5000	0.5000	0.3000	0.3000	0.3000	0.3000
4	0.3000	0.5000	0.5000	0.3000	0.3000	0.3000	0.3000
5	0.3000	0.5000	0.5000	0.3000	0.3000	0.3000	0.3000
6	0.3000	0.5000	0.5000	0.3000	0.3000	0.3000	0.3000

Table 4: The coin example for  $\mathbf{Y} = \{\langle HHH \rangle, \langle TTT \rangle, \langle HHH \rangle, \langle TTT \rangle\}$ , with  $p_1$  and  $p_2$  initialised to the same value. EM is stuck at a saddle point

Iteration	$\lambda$	$p_1$	$p_2$	$\tilde{p}_1$	$\tilde{p}_2$	$\tilde{p}_3$	$\tilde{p}_4$
0	0.3000	0.7001	0.7000	0.3001	0.2998	0.3001	0.2998
1	0.2999	0.5003	0.4999	0.3004	0.2995	0.3004	0.2995
2	0.2999	0.5008	0.4997	0.3013	0.2986	0.3013	0.2986
3	0.2999	0.5023	0.4990	0.3040	0.2959	0.3040	0.2959
4	0.3000	0.5068	0.4971	0.3122	0.2879	0.3122	0.2879
5	0.3000	0.5202	0.4913	0.3373	0.2645	0.3373	0.2645
6	0.3009	0.5605	0.4740	0.4157	0.2007	0.4157	0.2007
7	0.3082	0.6744	0.4223	0.6447	0.0739	0.6447	0.0739
8	0.3593	0.8972	0.2773	0.9500	0.0016	0.9500	0.0016
9	0.4758	0.9983	0.0477	0.9999	0.0000	0.9999	0.0000
10	0.4999	1.0000	0.0001	1.0000	0.0000	1.0000	0.0000
11	0.5000	1.0000	0.0000	1.0000	0.0000	1.0000	0.0000
Iteration	$\lambda$	$p_1$	$p_2$	$\tilde{p}_1$	$\tilde{p}_2$	$\tilde{p}_3$	$\tilde{p}_4$
0	0.3000	0.6999	0.7000	0.2999	0.3002	0.2999	0.3002
1	0.3001	0.4998	0.5001	0.2996	0.3005	0.2996	0.3005
2	0.3001	0.4993	0.5003	0.2987	0.3014	0.2987	0.3014
3	0.3001	0.4978	0.5010	0.2960	0.3041	0.2960	0.3041
4	0.3001	0.4933	0.5029	0.2880	0.3123	0.2880	0.3123
5	0.3002	0.4798	0.5087	0.2646	0.3374	0.2646	0.3374
6	0.3010	0.4396	0.5260	0.2008	0.4158	0.2008	0.4158
7	0.3083	0.3257	0.5777	0.0739	0.6448	0.0739	0.6448
8	0.3594	0.1029	0.7228	0.0016	0.9500	0.0016	0.9500
9	0.4758	0.0017	0.9523	0.0000	0.9999	0.0000	0.9999
10	0.4999	0.0000	0.9999	0.0000	1.0000	0.0000	1.0000
11	0.5000	0.0000	1.0000	0.0000	1.0000	0.0000	1.0000

Table 5: The coin example for  $\mathbf{Y} = \{\langle HHH \rangle, \langle TTT \rangle, \langle HHH \rangle, \langle TTT \rangle\}$ . If we initialise  $p_1$  and  $p_2$  to be a small amount away from the saddle point  $p_1 = p_2$ , the algorithm diverges from the saddle point and eventually reaches the global maximum