

PMML - RII / EM

PGM: Representation / Inference / Learning

Complete Data / Missing Data

Parameters / Structure

Complete Data / Missing Data

The EM idea: if we know the value of the missing data we could find MLE $\hat{\theta}$, maximizing the log-likelihood.

E: Try to complete (guess) the missing data based on observed data and current params $\theta^{(t)}$,
 M: Find new params $\theta^{(t+1)}$ maximizing the log-likelihood w. the filled-in data.

Inference

* We do it in a soft manner, computing and maximizing the expected log likelihood w.r.t. θ .

\uparrow
 w.r.t. $P(H|D, \theta^{(t)})$
 \uparrow
 expected sufficient statistics.

Formally:

1. Initialize $\theta^{(0)}$ to some values

2. Calculate

$$Q(\theta; \theta^{(t)}) = \sum_m \sum_{x_h^{(m)}} P(x_h^{(m)} | x_o^{(m)}; \theta^{(t)}) \cdot$$

3. Find $\theta^{(t+1)} = \underset{\theta}{\operatorname{argmax}} Q(\theta; \theta^{(t)})$

$$= \mathbb{E} \log P(x_h | x_o; \theta) \\ P(x_h | x_o; \theta^{(t)})$$

The expected complete-data log-likelihood w.r.t the missing data posterior (using current parameters).

EM - Proof of convergence

(2)

(show that $l(\theta^{(t+1)}; D) \geq l(\theta^{(t)}; D)$)

For any $q(x_h)$ (and specifically $q(x_h) = P(x_h | x_0; \theta^{(t)})$)

for a single sample
 $\leftarrow x_0, x_h \in D$

$$\begin{aligned}
 & \rightarrow \sum_{x_h} q(x_h) \log P(x_h, x_0; \theta) \\
 & = \sum_{x_h} q(x_h) \log [P(x_0; \theta) P(x_h | x_0; \theta)] \quad \leftarrow \text{chain rule} \\
 & = \sum_{x_h} q(x_h) \log P(x_0; \theta) + \sum_{x_h} q(x_h) \log P(x_h | x_0; \theta) \\
 & = \log P(x_0; \theta) + \sum_{x_h} q(x_h) \log \frac{P(x_h | x_0; \theta)}{q(x_h)} + \sum_{x_h} q(x_h) \log q(x_h) \\
 & = l(x_0; \theta) + \text{KL}(q(x_h) || P(x_h | x_0; \theta)) - H(q(x_h))
 \end{aligned}$$

is also possible

$Q(\theta^{(t+1)}; \theta^{(t)}) \geq Q(\theta^{(t)}; \theta^{(t)})$ by maximization.

$$\begin{aligned}
 0 & \leq Q(\theta^{(t+1)}; \theta^{(t)}) - Q(\theta^{(t)}; \theta^{(t)}) \\
 & = l(D; \theta^{(t+1)}) - l(D; \theta^{(t)}) \\
 & \quad - \text{KL}(P(H|D); \theta^{(t+1)}) || P(H|D; \theta^{(t)}) + \text{KL}(P(H|D; \theta^{(t)}) || P(H|D; \theta^{(t)})) \\
 & \quad - \cancel{H(q(H))} + \cancel{H(q(H))} \quad \leftarrow \text{same}
 \end{aligned}$$

$l(D; \theta^{(t+1)}) \geq l(D; \theta^{(t)}) + c$
 \uparrow
 ≥ 0 □

Note: $Q(\theta; \theta^{(t)}) = \sum_m \sum_{x_h \in D} P(x_h \in D | x_0 \in D; \theta^{(t)}) \log P(x_h \in D, x_0 \in D; \theta)$

(In $Q(\theta^{(t+1)}; \theta^{(t)})$ and $Q(\theta^{(t)}; \theta^{(t)})$, $q(H)$ is the same and θ changes.)

$$= E_q \log P(\theta; D^+) = E_q l(\theta; D^+)$$

Exercise: Incremental EM - CPD at a time ③

Instead of updating all params in each M step, we select one parameter x_i and update only its CPD: $\theta_{x_i | p_i}$. Does the likelihood increase in each iteration? i.e.

$$l(\theta^{(t+1)}; D) \geq l(\theta^{(t)}; D)$$

prove your answer.

Answer: Yes

Proof:

We saw that:

$$l(\theta^{(t+1)}; D) - l(\theta^{(t)}; D) \geq \mathbb{E}_{P(D; \theta^{(t)})} [l(\theta^{(t+1)}; D^*)] -$$

$$\mathbb{E}_{P(D; \theta^{(t)})} [l(\theta^{(t)}; D^*)]$$

$$Q(\theta^{(t+1)}; \theta^{(t)}) - Q(\theta^{(t)}; \theta^{(t)})$$

↳

We need to show that $Q(\theta^{(t+1)}; \theta^{(t)}) \geq Q(\theta^{(t)}; \theta^{(t)})$

In the move from $\theta^{(t)}$ to $\theta^{(t+1)}$ we maximized

$$Q(\theta; \theta^{(t)})$$

under constraint that all params except $\theta_{x_i | p_i}$ are unchanged $\Rightarrow Q(\theta^{(t+1)}; \theta^{(t)}) \geq Q(\theta^{(t)}; \theta^{(t)})$

□

[this is also true if we are maximizing

$$\operatorname{argmax}_{\theta_{x_i | p_i}} \mathbb{E}_{P(D; \theta^{(t)})} P(x_i | p_i; \theta_{x_i | p_i})$$

since other factors are unchanged]

(10.2) EM Algorithm using Expected Sufficient Statistics ④

Algorithm EM (G, θ^0, D): (for table CPDs)

for $t = 1 \dots$

$\bar{M}_t = \text{ESS}(G, \theta^t, D)$ // E

for $i = 1 \dots n$

for x_i, u_i in $\text{Val}(X_i \cup U_i)$

$\theta_{x_i, u_i}^{t+1} = \frac{\bar{M}_t[x_i, u_i]}{\bar{m}_t[u_i]}$ // M

return θ^t

Procedure ESS (G, θ, D):

$\bar{M} = 0$

for $m = 1 \dots M$

Perform inference* on (G, θ) using evidence $O[m]$

for $i = 1 \dots n$

for x_i, u_i in $\text{Val}(X_i \cup U_i)$

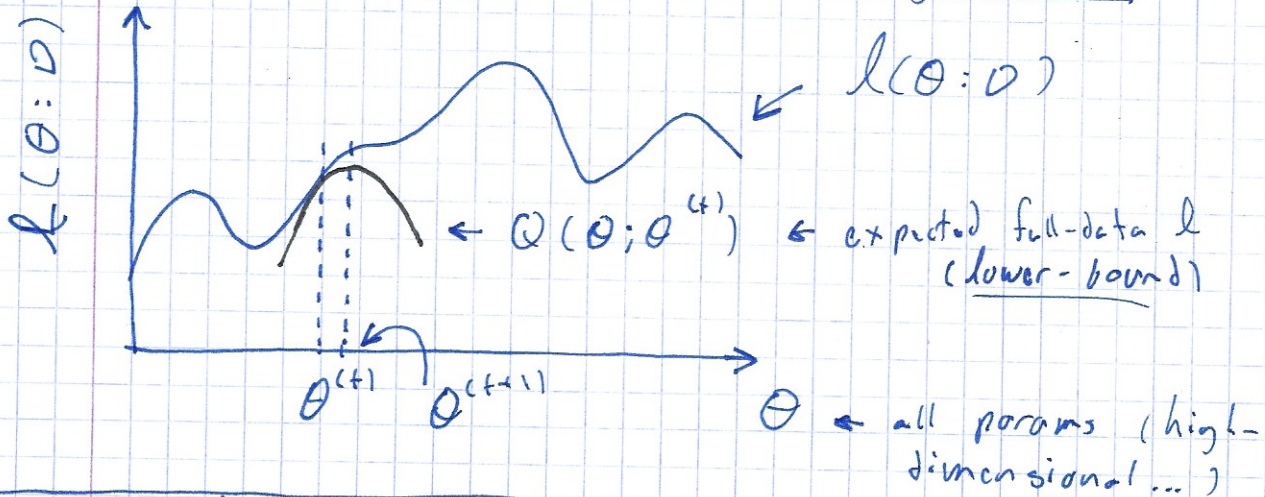
$\bar{M}[x_i, u_i] += P(x_i, u_i | O[m])$ ←

return \bar{M}

hidden → for normal type

* Inference e.g.
construct a c.t. w.
evidence, for calculating
all marginals $P(x_i, u_i | O)$

The EM likelihood hill-climbing process:



Fully observed:

$M[x, u] = \sum_m \mathbb{1}\{O[m] \langle x, u \rangle = (x, u)\}$

Partially observed

$\bar{M}[x, u] = \sum_m P(x, u | O[m], \theta)$

27/6

PMAI RII / Normalized Importance Sampling

PGM \rightarrow Inference \rightarrow Approximate \rightarrow sampling-based
 \rightarrow Importance Sampling \rightarrow NIS

Estimate the expectation of some $f(x)$ w.r.t. P by sampling from some Q . (and using weights based on Q and un-normalized version of P : $\tilde{P}(x) \propto P(x)$).

Derivation:

$$\begin{aligned} \mathbb{E}_P[f(x)] &= \mathbb{E}_Q\left[f(x) \frac{P(x)}{Q(x)}\right] = \frac{1}{Z} \mathbb{E}_Q\left[f(x) \frac{\tilde{P}(x)}{Q(x)}\right] \\ &= \frac{1}{Z} \mathbb{E}_Q[f(x) w(x)] \end{aligned}$$

$$\mathbb{E}_Q[w(x)] = \mathbb{E}_Q\left[\frac{\tilde{P}(x)}{Q(x)}\right] = \sum_x \tilde{P}(x) = Z \quad (P(x) = \frac{\tilde{P}(x)}{Z})$$

$$\mathbb{E}_P[f(x)] = \left(\frac{\mathbb{E}_Q\left[f(x) \frac{\tilde{P}(x)}{Q(x)}\right]}{\mathbb{E}_Q[w(x)]} \right) = \frac{\mathbb{E}_Q[f(x) w(x)]}{\mathbb{E}_Q[w(x)]}$$

$$\hat{\mathbb{E}}_D^{NIS}[f(x)] = \frac{\frac{1}{n} \sum_{m=1}^n f(x^{[m]}) w(x^{[m]})}{\frac{1}{n} \sum_{m=1}^n w(x^{[m]})}, \quad w(x^{[m]}) = \frac{\tilde{P}(x^{[m]})}{Q(x^{[m]})}$$

* "Sampling on a tree ..."

$$\tilde{P}(x|e) = \frac{P(x,e)}{P(e)} = \frac{P(x,e)}{\sum_x P(x,e)} \leftarrow \text{doesn't sum to 1 (joint w. some vals get)}$$

$$\hat{P}(x|e) = P(x,e)$$

$$w^{[m]} = \frac{P(x^{[m]}, e^{[m]})}{Q(x^{[m]}, e^{[m]} | e^{[m]})} \leftarrow \text{un-normalized}$$

$$\hat{P}(X=x|E=e) = \left(\sum_m w^{[m]} \mathbb{1}_{\{x^{[m]}=x\}} \right) / \left(\sum_m w^{[m]} \right) \leftarrow \text{normalized}$$

PMAI - RII - Gibbs Sampling

PGM \rightarrow Inference \rightarrow Approximate \rightarrow Sampling Based
 \rightarrow Gibbs Sampling

1. Initialize a sample w. $E=0$ and some values for the other variables.
2. for $t = 1 \dots T$ // for T samples
 $x(t) = x(t-1)$
 for x_i in $X \setminus E$
 sample $x_i^{(t)}$ from $P(x_i | x_{-i}^{(t)})$
 \uparrow
 \rightarrow is p_i
 \dots x_i
3. return $x[0], \dots, x[T]$

\rightarrow We saw in PS7 that $P(x_i | x_{-i})$ is easy.

Recitation 11: Inference and Learning Review

Teaching Assistant: Eitan Richardson

Prove that $P(X | \mathbf{e})$ is stationary w.r.t Gibbs process.

Proof:

A stationary distribution for a Markov chain \mathcal{T} :

$$\pi(X = x') = \sum_x \pi(X = x) \mathcal{T}(x \rightarrow x')$$

In our case: $\pi(X) = P(X | \mathbf{e})$ and $\mathcal{T}(x \rightarrow x') = P(x'_i | x_{-i}, \mathbf{e})$, for a random variable x_i .

We have:

$$\begin{aligned} \sum_{x_i} P(x_i, x_{-i} | \mathbf{e}) P(x'_i | x_{-i}, \mathbf{e}) &= \sum_{x_i} P(x_i | x_{-i}, \mathbf{e}) P(x_{-i} | \mathbf{e}) \frac{P(x'_i, x_{-i} | \mathbf{e})}{P(x_{-i} | \mathbf{e})} \\ &= \sum_{x_i} P(x_i | x_{-i}, \mathbf{e}) P(x'_i, x_{-i} | \mathbf{e}) \\ &= P(x'_i, x_{-i} | \mathbf{e}) \\ &= P(x' | \mathbf{e}) \end{aligned}$$

* for a chosen i . All other variables treated as constant.

¹Original LaTeX template courtesy of UC Berkeley.