

Expectation-Maximization (EM) algorithm

Yuzhen Ye
 School of Informatics and Computing
 Indiana University, Bloomington
 Spring 2018

Contents

- Introduce EM algorithm using the flipping coin experiment
- Formal definition of the EM algorithm
- Two toy examples
 - Coin toss with missing data
 - Coin toss with hidden data
- Applications of the EM algorithm
 - Motif finding
 - Baum-Welch algorithm
 - Binning of metagenomes

A coin-flipping experiment

a Maximum likelihood

	Coin A	Coin B	
HTTTHHTHTH	5H, 5T		$\hat{\theta}_A = \frac{24}{24+6} = 0.80$
HHHHTHHHHH	9H, 1T		
HTHHHHTHHH	8H, 2T		$\hat{\theta}_B = \frac{9}{9+11} = 0.45$
HTTTHHTHTT	4H, 6T		
THHHTHHHTH	7H, 3T	9H, 11T	
24H, 6T		9H, 11T	

5 sets, 10 tosses per set

θ , the probability of getting heads
 θ_A , the probability of coin A landing on head
 θ_B , the probability of coin B landing on head

Ref: What is the expectation maximization algorithm?
 Nature Biotechnology 26, 897 - 899 (2008)

When the identities of the coins are unknown

b Expectation maximization

Instead of picking up the single best guess, the EM algorithm computes probabilities for each possible completion of the missing data, using the current parameters

	Coin A	Coin B
HTTTHHTHTH	~2.2H, 2.2T	~2.8H, 2.8T
HHHHTHHHHH	~7.2H, 0.8T	~1.8H, 0.2T
HTHHHHTHHH	~5.9H, 1.5T	~2.1H, 0.5T
HTTTHHTHTT	~1.4H, 2.1T	~2.6H, 3.9T
THHHTHHHTH	~4.5H, 1.9T	~2.5H, 1.1T
	~21.3H, 8.6T	~11.7H, 8.4T

$\hat{\theta}_A^{(0)} = 0.80$
 $\hat{\theta}_B^{(0)} = 0.50$
 $\hat{\theta}_A^{(1)} = \frac{21.3}{21.3+8.6} = 0.71$
 $\hat{\theta}_B^{(1)} = \frac{11.7}{11.7+8.4} = 0.58$
 $\hat{\theta}_A^{(2)} = 0.80$
 $\hat{\theta}_B^{(2)} = 0.52$

E(H) for coin B

Main applications of the EM algorithm

- When the data indeed has missing values, due to problems with or limitations of the observation process
- When optimizing the likelihood function is analytically intractable but it can be simplified by assuming the existence of and values for additional but *missing (or hidden)* parameters.

The EM algorithm handles hidden data

Consider a model where, for observed data x and model parameters θ :

$$p(x|\theta) = \sum_z p(x, z|\theta).$$

z is the "hidden" variable that is marginalized out

Finding θ^* which maximizes $\sum_z p(x, z|\theta)$ is hard!

The EM algorithm reduces the difficult task of optimizing $\log P(x; \theta)$ into a sequence of simpler optimization subproblems. In each iteration, The EM algorithm receives parameters $\theta^{(l)}$, and returns new parameters $\theta^{(l+1)}$, s.t. $p(x|\theta^{(l+1)}) > p(x|\theta^{(l)})$.

The EM algorithm

In each iteration the EM algorithm does the following.

E step: Calculate

$$Q_t(\theta) = \sum_z P(z|x; \hat{\theta}^{(t)}) \log P(x, z; \theta)$$

M step: Find $\hat{\theta}^{(t+1)}$ which maximizes the Q function (Next iteration sets $\theta^{(t)} \leftarrow \hat{\theta}^{(t+1)}$ and repeats).

The EM update rule:

$$\hat{\theta}^{(t+1)} = \arg \max_{\theta} \sum_z P(z|x; \hat{\theta}^{(t)}) \log P(x, z; \theta)$$

Convergence of the EM algorithm

Compare the Q function and the g function

$$Q_t(\theta) = \sum_z P(z|x; \hat{\theta}^{(t)}) \log P(x, z; \theta)$$

$$g_t(\theta) = \sum_z P(z|x; \hat{\theta}^{(t)}) \log \frac{P(x, z; \theta)}{P(z|x; \hat{\theta}^{(t)})}$$

Fig. 1

Fig 1 demonstrates the convergence of the EM algorithm. Starting from initial parameters $\theta^{(0)}$, the E-step of the EM algorithm constructs a function g_t that lower-bounds the objective function $\log P(x; \theta)$ (i.e., $g_t \leq \log P(x; \theta)$); and $g_t(\hat{\theta}^{(t)}) = \log P(x; \hat{\theta}^{(t)})$. In the M-step, $\hat{\theta}^{(t+1)}$ is computed as the maximum of g_t . In the next E-step, a new lower-bound g_{t+1} is constructed; maximization of g_{t+1} in the next M-step gives $\hat{\theta}^{(t+2)}$, etc.

As the value of the lower-bound g_t matches the objective function at $\hat{\theta}^{(t)}$, it follows that

$$\log P(x; \hat{\theta}^{(t)}) = g_t(\hat{\theta}^{(t)}) \leq g_t(\hat{\theta}^{(t+1)}) = \log P(x; \hat{\theta}^{(t+1)}) \quad (2)$$

So the objective function monotonically increases during each iteration of expectation maximization!

Ref: Nature Biotechnology 26, 897 - 899 (2008)

The EM update rule

$$\hat{\theta}^{(t+1)} = \arg \max_{\theta} \sum_z P(z|x; \hat{\theta}^{(t)}) \log P(x, z; \theta)$$

The EM update rule maximizes the log likelihood of a dataset expanded to contain all possible completions of the unobserved variables, where each completion is weighted by the posterior probability!

Coin toss with missing data

- Given a coin with two possible outcomes: H (head) and T (tail), with probabilities θ and $1-\theta$, respectively.
- The coin is tossed twice, **but only the 1st outcome, T , is seen. So the data is $x = (T, *)$ (with incomplete data!)**
- We wish to apply the EM algorithm to get parameters that increase the likelihood of the data.
- Let the initial parameters be $\theta = 1/4$.

The EM algorithm at work

$$Q_t(\theta) = \sum_z P(z|x; \theta) \log P(x, z; \theta)$$

Inputs:
 Observation: $x=(T,*)$
 Hidden data: $z_1=(T,T)$ $z_2=(T,H)$
 Initial guess: $\theta = 1/4$

$$n_H \log \theta + n_T \log(1 - \theta) \text{ is maximized when } \theta = \frac{n_H}{n_H + n_T}$$

$$P(x; \theta) = P(z_1; \theta) + P(z_2; \theta) = (1 - \theta)^2 + (1 - \theta)\theta = 3/4$$

$$P(z_1|x; \theta) = P(x, z_1; \theta) / P(x; \theta) = (1 - \theta)^2 / (3/4) = 3/4$$

$$P(z_2|x; \theta) = 1 - P(z_1|x; \theta) = 1/4$$

$$n_H(z_1) = 0, n_T(z_1) = 2, n_H(z_2) = 1, \text{ and } n_T(z_2) = 1$$

$$n_H = 1/4 \times 1 = 1/4, n_T = 3/4 \times 2 + 1/4 \times 1 = 7/4, \theta = \frac{n_H}{n_H + n_T} = \frac{1/4}{1/4 + 7/4} = 1/8$$

The EM algorithm at work: continue

- Initial guess $\theta = 1/4$
- After one iteration $\theta = 1/8$
- ...
- The optimal parameter θ will never be reached by the EM algorithm!**

Coin toss with hidden data

Two coins A and B, with parameters $\theta = \{\theta_A, \theta_B\}$; compute θ that maximizes the log likelihood of the observed data $x = \{x_1, x_2, \dots, x_5\}$

E.g., initial parameter θ : $\theta_A = 0.60, \theta_B = 0.50$
 $P(z_1 = A | x; \theta^i) = P(z_1 = A | x_1; \theta^i)$ (x_1, x_2, \dots, x_5 are independent observations)

$$P(z_1 = A, x_1; \theta^i) = \frac{0.6^5 \times 0.4^5}{0.6^5 \times 0.4^5 + 0.5^5 \times 0.5^5} = 0.58$$

observation	Coin A		Coin B	
	nH	nT	nH	nT
x1: HTTTHHTHTH	5	5	0.58	0.42
x2: HHHHTHHHHH	9	1	0.84	0.16
x3: HTHHHHTHH	8	2	0.81	0.19
x4: HTTHTTHTT	4	6	0.25	0.75
x5: THHHHTHHHTH	8	2	0.81	0.19
			24.3H	8.4T
			9.7H	7.6T

New parameter: $\theta_A = 24.3 / (24.3 + 8.4) = 0.74, \theta_B = 9.7 / (9.7 + 7.6) = 0.56$

Motif finding problem

- Motif finding problem is not that different from the coin toss problem!
- Probabilistic approaches to motif finding
 - EM
 - Gibbs sampling (a generalized EM algorithm)
- There are also combinatorial approaches

Motif finding problem

- Given a set of DNA sequences:

```

cctgatagacgctatctggctatccacgtacgtaggctcctctgtgcaatctatgcgtttccaacct
agtactgggtgtacatttgatagctacgtacacccggcaacctgaaacaaacgctcagaaccagaagtgc
aaacgtacgtgcaccctctcttctctgtgctctggccaacgagggcgtatgtataagacgaaaatttt
agcctccgatgtaagtcatagctgtaactattacctgcacccctattacattctacgtacgtataca
ctgttatacaacgctcatggcgggtatgctgtttggtcgtctacgctcgatcgttaacgtacgtc
    
```

- Find the motif in each of the individual sequences

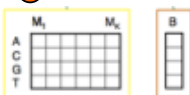
The MEME algorithm

- Collect all substrings with the same length w from the input sequences: $X = (X_1, \dots, X_n)$
- Treat sequences as bags of subsequences: a bag for motif, and a bag for background
- Need to figure out two models (one for motif, and one for the background), and assign each of the subsequences to one of the bags, such that the likelihood of the data (subsequences) is maximized
 - Difficult problem
 - Solved by the EM algorithm

Motif finding vs coin toss

tagacgctatc 0.3x(M) 0.7x(B)
 gctatccacgt 0.7x(M) 0.3x(B)
 gtaggtcctct 0.2x(M) 0.8x(B)

(M) Motif
 (B) Background model



Probability of a subsequence:
 $P(x|M)$, or $P(x|B)$



θ : the probability of getting heads
 θ_A : P(head) for coin A
 θ_B : P(head) for coin B

Probability of an observation sequence:
 $P(x|\theta) = \theta^{\#(\text{heads})} (1-\theta)^{\#(\text{tails})}$

Fitting a mixture model by EM

- A finite mixture model:
 - data $X = (X_1, \dots, X_n)$ arises from two or more groups with g models $\theta = (\theta_1, \dots, \theta_g)$.
- Indicator vectors $Z = (Z_1, \dots, Z_n)$, where $Z_i = (Z_{i1}, \dots, Z_{ig})$, and $Z_{ij} = 1$ if X_i is from group j , and $= 0$ otherwise.
- $P(Z_{ij} = 1 | \theta_j) = \lambda_j$. For any given i , all Z_{ij} are 0 except one j ;
- $g=2$: class 1 (the motif) and class 2 (the background) are given by position specific and a general multinomial distribution

The E- and M-step

- E-step: Since the log likelihood is the sum of over i and j of terms multiplying Z_{ij} , and these are independent across i , we need only consider the expectation of one such, given X_i . Using initial parameter values θ' and λ' , and the fact that the Z_{ij} are binary, we get

$$E(Z_{ij} | X, \theta', \lambda') = \lambda'_{ij} P(X_i | \theta'_{ij}) / \sum_k \lambda'_{ik} P(X_i | \theta'_{ik}) = Z'_{ij}$$

- M-step: The maximization over λ is independent of the rest and is readily achieved with

$$\lambda''_{ij} = \sum_i Z'_{ij} / n.$$

Baum-Welch algorithm for HMM parameter estimation

$$A_{kl} = \sum_{j=1}^n \frac{1}{p(x^j)} \sum_{i=1}^L p(s_{i-1}=k, s_i=l, x^j | \theta)$$

$$A_{kl} = \sum_{j=1}^n \frac{1}{p(x^j)} \sum_{i=1}^L f_k^j(i-1) a_{kl} e_l(x_i) b_l^j(i)$$

$$E_k(b) = \sum_{j=1}^n \frac{1}{p(x^j)} \sum_{i: x_i=b} f_k^j(i) f_k^j(i)$$

During each iteration, compute the expected transitions between any pair of states, and expected emissions from any state, using averaging process (E-step), which are then used to compute new parameters (M-step).

Application of EM algorithms in metagenomics: Binning

- AbundanceBin
 - Binning of short reads into bins (species)
 - A Novel Abundance-Based Algorithm for Binning Metagenomic Sequences Using l-tuples (JCB 2011, 18(3): 523-534. RECOMB 2010)
- MaxBin/MaxBin2
 - Binning of assembled metagenomic scaffolds using an EM algorithm (Microbiome, 2014 doi: 10.1186/2049-2618-2-26)

Pros and Cons

- Cons
 - Slow convergence
 - Converge to local optima
- Pros
 - The E-step and M-step are often easy to implement for many problems, thanks to the nice form of the complete-data likelihood function
 - Solutions to the M-steps often exist in the closed form
- Ref
 - On the convergence properties of the EM algorithm. CFJ WU, 1983
 - A gentle tutorial of the EM algorithm and its applications to parameter estimation for Gaussian mixture and hidden Markov models, JA Bilmes, 1998
 - What is the expectation maximization algorithm? 2008