



Tutorial on Expectation-Maximization Algorithm

Long Qiu

`qiul@comp.nus.edu.sg`

School of Computing
National University of Singapore

- Maximum Likelihood Estimation

- Maximum Likelihood Estimation
- Example #1 (from [1]): MLE with missing data

- Maximum Likelihood Estimation
- Example #1 (from [1]): MLE with missing data
- EM Algorithm

- Maximum Likelihood Estimation
- Example #1 (from [1]): MLE with missing data
- EM Algorithm
- Example #2 (from [1]): MLE with a hidden variable

- Maximum Likelihood Estimation
- Example #1 (from [1]): MLE with missing data
- EM Algorithm
- Example #2 (from [1]): MLE with a hidden variable
- Discussion

Maximum Likelihood Estimation

- Parameter Estimation:
Given

- a set of observed data $\{x_1, x_2, \dots\}$, and
- a proposed model $p(x \mid \theta)$

to find

- the value for one or more parameters θ that maximize the **likelihood**:

$$L(\theta; x) = p(x_1, x_2, \dots \mid \theta).$$

Example #1

A	B
0	0
0	0
0	0
0	?
0	1
1	0
1	1
1	1

Estimate parameters in joint distribution: θ

	\bar{A}	A
\bar{B}	$p(\bar{A}\bar{B})$	$p(A\bar{B})$
B	$p(\bar{A}B)$	$p(AB)$

Example #1

A	B
0	0
0	0
0	0
0	1
1	0
1	1
1	1

Estimate parameters in joint distribution: θ

	\bar{A}	A
\bar{B}	3/7	1/7
B	1/7	2/7

- Opt 1: Ignore the incomplete data

Example #1

A	B
0	0
0	0
0	0
0	0
0	1
1	0
1	1
1	1

Estimate parameters in joint distribution: θ

	\bar{A}	A
\bar{B}	4/8	1/8
B	1/8	2/8

- Opt 1: Ignore the incomplete data
- Opt 2: Fill in a best guessed value

Example #1

A	B
0	0
0	0
0	0
0	?
0	1
1	0
1	1
1	1

Estimate parameters in joint distribution: θ

	\bar{A}	A
\bar{B}	0.25	0.25
B	0.25	0.25

- Opt 1: Ignore the incomplete data
- Opt 2: Fill in a best guessed value
- Opt 3: Fill in with distribution

Example #1

A	B
0	0
0	0
0	0
0	0.5, 0 0.5, 1
0	1
1	0
1	1
1	1

Estimate parameters in joint distribution: θ

	\bar{A}	A
\bar{B}	0.25	0.25
B	0.25	0.25

- Opt 1: Ignore the incomplete data
- Opt 2: Fill in a best guessed value
- Opt 3: Fill in with distribution

Example #1

A	B
0	0
0	0
0	0
0	0.5, 0 0.5, 1
0	1
1	0
1	1
1	1

Estimate parameters in joint distribution: θ

	\bar{A}	A
\bar{B}	3.5/8	1/8
B	1.5/8	2/8

- Opt 1: Ignore the incomplete data
- Opt 2: Fill in a best guessed value
- Opt 3: Fill in with distribution

Example #1

A	B
0	0
0	0
0	0
0	0.7, 0 0.3, 1
0	1
1	0
1	1
1	1

Estimate parameters in joint distribution: θ

	\bar{A}	A
\bar{B}	3.5/8	1/8
B	1.5/8	2/8

- Opt 1: Ignore the incomplete data
- Opt 2: Fill in a best guessed value
- Opt 3: Fill in with distribution

Problem Setting

- observation: $x_i = (x_{i1}, x_{i2}, \dots, x_{iN}), i = 1 \dots k$
- modeling:
 - observable variables: $X = (X_1, X_2, \dots, X_N)$
 - hidden/latent variables: $Z = (Z_1, Z_2, \dots, Z_M)$
 - probability model: $p(x, z \mid \theta)$
- parameter estimation:
 - find the values for θ that maximize the *log likelihood*
$$l(\theta; x, z) \triangleq \log p(x, z \mid \theta)$$

Why observability matters?

- If Z is observable, we can maximize the *complete log likelihood*:

$$l_{\text{c}}(\theta; x, z) \triangleq \log p(x, z \mid \theta);$$

- If Z is hidden, we need to maximize the *incomplete log likelihood*:

$$l(\theta; x) \triangleq \log p(x \mid \theta) = \log \sum_z p(x, z \mid \theta).$$

Auxiliary Function $\mathcal{L}(q, \theta)$

an arbitrary distribution: $q(z \mid x)$

$$\begin{aligned} l(\theta; x) &= \log p(x \mid \theta) \\ &= \log \sum_z p(x, z \mid \theta) \\ &= \log \sum_z q(z \mid x) \frac{p(x, z \mid \theta)}{q(z \mid x)} \end{aligned}$$

Auxiliary Function $\mathcal{L}(q, \theta)$

an arbitrary distribution: $q(z \mid x)$

$$\begin{aligned} l(\theta; x) &= \log p(x \mid \theta) \\ &= \log \sum_z p(x, z \mid \theta) \\ &= \log \sum_z q(z \mid x) \frac{p(x, z \mid \theta)}{q(z \mid x)} \end{aligned}$$

Jensen's Inequality: For a real continuous concave function $f(x)$:

$$\sum_x p(x) f(x) \leq f \left(\sum_x p(x) x \right)$$

Auxiliary Function $\mathcal{L}(q, \theta)$

an arbitrary distribution: $q(z \mid x)$

$$\begin{aligned} l(\theta; x) &= \log p(x \mid \theta) \\ &= \log \sum_z p(x, z \mid \theta) \\ &= \log \sum_z q(z \mid x) \frac{p(x, z \mid \theta)}{q(z \mid x)} \\ &\geq \sum_z q(z \mid x) \log \frac{p(x, z \mid \theta)}{q(z \mid x)} \\ &\triangleq \mathcal{L}(q, \theta) \end{aligned}$$

Auxiliary Function $\mathcal{L}(q, \theta)$

an arbitrary distribution: $q(z \mid x)$

$$\begin{aligned} l(\theta; x) &= \log p(x \mid \theta) \\ &= \log \sum_z p(x, z \mid \theta) \\ &= \log \sum_z q(z \mid x) \frac{p(x, z \mid \theta)}{q(z \mid x)} \\ &\geq \sum_z q(z \mid x) \log \frac{p(x, z \mid \theta)}{q(z \mid x)} \\ &\triangleq \mathcal{L}(q, \theta) \end{aligned}$$

$\mathcal{L}(q, \theta)$ is a **lower bound** of $l(\theta; x)$

Expectation Step

With θ fixed, find q that maximizes $\mathcal{L}(q, \theta)$.

$$\mathcal{L}(q, \theta) = \sum_z q(z | x) \log \frac{p(x, z | \theta)}{q(z | x)}$$

Expectation Step

With θ fixed, find q that maximizes $\mathcal{L}(q, \theta)$.

$$\mathcal{L}(q, \theta) = \sum_z q(z \mid x) \log \frac{p(x, z \mid \theta)}{q(z \mid x)}$$

$$q^{(t+1)}(z \mid x) \Downarrow$$

$$\Downarrow p(z \mid x, \theta^{(t)})$$

Expectation Step

With θ fixed, find q that maximizes $\mathcal{L}(q, \theta)$.

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_z q(z | x) \log \frac{p(x, z | \theta)}{q(z | x)} \\ q^{(t+1)}(z | x) \\ &\Downarrow p(z | x, \theta^{(t)}) \\ \mathcal{L}(p(z | x, \theta^{(t)}), \theta^{(t)}) &= \sum_z p(z | x, \theta^{(t)}) \log \frac{p(x, z | \theta^{(t)})}{p(z | x, \theta^{(t)})} \\ &= \sum_z p(z | x, \theta^{(t)}) \log p(x | \theta^{(t)}) \\ &= \log p(x | \theta^{(t)}) \\ &= l(\theta^{(t)}; x)\end{aligned}$$

Maximization Step

With q fixed, find θ that maximizes $\mathcal{L}(q, \theta)$.

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_z q(z | x) \log \frac{p(x, z | \theta)}{q(z | x)} \\ &= \sum_z q(z | x) \log p(x, z | \theta) - \sum_z q(z | x) \log q(z | x) \\ &= \langle l_c(\theta; x, z) \rangle_q - \sum_z q(z | x) \log q(z | x)\end{aligned}$$

Maximization Step

With q fixed, find θ that maximizes $\mathcal{L}(q, \theta)$.

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_z q(z | x) \log \frac{p(x, z | \theta)}{q(z | x)} \\ &= \sum_z q(z | x) \log p(x, z | \theta) - \sum_z q(z | x) \log q(z | x) \\ &= \langle l_c(\theta; x, z) \rangle_q - \sum_z q(z | x) \log q(z | x)\end{aligned}$$

Just need to maximize $\langle l_c(\theta; x, z) \rangle_q$, which is no more difficult than maximizing $\log p(x, z | \theta)$, the complete log likelihood.

Maximization Step

With q fixed, find θ that maximizes $\mathcal{L}(q, \theta)$.

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_z q(z | x) \log \frac{p(x, z | \theta)}{q(z | x)} \\ &= \sum_z q(z | x) \log p(x, z | \theta) - \sum_z q(z | x) \log q(z | x) \\ &= \langle l_c(\theta; x, z) \rangle_q - \sum_z q(z | x) \log q(z | x)\end{aligned}$$

Just need to maximize $\langle l_c(\theta; x, z) \rangle_q$, which is no more difficult than maximizing $\log p(x, z | \theta)$, the complete log likelihood.

Now we have **improved** θ : 1) it leads to higher likelihood; and 2) q is ready to be further, better, estimated.

EM: an Iterative Process

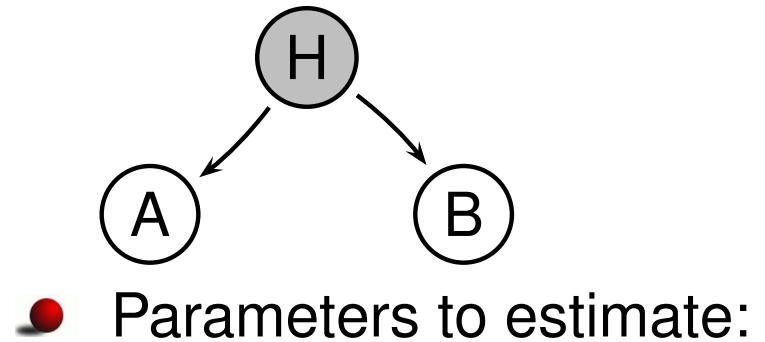
- Start with initial θ^0
- Repeat until convergence

E-Step: $p(z \mid x, \theta^{(t)}) \rightarrow q^{(t+1)} (= \operatorname{argmax}_q \mathcal{L}(q, \theta^{(t)}))$

M-Step: $\theta^{(t+1)} = \operatorname{argmax}_\theta \mathcal{L}(q^{(t+1)}, \theta)$

Example #2: MLE with hidden variables

A	B	#	$Pr(H^m \mid D^m, \theta)$
0	0	6	
0	1	1	
1	0	1	
1	1	4	



$$Pr(H) = ?$$

$$Pr(A \mid H) = ?$$

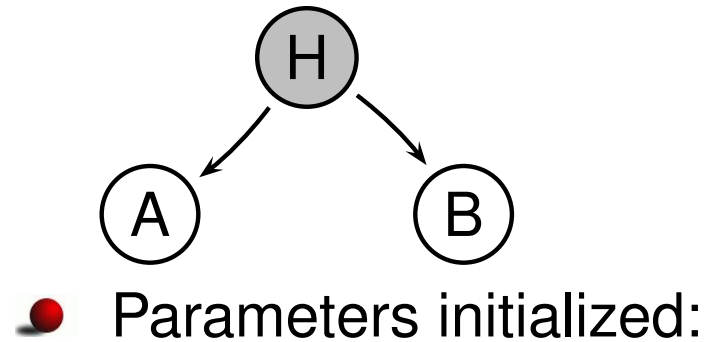
$$Pr(A \mid \overline{H}) = ?$$

$$Pr(B \mid H) = ?$$

$$Pr(B \mid \overline{H}) = ?$$

Example #2: MLE with hidden variables

A	B	#	$Pr(H^m \mid D^m, \theta)$
0	0	6	
0	1	1	
1	0	1	
1	1	4	



$$Pr(H) = 0.4$$

$$Pr(A \mid H) = 0.55$$

$$Pr(A \mid \overline{H}) = 0.61$$

$$Pr(B \mid H) = 0.43$$

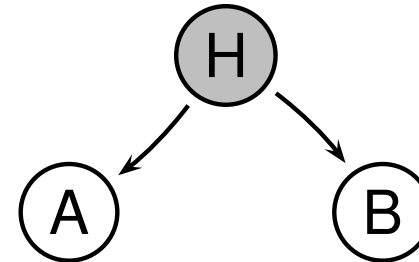
$$Pr(B \mid \overline{H}) = 0.52$$

Example #2: MLE with hidden variables

A	B	#	$Pr(H^m D^m, \theta)$
0	0	6	.48
0	1	1	.39
1	0	1	.42
1	1	4	.33

● Iteration 1: E-Step

$$\begin{aligned}
 Pr(H | A, B) &= \frac{Pr(A, B, H)}{Pr(A, B)} \\
 &= \frac{Pr(A, B | H)Pr(H)}{Pr(A, B)} \\
 &= \frac{Pr(A | H)Pr(B | H)Pr(H)}{Pr(A, B | H)Pr(H) + Pr(A, B | \bar{H})(1 - Pr(H))}
 \end{aligned}$$



● Parameters initialized:

$$Pr(H) = 0.4$$

$$Pr(A | H) = 0.55$$

$$Pr(A | \bar{H}) = 0.61$$

$$Pr(B | H) = 0.43$$

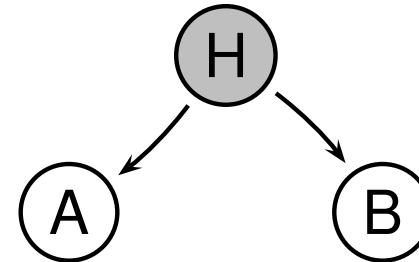
$$Pr(B | \bar{H}) = 0.52$$

Example #2: MLE with hidden variables

A	B	#	$Pr(H^m D^m, \theta)$
0	0	6	.48
0	1	1	.39
1	0	1	.42
1	1	4	.33

● Iteration 1: E-Step

$$\begin{aligned}
 Pr(H | A, B) &= \frac{Pr(A, B, H)}{Pr(A, B)} \\
 &= \frac{Pr(A, B | H)Pr(H)}{Pr(A, B)} \\
 &= \frac{Pr(A | H)Pr(B | H)Pr(H)}{Pr(A, B | H)Pr(H) + Pr(A, B | \bar{H})(1 - Pr(H))}
 \end{aligned}$$



● Iteration 1:
M-step
(parameters re-estimated):

$$Pr(H) = 0.42$$

$$Pr(A | H) = 0.35$$

$$Pr(A | \bar{H}) = 0.46$$

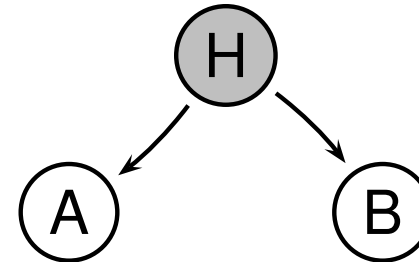
$$Pr(B | H) = 0.34$$

$$Pr(B | \bar{H}) = 0.47$$

Example #2: MLE with hidden variables

A	B	#	$Pr(H^m \mid D^m, \theta)$
0	0	6	.52
0	1	1	.39
1	0	1	.39
1	1	4	.28

● Iteration 2: E-Step



● Iteration 1:
M-step
(parameters re-estimated):

$$Pr(H) = 0.42$$

$$Pr(A \mid H) = 0.35$$

$$Pr(A \mid \overline{H}) = 0.46$$

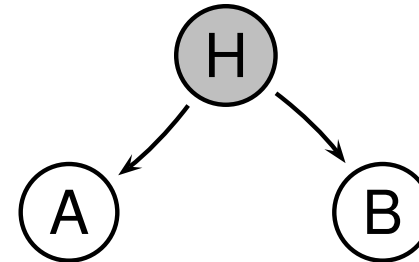
$$Pr(B \mid H) = 0.34$$

$$Pr(B \mid \overline{H}) = 0.47$$

Example #2: MLE with hidden variables

A	B	#	$Pr(H^m \mid D^m, \theta)$
0	0	6	.52
0	1	1	.39
1	0	1	.39
1	1	4	.28

● Iteration 2: E-Step



● Iteration 2:
M-step
(parameters re-estimated):

$$Pr(H) = 0.42$$

$$Pr(A \mid H) = 0.31$$

$$Pr(A \mid \overline{H}) = 0.50$$

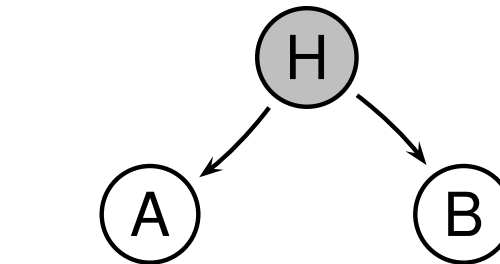
$$Pr(B \mid H) = 0.30$$

$$Pr(B \mid \overline{H}) = 0.50$$

Example #2: MLE with hidden variables

A	B	#	$Pr(H^m \mid D^m, \theta)$
0	0	6	.79
0	1	1	.31
1	0	1	.31
1	1	4	.05

● Iteration 5: E-Step



● Iteration 5:
M-step:

$$Pr(H) = 0.46$$

$$Pr(A \mid H) = 0.09$$

$$Pr(A \mid \overline{H}) = 0.69$$

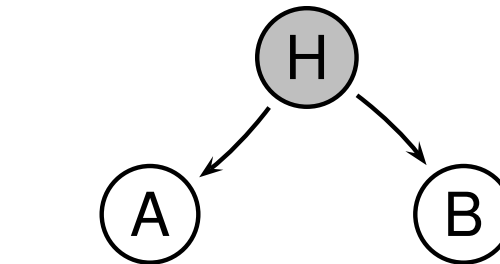
$$Pr(B \mid H) = 0.09$$

$$Pr(B \mid \overline{H}) = 0.69$$

Example #2: MLE with hidden variables

A	B	#	$Pr(H^m \mid D^m, \theta)$
0	0	6	.971
0	1	1	.183
1	0	1	.183
1	1	4	.001

● Iteration 10: E-Step



● Iteration 10:
M-step:

$$Pr(H) = 0.52$$

$$Pr(A \mid H) = 0.03$$

$$Pr(A \mid \overline{H}) = 0.83$$

$$Pr(B \mid H) = 0.03$$

$$Pr(B \mid \overline{H}) = 0.83$$

- Properties of EM:

- Properties of EM:
 - Each iteration increases the log likelihood;

- Properties of EM:

- Each iteration increases the log likelihood;
- Susceptible to local maximum.

Solution: Initialize $\theta^{(0)}$ randomly, or for multiple times.

- Properties of EM:
 - Each iteration increases the log likelihood;
 - Susceptible to local maximum.
Solution: Initialize $\theta^{(0)}$ randomly, or for multiple times.
- Commonly used for:

- Properties of EM:
 - Each iteration increases the log likelihood;
 - Susceptible to local maximum.
Solution: Initialize $\theta^{(0)}$ randomly, or for multiple times.
- Commonly used for:
 - mixture models;

- Properties of EM:
 - Each iteration increases the log likelihood;
 - Susceptible to local maximum.
Solution: Initialize $\theta^{(0)}$ randomly, or for multiple times.
- Commonly used for:
 - mixture models;
 - HMM;

- Properties of EM:
 - Each iteration increases the log likelihood;
 - Susceptible to local maximum.
Solution: Initialize $\theta^{(0)}$ randomly, or for multiple times.
- Commonly used for:
 - mixture models;
 - HMM;
 - models with hidden variables (why do we need hidden variables?), *etc.*

References

- [1] Leslie Pack Kaelbling, sma5504/MIT6.825 lecture notes (available on MIT OpenCourseWare:
ocw.mit.edu/OcwWeb/Electrical-Engineering-and-Computer-Science)
- [2] Michael I. Jordan, *An Introduction to Probabilistic Graphical Models* (Book Draft)