

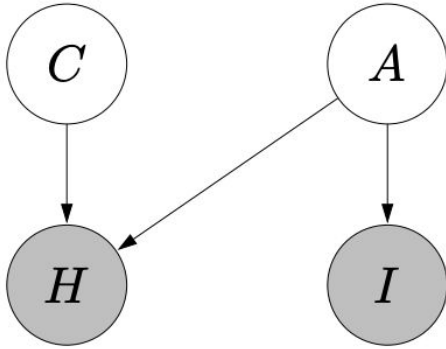
# CS221 Problem Workout

Week 8

# Outline

- **Bayesian networks: Learning**
  - Maximum likelihood
  - Smoothing
  - EM Algorithm
- Problem discussion

# Bayesian networks: Learning



Given local probability distributions, i.e.  $P(x \mid \text{parents}(x))$

Find conditional  $P(Q \mid E=e)$

**Inference**

Given observations / samples

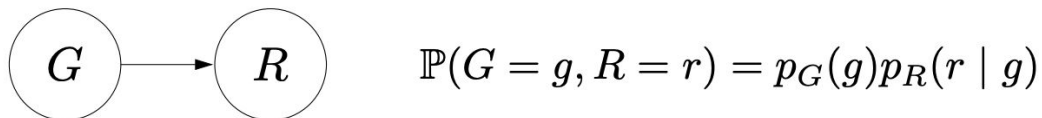
Find the local distributions, i.e.  $P(x \mid \text{Parents}(x))$

**Learning**

# Example

Variables:

- Genre  $G \in \{\text{drama}, \text{comedy}\}$
- Rating  $R \in \{1, 2, 3, 4, 5\}$



$$\mathcal{D}_{\text{train}} = \{(d, 4), (d, 4), (d, 5), (c, 1), (c, 5)\}$$

Parameters:  $\theta = (p_G, p_R)$

Example borrowed from lecture slides

# Outline

- Bayesian networks: Learning
  - **Maximum likelihood**
  - Smoothing
  - EM Algorithm
- Problem discussion

# Maximum likelihood

**Input:** training examples  $\mathcal{D}_{\text{train}}$  of full assignments

**Output:** parameters  $\theta = \{p_d : d \in D\}$



## Algorithm: count and normalize

### Count:

For each  $x \in \mathcal{D}_{\text{train}}$ :

For each variable  $x_i$ :

Increment count $_{d_i}(x_{\text{Parents}(i)}, x_i)$

### Normalize:

For each  $d$  and local assignment  $x_{\text{Parents}(i)}$ :

Set  $p_d(x_i \mid x_{\text{Parents}(i)}) \propto \text{count}_d(x_{\text{Parents}(i)}, x_i)$

Slide borrowed from lecture slides

# Outline

- Bayesian networks: Learning
  - Maximum likelihood
  - **Smoothing**
  - EM Algorithm
- Problem discussion

# Smoothing

Why?

- What if count is 0? Should P be 0?

How to solve?

- Initialize all counts with a non-zero constant  $\lambda$

Observations

- Larger  $\lambda$  -> more uniform distributions, less influenced by data
- Smaller  $\lambda$  -> more influenced by data
- Infinite data -> effect of  $\lambda$  vanishes



# Final algorithm

**Input:** training examples  $\mathcal{D}_{\text{train}}$  of full assignments

**Output:** parameters  $\theta = \{p_d : d \in D\}$



## Algorithm: count and normalize

### Count:

For each  $x \in \mathcal{D}_{\text{train}}$ :

For each variable  $x_i$ :

Increment count $_{d_i}(x_{\text{Parents}(i)}, x_i)$

### Normalize:

For each  $d$  and local assignment  $x_{\text{Parents}(i)}$ :

Set  $p_d(x_i \mid x_{\text{Parents}(i)}) \propto \text{count}_d(x_{\text{Parents}(i)}, x_i) + \lambda$

Slide modified from lecture slides

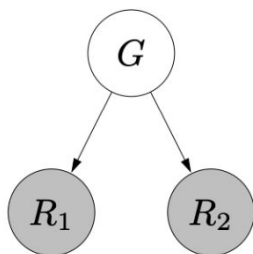
# Outline

- Bayesian networks: Learning
  - Maximum likelihood
  - Smoothing
  - **EM Algorithm**
- Problem discussion

# EM Algorithm

Variables:  $H$  is hidden,  $E = e$  is observed

Example:



$$H = G \quad E = (R_1, R_2) \quad e = (1, 2)$$
$$\theta = (p_G, p_R)$$

Maximum marginal likelihood objective:

$$\begin{aligned} & \max_{\theta} \prod_{e \in \mathcal{D}_{\text{train}}} \mathbb{P}(E = e; \theta) \\ &= \max_{\theta} \prod_{e \in \mathcal{D}_{\text{train}}} \sum_h \mathbb{P}(H = h, E = e; \theta) \end{aligned}$$

Slide borrowed from lecture slides

# EM Algorithm

Initialize  $\theta$  randomly

Until convergence:

## # E Step

for each  $e$  in Data:

for each  $h$ :

$q(h; e) = P(H = h \mid E = e; \theta) \dots$  inference

# Update table from  $\{e, \text{count}(e)\}$  to  $\{(h,e), (q(h; e) \times \text{count}(e))\}$

# Now no variables are hidden

## # M step

update( $\theta$ ) using Table  $\{(h,e), (q(h; e) \times \text{count}(e))\} \dots$  MLE

# Summary

- Given data learn the parameters of bayesian net

## MLE

$$p \propto \text{count}(x_i; \text{parents}(x_i))$$

## Smoothing

$$p \propto \text{count}(x_i; \text{parents}(x_i)) + \lambda$$

## EM

Data is *incomplete*

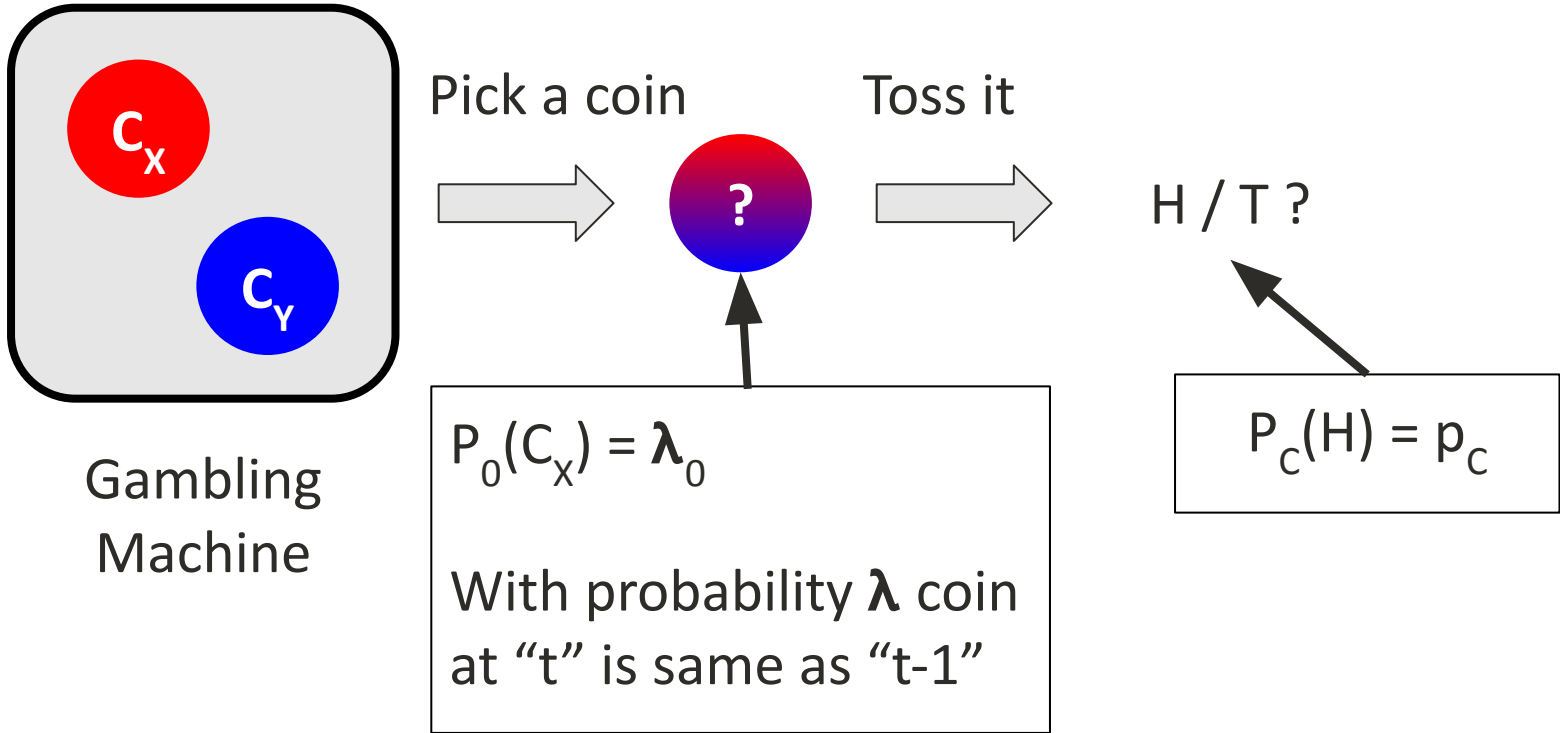
*E step*: compute counts

*M step*: MLE

# Outline

- Bayesian networks: Learning
  - Maximum likelihood
  - Smoothing
  - EM Algorithm
- **Problem discussion**

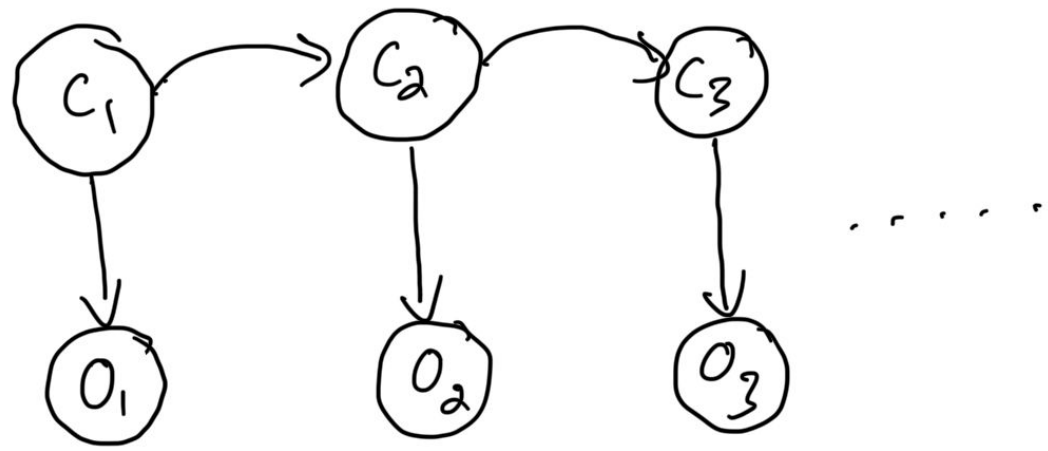
# Problem: P2, Winter 2021 Exam 2



# How does the bayesian net look?

$$P(C_1) = \begin{cases} p_0 & \text{if } C_1 = X \\ 1-p_0 & \text{else} \end{cases}$$

$$P(C_{i+1} | C_i) = \begin{cases} p & \text{if } C_{i+1} = C_i \\ 1-p & \text{else} \end{cases}$$



$$P(O_i | C_i) = \begin{cases} p_{C_i} & \text{if } O_i = H \\ 1-p_{C_i} & \text{else} \end{cases}$$



# Learning using EM algorithm

- Data = {H, H, T}
- $\lambda_0$  and  $\lambda$  are given. To find:  $p_x$  and  $p_y$
  
- Why do we need EM?
  - What is not observed?
  - $C_i$  is not observed
  
- How do we use EM?
  - Compute  $q(c_i)$  using  $p'_x$  and  $p'_y$
  - Use ML to update  $p'_x$  and  $p'_y$

# Given $q$ 's compute update

	T=1	T=2	T=3
X	0.1	0.5	0.3
Y	0.9	0.5	0.7

Data = {H, H, T}

Compute:

$C_i$	$O_i$	Count
?	?	?

# Given q's compute update

	T=1	T=2	T=3
X	0.1	0.5	0.3
Y	0.9	0.5	0.7

Data = {H, H, T}

$C_i$	$O_i$	Count
X	H	0.1
X	H	0.5
X	T	0.3

$$\begin{aligned} P(H | X) &= p'_X \\ &= (0.1 + 0.5) / (0.1 + 0.5 + 0.3) \dots \text{MLE} \end{aligned}$$



Thank You