

Incomplete dataReview

How to minimize  $f(\theta)$ ?

1) gradient descent  $\vec{\theta} \leftarrow \vec{\theta} - \eta \left( \frac{\partial f}{\partial \vec{\theta}} \right)$

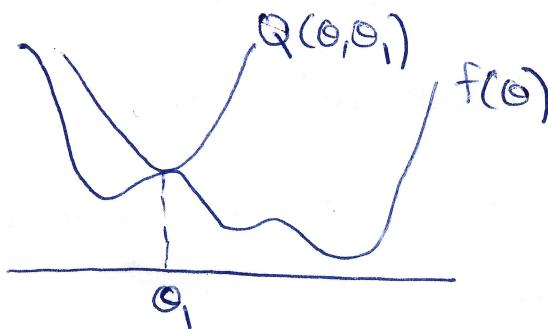
2) Newton's method  $\vec{\theta} \leftarrow \vec{\theta} - H^{-1} \left( \frac{\partial f}{\partial \vec{\theta}} \right)$

3) Auxiliary fn:  $Q(\theta | \theta^*, \vec{\theta})$  such that

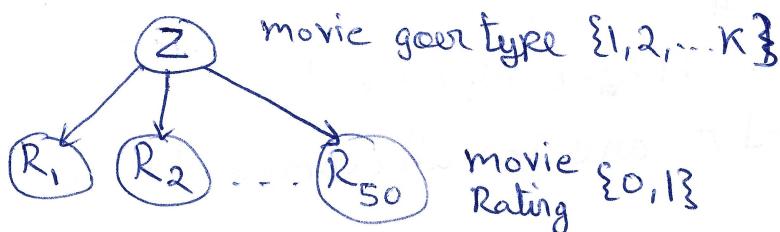
$Q(\theta, \theta^*) = f(\theta)$  and  $Q(\theta^*, \theta) \geq f(\theta)$

choose  $\theta_{\text{new}} = \underset{\theta}{\operatorname{arg\,min}} Q(\theta, \theta_{\text{old}})$

Then  $f(\theta_{\text{new}}) \leq f(\theta_{\text{old}})$

Incomplete data

Ex: movie recommender system



## ML estimation for incomplete data

(2)

Examples  $1, 2, \dots, T$

Hidden Nodes  $H^{(t)}$

Visible Nodes  $V^{(t)}$

How to choose CPTs to maximize  $\lambda = \sum_t \log P(V^{(t)})$ ?

Today: EM algorithm

- Statement & Induction
- Formal derivation.

Next Week: Many concrete examples.

## Expectation Maximization (EM) algorithm

- To maximize log likelihood for incomplete data
  - Initialize CPTs to some random (non zero) values
  - Iterate until convergence.

### A) E-step (Inference)

Compute posterior probabilities.

At root nodes  $P(X_i=x | V^{(t)})$

At nodes with parents  $P(X_i=x, pa_i=\pi | V^{(t)})$

for all values of  $x$  and  $\pi$  and for all examples  $t=1, 2, \dots, T$

### B) M-step (Learning)

Root node  $P(X_i=x) \leftarrow \frac{1}{T} \sum_{t=1}^T P(X_i=x | V^{(t)})$

Nodes with parents  $P(X_i=x | pa_i=\pi) \leftarrow \frac{\sum_t P(X_i=x, pa_i=\pi | V^{(t)})}{\sum_t P(pa_i=\pi | V^{(t)})}$

(3)

## Intuition for updates

By analogy to complete data case.

$$\begin{aligned} \text{Node w/ parents } P_{ML}(X_i=x | \text{pa}_i=\pi) &= \frac{\text{count}(X_i=x, \text{pa}_i=\pi)}{\text{count}(\text{pa}_i=\pi)} \\ &= \frac{\sum_{t=1}^T I(X_i^{(t)}=x) I(\text{pa}_i^{(t)}=\pi)}{\sum_{t=1}^T I(\text{pa}_i^{(t)}=\pi)} \end{aligned}$$

$$\text{Root nodes: } P_{ML}(X_i=x) = \text{count}(X_i=x) = \frac{1}{T} \sum_{t=1}^T I(X_i^{(t)}=x)$$

For incomplete data, we must "fill in" missing values.

Expected statistics under  $P(H|V^{(t)})$  from current CPTs substitute for observed ~~data~~ statistics when data is incomplete.

RHS of EM updates reduces to complete data counts at nodes where child & parents are observed.

## Properties of EM algorithm

1) No learning state

2) Monotonic convergence: each iteration improves  $\lambda = \sum_t \log P(V^{(t)})$

## Derivation of EM

- Key inequality: let  $P(x)$  &  $\tilde{P}(x)$  be different distributions over

$$X = \{X_1, X_2, \dots, X_n\}$$

Let  $V$  be subset of observed nodes.

(4)

$$\log \tilde{P}(V) = \log \frac{\tilde{P}(H=h|V)}{\tilde{P}(H=h|V)} \text{ for any setting } h \text{ of unobserved nodes.}$$

$$= \sum_h P(h|V) \log \frac{\tilde{P}(H=h|V)}{\tilde{P}(H=h|V)}$$

$$= \sum_h P(h|V) \log \frac{\tilde{P}(H=h|V)}{\tilde{P}(H=h|V)} + \sum_n P(n|V) \log P(n|V) - \sum_n P(n|V) \log \tilde{P}(n|V)$$

$$= \sum_n P(n|V) \log \frac{\tilde{P}(n|V)}{P(n|V)} + \sum_n P(n|V) \log \frac{P(n|V)}{\tilde{P}(n|V)}$$

↳ This is KL distance between posteriors  $P$  &  $\tilde{P}$

$$\geq \sum_n P(n|V) \log \frac{\tilde{P}(n|V)}{P(n|V)}$$

For ML estimation,

Imagine  $P(X)$  is from BN with current (old) CPTs (like  $\vec{\Theta}$ )

Imagine  $\tilde{P}(X)$  is from BN with updated (new) CPTs (like  $\vec{\Theta}'$ )

Formal statement of EM algorithm

E-step: Compute auxiliary function.

$$\textcircled{a} Q(\tilde{P}, P) = \sum_{t=1}^T \left[ \sum_n P(n|V^{(t)}) \log \frac{\tilde{P}(n|V^{(t)})}{P(n|V^{(t)})} \right]$$

To validate:

$$\begin{aligned} \text{(ii)} Q(P, P) &= \sum_t \sum_n P(n|V^{(t)}) \log \frac{P(n|V^{(t)})}{P(n|V^{(t)})} \\ &= \sum_t \sum_n P(n|V^{(t)}) \log P(V^{(t)}) \\ &= \sum_t \log P(V^{(t)}) = N_{\text{old}} \end{aligned}$$

(5)

$$(ii) Q(\tilde{p}, p) \leq \sum_t \log \tilde{P}(\hat{v}^{(t)}) = \lambda(\tilde{p}) - \lambda_{\text{new}} \\ (\text{by key inequality})$$

M step Maximize  $Q(\tilde{p}, p)$

choose CPTs  $\tilde{P}(x_i=x | p_{a_i}=\pi)$  to maximize it.

Suppose we can choose CPTs in this way

$$\text{Then } \lambda_{\text{new}} = \sum_t \log \tilde{P}(v^{(t)})$$

$$\geq \sum_t \sum_h P(h|v^{(t)}) \log \frac{\tilde{P}(h, v^{(t)})}{P(h|v^{(t)})}$$

$$\geq \sum_t \sum_h P(h|v^{(t)}) \log \frac{P(h, v^{(t)})}{P(h|v^{(t)})}$$

$$= \sum_t \sum_h P(h|v^{(t)}) \log P(v^{(t)})$$

$$= \sum_t \log P(v^{(t)})$$

$$= \lambda_{\text{old}}$$

b/c update rule  
maximizes above RHS  
w.r.t  $\tilde{p}$

Derive M-step updates

$$\text{Maximize } Q(\tilde{p}, p) = \sum_t \sum_h P(h|v^{(t)}) \log \frac{\tilde{P}(h, v^{(t)})}{P(h|v^{(t)})}$$

$$- \sum_t \sum_h P(h|v^{(t)}) \log P(h|v^{(t)})$$

$$= \sum_t \sum_h P(h|v^{(t)}) \log \left. \prod_{i=1}^n \tilde{P}(x_i=x | p_{a_i}=\pi) \right|_{H=h, V=v^{(t)}} - \dots$$

$$= \sum_{i=1}^n \sum_h \sum_t P(h|v^{(t)}) \log \tilde{P}(x_i=x | p_{a_i}=\pi)$$

(6)

$$= \sum_{i=1}^n \sum_x \sum_{\pi} \sum_t P(x_i=x, pa_i=\pi / V^{(t)}) \log \tilde{P}(x_i=x | pa_i=\pi)$$

$$\text{Maximize } Q(\tilde{p}, p) = \sum_i \sum_{\pi} \sum_x \left\{ \sum_t P(x_i=x, pa_i=\pi / V^{(t)}) \right\} \log \tilde{P}(x_i=x | pa_i=\pi)$$

Recall complete data,

$$\begin{aligned} \text{maximize } \lambda_{\text{complete}} &= \sum_i \sum_{\pi} \sum_x \text{count}(x_i=x, pa_i=\pi) \log \tilde{P}(x_i=x | pa_i=\pi) \\ &= \sum_i \sum_{\pi} \sum_x \left\{ \sum_{t=1}^T I(x_i^{(t)}, x) I(pa_i^{(t)}, \pi) \right\} \log \tilde{P}(x_i=x | pa_i=\pi) \end{aligned}$$

$$P_{m_2}(x_i=x | pa_i=\pi) = \frac{\sum_t I(x_i^{(t)}, x) I(pa_i^{(t)}, \pi)}{\sum_t I(pa_i^{(t)}, \pi)}$$

To maximize  $Q(\tilde{p}, p)$

$$\tilde{P}(x_i=x | pa_i=\pi) \leftarrow \frac{\sum_t P(x_i=x, pa_i=\pi / V^{(t)})}{\sum_t P(pa_i=\pi / V^{(t)})}$$

Matches update rule for EM