

10/21

* ML estimation for complete data

$$P_{ML}(X_i = x \mid p_{a_i} = \pi) = \frac{\text{count}(X_i = x, p_{a_i} = \pi)}{\text{count}(p_{a_i} = \pi)}$$

* ML estimation for incomplete dataExamples $t = 1, 2, \dots, T$ Hidden nodes $H^{(t)}$ Visible nodes $V^{(t)}$

Review

* EM algorithm

E-step : compute posterior probabilities

$$P(X_i = x, p_{a_i} = \pi \mid V^{(t)}) \quad \text{inference}$$

M-step : update CPTs

$$P(X_i = x \mid p_{a_i} = \pi) \leftarrow \frac{\sum_{\pi} P(X_i = x, p_{a_i} = \pi \mid V^{(t)})}{\sum_{\pi} P(p_{a_i} = \pi \mid V^{(t)})}$$

Iterate until convergence. Note that RHS depends on current CPT estimates.

* Properties

- 1) no learning rate or tuning parameters
- 2) monotonic convergence

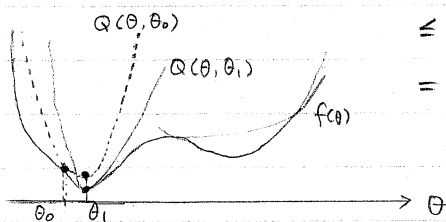
- each iteration improves log-likelihood. $\mathcal{L} = \sum_{\mathbf{x}} \log P(V^{(t)})$ Detour - numerical optimizationHow to minimize $f(\vec{\theta})$?

- 1) gradient descent: $\vec{\theta}_{t+1} = \vec{\theta}_t - \eta \nabla \frac{\partial f}{\partial \vec{\theta}}$ learning rate $\eta > 0$
- 2) Newton's method: $\vec{\theta}_{t+1} = \vec{\theta}_t - H^{-1} \frac{\partial f}{\partial \vec{\theta}}$ (not always converge; expensive to compute H)
- 3) Auxiliary function: $Q(\vec{\theta}, \vec{\theta}')$

Suppose $Q(\vec{\theta}, \vec{\theta}')$ satisfies: (i) $Q(\vec{\theta}, \vec{\theta}) = f(\vec{\theta})$ for all $\vec{\theta}$.(ii) $Q(\vec{\theta}, \vec{\theta}') \geq f(\vec{\theta})$ for all $\vec{\theta}, \vec{\theta}'$ Consider update rule: $\vec{\theta}_{t+1} = \arg \min_{\vec{\theta}} Q(\vec{\theta}, \vec{\theta}_t)$ It follows that: $f(\vec{\theta}_{t+1}) \leq Q(\vec{\theta}_{t+1}, \vec{\theta}_t)$ by property (ii)

$$\leq Q(\vec{\theta}_t, \vec{\theta}_t) \quad \text{because } \vec{\theta}_{t+1} = \arg \min_{\vec{\theta}} Q(\vec{\theta}, \vec{\theta}_t)$$

$$= f(\vec{\theta}_t)$$



Properties: - no learning rate

- monotonic improvement

- convergence to local stationary point (local minimum in general)

How to derive auxiliary function for ML estimation?

* Key inequality

Let $P(X)$ and $\tilde{P}(X)$ be different distributions over $X = \{X_1, X_2, \dots, X_n\}$

$$\begin{aligned} \log \tilde{P}(V) &= \log \left[\frac{\tilde{P}(h, v)}{\tilde{P}(h|v)} \right] \text{ for any instantiation } h \in H \text{ of hidden nodes} \\ &= \sum_h P(h|v) \log \left[\frac{\tilde{P}(h, v)}{\tilde{P}(h|v)} \right] \\ &= \sum_h P(h|v) \{ \log \tilde{P}(h, v) - \log \tilde{P}(h|v) + \log P(h|v) - \log P(h|v) \} \\ &= \sum_h P(h|v) \{ \log \tilde{P}(h, v) - \log P(h|v) + \log \frac{P(h|v)}{\tilde{P}(h|v)} \} \\ &= \sum_h P(h|v) \{ \log \tilde{P}(h, v) - \log P(h|v) \} + \text{KL}(P(h|v), \tilde{P}(h|v)) \\ \log \tilde{P}(V) &\geq \sum_h P(h|v) \{ \log \tilde{P}(h, v) - \log P(h|v) \} \end{aligned}$$

* Relation to EM algorithm

Imagine $P(X) = P_{\text{OLD}}(X; \theta)$ with old CPTs θ

Imagine $\tilde{P}(X) = P_{\text{NEW}}(X; \theta')$ with new CPTs θ'

How to derive update rule $\theta \rightarrow \theta'$, $P(X) \rightarrow \tilde{P}(X)$?

* Formal statement of EM

(i) E-step

Compute auxiliary function

$$Q(\theta, \theta') = \sum_v \sum_h \underbrace{P(h|v^{ct})}_{\text{old and new CPTs}} \log \tilde{P}(h, v^{ct}) - \sum_v \sum_h \underbrace{P(h|v^{ct})}_{\text{expected value of } \log \tilde{P}(h, v^{ct})} \log P(h|v^{ct})$$

(ii) M-step

Maximize $\sum_v \sum_h P(h|v^{ct}) \log \tilde{P}(h, v^{ct})$ in terms of new CPTs $\tilde{P}(X_i = x | pa_i = \pi)$

* Convergence proof

Suppose we choose new CPTs in this way.

$$\begin{aligned} \mathcal{L}_{\text{new}} &= \sum_v \log \tilde{P}(V^{ct}) \\ &\geq \sum_v \left\{ \sum_h P(h|V^{ct}) \log \tilde{P}(h, V^{ct}) - \sum_h P(h|V^{ct}) \log P(h|V^{ct}) \right\} \text{ (key inequality)} \\ &\geq \sum_v \left\{ \sum_h P(h|V^{ct}) \log P(h, V^{ct}) - \sum_h P(h|V^{ct}) \log P(h|V^{ct}) \right\} \text{ because of} \\ &\hspace{15em} \text{how } \hat{P} \text{ is chosen in the M-step} \\ &= \sum_v \sum_h P(h|V^{ct}) \log \left[\frac{P(h, V^{ct})}{P(h|V^{ct})} \right] \\ &= \sum_v \left[\sum_h P(h|V^{ct}) \right] \log P(V^{ct}) \\ &= \sum_v \log P(V^{ct}) = \mathcal{L}_{\text{old}} \end{aligned}$$

$$\therefore \boxed{\mathcal{L}_{\text{new}} \geq \mathcal{L}_{\text{old}}}$$

- M-step guarantees monotonic improvement
- Stronger guarantee than gradient ascent.

* Derivation of M-step for discrete BNs with lookup CPTs.

$$\begin{aligned}
 & \text{maximize } \sum_{\mathbf{h}} \sum_{\mathbf{h}'} P(\mathbf{h} | V^{(t)}) \log \tilde{P}(\mathbf{h}, V^{(t)}) \\
 & = \sum_{\mathbf{h}} \sum_{\mathbf{h}'} P(\mathbf{h} | V^{(t)}) \log \prod_i \tilde{P}(X_i | \text{pa}_i) \Big|_{\mathbf{x} = \mathbf{h}, V^{(t)}} \\
 \text{nodes in BN} & = \sum_{\mathbf{h}} \sum_{\mathbf{h}'} \sum_{\mathbf{x}} P(\mathbf{h} | V^{(t)}) \log \tilde{P}(X_i | \text{pa}_i) \Big|_{\mathbf{x} = \mathbf{h}, V^{(t)}} \\
 & = \sum_{\mathbf{h}} \sum_{\mathbf{h}'} \sum_{\pi} \sum_{\mathbf{x}} P(\mathbf{h} | V^{(t)}) \log \tilde{P}(X_i = \mathbf{x} | \text{pa}_i = \pi) \\
 & \quad \text{sum over possible } \mathbf{x} \text{ of node } X_i \\
 & \quad \text{sum over possible } \pi \text{ of parent configurations.} \\
 & = \sum_i \sum_{\pi} \sum_{\mathbf{x}} \underbrace{\left[\sum_{\mathbf{h}} P(X_i = \mathbf{x}, \text{pa}_i = \pi | V^{(t)}) \right]}_{\text{expected count}} \log \tilde{P}(X_i = \mathbf{x} | \text{pa}_i = \pi) \quad \text{: regrouping}
 \end{aligned}$$

Just like ML for complete data case, with expected count replacing count $(X_i = \mathbf{x}, \text{pa}_i = \pi)$.

Solution (M-step) of EM:

$$\tilde{P}(X_i = \mathbf{x} | \text{pa}_i = \pi) = \frac{\sum_{\mathbf{h}} P(X_i = \mathbf{x}, \text{pa}_i = \pi | V^{(t)})}{\sum_{\mathbf{x}'} \sum_{\mathbf{h}} P(X_i = \mathbf{x}', \text{pa}_i = \pi | V^{(t)})}$$

Denominator simplifies to: $\sum_{\mathbf{h}} P(\text{pa}_i = \pi | V^{(t)})$