

Expectation Maximization with Coin Toss example

Marijo Simunovic

October 18, 2017

1 Coin toss problem description

This document contains detailed solution for the problem described in [2]. In summary, someone choose randomly one of two biased coins five times. Each time, selected coin has been tossed and outcomes are recorded. We are given the outcomes of this five events but not the identities of coins whom each event belongs. Also, coin biases are unknown. Our goal is to somehow infer these parameters from the given data. In order to solve the problem we resort the Expectation Maximization (EM) algorithm.

2 Theoretical background of EM

Given the statistical model which generates a set \mathbf{X} of observed data, a set of unobserved latent data or missing values \mathbf{Z} , and a vector of unknown parameters Θ , along with a likelihood function $L(\Theta; \mathbf{X}, \mathbf{Z}) = p(\mathbf{X}, \mathbf{Z}|\Theta)$, the maximum likelihood estimate (MLE) of the unknown parameters is determined by the marginal likelihood of the observed data

$$L(\Theta; \mathbf{X}) = p(\mathbf{X}|\Theta) = \int p(\mathbf{X}, \mathbf{Z}|\Theta) d\mathbf{Z} \quad (1)$$

The EM algorithm seeks to find the MLE of the marginal likelihood by iteratively applying these two steps:

Expectation step (E step): Calculate the expected value of the log likelihood function, with respect to the conditional distribution of \mathbf{Z} given \mathbf{X} under the current estimate of the parameters $\Theta^{(t)}$:

$$Q(\Theta|\Theta^{(t)}) = E_{\mathbf{Z}|\mathbf{X}, \Theta^{(t)}}[\log L(\Theta; \mathbf{X}, \mathbf{Z})] \quad (2)$$

Maximization step (M step): Find the parameter that maximizes this quantity:

$$\Theta^{(t+1)} = \arg \min_{\Theta} Q(\Theta|\Theta^{(t)}) \quad (3)$$

This is a word by word definition copied from [1]¹. If you don't understand what the above means at first, don't worry. It is merely due to the compactness of the mathematical language. The first term means: Set up the model for the likelihood function pretending that you know what is the value of Θ . In the actual solution, we randomly instantiate its values to $\Theta_i \in [0, 1]$ where i

¹I choose to take this particular definition (out of many out there) because it is the first one I actually understood. I don't think that this is a better or worse than the others, it's just that my brain was finally ready.

is the i -th component of the vector Θ . Once we have defined this equation (model) we take the expectation over it to find the probability of each possible value of Z , given Θ . Then we compute a better estimate for the parameters Θ using these probabilities. We iterate these procedure until convergence.

2.1 Math Mode

Vector \mathbf{z} corresponds to the event of choosing one of the two given coins. It is bivariate, the values of each of the two possible outcomes are mutually exclusive.

$$\mathbf{z}_n = \begin{bmatrix} z_{n1} \\ z_{n2} \end{bmatrix} \in \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\} \quad (4)$$

Probability of the event \mathbf{z}_n (which coin we selected in the n -th event) is given by

$$p(\mathbf{z}_n) = \prod_{k=1}^2 \pi_k^{z_{nk}} \quad (5)$$

with requirement that $\sum_k \pi_k = 1$.

Each single coin toss is independent on the other tosses. It only depends on the coin we choose previously. Probability of the single coin toss outcome given coin k is

$$p(x_n^j | \mathbf{z}_n, \Theta) = \prod_{k=1}^2 [\Theta_k^{x_n^j} (1 - \Theta_k^{1-x_n^j})]^{z_{nk}} \quad (6)$$

where $x_n^j = 1$ if j -th outcome was head and $x_n^j = 0$ if j -th outcome was tail. We can now determine the probability of each of our 5 events

$$p(\{x_n^1 \cdots x_n^{10}\} | \mathbf{z}_n, \Theta) = \prod_{j=1}^{10} p(x_n^j | \mathbf{z}_n, \Theta) \quad (7)$$

Joint probability of the coin toss outcomes and coin selections is given by:

$$p(X_1, \dots, X_5, z_1, \dots, z_5 | \Theta) \quad (8)$$

$$= p(\{x_1^1 \cdots x_1^{10}\}, \dots, \{x_5^1 \cdots x_5^{10}\}, z_1, \dots, z_5 | \Theta) \quad (9)$$

$$= p(\{x_1^1 \cdots x_1^{10}\}, \dots, \{x_5^1 \cdots x_5^{10}\} | z_1, \dots, z_5, \Theta) p(z_1, \dots, z_5) \quad (10)$$

$$= \prod_{n=1}^5 p(\{x_n^1 \cdots x_n^{10}\} | \mathbf{z}_n, \Theta) \prod_{n=1}^5 p(\mathbf{z}_n) \quad (11)$$

Finally, combining the above we get the equation

$$p(X_1, \dots, X_5, z_1, \dots, z_5 | \Theta) = \prod_{n=1}^5 \prod_{j=1}^{10} \prod_{k=1}^2 [\Theta_k^{x_n^j} (1 - \Theta_k^{1-x_n^j})]^{z_{nk}} \prod_{n=1}^5 \prod_{k=1}^2 \pi_k^{z_{nk}} \quad (12)$$

Taking the log of the expression above we get rid of the products in exchange of the summations.

$$\log(P) = \sum_{n=1}^5 \sum_{k=1}^2 z_{nk} \left[\sum_{j=1}^{10} \log(\Theta_k^{x_n^j} (1 - \Theta_k^{1-x_n^j})) + \log(\pi_k) \right] \quad (13)$$

The final step is taking the expectation over the z random variable. In this case, we are using the conditional expectation, as we do not know what was the outcome on each of the step. Using the linearity of the expectation as

$$E_{Z|X}[\log(P)] = \sum_{n=1}^5 \sum_{k=1}^2 E_{Z|X}[z_{nk}] \left[\sum_{j=1}^{10} \log(\Theta_k^{x_n^j} (1 - \Theta_k^{1-x_n^j})) + \log(\pi_k) \right] \quad (14)$$

Conditional expectation gives as the expected value of the random experiment in light of the new information on which we are conditioning. Consider throwing an ordinary dice. The expected value of the outcome is $E(\text{number on dice}) = (1 + 2 + 3 + 4 + 5 + 6)/6 = 3.5$. But if I told you the parity of the outcome, in that case we would calculate the expectation $E(\text{odd}) = (1+3+5)/3 = 3$ and $E(\text{even}) = (2+4+6)/3 = 4$. We see that it is not enough to express the expectation with the one value only. We can use similar analogy in our experiment. As we have seen the outcomes of each of the five events, we will assign bigger value of Θ to the ones that have had more heads. If there isn't any significant difference in the outcomes we will have a hard time trying to figure out which event belongs to which coin.

$$E_{Z|X}[z_{nk}] = \sum_z z_{nk} p(z_n | \{x_n^1 \cdots x_n^{10}\}; \Theta) \quad (15)$$

Using the Bayes formula, the posterior of z_n given values of the n -th event is

$$p(z_n | \{x_n^1 \cdots x_n^{10}\}; \Theta) = \frac{p(\{x_n^1 \cdots x_n^{10}\} | z_n; \Theta) p(z_n)}{p(\{x_n^1 \cdots x_n^{10}\}; \Theta)} \quad (16)$$

We already know the expressions for the terms in numerator, and the expression in denominator is given by marginalizing over all possible outcomes of z (in our case there are only two, as already mentioned).

$$p(z_n | x_n; \Theta) = \frac{\prod_{k=1}^2 [\pi_k p(x_n; \Theta_k)]^{z_{nk}}}{\sum_{z_m} \prod_{m=1}^2 [\pi_m p(x_n; \Theta_m)]^{z_{nm}}} \quad (17)$$

Going back to our expectation formula

$$E_{Z|X}[z_{nk}] = \frac{\sum_{z_n} z_{nk} \prod_{k=1}^2 [\pi_k p(x_n; \Theta_k)]^{z_{nk}}}{\sum_{z_m} \prod_{m=1}^2 [\pi_m p(x_n; \Theta_m)]^{z_{nm}}} \quad (18)$$

We see that in the numerator those events where $z_{nk} = 0$ will be eliminated, while in the denominator we are still summing over all possibilities.

$$E_{Z|X}[z_{nk}] = \frac{\pi_k p(x_n; \Theta)}{\sum_{m=1}^2 \pi_m p(x_n; \Theta_m)} \quad (19)$$

As the probability of the choosing both coins is equal, we know that the $\pi_k = 0.5$ in both cases. Thus, we can eliminate it from our expression. Also, we fix the Θ_k values (first iteration uses randomly initiated values, while the subsequent steps use values calculated in the previous step) in the E step of the EM algorithm.

In the M step, we fix the expectation value calculated in the E step vary Θ_k values. To find the values which maximize the expectation of log equation, we use the derivation. Second term, after the summation sign becomes constant, so it does not have the effect on the maximization procedure.

$$\max L(\Theta) = E_{Z|X}[p(X_1, \dots, X_5, z_1, \dots, z_5 | \Theta)] \quad (20)$$

$$= \sum_{n=1}^5 \sum_{j=1}^{10} \sum_{k=1}^2 E_{Z|X}[z_{nk}] \log(\Theta_k^{x_n^j} (1 - \Theta_k^{1-x_n^j})) + \text{CONST.} \quad (21)$$

$$= \sum_{n=1}^5 \sum_{j=1}^{10} \sum_{k=1}^2 E_{Z|X}[z_{nk}] (x_n^j \log(\Theta_k) + (1 - x_n^j) \log(1 - \Theta_k)) + \text{CONST.} \quad (22)$$

For a single coin k

$$\frac{dL(\Theta)}{d\Theta_k} = \sum_{n=1}^5 \sum_{j=1}^{10} E_{Z|X}[z_{nk}] \left(\frac{x_n^j}{\Theta_k} - \frac{1 - x_n^j}{1 - \Theta_k} \right) = 0 \quad (23)$$

Multiplying with both denominators $\Theta_k(1 - \Theta_k)$ we get

$$\frac{dL(\Theta)}{d\Theta_k} = \sum_{n=1}^5 \sum_{j=1}^{10} E_{Z|X}[z_{nk}] (x_n^j - \Theta_k) = 0 \quad (24)$$

and finally, since Θ 's index does not appear in the summation indexes, we can simply divide to obtain its (maximal) value

$$\Theta_k = \frac{\sum_{n=1}^5 \sum_{j=1}^{10} E_{Z|X}[z_{nk}] x_n^j}{10 \sum_{n=1}^5 E_{Z|X}[z_{nk}]} \quad (25)$$

We see that the result for each coin is weighted expectation of a single coin over a sum of the expectations of both coins.

EM is sort of like moving a heavy table without anyone's help. You push a bit from one side and then from the other. Bit by bit you are reaching the goal.

References

- [1] Expectation Maximization algorithm wiki page, <https://en.wikipedia.org/wiki/Expectation>
- [2] Chuong B Do, Serafim Batzoglou, *What is the expectation maximization algorithm*, Nature Biotechnology 26, 897-899, doi:10.1038/nbt14052008, <http://www.nature.com/nbt/journal/v26/n8/full/nbt1406.html>