

Learning with Partially Observed Data

Lecture 12 – May 4, 2011
CSE 515, Statistical Methods, Spring 2011

Instructor: Su-In Lee
University of Washington, Seattle

Model Selection

- So far, we focused on single model
 - Given $D = \{X[1], \dots, X[M]\}$, find best scoring model $\tilde{G} = \arg \max_G P(G|D)$
 - Use it to predict next example $P(X[M+1]|D, \tilde{G})$
- Implicit assumption
 - Making predictions based on the Bayesian estimation rule:
$$\rightarrow P(X[M+1]|D) = \sum_G P(X[M+1]|D, G) P(G|D)$$
 - Best scoring model dominates the weighted sum
$$P(X[M+1]|D) \approx P(X[M+1]|D, \tilde{G})$$
 - Valid with many data instances (very large M)
- Pros:
 - We get a single structure
 - Allows for efficient use in our prediction tasks
- Cons:
 - Committing to the independencies of a particular structure
 - Other structures with similar score might be probable given D

Model Selection

- Density estimation
 - Picking one structure may suffice if its distribution $P(\mathbf{X}[M+1] | D, G)$ is similar for different high-scoring structures.
- Structure discovery
 - Several networks with similar scores \rightarrow one or several of them might be close to the "true" structure, but we cannot distinguish between them given the data D .
 - Drawing a conclusion about the structure from one of the networks can be wrong
 - Thus, instead of picking one of the high-scoring structures, we should focus on estimating the "confidence" of the structural properties we are interested in.
 - Define features $f(G)$ (e.g., edge, sub-structure, d-sep property)
 - Compute $P(f | D) = \sum_G f(G) P(G | D) \rightarrow p(f, G | D)$
 - Requires summing over exponentially many structures
 - We can reduce the computation assuming a certain ordering

Model Averaging Given an Order

- Assumptions
 - Known total order of variables α ($X_1, X_2, \dots, X_{i-1}, X_i, X_{i+1}, \dots, X_{n-1}, X_n$)
 - Maximum in-degree for variables d

■ Marginal likelihood

Using decomposability assumption on prior $P(G | \alpha)$

$$\begin{aligned}
 P(D | \alpha) &= \sum_{G \in G_{d, \alpha}} P(D | G) P(G | \alpha) \\
 &= \sum_{G \in G_{d, \alpha}} \prod_i \exp\{FamScore_B(X_i | Pa_{X_i}^G : D)\} \\
 &= \prod_i \left[\sum_{U_i \in \{U \mid U \prec X_i \in \alpha, U \prec G\}} \exp\{FamScore_B(X_i | U_i : D)\} \right] \\
 &= (f_1^{g^1} + \dots + f_1^{g^m}) \dots (f_n^{g^1} + \dots + f_n^{g^m})
 \end{aligned}$$

Cost per family: $O(n^d)$
 Total cost: $O(n^{d+1})$

Since given ordering α , parent choices are independent

Model Averaging Given an Order

- Posterior probability of a **general feature f**

$$P(f | \alpha, D) = \frac{P(f, D | \alpha)}{P(D | \alpha)} = \frac{\sum_{G \in \mathcal{G}_{f, \alpha}} f(G) P(D | G) P(G | \alpha)}{\prod_i \sum_{U \in \{U: U < X_i \in \alpha, |U| < d\}} \exp\{FamScore_B(X_i | U_i : D)\}}$$

- f**: particular choice of parents **U** for X_i

$$P(Pa_{X_i}^G = U | D, \alpha) = \frac{\exp\{FamScore_B(X_i | U : D)\}}{\sum_{U \in \{U: U < X_i \in \alpha, |U| < d\}} \exp\{FamScore_B(X_i | U_i : D)\}}$$

All terms cancel out

- f**: existence of a particular edge between $X_j \rightarrow X_i$

$$P(X_j \in Pa_{X_i}^G | D, \alpha) = \frac{\sum_{U \in \{U: X_j \in U \text{ and } U < X_i \in \alpha, |U| < d\}} \exp\{FamScore_B(X_i | U_i : D)\}}{\sum_{U \in \{U: U < X_i \in \alpha, |U| < d\}} \exp\{FamScore_B(X_i | U_i : D)\}}$$

5

Model Averaging

- We cannot assume that order is known α
- Solution: Sample from posterior distribution of $P(G|D)$
 - If we manage to sample graphs G_1, \dots, G_K from $P(G|D)$
 - Estimate feature probability by $P(f | D) \approx \frac{1}{K} \sum_{i=1}^K f(G_i)$
 - Sampling can be done by **MCMC** (Markov chain Monte Carlo)
 - Next week $G_1, \dots, G_K \leftarrow P(G|D)$

6

Notes on Learning Local Structures

- Beyond table CPDs
- Define score with local structures
 - Example: in tree CPDs, score decomposes by leaves (not by X_i and a particular value on Par X_i)
- Prior may need to be extended
 - Example: in tree CPDs, penalty for tree structure per CPD (depth of the tree)
- Extend search operators to local structure
 - Example: in tree CPDs, we need to search for tree structure
 - Can be done by local encapsulated search or by defining new global operations

7

Structure Search: Summary

- Discrete optimization problem
- In general, NP-Hard
 - Need to resort to heuristic search
 - In practice, search is relatively fast (~ 100 vars in ~ 10 min):
 - Decomposability ←
 - Sufficient statistics ←
- In some cases, we can reduce the search problem to an easy optimization problem
 - Example: learning trees, a fixed ordering α

8

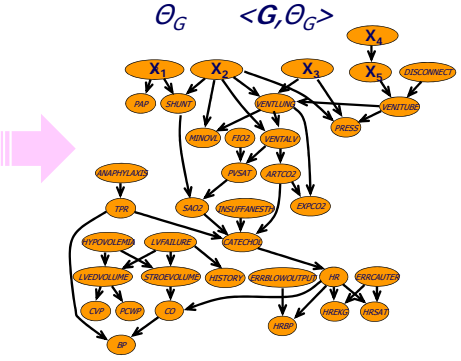
Let's turn to the main topic for today...

LEARNING WITH PARTIALLY OBSERVED DATA

Training Data D

Training instance

X_1	3	10	-1	1	100	298	...
X_2	1	11	1	0	072	365	...
X_3	0	01	0	108	2123	...	
X_4	12	5	-2	30	1345	...	
:	:	:	:	:	:	:	:
X_{N-1}	1	32	3	6	...		
X_N	0	74	-4	7	...		



- Until now, we assumed that the training data is **fully observed**
 - Each instance assigns values to all the variables in our domain

Incomplete Data

- In reality, this assumption might not be true.

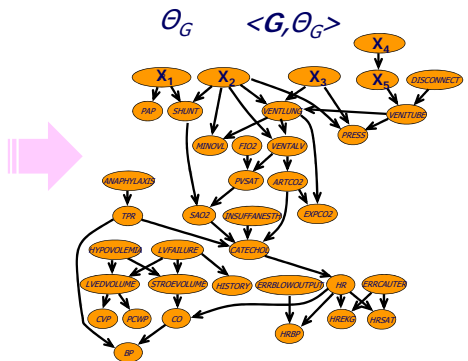
D Training instance

X_1	3	1	?	-1	1	1	?	0	2	9	8	...
X_2	1	?	1	1	0	0	7	2	3	6	5	...
X_3	0	0	1	?	1	0	8	2	?	2	3	...
X_4	1	2	5	-2	?	0	1	3	4	5	...	

Lung cancer?

X_{N-1} ??????????????????

X_N 0 7 4 ? ?



- Missing values, Hidden variables
- Challenges
 - Foundational** – is the learning task well defined?
 - Computational** – how can we learn with missing data?

Treating Missing Data

- How should we treat missing data?
 - Based on data missing mechanism
- Case I:** A coin is tossed on a table, occasionally it drops and measurements are not taken (random missing)
 - Sample sequence: H, T, ?, ?, T, ?, H
 - Treat missing data by ignoring it
- Case II:** A coin is tossed, but only heads are reported (deliberate missing values)
 - Sample sequence: H, ?, ?, ?, H, ?, H
 - Treat missing data by filling it with Tails

↓

We need to consider the data missing mechanism

Modeling Data Missing Mechanism

- Let's try to model the data missing mechanism

- $X = \{X_1, \dots, X_n\}$ are random variables

- $O_x = \{O_{x_1}, \dots, O_{x_n}\}$ are *observability variables*
 - Always observed

- $Y = \{Y_1, \dots, Y_n\}$ new random variables

- $\text{Val}(Y_i) = \text{Val}(X_i) \cup \{?\}$
- Y_i is a deterministic function of X_i and O_{x_i} :

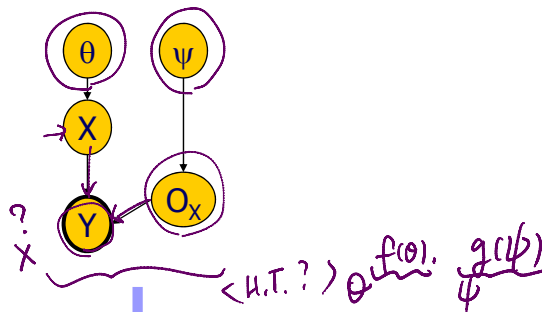
$$Y_i = \begin{cases} X_i & O_{x_i} = o^1 \\ ? & O_{x_i} = o^0 \end{cases}$$

13

Modeling Missing Data Mechanism

Case I

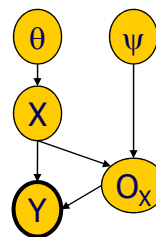
(random missing values)



$$\begin{aligned} P(Y=H) &= \theta\psi & M_H & \hat{\theta} = \frac{\partial \log L}{\partial \theta} \\ P(Y=T) &= (1-\theta)\psi & M_T & \\ P(Y=?) &= (1-\psi) & M_? & \\ L(D; \theta, \psi) &= \theta^{M_H} \cdot (1-\theta)^{M_T} \cdot \psi^{M_H+M_T} \cdot (1-\psi)^{M_?} & = f(\theta, \psi) & \leftarrow \end{aligned}$$

Case II

(deliberate missing values)



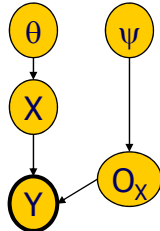
MLE

$$\begin{aligned} \hat{\theta} &= \frac{M_H}{M_H + M_T} \\ \hat{\psi} &= \frac{M_H + M_T}{M_H + M_T + M_?} \end{aligned}$$

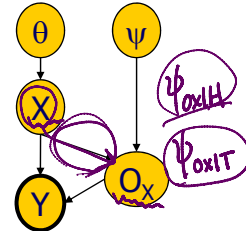
14

Modeling Missing Data Mechanism

Case I
(random missing values)



Case II
(deliberate missing values)



MLE ?

$\hat{\theta} =$
 $\hat{\psi} =$

$f(\theta)g(\psi)$

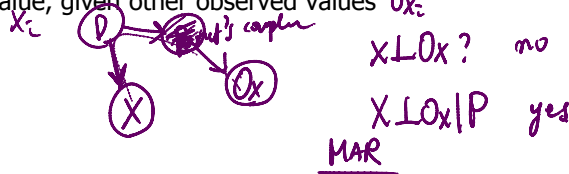
$$\begin{aligned} P(Y=H) &= \theta \psi_{O_x|H} \\ P(Y=T) &= (1-\theta) \psi_{O_x|T} \\ P(Y=?) &= \theta(1-\psi_{O_x|H}) + (1-\theta)(1-\psi_{O_x|T}) \end{aligned}$$

$$L(D; \theta, \psi) = \theta^{M_H} \cdot (1-\theta)^{M_T} \cdot \psi^{M_H} \cdot \psi^{M_T} \cdot (\theta(1-\psi_{O_x|H}) + (1-\theta)(1-\psi_{O_x|T}))^{M_?}$$

Decoupling of Observation Mechanism

When can we ignore the missing data mechanism and focus only on the likelihood? $L = f(\theta)g(\psi)$

- Missing Completely at Random (MCAR)
 - For every X_i , $\text{Ind}(X_i; O_{X_i})$, a very strong assumption
 - Sufficient but not necessary for the decomposition of the likelihood
- Missing at Random (MAR) is sufficient
 - The probability that the value of X_i is missing is independent of its actual value, given other observed values



- In both cases, the likelihood decomposes $f \theta$
 - When there are missing values in D , try to model such that MAR holds. $\psi \theta$

Incomplete Data

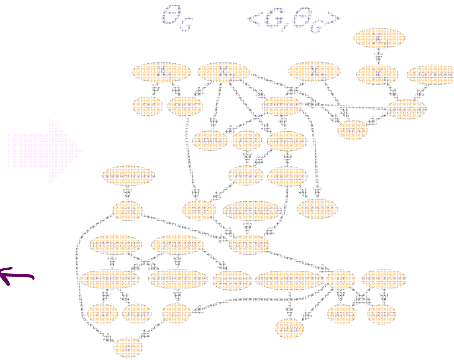
- In reality, this assumption might not be true.

D

X_1	3 1 ? -1 1 1 ? 0 2 9 8 ...
X_2	1 ? 1 1 0 0 7 2 3 6 5 ...
X_3	0 0 1 ? 1 0 8 2 ? 2 3 ...
X_4	1 2 5 -2 ? 0 1 3 4 5 ...

Lung cancer?

X_{N-1}	????????????????
X_N	0 7 4 ? 7



- Missing values, Hidden variables

■ Challenges

- Foundational – is the learning task well defined?
- Computational – how can we learn with missing data?

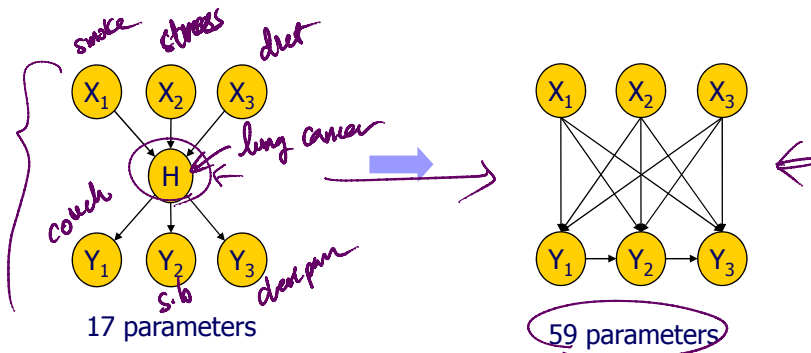
Hidden (Latent) Variables

- Attempt to learn a model with hidden variables

- In this case, MCAR always holds (variable is always missing)

MCAR
MAR $O_H = F \otimes$

- Why should we care about unobserved variables?

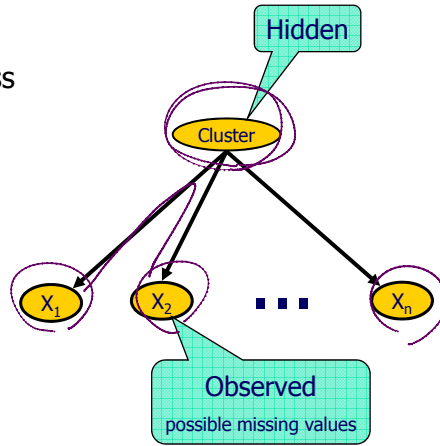


Hidden (Latent) Variables

- Hidden variables also appear in **clustering**
- **Naïve Bayes** model:
 - Class variable is hidden
 - Observed attributes are independent given the class

D

X_1	3	1	0	-1	1	1	0	0	2	9	8	...	
X_2	1	1	1	1	0	0	7	2	3	6	5	...	
X_3	0	0	1	0	1	0	8	2	1	2	3	...	
X_4	1	2	5	-2	3	0	1	3	4	5	...		
:	:	:	:	:	:	:	:	:	:	:	:	:	
X_{N-1}	1	3	2	3	6	
X_N	0	7	4	-4	7	
H	1	2	1	1	3	3	1	1	2	2	1	1	...



19

How do missing data affect the likelihood function?

20

Likelihood for Complete Data

Input Data:

x	y
x^0	y^0
x^0	y^1
x^1	y^0

Likelihood:

$p(D; \theta)$

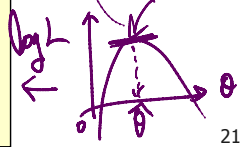
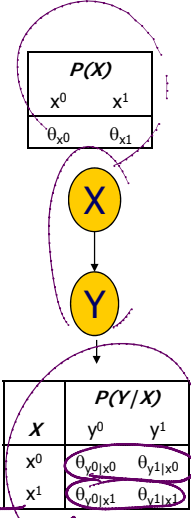
$$L(D; \theta) = P(x[1], y[1]) \cdot P(x[2], y[2]) \cdot P(x[3], y[3])$$

$$= P(x^0, y^0) \cdot P(x^0, y^1) \cdot P(x^1, y^0)$$

$$= \theta_{x^0} \cdot \theta_{y^0|x^0} \cdot \theta_{x^0} \cdot \theta_{y^1|x^0} \cdot \theta_{x^1} \cdot \theta_{y^0|x^1}$$

$$= (\theta_{x^0} \cdot \theta_{x^1}) \cdot (\theta_{y^0|x^0} \cdot \theta_{y^1|x^0}) \cdot (\theta_{y^0|x^1})$$

- Likelihood decomposes by variables
- Likelihood decomposes within CPDs
- Likelihood function is log-concave → unique global maximum that has a simple analytic closed form.



Likelihood for Incomplete Data

Input Data:

x	y
?	y^0
x^0	y^1
?	y^0

Likelihood:

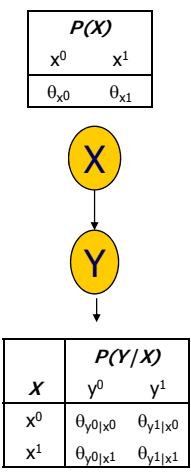
$$L(D; \theta) = P(y^0) \cdot P(x^0, y^1) \cdot P(y^0)$$

$$= \sum_{x \in X} P(x, y^0) \cdot P(x^0, y^1) \cdot \left(\sum_{x \in X} P(x, y^0) \right)$$

$$= (\theta_{x^0} \cdot \theta_{y^0|x^0} + \theta_{x^1} \cdot \theta_{y^0|x^1}) \cdot \theta_{x^0} \cdot \theta_{y^1|x^0} \cdot (\theta_{x^0} \cdot \theta_{y^0|x^0} + \theta_{x^1} \cdot \theta_{y^0|x^1})$$

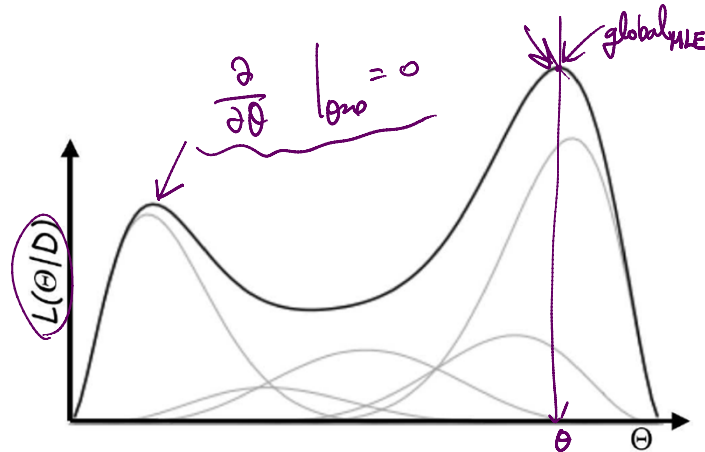
$$= (\theta_{x^0} \cdot \theta_{y^0|x^0} + \theta_{x^1} \cdot \theta_{y^0|x^1})^2 \cdot \theta_{x^0} \cdot \theta_{y^1|x^0}$$

- Likelihood does not decompose by variables
- Likelihood does not decompose within CPDs
- Computing likelihood per instance requires inference!



Likelihood with Missing Data

- **Multimodal likelihood function** with incomplete data
 - Likelihood function is not log-concave → local maxima cannot be obtained by a simple analytic closed form

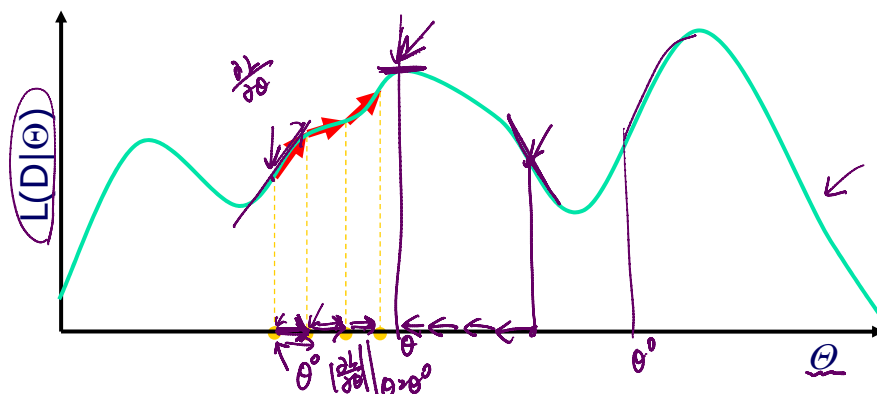


CSE 515 – Statistical Methods – Spring 2011

23

MLE from Incomplete Data

- Take steps proportional to the positive of the gradient.



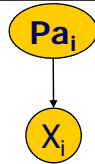
Gradient Ascent:

- Follow gradient of likelihood w.r.t. to parameters
- Add line search and conjugate gradient methods to get fast convergence

24

Gradient Ascent

Pa_i	x_i	
	H	T
(H, ..., H)	$\Theta_{H (H, \dots, H)}$	$\Theta_{T (H, \dots, H)}$
(H, ..., T)	$\Theta_{H (H, \dots, T)}$	$\Theta_{T (H, \dots, T)}$



■ Theorem:

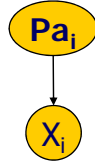
$$\frac{\partial \log P(D | \Theta)}{\partial \theta_{x_i, pa_i}} = \frac{1}{\theta_{x_i, pa_i}} \sum P(x_i, pa_i | o[m], \Theta)$$

■ Proof:

$$\begin{aligned} \frac{\partial \log P(D | \Theta)}{\partial \theta_{x_i, pa_i}} &= \sum_m \frac{\partial \log P(o[m] | \Theta)}{\partial \theta_{x_i, pa_i}} \\ &= \sum_m \frac{1}{P(o[m] | \Theta)} \frac{\partial P(o[m] | \Theta)}{\partial \theta_{x_i, pa_i}} \end{aligned}$$

How do we compute? $\frac{\partial P(o[m] | \Theta)}{\partial \theta_{x_i, pa_i}}$

Gradient Ascent



$$\frac{\partial P(o[m] | \Theta)}{\partial \theta_{x_i, pa_i}} = \sum_{\mathbf{X}: \mathbf{X} < \mathbf{O} > o[m]} \frac{\partial P(\mathbf{X} | \Theta)}{\partial \theta_{x_i, pa_i}}$$

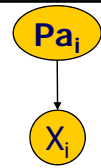
chain rule $\rightarrow \frac{\partial P(\mathbf{X} | \Theta)}{\partial \theta_{x_i, pa_i}}$

$$= \sum_{\mathbf{X}: \mathbf{X} < \mathbf{O} > o[m], X_{x_i, pa_i} = x_i, pa_i} \frac{P(\mathbf{X} | \Theta)}{\theta_{x_i, pa_i}}$$

$$= \frac{1}{\theta_{x_i, pa_i}} P(x_i, pa_i, o[m] | \Theta)$$

$$\begin{aligned} \frac{\partial \log P(D | \Theta)}{\partial \theta_{x_i, pa_i}} &= \sum_m \frac{\partial \log P(o[m] | \Theta)}{\partial \theta_{x_i, pa_i}} \\ &= \sum_m \frac{1}{P(o[m] | \Theta)} \frac{\partial P(o[m] | \Theta)}{\partial \theta_{x_i, pa_i}} \end{aligned}$$

Gradient Ascent



$$\begin{aligned} \frac{\partial \log P(D | \Theta)}{\partial \theta_{x_i, pa_i}} &= \sum_m \frac{1}{P(o[m] | \Theta)} \frac{\partial P(o[m] | \Theta)}{\partial \theta_{x_i, pa_i}} \\ &= \sum_m \frac{1}{P(o[m] | \Theta)} \frac{P(x_i, pa_i, o[m] | \Theta)}{\theta_{x_i, pa_i}} \\ &= \sum_m \frac{P(x_i, pa_i | o[m], \Theta)}{\theta_{x_i, pa_i}} \quad \frac{1}{\theta_{x_i, pa_i}} \sum_m P(x_i, pa_i | o[m], \Theta) \end{aligned}$$

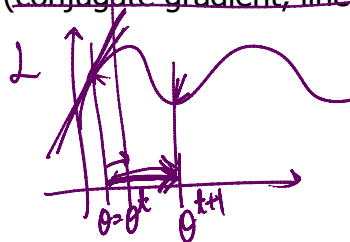
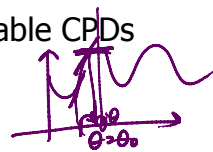
- Requires computation: $P(x_i, pa_i | o[m], \Theta)$ for all j, m
- Can be done with clique-tree algorithm, since X_i, Pa_i are in the same clique

29

Gradient Ascent Summary

- Pros
 - Flexible, can be extended to non table CPDs
- Cons
 - Need to project gradient onto space of legal parameters
 - For reasonable convergence, need to combine with advanced methods (conjugate gradient, line search)

$$L = f(\theta)$$



30

Expectation Maximization (EM)

- Tailored algorithm for optimizing likelihood functions

- **Intuition**

- Parameter estimation is easy given complete data
- Computing probability of missing data is "easy" (=inference) given parameters



$$\langle \theta, H \rangle \quad \sum_x P(x) p(y|x)$$

- **Strategy**

- Pick a starting point for parameters (θ^0)
- "Complete" the data using current parameters
- Estimate parameters relative to data completion
- Iterate
- Procedure guaranteed to improve at each iteration

Expectation Maximization

$$P(X_i | \theta^{(t)}, \theta)$$



Expectation Maximization (EM)

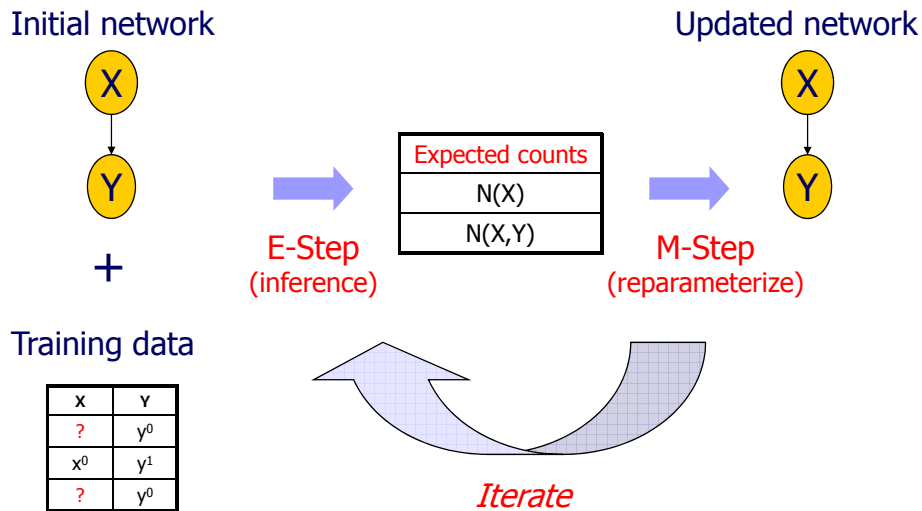
- Initialize parameters to θ^0
- Iterate E-step and M-step
- In the t-th iteration, we do
- **Expectation (E-step):**
 - Let $o[m]$ be the observed data in the m-th training instance.
 - For each m and each family X_i, \mathbf{Pa}_i , compute $P(X_i, \mathbf{Pa}_i | o[m], \theta^{(t)})$
 - Compute the *expected sufficient statistics* for each values x, \mathbf{u} on X_i, \mathbf{Pa}_i , respectively.

$$\bar{M}_{\theta^{(t)}}[X_i = x, \mathbf{Pa}_i = \mathbf{u}] = \sum_m P(X_i = x, \mathbf{Pa}_i = \mathbf{u} | o[m], \theta^{(t)})$$

- **Maximization (M-step):**
 - Treat the expected sufficient statistics as observed and set the parameters to the MLE with respect to the ESS

$$\theta_{X_i=x, \mathbf{Pa}_i=\mathbf{u}}^{(t+1)} = \frac{\bar{M}_{\theta^{(t)}}[X_i = x, \mathbf{Pa}_i = \mathbf{u}]}{\bar{M}_{\theta^{(t)}}[\mathbf{Pa}_i = \mathbf{u}]}$$

Expectation Maximization (EM)



33

Expectation Maximization (EM)

- **Formal Guarantees:**
 - $L(D:\Theta^{(t+1)}) \geq L(D:\Theta^{(t)})$
 - Each iteration improves the likelihood
 - If $\Theta^{(t+1)} = \Theta^{(t)}$, then $\Theta^{(t)}$ is a stationary point of $L(D:\Theta)$
 - Usually, this means a local maximum
- **Main cost:**
 - Computations of expected counts in E-Step
 - Requires inference for each instance in training set
 - Exactly the same as in gradient ascent!
- **Reading material on EM**
 - Please read Andrew Ng's lecture note

34

EM – Practical Considerations

- **Initial parameters**
 - Highly sensitive to starting parameters
 - Choose randomly
 - Choose by guessing from another source
- **Stopping criteria**
 - Small change in data likelihood
 - Small change in parameters
- **Avoiding bad local maxima**
 - Multiple restarts
 - Early pruning of unpromising starting points

35

Acknowledgement

- These lecture notes were generated based on the slides from Prof Eran Segal.

CSE 515 – Statistical Methods – Spring 2011

36