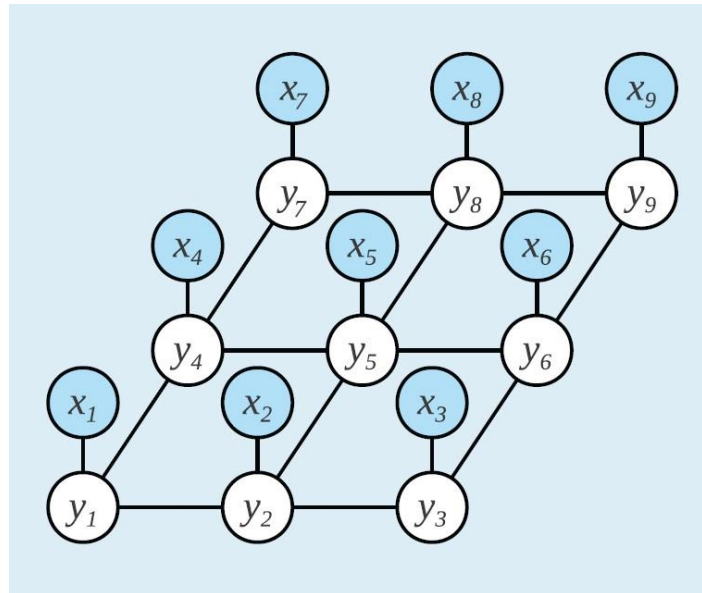


# Probabilistic Graphical Models in Bioinformatics

Lecture 9: score-based structure learning; partially observed data



# Score-based structure learning

# Learning general Bayesian networks

- **Theorem:** Finding the maximum-score network with at most  $d \geq 2$  parents for each variable is *NP-hard*.
- We will use heuristic algorithm to search the space of graphs and return a high-scoring one.
- Local operators:
  - Edge addition
  - Edge deletion
  - Edge reversal
- Search techniques
  - Greedy hill climbing
  - Simulated annealing
  - ...

# Greedy hill climbing

- The search procedure
  - Pick an initial network structure  $g$ 
    - Empty network
    - A random choice
    - The best tree
  - At each iteration
    - We consider all legal networks in neighbors of  $g$  using local operators
    - Apply the change that leads to the best improvement
  - Stopping condition: when no modification improves the score
- Two issues:
  - Sticking in a local maxima
  - Reaching a *plateau* (a large set of neighboring networks that have the same score).

# Score decomposition and search

$$\text{score}(\mathcal{G} : \mathcal{D}) = \sum_i \text{FamScore}(X_i \mid \text{Pa}_{X_i}^{\mathcal{G}} : \mathcal{D}),$$

We define the delta score as:

$$\delta(\mathcal{G} : o) = \text{score}(o(\mathcal{G}) : \mathcal{D}) - \text{score}(\mathcal{G} : \mathcal{D})$$

**Proposition 18.5** *Let  $\mathcal{G}$  be a network structure and score be a decomposable score.*

- *If  $o$  is “Add  $X \rightarrow Y$ ,” and  $X \rightarrow Y \notin \mathcal{G}$ , then*

$$\delta(\mathcal{G} : o) = \text{FamScore}(Y, \text{Pa}_Y^{\mathcal{G}} \cup \{X\} : \mathcal{D}) - \text{FamScore}(Y, \text{Pa}_Y^{\mathcal{G}} : \mathcal{D}).$$

- *If  $o$  is “Delete  $X \rightarrow Y$ ” and  $X \rightarrow Y \in \mathcal{G}$ , then*

$$\delta(\mathcal{G} : o) = \text{FamScore}(Y, \text{Pa}_Y^{\mathcal{G}} - \{X\} : \mathcal{D}) - \text{FamScore}(Y, \text{Pa}_Y^{\mathcal{G}} : \mathcal{D}).$$

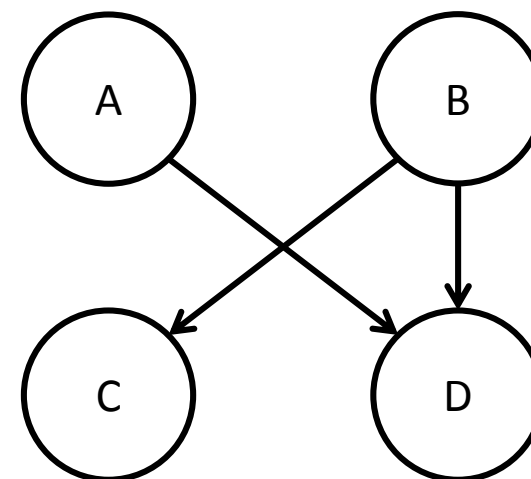
- *If  $o$  is “Reverse  $X \rightarrow Y$ ” and  $X \rightarrow Y \in \mathcal{G}$ , then*

$$\begin{aligned} \delta(\mathcal{G} : o) = & \text{FamScore}(X, \text{Pa}_X^{\mathcal{G}} \cup \{Y\} : \mathcal{D}) + \text{FamScore}(Y, \text{Pa}_Y^{\mathcal{G}} - \{X\} : \mathcal{D}) \\ & - \text{FamScore}(X, \text{Pa}_X^{\mathcal{G}} : \mathcal{D}) - \text{FamScore}(Y, \text{Pa}_Y^{\mathcal{G}} : \mathcal{D}). \end{aligned}$$

# Score decomposition and search-example

**Question:** compute the delta score for the following operations

1. Add  $A \rightarrow B$
2. Add  $C \rightarrow D$
3. Remove  $A \rightarrow D$
4. Remove  $B \rightarrow C$
5. Reverse  $B \rightarrow C$
6. Reverse  $B \rightarrow D$



# Structure learning methods

- Constraint-based structure learning
- Score-based structure learning
- **Bayesian model averaging methods**

# Bayesian model averaging

- Fully Bayesian inference
  - We consider the structure as a random variable
- Bayesian prediction

$$P(\xi[M+1] \mid \mathcal{D}) = \sum_{\mathcal{G}} P(\xi[M+1] \mid \mathcal{D}, \mathcal{G}) P(\mathcal{G} \mid \mathcal{D}),$$

where 
$$P(\mathcal{G} \mid \mathcal{D}) = \frac{P(\mathcal{G})P(\mathcal{D} \mid \mathcal{G})}{P(\mathcal{D})}.$$

- How to approximate  $P(g \mid D)$ ?
  - often by Markov chain Monte Carlo (MCMC) over structures

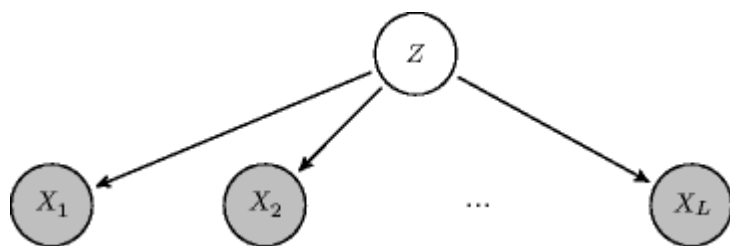


# Partially observed data

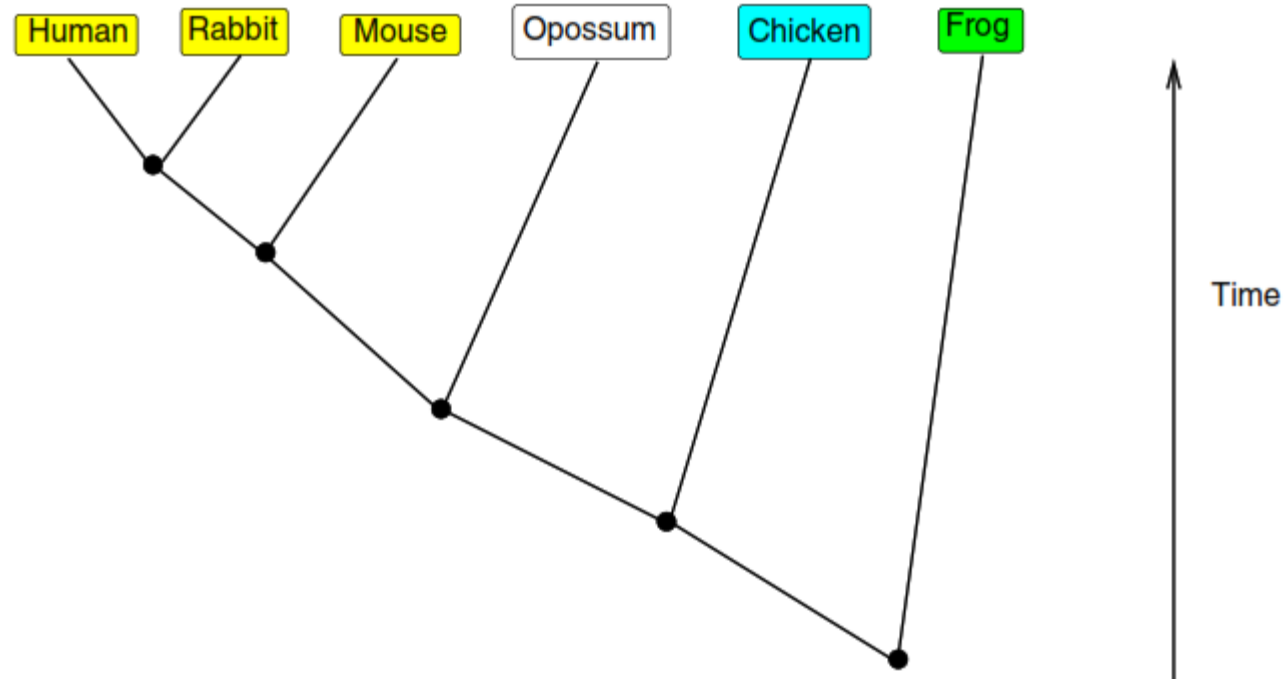
	Known structure (Parameter estimation)	Unknown structure (Structure learning )
<b>Fully observable</b>	<p><i>MLE</i></p> <p><i>Bayesian methods</i></p>	<p><i>Constraint-based methods</i></p> <p><i>Score-based methods e.g. hill climbing</i></p>
<b>Partially observable</b>		

# Incomplete data

- In real-world applications of learning, we rarely have fully observed data.
- Incomplete data
  - Hidden variables
  - Missing values

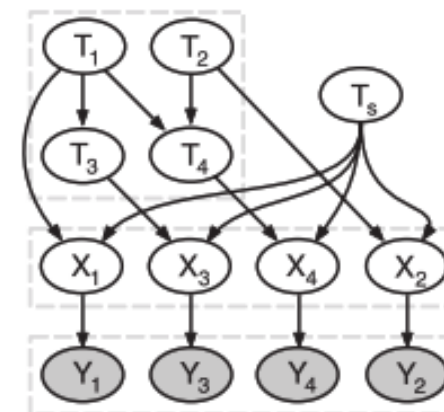
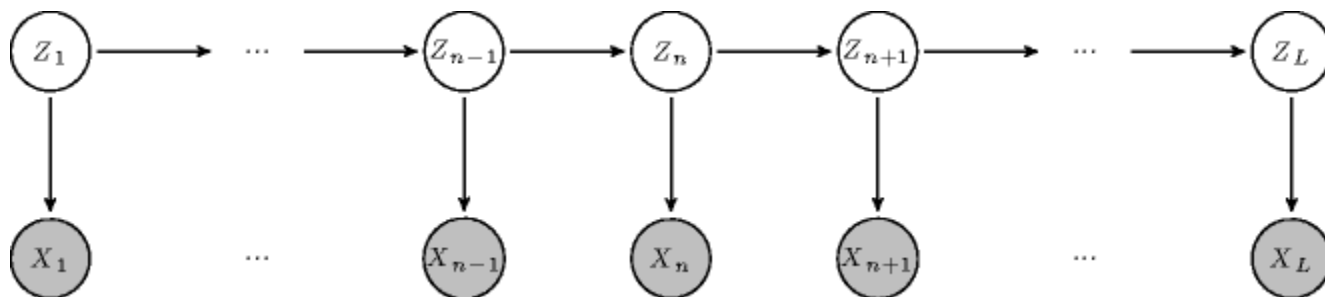


Clustering



Phylogenetic trees

M M M M M M M M M M M M M M M M X X M M M M M M M M M M M M M M Y M M M M Y M M M M M  
 Seq1 = CTRPNNNTRKSIRPQIGPGQAFYATGD-IGDI-RQAHC  
 Seq2 = CGRPNNHRIKGLR--IGPGRAFFAMGAIRGGEIRQAHC



Cancer progression

# The likelihood function

- Question: write the likelihood function for  $X \rightarrow Y$ 
  - Based on counts  $M[x^0, y^0], \dots$

$$L(\theta_X, \theta_{Y|x^0}, \theta_{Y|x^1} : \mathcal{D}) = \theta_{x^1}^{M[x^1]} \theta_{x^0}^{M[x^0]} \cdot \theta_{y^1|x^0}^{M[x^0, y^1]} \theta_{y^0|x^0}^{M[x^0, y^0]} \cdot \theta_{y^1|x^1}^{M[x^1, y^1]} \theta_{y^0|x^1}^{M[x^1, y^0]}.$$

- Data:
  - $M[x^1, y^1] = 13, M[x^1, y^0] = 16, M[x^0, y^1] = 10, M[x^0, y^0] = 4$

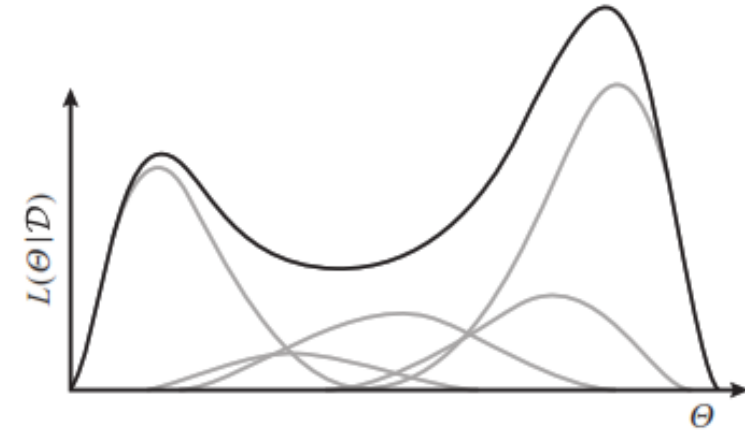
$$\theta_{x^1}^{29} (1 - \theta_{x^1})^{14} \cdot \theta_{y^1|x^0}^{10} (1 - \theta_{y^1|x^0})^4 \cdot \theta_{y^1|x^1}^{13} (1 - \theta_{y^1|x^1})^{16}.$$

- The function is well-behaved ; it is log-concave and has a unique global maximum.
- Now assume the first observation was  $X[1] = x^0, Y[1] = y^1$ . Consider we observed only  $Y[1] = y^1$ .
- Question: write down modified likelihood.

$$\theta_{x^1}^{29} (1 - \theta_{x^1})^{13} \cdot \theta_{y^1|x^0}^9 (1 - \theta_{y^1|x^0})^4 \cdot \theta_{y^1|x^1}^{13} (1 - \theta_{y^1|x^1})^{16} [\theta_{x^1} \theta_{y^1|x^1} + (1 - \theta_{x^1}) \theta_{y^1|x^0}].$$

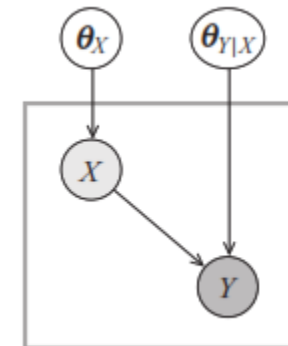
# A multimodal likelihood function with incomplete data

- Complete data likelihood defines a unimodal function.
- Their sum can be a multimodal function.
  - Mixture of peaks



$$\theta_{x^1}^{29} (1 - \theta_{x^1})^{13} \cdot \theta_{y^1|x^0}^9 (1 - \theta_{y^1|x^0})^4 \cdot \theta_{y^1|x^1}^{13} (1 - \theta_{y^1|x^1})^{16} [\theta_{x^1} \theta_{y^1|x^1} + (1 - \theta_{x^1}) \theta_{y^1|x^0}]$$

- Are parameters independence given incomplete data?
  - We lose the property of *parameter independence*
  - The likelihood function is not decomposable!



# On global decomposability of likelihood function

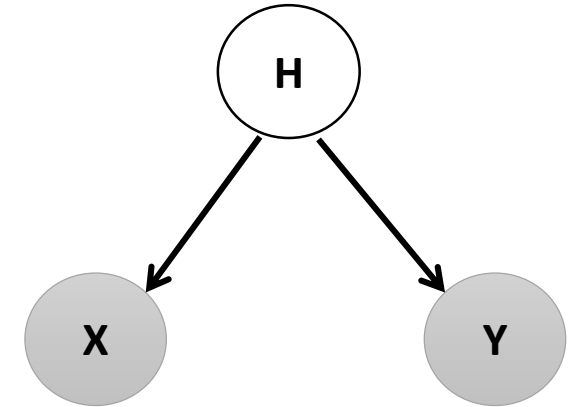
- The probability of observing  $x$  and  $y$

$$P(x, y) = \sum_h P(h)P(x | h)P(y | h).$$

- Likelihood function

$$L(\theta : \mathcal{D}) = \prod_{x,y} \left( \sum_h P(h)P(x | h)P(y | h) \right)^{M[x,y]}$$

*Can we write the likelihood as the product of local likelihoods?*



- Likelihood function in the general case:

$$L(\theta : \mathcal{D}) = P(\mathcal{D} | \theta) = \sum_{\mathcal{H}} P(\mathcal{D}, \mathcal{H} | \theta).$$

exponential numbers of modes  
(worst case)

a unimodal function

Likelihood function for IID instances:

$$L(\theta : \mathcal{D}) = \prod_m P(o[m] | \theta) = \prod_m \sum_{h[m]} P(o[m], h[m] | \theta).$$

**The learning problem becomes substantially more complex!**

# Parameter estimation for the partially observed case

- We cover methods for maximum likelihood estimation.
- Problem definition:
  - **Input**
    - a network structure  $g$  and the form of the CPDs
    - A data set  $D$  that consists of  $M$  partial instances
  - **Goal:** Find the values  $\hat{\theta}$  that maximize the log-likelihood function:  $\hat{\theta} = \arg \max_{\theta} l(\theta; D)$
  - The problem requires optimizing a highly nonlinear and multimodal function over a high-dimensional space.
- Two main classes of methods for performing this optimization
  - A generic nonconvex optimization algorithm such as *gradient ascent*
  - *Expectation maximization*: a more specialized approach for optimizing likelihood functions.



# Gradient ascent

- Review A.5.1 and A.5.2 from the text book

---

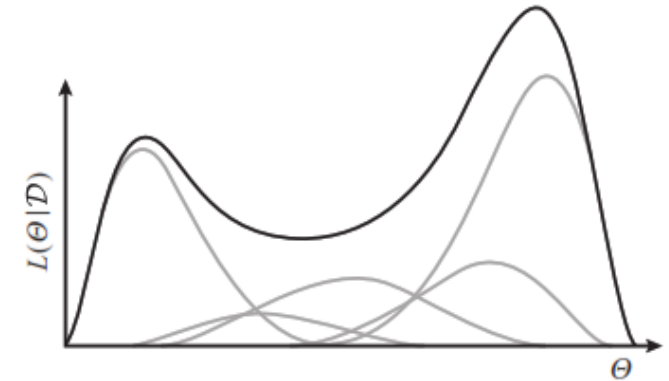
**Algorithm A.10 Simple gradient ascent algorithm**

---

```
Procedure Gradient-Ascent (  
   $\theta^1$ , // Initial starting point  
   $f_{\text{obj}}$ , // Function to be optimized  
   $\delta$  // Convergence threshold  
)  
1   $t \leftarrow 1$   
2  do  
3     $\theta^{t+1} \leftarrow \theta^t + \eta \nabla f_{\text{obj}}(\theta^t)$   
4     $t \leftarrow t + 1$   
5  while  $\|\theta^t - \theta^{t-1}\| > \delta$   
6  return ( $\theta^t$ )
```

---

$$\nabla f = \left\langle \frac{\partial f}{\partial \theta_1}, \dots, \frac{\partial f}{\partial \theta_n} \right\rangle.$$



- Question:** consider  $X \sim \text{Poisson}(\lambda)$  with  $P(x) = \frac{\lambda^x e^{-\lambda}}{x!}$ . Suppose  $D = \{4, 6, 3, 7\}$ ,  $\lambda^1 = 10$ , and  $\eta = 1$ .
  - Write the gradient-ascent updating formula for  $\lambda$ .
  - Compute  $\lambda^4$

# Gradient ascent for Bayesian networks

- Let  $D = \{o[1], \dots, o[M]\}$  be a partially observed data set and  $X$  be a variable and  $U$  its parent in  $g$ .  
Then

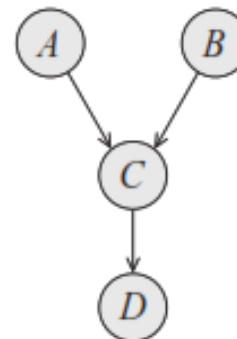
$$\frac{\partial \ell(\boldsymbol{\theta} : \mathcal{D})}{\partial P(x | \mathbf{u})} = \frac{1}{P(x | \mathbf{u})} \sum_{m=1}^M P(x, \mathbf{u} | o[m], \boldsymbol{\theta}).$$

- Given  $\theta$  is a parameter in a CPD, we can use *chain rule of derivative*

$$\frac{\partial \ell(\boldsymbol{\theta} : \mathcal{D})}{\partial \theta} = \sum_{x, \mathbf{u}} \frac{\partial \ell(\boldsymbol{\theta} : \mathcal{D})}{\partial P(x | \mathbf{u})} \frac{\partial P(x | \mathbf{u})}{\partial \theta},$$

# Example

- given  $o = \langle a^1, ?, ?, d^0 \rangle$  for the following network



$$\frac{\partial \ell(\theta : \mathcal{D})}{\partial P(x | \mathbf{u})} = \frac{1}{P(x | \mathbf{u})} \sum_{m=1}^M P(x, \mathbf{u} | \mathbf{o}[m], \theta).$$

$$\begin{aligned} \theta_{a^1} &= 0.3 \\ \theta_{b^1} &= 0.9 \\ \theta_{c^1|a^0,b^0} &= 0.83 \\ \theta_{c^1|a^0,b^1} &= 0.09 \\ \theta_{c^1|a^1,b^0} &= 0.6 \\ \theta_{c^1|a^1,b^1} &= 0.2 \\ \theta_{d^1|c^0} &= 0.1 \\ \theta_{d^1|c^1} &= 0.8. \end{aligned}$$

*Probabilities of the consistent cases with  $o$*

$$\begin{aligned} P(\langle a^1, b^1, c^1, d^0 \rangle) &= 0.3 \cdot 0.9 \cdot 0.2 \cdot 0.2 = 0.0108 \\ P(\langle a^1, b^1, c^0, d^0 \rangle) &= 0.3 \cdot 0.9 \cdot 0.8 \cdot 0.9 = 0.1944 \\ P(\langle a^1, b^0, c^1, d^0 \rangle) &= 0.3 \cdot 0.1 \cdot 0.6 \cdot 0.2 = 0.0036 \\ P(\langle a^1, b^0, c^0, d^0 \rangle) &= 0.3 \cdot 0.1 \cdot 0.4 \cdot 0.9 = 0.0108. \end{aligned}$$

$$\begin{aligned} P(\langle a^1, b^1, c^1, d^0 \rangle | \mathbf{o}) &= 0.0492 \\ P(\langle a^1, b^1, c^0, d^0 \rangle | \mathbf{o}) &= 0.8852 \\ P(\langle a^1, b^0, c^1, d^0 \rangle | \mathbf{o}) &= 0.0164 \\ P(\langle a^1, b^0, c^0, d^0 \rangle | \mathbf{o}) &= 0.0492. \end{aligned}$$

$$\begin{aligned} \frac{\partial \log P(\mathbf{o})}{\partial P(d^1 | c^0)} &= \frac{P(d^1, c^0 | \mathbf{o})}{P(d^1 | c^0)} = \frac{0}{0.1} = 0 \\ \frac{\partial \log P(\mathbf{o})}{\partial P(d^0 | c^0)} &= \frac{P(d^0, c^0 | \mathbf{o})}{P(d^0 | c^0)} = \frac{0.8852 + 0.0492}{0.9} = 1.0382 \\ \frac{\partial \log P(\mathbf{o})}{\partial P(d^1 | c^1)} &= \frac{P(d^1, c^1 | \mathbf{o})}{P(d^1 | c^1)} = \frac{0}{0.8} = 0 \\ \frac{\partial \log P(\mathbf{o})}{\partial P(d^0 | c^1)} &= \frac{P(d^0, c^1 | \mathbf{o})}{P(d^0 | c^1)} = \frac{0.0492 + 0.0164}{0.2} = 0.328. \end{aligned}$$

1. Inference for each instance  $P(X[m], U[m] | \mathbf{o}[m], \theta)$  at each iteration
2. Ensure that parameters describe a legal probability distribution
  - Reparametrization
  - Lagrange multipliers