

EM algorithm

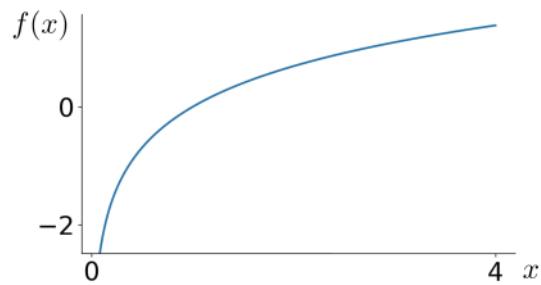
Saturday, June 23, 2018 4:55 PM



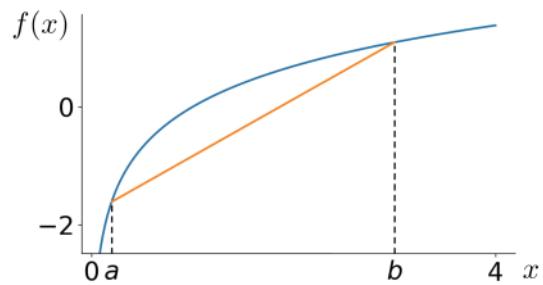
em1

General form of Expectation Maximization

Concave functions

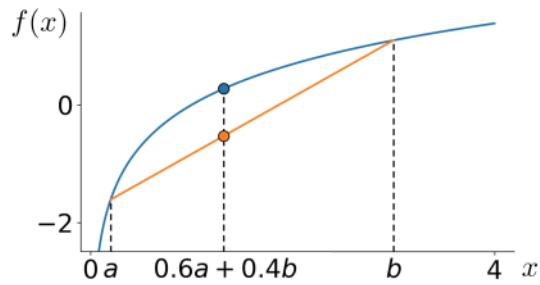


Concave functions



Concave functions

Concave function



Def.: $f(x)$ is concave if

for any $a, b, \alpha : f(\alpha a + (1 - \alpha)b) \geq \alpha f(a) + (1 - \alpha)f(b)$

$$0 \leq \alpha \leq 1$$

Jensen's inequality

2. Probability

$$f(\underbrace{\alpha a + (1 - \alpha)b}_{\text{expected value}}) \geq \alpha f(a) + (1 - \alpha)f(b)$$

Jensen's inequality

If $f(\alpha a + (1 - \alpha)b) \geq \alpha f(a) + (1 - \alpha)f(b)$

Then $\alpha_1 + \alpha_2 + \alpha_3 = 1; \alpha_k \geq 0.$

$$f(\alpha_1 a_1 + \alpha_2 a_2 + \alpha_3 a_3) \geq \alpha_1 f(a_1) + \alpha_2 f(a_2) + \alpha_3 f(a_3)$$

Jensen's inequality

If $f(\alpha a + (1 - \alpha)b) \geq \alpha f(a) + (1 - \alpha)f(b)$

Then $\alpha_1 + \alpha_2 + \alpha_3 = 1; \alpha_k \geq 0.$

$$f(\alpha_1 a_1 + \alpha_2 a_2 + \alpha_3 a_3) \geq \alpha_1 f(a_1) + \alpha_2 f(a_2) + \alpha_3 f(a_3)$$

$$p(t = a_1) = \alpha_1,$$

$$p(t = a_2) = \alpha_2,$$

$$p(t = a_3) = \alpha_3$$

Jensen's inequality

If $f(\alpha a + (1 - \alpha)b) \geq \alpha f(a) + (1 - \alpha)f(b)$

Then $\alpha_1 + \alpha_2 + \alpha_3 = 1; \alpha_k \geq 0.$

$$f(\underbrace{\alpha_1 a_1 + \alpha_2 a_2 + \alpha_3 a_3}_{\mathbb{E}_{p(t)} t}) \geq \underbrace{\alpha_1 f(a_1) + \alpha_2 f(a_2) + \alpha_3 f(a_3)}_{\mathbb{E}_{p(t)} f(t)}$$

$$p(t = a_1) = \alpha_1,$$

$$p(t = a_2) = \alpha_2,$$

$$p(t = a_3) = \alpha_3$$

\log

$| \# x^2$

Jensen's inequality

If $f(\alpha a + (1 - \alpha)b) \geq \alpha f(a) + (1 - \alpha)f(b)$

Then

Jensen's inequality:

$$f(\underbrace{\mathbb{E}_{p(t)} t}) \geq \underbrace{\mathbb{E}_{p(t)} f(t)}$$

Ininitely many points

Random variable x follows standard Gaussian distribution: $x \sim N(0, 1)$. A random variable y is a deterministic function of x : $y = 1 + x^2$. Apply Jensen's inequality to $\log(y)$ (logarithm is a concave function of its argument y) and choose all the correct inequalities from the options below:

$\mathbb{E} \log(1 + x^2) \geq \log(\mathbb{E}(1 + x^2))$

Un-selected is correct

Jensen's inequality is not applicable, since $\log(1 + x^2)$ is not a concave function of x

Un-selected is correct

$\mathbb{E} \log(1 + x^2) \leq \log(\mathbb{E}(1 + x^2))$

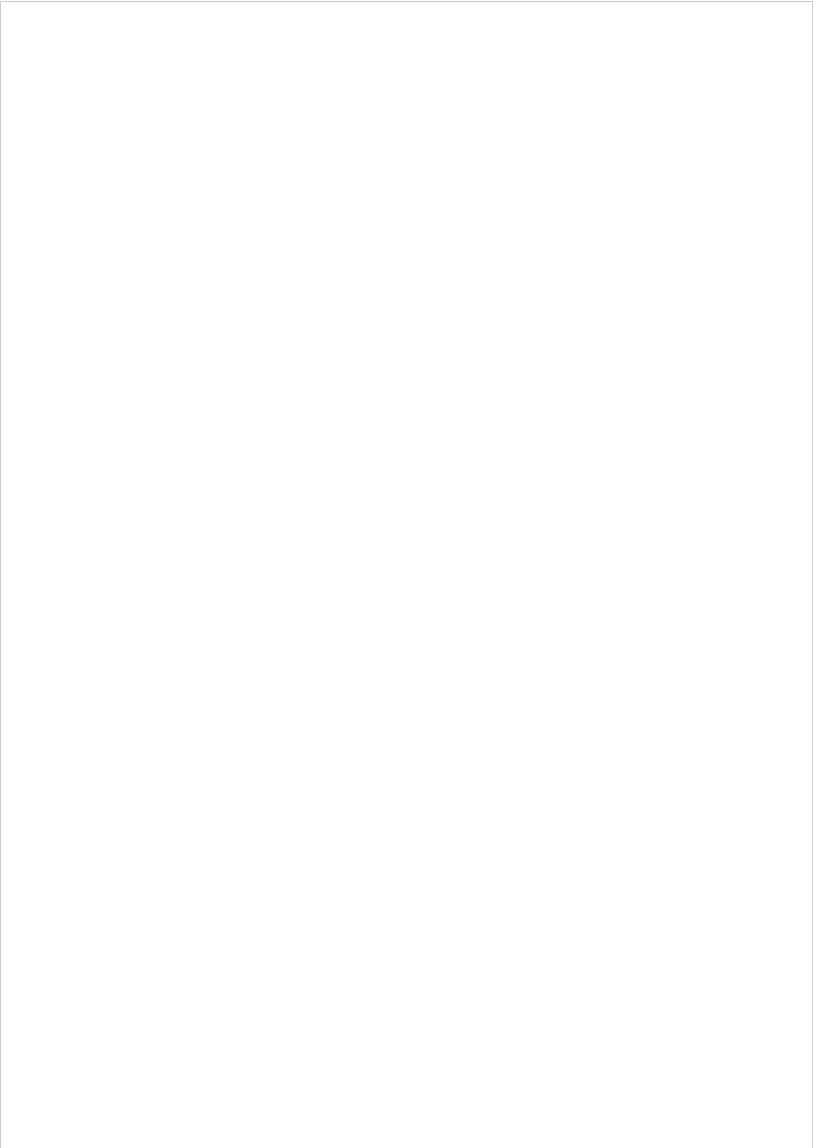
Correct

Here the order is changed, but the inequality sign is reversed as well, so it doesn't change anything. Also, instead of y we write $1 + x^2$, but they are identical.

$\text{lo}\sigma(\mathbb{E}_H) > \mathbb{E} \text{lo}\sigma(u)$

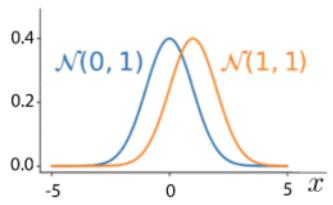
Kullback–Leibler divergence

Differences between
two probability distribution

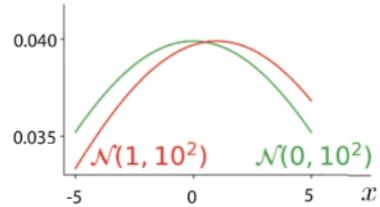


Kullback–Leibler divergence

Parameters difference: 1

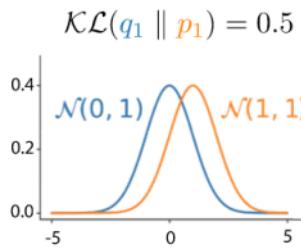


Parameters difference: 1

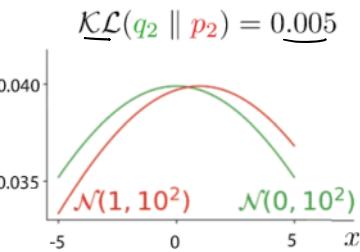


Kullback–Leibler divergence

Parameters difference: 1



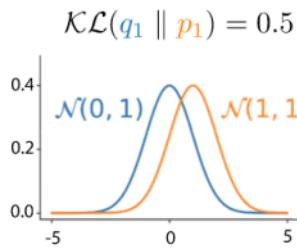
Parameters difference: 1



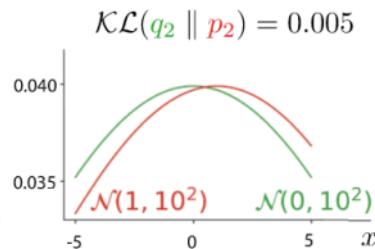
the second
set is closer
to each other
with this measure

Kullback–Leibler divergence

Parameters difference: 1



Parameters difference: 1



$$\underline{\mathcal{KL}(q \parallel p)} = \int \underbrace{q(x)}_{\text{the expected value of}} \log \frac{q(x)}{\underbrace{p(x)}_{\text{the difference between the two distribution}}} dx$$

the expected value of
the difference between the
two distribution

Property: It's non-symmetric. if change p, q, things will
be different

Kullback–Leibler divergence

$$\mathcal{KL}(q \parallel p) = \int q(x) \log \frac{q(x)}{p(x)} dx$$

Kullback–Leibler divergence

$$\mathcal{KL}(q \parallel p) = \int q(x) \log \frac{q(x)}{p(x)} dx$$

1. $\underbrace{\mathcal{KL}(q \parallel p) \neq \mathcal{KL}(p \parallel q)}$ 

Kullback–Leibler divergence

$$\mathcal{KL}(q \parallel p) = \int q(x) \log \frac{q(x)}{p(x)} dx$$

1. $\mathcal{KL}(q \parallel p) \neq \mathcal{KL}(p \parallel q)$
2. $\mathcal{KL}(q \parallel q) = 0$ $\int q(x) \log 1 dx = 0$

Kullback–Leibler divergence

$$\mathcal{KL}(q \parallel p) = \int q(x) \log \frac{q(x)}{p(x)} dx$$

1. $\mathcal{KL}(q \parallel p) \neq \mathcal{KL}(p \parallel q)$
2. $\mathcal{KL}(q \parallel \textcolor{red}{q}) = 0$
3. $\mathcal{KL}(q \parallel p) \geq 0$

Kullback–Leibler divergence

$$\mathcal{KL}(q \parallel p) = \int q(x) \log \frac{q(x)}{p(x)} dx$$

1. $\mathcal{KL}(q \parallel p) \neq \mathcal{KL}(p \parallel q)$
2. $\mathcal{KL}(q \parallel \textcolor{red}{q}) = 0$
3. $\mathcal{KL}(q \parallel p) \geq 0$

Proof: $-\mathcal{KL}(q \parallel p) = \mathbb{E}_q \left(-\log \frac{q}{p} \right)$

Kullback–Leibler divergence

$$\mathcal{KL}(q \parallel p) = \int q(x) \log \frac{q(x)}{p(x)} dx$$

1. $\mathcal{KL}(q \parallel p) \neq \mathcal{KL}(p \parallel q)$
2. $\mathcal{KL}(q \parallel \textcolor{red}{q}) = 0$
3. $\mathcal{KL}(q \parallel p) \geq 0$

Proof: $-\mathcal{KL}(q \parallel p) = \mathbb{E}_q \left(-\log \frac{q}{p} \right) = \mathbb{E}_q \underbrace{\left(\log \frac{p}{q} \right)}_{\sim}$

Kullback–Leibler divergence

$$\mathcal{KL}(q \parallel p) = \int q(x) \log \frac{q(x)}{p(x)} dx$$

1. $\mathcal{KL}(q \parallel p) \neq \mathcal{KL}(p \parallel q)$
2. $\mathcal{KL}(q \parallel \textcolor{red}{q}) = 0$
3. $\mathcal{KL}(q \parallel p) \geq 0$

Proof: $-\mathcal{KL}(q \parallel p) = \mathbb{E}_q \left(-\log \frac{q}{p} \right) = \mathbb{E}_q \left(\log \frac{p}{q} \right)$

$\leq \log(\mathbb{E}_q \frac{p}{q})$

$\xrightarrow{\text{Jensen Inequality}}$

\downarrow
concave function

Kullback–Leibler divergence

$$\mathcal{KL}(q \parallel p) = \int q(x) \log \frac{q(x)}{p(x)} dx$$

1. $\mathcal{KL}(q \parallel p) \neq \mathcal{KL}(p \parallel q)$

2. $\mathcal{KL}(q \parallel q) = 0$

~~3.~~ 3. $\mathcal{KL}(q \parallel p) \geq 0$

Proof: $-\mathcal{KL}(q \parallel p) = \mathbb{E}_q \left(-\log \frac{q}{p} \right) = \mathbb{E}_q \left(\log \frac{p}{q} \right)$

$$\leq \log(\mathbb{E}_q \frac{p}{q}) = \log \int q(x) \frac{p(x)}{q(x)} dx = 0$$

$$\begin{aligned}
 q(x) &= N(\mu, \sigma_2^2) & p(x) &= N(\mu, \sigma_1^2) \\
 \mathcal{KL}(q \parallel p) &= \int q(x) \log \frac{q(x)}{p(x)} dx \\
 &= \int \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp \left\{ \frac{(x-\mu)^2}{2\sigma_2^2} \right\} \log \frac{\frac{1}{\sqrt{2\pi\sigma_2^2}} \exp \left\{ \frac{(x-\mu)^2}{2\sigma_2^2} \right\}}{\frac{1}{\sqrt{2\pi\sigma_1^2}} \exp \left\{ \frac{(x-\mu)^2}{2\sigma_1^2} \right\}} dx \\
 &= \int \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp \left\{ \frac{(x-\mu)^2}{2\sigma_2^2} \right\} \log \left[\frac{\sigma_1}{\sigma_2} \exp \left\{ \frac{(x-\mu)^2}{2\sigma_1^2 \sigma_2^2} \right\} \right] dx \\
 &= \int \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp \left\{ \frac{(x-\mu)^2}{2\sigma_2^2} \right\} \left[\log \frac{\sigma_1}{\sigma_2} + \frac{(x-\mu)^2 (\sigma_1^2 + \sigma_2^2) - (\sigma_1 - \sigma_2)^2}{2\sigma_1^2 \sigma_2^2} \right] dx \\
 &= \underbrace{\log \frac{\sigma_1}{\sigma_2} \times 1}_{\text{Variance} \leq \sigma^2} + \underbrace{\int \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp \left\{ \frac{(x-\mu)^2}{2\sigma_2^2} \right\} (x-\mu)^2 dx}_{\times \frac{\sigma_1^2 - \sigma_2^2}{2\sigma_1^2 \sigma_2^2}} \\
 &\approx \log \frac{\sigma_1}{\sigma_2} + \frac{\sigma_1^2 - \sigma_2^2}{2\sigma_1^2} = \log \frac{\sigma_1}{\sigma_2} + \frac{1}{2} - \frac{\sigma_2^2}{2\sigma_1^2}
 \end{aligned}$$

Kullback–Leibler divergence

$$\mathcal{KL}(q \parallel p) = \int q(x) \log \frac{q(x)}{p(x)} dx$$

Summary

A way to compare distributions
not a proper distance

1. $\mathcal{KL}(q \parallel p) \neq \mathcal{KL}(p \parallel q)$
2. $\mathcal{KL}(q \parallel q) = 0$
3. $\mathcal{KL}(q \parallel p) \geq 0$

\mathcal{KL} 不对称
 $\mathcal{KL}(q \parallel p) \neq \mathcal{KL}(p \parallel q)$

General form of Expectation Maximization

General form of Expectation Maximization



General form of Expectation Maximization



$$p(x_i \mid \theta) = \sum_{c=1}^3 p(x_i \mid t_i = c, \theta) p(t_i = c \mid \theta)$$

General form of Expectation Maximization

$$\max_{\theta} \quad p(X \mid \theta)$$

General form of Expectation Maximization

$$\max_{\theta} \log p(X | \theta)$$

General form of Expectation Maximization

$$\max_{\theta} \log p(X | \theta) = \log \prod_{i=1}^N p(x_i | \theta)$$

General form of Expectation Maximization

$$\begin{aligned}\max_{\theta} \log p(X \mid \theta) &= \log \prod_{i=1}^N p(x_i \mid \theta) \\ &= \sum_{i=1}^N \log p(x_i \mid \theta)\end{aligned}$$

General form of Expectation Maximization

$$\log p(X \mid \theta) = \sum_{i=1}^N \log p(x_i \mid \theta)$$

General form of Expectation Maximization

$$\begin{aligned}\log p(X \mid \theta) &= \sum_{i=1}^N \log p(x_i \mid \theta) \\ &= \sum_{i=1}^N \log \sum_{\substack{c=1 \\ z}}^3 p(x_i, t_i = c \mid \theta)\end{aligned}$$

General form of Expectation Maximization

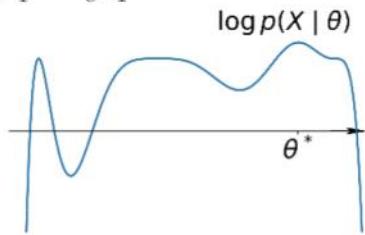
$$\begin{aligned}\log p(X \mid \theta) &= \sum_{i=1}^N \log p(x_i \mid \theta) \\ &= \sum_{i=1}^N \log \sum_{c=1}^3 p(x_i, t_i = c \mid \theta) \geq \mathcal{L}(\theta)\end{aligned}$$

Jensen's inequality



General form of Expectation Maximization

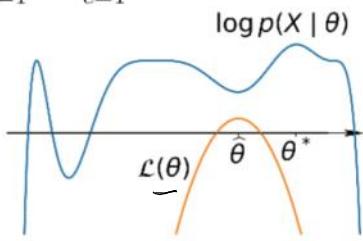
$$\begin{aligned}\log p(X \mid \theta) &= \sum_{i=1}^N \log p(x_i \mid \theta) \\ &= \sum_{i=1}^N \log \sum_{c=1}^3 p(x_i, t_i = c \mid \theta) \geq \mathcal{L}(\theta)\end{aligned}$$



General form of Expectation Maximization

$$\log p(X | \theta) = \sum_{i=1}^N \log p(x_i | \theta)$$

$$= \sum_{i=1}^N \log \sum_{c=1}^3 p(x_i, t_i = c | \theta) \geq \underline{\mathcal{L}(\theta)}$$



instead of
one lower bound
 \Rightarrow choose a good
family of lower bound.

General form of Expectation Maximization

$$\begin{aligned}\log p(X \mid \theta) &= \sum_{i=1}^N \log p(x_i \mid \theta) \\ &= \sum_{i=1}^N \log \sum_{c=1}^3 p(x_i, t_i = c \mid \theta)\end{aligned}$$

General form of Expectation Maximization

$$\begin{aligned}\log p(X \mid \theta) &= \sum_{i=1}^N \log p(x_i \mid \theta) \\ &= \sum_{i=1}^N \log \sum_{c=1}^3 \frac{q(t_i = c)}{q(t_i = c)} p(x_i, t_i = c \mid \theta)\end{aligned}$$

General form of Expectation Maximization

$$\begin{aligned}\log p(X \mid \theta) &= \sum_{i=1}^N \log p(x_i \mid \theta) \\ &= \sum_{i=1}^N \log \sum_{c=1}^3 \frac{q(t_i = c)}{q(t_i = c)} p(x_i, t_i = c \mid \theta)\end{aligned}$$

treat q as weights
of Jensen inequality

Jensen's inequality

$$\log \left(\sum_c \alpha_c v_c \right) \geq \sum_c \alpha_c \log(v_c)$$

General form of Expectation Maximization

$$\begin{aligned}\log p(X | \theta) &= \sum_{i=1}^N \log p(x_i | \theta) \\ &= \sum_{i=1}^N \log \sum_{c=1}^3 \frac{q(t_i = c)}{q(t_i = c)} p(x_i, t_i = c | \theta) \\ &\geq \sum_{i=1}^N \sum_{c=1}^3 \underbrace{q(t_i = c)}_{\text{Jensen's inequality}} \log \frac{p(x_i, t_i = c | \theta)}{q(t_i = c)}\end{aligned}$$

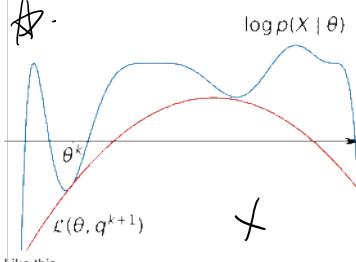
$$\log \left(\sum_c \alpha_c v_c \right) \geq \sum_c \alpha_c \log(v_c)$$

we get

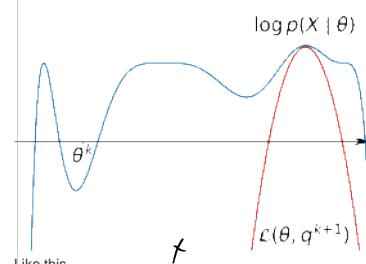
How do you think the optimal lower bound $L(\theta, q^{k+1})$ will look like? Assume that q_k is the solution of the following optimization problem:

$$\max_q L(\theta, q)$$

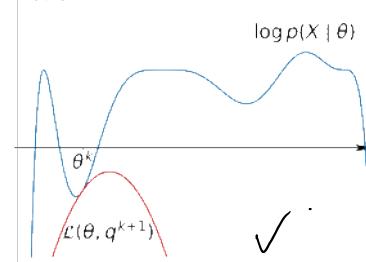
Like this



Like this



Like this



Correct

From <https://www.coursera.org/learn/bayesian-methods-in-machine-learning/lecture/Fm3mY/expectation-maximization-algorithm>

$$\log \sum_v q(t_i = v) \frac{p}{q_i}$$

$\left(\log q_1 \cdot \frac{p}{q_1} + \log q_2 \cdot \frac{p_2}{q_2} + \log q_3 \cdot \frac{p_3}{q_3} \right)$

General form of Expectation Maximization

$$\begin{aligned}\log p(X | \theta) &= \sum_{i=1}^N \log p(x_i | \theta) \\ &= \sum_{i=1}^N \log \sum_{c=1}^3 \frac{q(t_i = c)}{q(t_i = c)} p(x_i, t_i = c | \theta) \\ &\geq \sum_{i=1}^N \sum_{c=1}^3 \underbrace{q(t_i = c)}_{\text{Let } q \text{ be the distribution}} \log \frac{p(x_i, t_i = c | \theta)}{q(t_i = c)} \\ &= \mathcal{L}(\theta, q)\end{aligned}$$

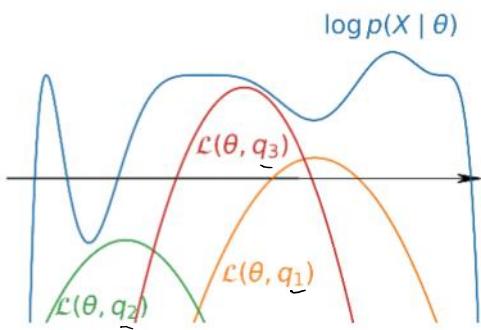
$$= \mathcal{L}(\theta, q)$$

or $\underbrace{\dots}_{\text{why this?}}$

$$\log p(X | \theta) = \sum_i (\log q_i) \cdot \frac{p_i}{q_i}$$

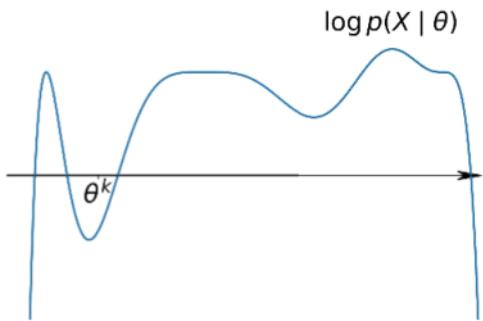
General form of Expectation Maximization

$$\log p(X | \theta) \geq \mathcal{L}(\theta, \underline{q}) \text{ for any } \underline{q}$$



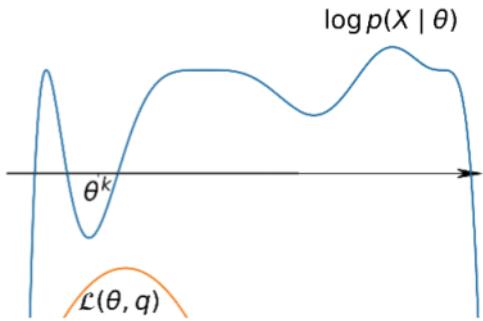
General form of Expectation Maximization

$\log p(X | \theta) \geq \mathcal{L}(\theta, q)$ for any q



General form of Expectation Maximization

$\log p(X | \theta) \geq \mathcal{L}(\theta, q)$ for any q

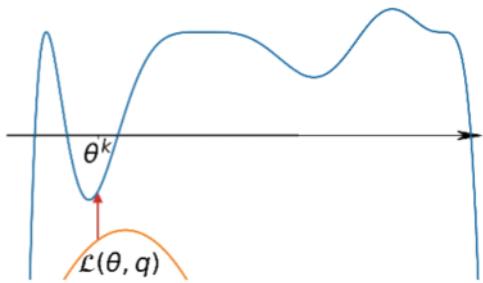


General form of Expectation Maximization

$\log p(X | \theta) \geq \mathcal{L}(\theta, q)$ for any q

$$q^{k+1} = \arg \max_q \mathcal{L}(\theta^k, q)$$

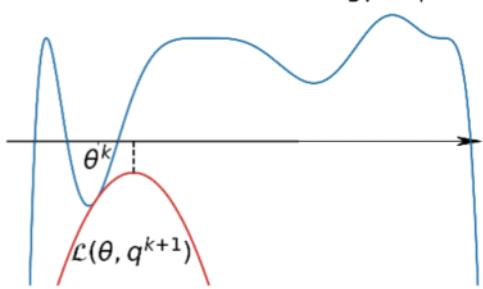
$\log p(X | \theta)$



General form of Expectation Maximization

$\log p(X | \theta) \geq \mathcal{L}(\theta, q)$ for any q

$$q^{k+1} = \arg \max_q \mathcal{L}(\theta^k, q)$$

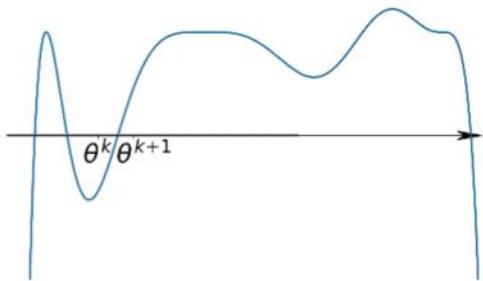


General form of Expectation Maximization

$\log p(X | \theta) \geq \mathcal{L}(\theta, q)$ for any q

$$q^{k+1} = \arg \max_q \mathcal{L}(\theta^k, q)$$

$\log p(X | \theta)$

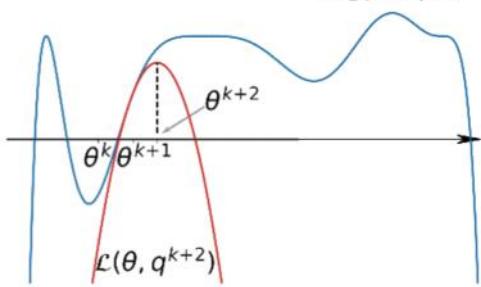


General form of Expectation Maximization

$$\log p(X | \theta) \geq \mathcal{L}(\theta, q) \text{ for any } q$$

$$q^{k+1} = \arg \max_q \mathcal{L}(\theta^k, q)$$

$\int \cdot$ variational
(lower bound)



Summary of Expectation Maximization

$$\log p(X | \theta) \geq \mathcal{L}(\theta, q) \text{ for any } q$$

Variational lower bound

E-step

$$q^{k+1} = \arg \max_q \mathcal{L}(\theta^k, q)$$

M-step

$$\theta^{k+1} = \arg \max_{\theta} \mathcal{L}(\theta, q^{k+1})$$

fix θ , find $(q_i)_c$ (distribution)

fix $\theta - q$, find θ

em-e

E-step details

$\text{GAP} = \log p(x|\theta) - \mathcal{L}(\theta, q)$

$$\sum_{t=1}^T \sum_{c=1}^3 q(t_i=c) \log \frac{p(x_i, t_i=c | \theta)}{q(t_i=c)}$$

$\text{Want min}_q \text{GAP}$

$$= \sum_{t=1}^T \sum_{c=1}^3 q(t_i=c) \log \left[\log p(x_i | \theta) - \log \frac{p(x_i, t_i=c | \theta)}{q(t_i=c)} \right]$$

$$= \sum_{t=1}^T \sum_{c=1}^3 q(t_i=c) \log \frac{p(x_i | \theta) q(t_i=c)}{p(x_i, t_i=c | \theta)}$$

Want min_q GAP
 minimize the sum of KL divergence of KL properties of KL

$$\begin{aligned}
 & \min_q \sum_{i=1}^N KL(q(t_i) || p(t_i | x_i, \theta)) \\
 & \Rightarrow q(t_i) = p(t_i | x_i, \theta)
 \end{aligned}$$

 \begin{aligned}
 & = \sum_{i=1}^N \sum_{c=1}^C q(t_i=c) \log \frac{p(x_i | \theta) q(t_i=c)}{p(x_i, t_i=c | \theta)} \\
 & = \sum_{i=1}^N \sum_{c=1}^C q(t_i=c) \log \frac{q(t_i=c)}{p(t_i=c | x_i, \theta)} \\
 & \quad \underbrace{p(t_i=c | x_i, \theta) p(x_i | \theta)}_{KL(q(t_i) || p(t_i | x_i, \theta))}
 \end{aligned}

E-step details

$$\log p(X | \theta) \geq \mathcal{L}(\theta, q)$$

E-step details

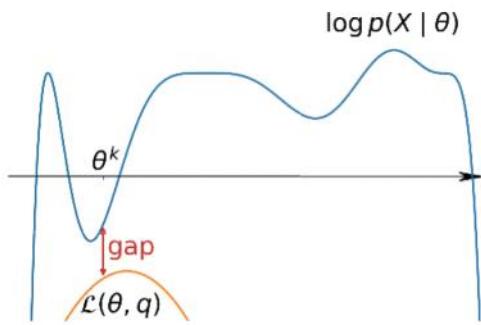
$$\log p(X \mid \theta) \geq \mathcal{L}(\theta, q)$$

E-step: $\max_q \mathcal{L}(\theta^k, q)$

E-step details

$$\log p(X \mid \theta) \geq \mathcal{L}(\theta, q)$$

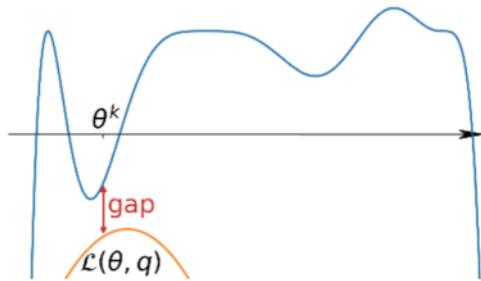
E-step: $\max_q \mathcal{L}(\theta^k, q)$



E-step summary

$$\log p(X \mid \theta) - \mathcal{L}(\theta, q) = \sum_i \mathcal{KL}(q(t_i) \parallel p(t_i \mid x_i, \theta))$$

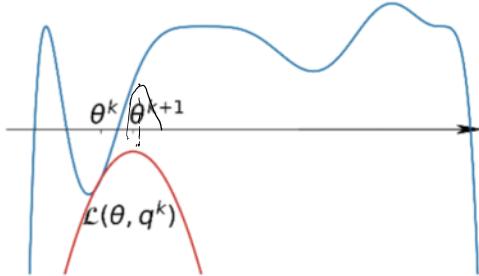
$$\text{E-step: } \arg \max_{q(t_i)} \mathcal{L}(\theta^k, q) = p(t_i \mid x_i, \theta)$$



E-step summary

$$\log p(X | \theta) - \mathcal{L}(\theta, q) = \sum_i \mathcal{KL}(q(t_i) \| p(t_i | x_i, \theta))$$

E-step: $\arg \max_{q(t_i)} \mathcal{L}(\theta^k, q) = p(t_i | x_i, \theta)$



M-step

$$\begin{aligned} f(\theta, q) &= \sum_i \sum_c q(t_i=c) \log \frac{p(x_i, t_i=c | \theta)}{p(t_i=c)} \\ &= \sum_i \sum_c q(t_i=c) \log p(x_i, t_i=c | \theta) \\ &\quad - \sum_i \sum_c q(t_i=c) \log q(t_i=c) \\ &= \mathbb{E}_q \underbrace{\log(p(x, t | \theta))}_{\text{concave function wrt. } \theta.} + \text{constant} \xrightarrow{\text{No. } \theta.} \end{aligned}$$

E-step

$$q^{k+1} = \arg \min_q \mathcal{KL}[q(t) \| p(t | X, \theta^k)]$$

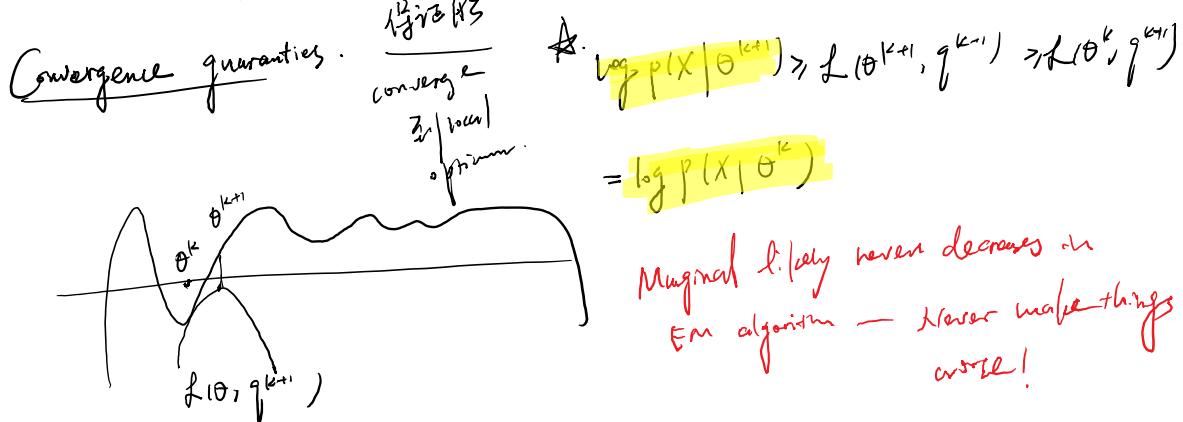
$$\Rightarrow q^{k+1}(t_i) = p(t_i | x_i, \theta^k)$$

$$q^{k+1} = \arg \min_q L(\theta^k, q)$$

$$\Leftrightarrow q^{k+1}(t_i) = p(t_i | x_i, \theta^k)$$

M-step

$$\theta^{k+1} = \arg \max_{\theta} E_{q^{k+1}} \log p(X, T | \theta)$$



Can we compare value of $\mathcal{L}(\theta^k, q^{k+1})$ with value of log-likelihood $\log p(X | \theta^k)$ in the general case?

- $\mathcal{L}(\theta^k, q^{k+1}) \geq \log p(X | \theta^k)$
- $\mathcal{L}(\theta^k, q^{k+1}) = \log p(X | \theta^k)$

Correct

Yes, since for any q the following inequality holds true $\mathcal{L}(\theta^k, q) \leq \log p(X | \theta^k)$, but for q^{k+1} is the maximizer of the lower bound $\mathcal{L}(\theta^k, q^{k+1}) = \max_q \mathcal{L}(\theta^k, q)$ and in the best case it became as large as $\log p(X | \theta^k)$

- $\mathcal{L}(\theta^k, q^{k+1}) \leq \log p(X | \theta^k)$

- No, the inequality sign depends on the particular model and the data at hand

Always improving
A debugging tool. If decreasing
 \Rightarrow wrong!

EM example E-step

data, discrete distribution, 1, 2, 3

know: the data is generated from a mixture of distributions



$$p(x_i) = \gamma p_1(x_i) + (1-\gamma) p_2(x_i)$$

	1	2	3
p_1	α	$1-\alpha$	0
p_2	0	β	β

Task: estimate α and β . with the EM algorithm

$$\text{Initiate: } \alpha_0 = \beta_0 = \gamma_0 = 0.5$$

$(x_i) \leftarrow (t_i)$ Define latent variable.

$$\text{if } t_i = 1 \text{ then } \gamma$$

$$p(x_i | t_i = 2) = P_2(x_i)$$

$(x_i) \leftarrow (t_i)$ Define latent variable

$$p(t_i=1) = \gamma, \quad p(x_i | t_i=1) = P_1(x_i)$$

[E-step]: Task: find the posterior distribution of $q(t_i)$

$$q(t_i=c) = p(t_i=c | x_i)$$

$$p(t_i=1 | x_i=1) = \frac{p(x_i=1 | t_i=1) \times p(t_i=1)}{p(x_i=1 | t_i=1) p(t_i=1) + p(x_i=1 | t_i=2) p(t_i=2)}$$

$$= \frac{\alpha \gamma}{\alpha \gamma + (1-\alpha)(1-\gamma)} = 1 \quad \text{if observe } x_i=1, \text{ it's certainly generated by } t_i=1$$

$$p(t_i=1 | x_i=2) = \frac{p(x_i=2 | t_i=1) \times p(t_i=1)}{p(x_i=2 | t_i=1) p(t_i=1) + p(x_i=2 | t_i=2) p(t_i=2)}$$

$$= \frac{(1-\alpha)\gamma}{(1-\alpha)\gamma + (1-\beta)(1-\gamma)} = \frac{0.5 \times 0.5}{0.5 \times 0.5 + 0.5 \times 0.5} = 0.5$$

$$q(t_i=1) = p(t_i=1 | x_i) = \begin{cases} 1 & x_i=1 \\ 0.5 & x_i=2 \\ 0 & x_i=3 \end{cases}$$

$$q(t_i=2) = 1 - q(t_i=1)$$

[M-step] Want to maximize the expected value

$$\begin{aligned} & \underset{\alpha, \beta, \gamma}{\operatorname{max}} \sum_{i=1}^N \mathbb{E}_{q(t_i)} [\ell \cdot q p(x_i | t_i) p(t_i)] \\ &= \sum_{i=1}^N q(t_i=1) \log p(x_i | t_i) p(t_i) + \sum_{i=1}^N q(t_i=2) \log p(x_i | t_i) p(t_i) \end{aligned}$$

Assume data: $N_1 = 30, N_2 = 20, N_3 = 60$ [Observation. # of values in each category]

$$= 30 \times p(t_i=1 | x_i=1) \log \alpha \gamma + 20 \times 0.5 \times \log(1-\alpha)\gamma$$

$$+ 60 \times 0.5 \times \log(1-\gamma)$$

$$+ 20 \times p(t_i=2 | x_i=1) \log \dots$$

$$+ 20 \times 0.5 \log(1-\beta)(1-\gamma)$$

$$+ 60 \times 1 \times \log \beta(1-\gamma)$$

$$\propto 30 \log \alpha + 10 \log(1-\alpha) + \text{const}(\alpha)$$

$$\frac{\partial}{\partial \alpha} = 30 \frac{1}{\alpha} + 10 \frac{1 \times (-1)}{1-\alpha} = 0 \quad \alpha = \frac{30}{40} = \frac{3}{4}$$

$$\beta = \frac{6}{7} \quad \gamma = \frac{4}{11}$$

$$\frac{\partial}{\partial \alpha} = 30 \frac{1}{\alpha} + 10 \frac{1 \times 11}{1-\alpha} - 0 \quad \alpha = \frac{1}{40} = \frac{2}{7}$$

$$\beta = \frac{6}{7} \quad \gamma = \frac{4}{11}$$

Summary of the Expectation-Maximization Algorithm

 em-summary

Summary of Expectation Maximization

Summary of Expectation Maximization

- Method for training Latent Variable Models



Summary of Expectation Maximization

- Method for training Latent Variable Models
- Handles missing data

	High school grade	University grade	IQ score	Phone Interview
<i>John</i>	4.0	4.0	120	3/4
<i>Helen</i>	3.7	3.6	N/A	4/4
<i>Jack</i>	3.2	N/A	112	2/4
<i>Emma</i>	2.9	3.2	N/A	3/4

Summary of Expectation Maximization

- Method for training Latent Variable Models
- Handles missing data
- Sequence of simple task instead of one hard

Summary of Expectation Maximization

- Method for training Latent Variable Models
- Handles missing data
- Sequence of simple task instead of one hard
- Guarantees to converge

Summary of Expectation Maximization

- Method for training Latent Variable Models
- Handles missing data
- Sequence of simple task instead of one hard
- Guarantees to converge
- Helps with complicated parameter constraints

$$\Sigma_c \succ 0$$

Would you use EM algorithm to fit a Gaussian to one-dimensional data (i.e. estimate the parameters μ and σ)? What about multidimensional data (estimating the mean vector μ and the covariance matrix σ)?

- Yes, both for one-dimensional and multidimensional data
 Only for multi-dimensional data

Correct

We can treat missing values as latent variables and still estimate the Gaussian parameters. But with one-dimensional data, if a data point has missing values, it means that we don't know anything about it (its only dimension is missing) and not even the smartest latent variable model can extract information from a point like this. So the only thing that is left is to throw away points with missing data.

Note that we also don't need EM to estimate the mean vector (e.g. we need it only for the covariance matrix) in the multi-dimensional case: since each coordinate of the mean vector can be treated independently, we can treat each coordinate as one-dimensional case and just throw away missing values.

Summary of Expectation Maximization

- Method for training Latent Variable Models
- Handles missing data
- Sequence of simple task instead of one hard
- Guarantees to converge
- Helps with complicated parameter constraints can solve M-step analytically.
- Numerous extensions:
 - Variational E-step: restrict the set of possible q Approximation.
(week 3 and 5)
 - Sampling on M-step (week 4) for complicated distributions

Summary of Expectation Maximization

- Method for training Latent Variable Models
- Handles missing data
- Sequence of simple task instead of one hard
- Guarantees to converge
- Helps with complicated parameter constraints
- Numerous extensions

Cons

- Only local maximum (or saddle point)
- Requires math :)