

# The Expectation-Maximization Algorithm

April 25, 2006

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>The EM algorithm</b>	<b>2</b>
2.1	MAP version . . . . .	2
2.2	ML version . . . . .	3
2.3	Relation between MAP and ML version . . . . .	3
<b>3</b>	<b>Example</b>	<b>3</b>
3.1	Method 1 . . . . .	4
3.2	Method 2 . . . . .	5

## 1 Introduction

### Problem: estimation in the presence of nuisance parameters

Consider the following problem: we transmit a sequence  $\mathbf{a}$  of  $N$  data symbols, belonging to a constellation  $\Omega$  (for instance  $\Omega = \{-1, +1\}$ ) over a channel. The symbols are drawn iid and uniformly from  $\Omega$ . At the receiver we have the following observation

$$\mathbf{r} = h\mathbf{a} + \mathbf{n}$$

where  $h$  is a channel gain,  $h \in \mathbb{R}$ , and  $\mathbf{n}$  is a vector of  $N$  iid normal RVs, zero-mean, variance  $\sigma^2$ . Our goal is to detect the data sequence  $\mathbf{a}$ .

### Solution - attempt 1

The optimal way to proceed (as we know from detection theory) is to determine the MAP estimate of  $\mathbf{a}$

$$\begin{aligned}\hat{\mathbf{a}}_{MAP} &= \arg \max_{\mathbf{a}} p(\mathbf{a} | \mathbf{r}) \\ &= \arg \max_{\mathbf{a}} \int p(\mathbf{a}, h | \mathbf{r}) dh \\ &= \arg \max_{\mathbf{a}} \int p(\mathbf{r} | \mathbf{a}, h) p(\mathbf{a}) p(h) dh \\ &= \arg \max_{\mathbf{a}} \int \prod_{k=1}^N p(r_k | a_k, h) p(h) dh\end{aligned}$$

and at this point we are stuck, since this integral is generally very hard, and  $p(h)$  is perhaps not known.

### Solution - attempt 2

A common way receivers are designed is as follows: we first estimate  $h$  (leading to  $\hat{h}$ ), and then determine  $\mathbf{a}$ , assuming  $h$  to be known:

$$\begin{aligned}\hat{\mathbf{a}} &= \arg \max_{\mathbf{a}} p(\mathbf{a} | \mathbf{r}, \hat{h}) \\ &= \arg \max_{\mathbf{a}} p(\mathbf{r} | \mathbf{a}, \hat{h}) \\ &= \arg \max_{\mathbf{a}} \prod_{k=1}^N p(r_k | a_k, \hat{h}) \\ &= \arg \max_{\mathbf{a}} \sum_{k=1}^N \log p(r_k | a_k, \hat{h})\end{aligned}$$

so that

$$\hat{a}_k = \arg \max_{a_k \in \Omega} \log p(r_k | a_k, \hat{h}).$$

So how do we find  $\hat{h}$ ? Suppose  $h$  has a zero-mean Gaussian prior (which may or may not be known). Then

$$\hat{h}_{MAP} = \arg \max_h \sum_{\mathbf{a}} p(\mathbf{r} | \mathbf{a}, h) p(h)$$

and

$$\hat{h}_{ML} = \arg \max_h \sum_{\mathbf{a}} p(\mathbf{r} | \mathbf{a}, h)$$

As this summation is over  $|\Omega|^N$  sequences, it is again very hard to compute.

## 2 The EM algorithm

The Expectation-Maximization (EM) algorithm is a technique that solves ML and MAP problems iteratively. To obtain an estimate of a parameter  $\theta$ , the EM algorithm generates a sequence of estimate  $\hat{\theta}^{(1)}, \hat{\theta}^{(2)}, \dots$ , starting from a well-chosen initial estimate  $\hat{\theta}^{(0)}$ . Many variations exist. We will focus on the most common ones.

### 2.1 MAP version

The MAP estimate of a parameter  $\theta$  from an observation  $x$  is obtained by maximizing the a posteriori distribution:

$$\begin{aligned}\hat{\theta}_{MAP} &= \arg \max_{\theta} p(\theta | x) \\ &= \arg \max_{\theta} p(x | \theta) p(\theta)\end{aligned}$$

When the solution is hard to obtain, we can introduce so-called missing or unobserved data (say  $z$ ), in such a way that if  $z$  were known, the estimation problem would be easier to solve ( $p(\theta | x, z)$  is easy to maximize w.r.t.  $\theta$ ). Once we have introduced the missing data, we can execute the EM algorithm. Starting from an initial estimate of  $\theta$ ,  $\hat{\theta}^{(0)}$ , the EM algorithm iterates between the E-step and the M-step:

**E-step:**

$$Q_{MAP}(\theta | \hat{\theta}^{(n)}) = \int p(z | x, \hat{\theta}^{(n)}) \log p(z, x, \theta) dz. \quad (1)$$

**M-step:**

$$\hat{\theta}^{(n)} = \arg \max_{\theta} Q_{MAP} \left( \theta | \hat{\theta}^{(n)} \right)$$

Hence, the EM algorithm generates a sequence of estimates,  $\hat{\theta}^{(0)}, \hat{\theta}^{(1)}, \dots$ . When this sequence converges, we name  $\hat{\theta}^{(+\infty)}$  a *solution* of the EM algorithm. The missing data  $z$  should be chosen such that the E-step is easy to compute.

**Convergence**

- Under some regularity conditions, one can show that the EM algorithm converges to an extremum or a saddle point of the a posteriori distribution  $p(\theta | x)$
- The a posteriori probability of successive estimates is non-decreasing:

$$\log p \left( \hat{\theta}^{(n+1)} | x \right) \geq \log p \left( \hat{\theta}^{(n)} | x \right). \quad (2)$$

## 2.2 ML version

The ML estimate of a parameter  $\theta$  from an observation  $x$  is obtained by maximizing the likelihood function:

$$\hat{\theta}_{ML} = \arg \max_{\theta} p(x | \theta)$$

The ML version is very similar, but uses a slightly different  $Q$ -function.

**E-step:**

$$Q_{ML} \left( \theta | \hat{\theta}^{(n)} \right) = \int p \left( z | x, \hat{\theta}^{(n)} \right) \log p(z, x | \theta) dz. \quad (3)$$

**M-step:**

$$\hat{\theta}^{(n)} = \arg \max_{\theta} Q_{ML} \left( \theta | \hat{\theta}^{(n)} \right)$$

**Convergence**

Similar convergence properties: extremum of likelihood, non-decreasing likelihood of successive estimates.

## 2.3 Relation between MAP and ML version

Since  $p(z, x, \theta) = p(z, x | \theta) p(\theta)$ , we can re-write (1) as

$$\begin{aligned} Q_{MAP} \left( \theta | \hat{\theta}^{(n)} \right) &= \int p \left( z | x, \hat{\theta}^{(n)} \right) \log (p(z, x | \theta) p(\theta)) dz \\ &= \int p \left( z | x, \hat{\theta}^{(n)} \right) \log (p(z, x | \theta)) dz + \int p \left( z | x, \hat{\theta}^{(n)} \right) \log (p(\theta)) dz \\ &= Q_{ML} \left( \theta | \hat{\theta}^{(n)} \right) + \log (p(\theta)) \end{aligned}$$

## 3 Example

As we will see, it is important to simplify the E-step as much as possible. This includes removing terms and factors that do not depend on  $\theta$ . For this reason we will often (ab)use the notation ' $\propto$ '. We will assume  $h \sim \mathcal{N}(0, \sigma_h^2)$ . This distribution may or may not be known to the receiver.

### 3.1 Method 1

Let us return to our example. We set  $\mathbf{a} \leftrightarrow \theta$  and  $h \leftrightarrow z$ . Since  $\mathbf{a}$  is uniformly distributed over  $\Omega^N$ , we can use the ML-version of the EM algorithm:

$$Q_{ML}(\mathbf{a}|\hat{\mathbf{a}}^{(n)}) = \int p(h|\mathbf{r}, \hat{\mathbf{a}}^{(n)}) \log p(h, \mathbf{r}|\mathbf{a}) dh$$

#### E-step

Since  $\mathbf{a}$  and  $h$  are independent,  $p(h, \mathbf{r}|\mathbf{a}) = p(\mathbf{r}|h, \mathbf{a})p(h)$  so that

$$\begin{aligned} Q_{ML}(\mathbf{a}|\hat{\mathbf{a}}^{(n)}) &= \int p(h|\mathbf{r}, \hat{\mathbf{a}}^{(n)}) \log p(\mathbf{r}|h, \mathbf{a}) dh \\ &+ \int p(h|\mathbf{r}, \hat{\mathbf{a}}^{(n)}) \log p(h) dh \end{aligned}$$

As the latter term is independent of  $\mathbf{a}$ , it will not affect the M-step, so it can be dropped. This gives us

$$Q_{ML}(\mathbf{a}|\hat{\mathbf{a}}^{(n)}) = \int p(h|\mathbf{r}, \hat{\mathbf{a}}^{(n)}) \log p(\mathbf{r}|h, \mathbf{a}) dh. \quad (4)$$

Let us first focus on  $\log p(\mathbf{r}|h, \mathbf{a})$ . Since the noise is iid Gaussian,

$$\begin{aligned} \log p(\mathbf{r}|h, \mathbf{a}) &\propto -\frac{1}{2\sigma^2} \sum_{k=1}^N (r_k - ha_k)^2 \\ &= -\frac{1}{2\sigma^2} \sum_{k=1}^N (r_k^2 + h^2 a_k^2 - 2hr_k a_k) \\ &\propto -\sum_{k=1}^N (h^2 a_k^2 - 2hr_k a_k). \end{aligned} \quad (5)$$

Now, we look into  $p(h|\mathbf{r}, \hat{\mathbf{a}}^{(n)})$ . Using Bayes' rule, we know that

$$\begin{aligned} p(h|\mathbf{r}, \hat{\mathbf{a}}^{(n)}) &= p(h|\mathbf{r}, \hat{\mathbf{a}}^{(n)}) p(\mathbf{r}|\hat{\mathbf{a}}^{(n)}) \\ &\text{and} \\ &= p(\mathbf{r}|h, \hat{\mathbf{a}}^{(n)}) p(h|\hat{\mathbf{a}}^{(n)}) \\ &= p(\mathbf{r}|h, \hat{\mathbf{a}}^{(n)}) p(h) \end{aligned}$$

so that

$$p(h|\mathbf{r}, \hat{\mathbf{a}}^{(n)}) = Cp(\mathbf{r}|h, \hat{\mathbf{a}}^{(n)}) p(h)$$

where  $C$  is a constant (independent of  $h$ ). Taking into account (5), it is clear that we only require the first and second order moments of  $p(h|\mathbf{r}, \hat{\mathbf{a}}^{(n)})$  to evaluate  $Q_{ML}(\mathbf{a}|\hat{\mathbf{a}}^{(n)})$ . We find that  $p(h|\mathbf{r}, \hat{\mathbf{a}}^{(n)})$  is Gaussian with mean  $M$  and variance  $\Sigma^2$ . Let us find  $M$  and  $\Sigma^2$ :

$$\begin{aligned} p(h|\mathbf{r}, \hat{\mathbf{a}}^{(n)}) &= Cp(\mathbf{r}|h, \hat{\mathbf{a}}^{(n)}) p(h) \\ &= C' \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{r} - h\hat{\mathbf{a}}^{(n)}\|^2 - \frac{1}{2\sigma_h^2} h^2\right) \\ &= C'' \exp\left(-\frac{1}{2\sigma^2} \left(h^2 \|\hat{\mathbf{a}}^{(n)}\|^2 - 2h\mathbf{r}^T \hat{\mathbf{a}}^{(n)}\right) - \frac{1}{2\sigma_h^2} h^2\right) \\ &= C'' \exp\left(-\frac{1}{2} (h^2) \left(\frac{\|\hat{\mathbf{a}}^{(n)}\|^2}{\sigma^2} + \frac{1}{\sigma_h^2}\right) + \frac{1}{\sigma^2} h\mathbf{r}^T \hat{\mathbf{a}}^{(n)}\right) \end{aligned}$$

which implies that  $p(h|\mathbf{r}, \hat{\mathbf{a}}^{(n)})$  is a Gaussian distribution with variance  $\Sigma^2 = 1 / \left( \left( \frac{\|\hat{\mathbf{a}}^{(n)}\|^2}{\sigma^2} + \frac{1}{\sigma_h^2} \right) \right)$  so that

$$p(h|\mathbf{r}, \hat{\mathbf{a}}^{(n)}) = C'' \exp \left( -\frac{h^2}{2\Sigma^2} + \frac{h}{\Sigma^2} \left( \frac{\Sigma^2}{\sigma^2} \mathbf{r}^T \hat{\mathbf{a}}^{(n)} \right) \right)$$

from which we see that the mean of  $p(h|\mathbf{r}, \hat{\mathbf{a}}^{(n)})$  is given by

$$\begin{aligned} M &= \frac{\Sigma^2}{\sigma^2} \mathbf{r}^T \hat{\mathbf{a}}^{(n)} \\ &= \frac{\sigma_h^2}{\|\hat{\mathbf{a}}^{(n)}\|^2 \sigma_h^2 + \sigma^2} \mathbf{r}^T \hat{\mathbf{a}}^{(n)} \end{aligned}$$

Note that when  $\sigma_h^2 \rightarrow +\infty$ , then  $M \rightarrow \mathbf{r}^T \hat{\mathbf{a}}^{(n)} / \|\hat{\mathbf{a}}^{(n)}\|^2$ , as we would expect. We find that

$$\begin{aligned} \int p(h|\mathbf{r}, \hat{\mathbf{a}}^{(n)}) h dh &= M \\ &= \tilde{h} \end{aligned}$$

and

$$\begin{aligned} \int p(h|\mathbf{r}, \hat{\mathbf{a}}^{(n)}) h^2 dh &= M^2 + \Sigma^2 \\ &= \widetilde{h^2} \end{aligned}$$

Using  $M$  and  $\Sigma^2$  to compute (4):

$$Q_{ML}(\mathbf{a}|\hat{\mathbf{a}}^{(n)}) = -\sum_{k=1}^N \left( \widetilde{h^2} a_k^2 - 2\tilde{h} r_k a_k \right).$$

### M-step

To maximize this function w.r.t.  $\mathbf{a}$ , we find that the  $k$ -th component is given by

$$\hat{a}_k^{(n+1)} = \arg \max_{a \in \Omega} \left( 2\tilde{h} r_k a_k - \widetilde{h^2} a_k^2 \right)$$

### Comments

- We see that this approach requires knowledge of  $p(h|\mathbf{r}, \hat{\mathbf{a}}^{(n)})$ . When  $p(h)$  is known, this distribution can be found. When  $p(h)$  is Gaussian, the EM algorithm leads to an elegant solution.
- When  $\Omega = \{-1, +1\}$ , we see that  $Q_{ML}(\mathbf{a}|\hat{\mathbf{a}}^{(n)}) = \sum_{k=1}^N (\tilde{h} r_k a_k)$ , in which case the M-step becomes ML detection of  $a_k$ , given a channel estimate  $\hat{h} = \tilde{h}$ .

## 3.2 Method 2

As we mentioned earlier, it is sometimes more convenient to estimate  $h$ , considering  $\mathbf{a}$  as a nuisance parameter. We set  $\mathbf{a} \leftrightarrow z$  and  $h \leftrightarrow \theta$ . We now have a choice to perform MAP estimation or ML estimation, depending on whether or not we know  $p(h)$ . We focus on MAP. By replacing the prior with a constant (or letting  $\sigma_h^2 \rightarrow \infty$ ), we find ML.

### E-step

$$Q_{MAP} \left( h | \hat{h}^{(n)} \right) = \int p \left( \mathbf{a} | \mathbf{r}, \hat{h}^{(n)} \right) \log p \left( \mathbf{a}, \mathbf{r}, h \right) dz \quad (6)$$

Now,  $p \left( \mathbf{a}, \mathbf{r}, h \right) = p \left( \mathbf{r} | \mathbf{a}, h \right) p \left( \mathbf{a} \right) p \left( h \right)$ , so that

$$\begin{aligned} \log p \left( \mathbf{a}, \mathbf{r}, h \right) &\propto -\frac{1}{2\sigma_h^2} h^2 - \frac{1}{2\sigma^2} \|\mathbf{r} - h\mathbf{a}\|^2 \\ &= -\frac{1}{2\sigma_h^2} h^2 - \frac{1}{2\sigma^2} h^2 \sum_k a_k^2 + \frac{1}{\sigma^2} h \sum_k r_k a_k. \end{aligned}$$

We see that to evaluate (6), we require the first and second order moments of the marginals  $p \left( a_k | \mathbf{r}, \hat{h}^{(n)} \right)$ :

$$\begin{aligned} Q_{MAP} \left( h | \hat{h}^{(n)} \right) &= -\frac{1}{2\sigma_h^2} h^2 - \frac{1}{2\sigma^2} \left\{ h^2 \sum_k \int a_k^2 p \left( \mathbf{a} | \mathbf{r}, \hat{h}^{(n)} \right) da_k - 2h \sum_k r_k \int a_k p \left( \mathbf{a} | \mathbf{r}, \hat{h}^{(n)} \right) da_k \right\} \\ &= -\frac{1}{2\sigma_h^2} h^2 - \frac{1}{2\sigma^2} \left\{ h^2 \sum_k \int a_k^2 p \left( a_k | \mathbf{r}, \hat{h}^{(n)} \right) da_k - 2h \sum_k r_k \int a_k p \left( a_k | \mathbf{r}, \hat{h}^{(n)} \right) da_k \right\} \\ &= -\frac{1}{2\sigma_h^2} h^2 - \frac{1}{2\sigma^2} \left\{ h^2 \sum_{k=1}^N \sum_{a \in \Omega} a^2 \times p \left( a_k = a | \mathbf{r}, \hat{h}^{(n)} \right) - 2h \sum_{k=1}^N r_k \sum_{a \in \Omega} a \times p \left( a_k = a | \mathbf{r}, \hat{h}^{(n)} \right) \right\} \\ &= -\frac{1}{2\sigma_h^2} h^2 - \frac{1}{2\sigma^2} \left\{ h^2 \sum_k \widetilde{a_k^2} - 2h \sum_k r_k \widetilde{a_k} \right\} \end{aligned}$$

So we need to determine  $p \left( a_k | \mathbf{r}, \hat{h}^{(n)} \right)$ .

$$\begin{aligned} p \left( a_k | \mathbf{r}, \hat{h}^{(n)} \right) &= \frac{1}{p \left( \mathbf{r} | \hat{h}^{(n)} \right)} p \left( \mathbf{r} | a_k, \hat{h}^{(n)} \right) p \left( a_k \right) \\ &= C \exp \left( -\frac{1}{2\sigma^2} (r_k - h a_k)^2 \right) \end{aligned}$$

for some constant  $C$ . This constant can be found by determining  $f(a_k = a) = \exp \left( -\frac{1}{2\sigma^2} (r_k - h a)^2 \right)$ ,  $\forall a \in \Omega$ , so that  $C = 1 / \sum_{a \in \Omega} f(a_k = a)$ .

### M-step

We now need to maximize

$$\begin{aligned} Q_{MAP} \left( h | \hat{h}^{(n)} \right) &= -\frac{1}{2\sigma_h^2} h^2 - \frac{1}{2\sigma^2} \left\{ h^2 \sum_k \widetilde{a_k^2} - 2h \sum_k r_k \widetilde{a_k} \right\} \\ &= h^2 \underbrace{\left( -\frac{1}{2\sigma^2} \left( \sum_k \widetilde{a_k^2} + \frac{\sigma^2}{\sigma_h^2} \right) \right)}_{\alpha} - 2h \underbrace{\left( -\frac{1}{2\sigma^2} \sum_k r_k \widetilde{a_k} \right)}_{\beta} \\ &= h^2 \alpha - 2h \beta \end{aligned}$$

w.r.t.  $h$ , giving

$$\begin{aligned} \hat{h}^{(n+1)} &= \frac{\beta}{\alpha} \\ &= \frac{\sum_k r_k \widetilde{a_k}}{\sum_k \widetilde{a_k^2} + \frac{\sigma^2}{\sigma_h^2}} \end{aligned}$$

## Comments

- When  $p(h)$  is unknown, we find the ML-version of the EM algorithm by simply letting  $\sigma_h^2 \rightarrow +\infty$ :

$$Q_{ML} \left( h | \hat{h}^{(n)} \right) = -h^2 \sum_k \widetilde{a_k^2} + 2h \sum_k r_k \tilde{a}_k$$

and

$$\hat{h}^{(n+1)} = \frac{\sum_k r_k \tilde{a}_k}{\sum_k \widetilde{a_k^2}}.$$

- When  $\Omega = \{-1, +1\}$ ,  $a_k^2 = 1$  so that

$$\begin{aligned} Q_{MAP} \left( h | \hat{h}^{(n)} \right) &= -\frac{1}{2\sigma_h^2} h^2 - \frac{1}{2\sigma^2} \left\{ h^2 N - 2h \sum_k r_k \tilde{a}_k \right\} \\ \hat{h}^{(n+1)} &= \frac{\sum_k r_k \tilde{a}_k}{N + \frac{\sigma^2}{\sigma_h^2}} \end{aligned}$$

and

$$\begin{aligned} Q_{ML} \left( h | \hat{h}^{(n)} \right) &= -h^2 N + 2h \sum_k r_k \tilde{a}_k \\ \hat{h}^{(n+1)} &= \frac{\sum_k r_k \tilde{a}_k}{N} \end{aligned}$$

- After a few iterations (say  $K$ ), we can make final decisions w.r.t. the data symbols

$$\hat{a}_k = \arg \max_{a \in \Omega} p \left( a_k = a | \mathbf{r}, \hat{h}^{(K)} \right)$$