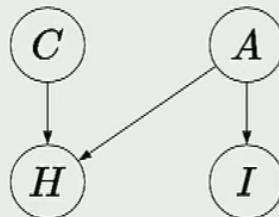


## 15. Bayesian Network 3 - Maximum Likelihood

### Review: Bayesian network



$$\begin{aligned} \mathbb{P}(C = c, A = a, H = h, I = i) \\ = p(c)p(a)p(h | c, a)p(i | a) \end{aligned}$$



#### Definition: Bayesian network

Let  $X = (X_1, \dots, X_n)$  be random variables.

A **Bayesian network** is a directed acyclic graph (DAG) that specifies a **joint distribution** over  $X$  as a product of **local conditional distributions**, one for each node:

$$\mathbb{P}(X_1 = x_1, \dots, X_n = x_n) = \prod_{i=1}^n p(x_i | x_{\text{Parents}(i)})$$

Stanford

### Review: probabilistic inference

Bayesian network:

$$\mathbb{P}(X = x) = \prod_{i=1}^n p(x_i | x_{\text{Parents}(i)})$$

Probabilistic inference:

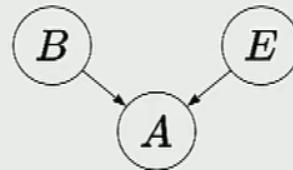
$$\mathbb{P}(Q | E = e)$$

Algorithms:

- Forward-backward: HMMs, exact
- Particle filtering: HMMs, approximate
- Gibbs sampling: general, approximate

We are gonna talk about learning

## Where do parameters come from?



b	$p(b)$
1	?
0	?

e	$p(e)$
1	?
0	?

b	e	a	$p(a   b, e)$
0	0	0	?
0	0	1	?
0	1	0	?
0	1	1	?
1	0	0	?
1	0	1	?
1	1	0	?
1	1	1	?

Stanford

So far we have just assumed that someone hands you these local conditional distribution which has numbers filled out In real world we have to figure them out from data

# Roadmap

**Supervised learning**

Laplace smoothing

Unsupervised learning with EM

Supervised learning

---

# Learning task

**Training data**

$\mathcal{D}_{\text{train}}$  (an example is an assignment to  $X$ )



**Parameters**

$\theta$  (local conditional probabilities)



[cs221.stanford.edu/q](http://cs221.stanford.edu/q)

Question

Which is computationally more expensive for Bayesian networks?

probabilistic inference given the parameters

learning the parameters given fully labeled data

It turns out that probabilistic inference is more expensive So in fully supervised learning it should be easier

## Example: one variable

Setup:

- One variable  $R$  representing the rating of a movie  $\{1, 2, 3, 4, 5\}$

$$\textcircled{R} \quad \mathbb{P}(R = r) = p(r)$$

Parameters:

$$\theta = (p(1), p(2), p(3), p(4), p(5))$$

Training data:

$$\mathcal{D}_{\text{train}} = \{1, 3, 4, 4, 4, 4, 4, 4, 5, 5, 5\}$$

Fully observed

## Example: one variable

Learning:

$$\mathcal{D}_{\text{train}} \Rightarrow \theta$$

Intuition:  $p(r) \propto$  number of occurrences of  $r$  in  $\mathcal{D}_{\text{train}}$

Example:

$$\mathcal{D}_{\text{train}} = \{1, 3, 4, 4, 4, 4, 4, 4, 5, 5, 5\}$$



$r$	$\text{count}(r)$
1	1
2	0
3	1
4	5
5	3

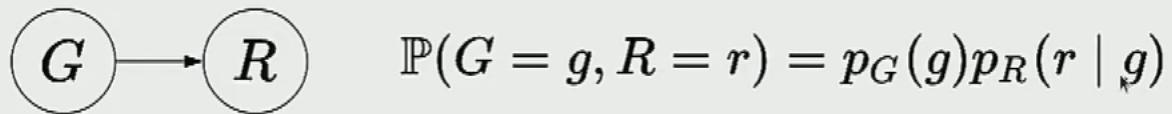
Stanford

count and divide, then normalize

## Example: two variables

Variables:

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

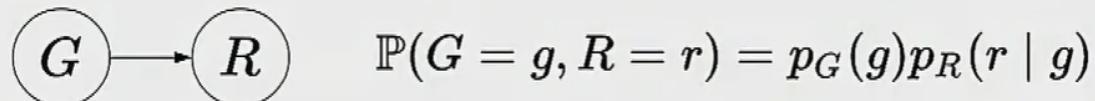


$$\mathcal{D}_{\text{train}} = \{(d, 4), (d, 4), (d, 5), (c, 1), (c, 5)\}$$

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

We first access  $p_G$  and  $p_R$  separately

## Example: two variables



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

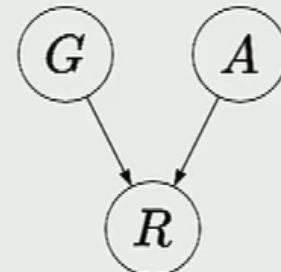
Intuitive strategy: Estimate each local conditional distribution ( $p_G$  and  $p_R$ ) separately

$\theta:$	$\begin{array}{ c c } \hline g & p_G(g) \\ \hline d & 3/5 \\ c & 2/5 \\ \hline \end{array}$	$\begin{array}{ c c c } \hline g & r & p_R(r   g) \\ \hline d & 4 & 2/3 \\ d & 5 & 1/3 \\ c & 1 & 1/2 \\ c & 5 & 1/2 \\ \hline \end{array}$
-----------	---------------------------------------------------------------------------------------------	---------------------------------------------------------------------------------------------------------------------------------------------

# Example: v-structure

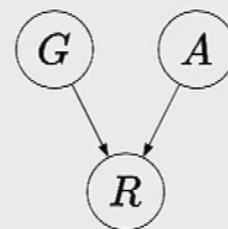
Variables:

- Genre  $G \in \{\text{drama, comedy}\}$
- Won award  $A \in \{0, 1\}$
- Rating  $R \in \{1, 2, 3, 4, 5\}$



$$\mathbb{P}(G = g, A = a, R = r) = p_G(g)p_A(a)p_R(r | g, a)$$

# Example: v-structure



$$\mathcal{D}_{\text{train}} = \{(d, 0, 3), (d, 1, 5), (d, 0, 1), (c, 0, 5), (c, 1, 4)\}$$

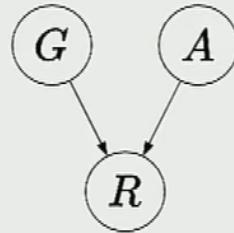
Parameters:  $\theta = (p_G, p_A, p_R)$

$\theta:$	$g \quad p_G(g)$	$a \quad p_A(a)$	$g \quad a \quad r \quad \text{count}_R(g, a, r)$
	d    3/5	0    3/5	d    0    1    1
	c    2/5	1    2/5	d    0    3    1 d    1    5    1 c    0    5    1 c    1    4    1

Be careful when you normalize  $P(r | g, a)$

When normalize, we focus on  $r$  only. For every possible unique setting  $(g, a)$ , I have a different distribution. So  $(d, 0)$  is a distribution over 1 or 3.

## Example: v-structure



$$\mathcal{D}_{\text{train}} = \{(d, 0, 3), (d, 1, 5), (d, 0, 1), (c, 0, 5), (c, 1, 4)\}$$

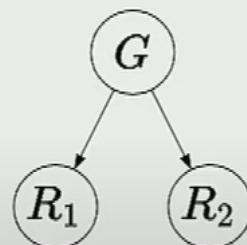
Parameters:  $\theta = (p_G, p_A, p_R)$

$\theta:$	$g \quad p_G(g)$	$a \quad p_A(a)$	$g \quad a \quad r \quad p_R(r   g, a)$
	d    3/5	0    3/5	d    0    1    1/2
	c    2/5	1    2/5	d    0    3    1/2
			d    1    5    1
			c    0    5    1
			c    1    4    1

## Example: inverted-v structure

Variables:

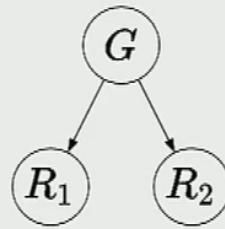
- Genre  $G \in \{\text{drama, comedy}\}$
- Jim's rating  $R_1 \in \{1, 2, 3, 4, 5\}$
- Martha's rating  $R_2 \in \{1, 2, 3, 4, 5\}$



$$\mathbb{P}(G = g, R_1 = r_1, R_2 = r_2) = p_G(g)p_{R_1}(r_1 | g)p_{R_2}(r_2 | g)$$

For R1 and R2, we are normalize based on g

## Example: inverted-v structure



$$\mathcal{D}_{\text{train}} = \{(d, 4, 5), (d, 4, 4), (d, 5, 3), (c, 1, 2), (c, 5, 4)\}$$

Parameters:  $\theta = (p_G, p_{R_1}, p_{R_2})$

$\theta:$

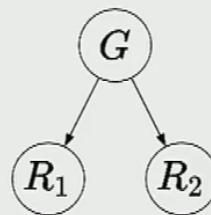
$g$	$p_G(g)$
d	3/5
c	2/5

$g$	$r_1$	$p_{R_1}(r   g)$
d	4	2/3
d	5	1/3
c	1	1/2
c	5	1/2

$g$	$r_2$	$p_{R_2}(r   g)$
d	3	1/3
d	4	1/3
d	5	1/3
c	2	1/2
c	4	1/2

We now want a pR to generally represent all generally rating, cuz we don't want a single distribution for each rating So we can combine R1 and R2

## Example: inverted-v structure



$$\mathcal{D}_{\text{train}} = \{(d, 4, 5), (d, 4, 4), (d, 5, 3), (c, 1, 2), (c, 5, 4)\}$$

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

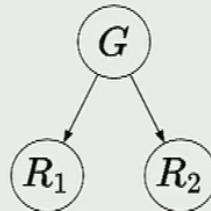
$\theta:$

$g$	$p_G(g)$
d	3/5
c	2/5

$g$	$r$	$\text{count}_R(g, r)$
d	3	1
d	4	3
d	5	2
c	1	1
c	2	1
c	4	1
c	5	1

Stanford

## Example: inverted-v structure



$$\mathcal{D}_{\text{train}} = \{(d, 4, 5), (d, 4, 4), (d, 5, 3), (c, 1, 2), (c, 5, 4)\}$$

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

$\theta:$

$g$	$p_G(g)$
d	3/5
c	2/5

$g$	$r$	$p_R(r   g)$
d	3	1/6
d	4	3/6
d	5	2/6
c	1	1/4
c	2	1/4
c	4	1/4
c	5	1/4

Stanford

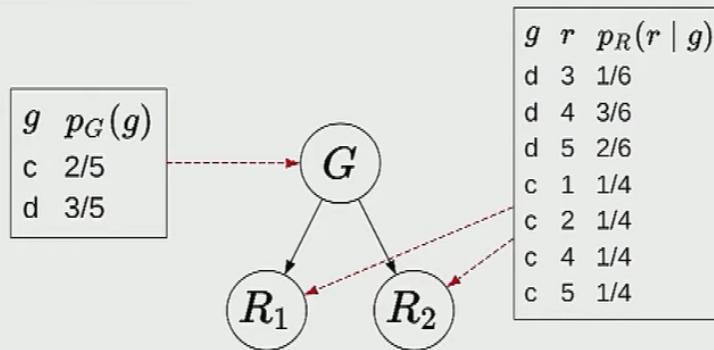
We are merging previous table R1 and R2 together Note that we are assuming all data in dataset are independent Also condition on  $g$ ,  $R_1$  and  $R_2$  are independent

## Parameter sharing



### Key idea: parameter sharing

The local conditional distributions of different variables use the same parameters.



Impact: more reliable estimates, less expressive model

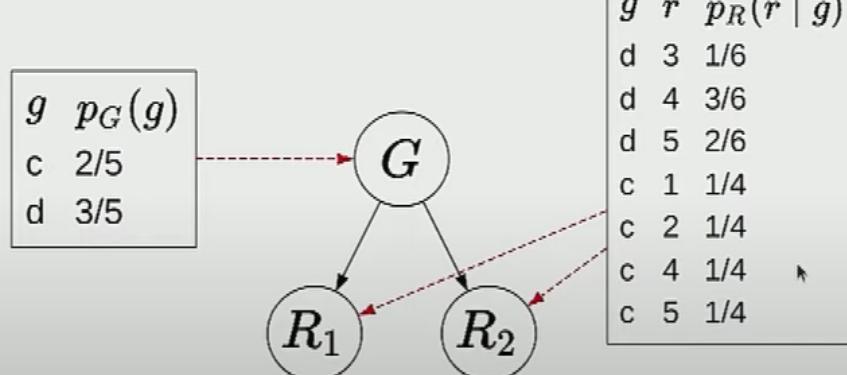
Stanford

## Parameter sharing



### Key idea: parameter sharing

The local conditional distributions of different variables use the same parameters.



Impact: more reliable estimates, less expressive model

Each node is gonna be powered by some table, So R1 and R2 now both powered by their joint table

## Laplace smoothing

## Unsupervised learning with EM

