

Example: one variable

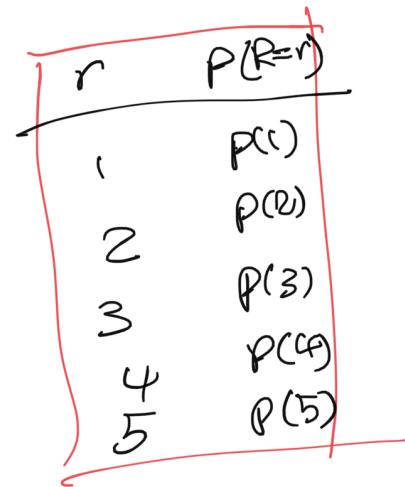
Setup:

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

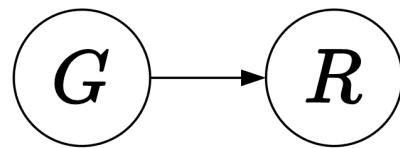
$$R \quad \mathbb{P}(R = r) = p(r)$$

Parameters:

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



Example: two variables



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

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

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

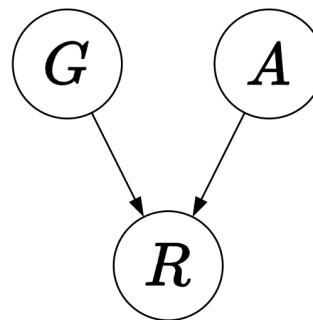
θ :

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

Red annotations: A red bracket groups the first two rows (d, 4) and (d, 5). A red bracket groups the last two rows (c, 1) and (c, 5). Red arrows point from the labels "g" and "r" to the first column and second column respectively.

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

Example: v-structure



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

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

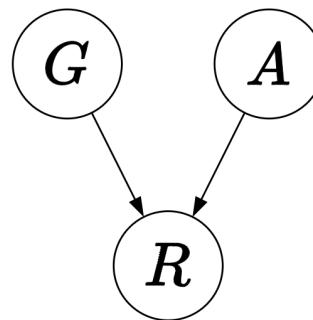
$\theta:$

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

a	$p_A(a)$
0	3/5
1	2/5

g	a	r	$\text{count}_R(g, a, r)$
d	0	1	1
d	0	3	1
d	1	5	1
c	0	5	1
c	1	4	1

Example: v-structure



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

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

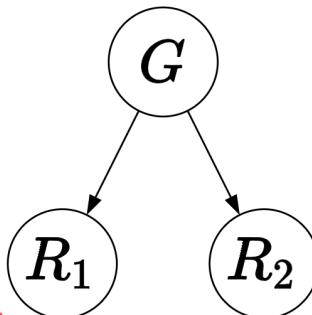
$\theta:$

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

a	$p_A(a)$
0	3/5
1	2/5

g	a	r	$\text{count}_R(g, a, r)$
d	0	1	1
d	0	3	1
d	1	5	1
c	0	5	1
c	1	4	1

Example: inverted-v structure



$$\underline{P(R_1|G)} = \underline{P(R_2|G)}$$

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

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

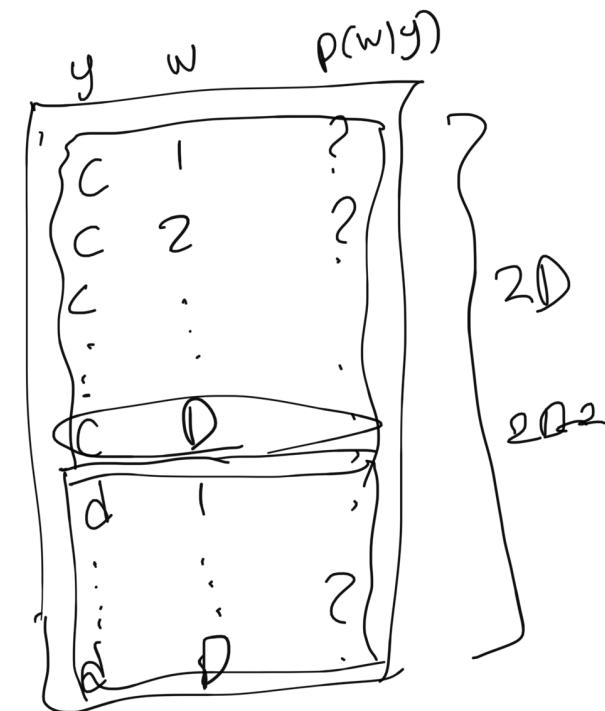
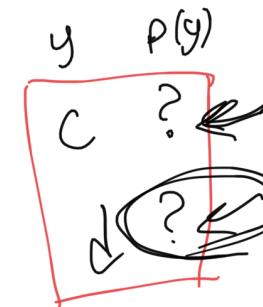


answer in chat

Question

If Y can take on 2 values and each W_j can take on D values, how many parameters are there?

activate deactivate reset report



General case

Bayesian network: variables X_1, \dots, X_n

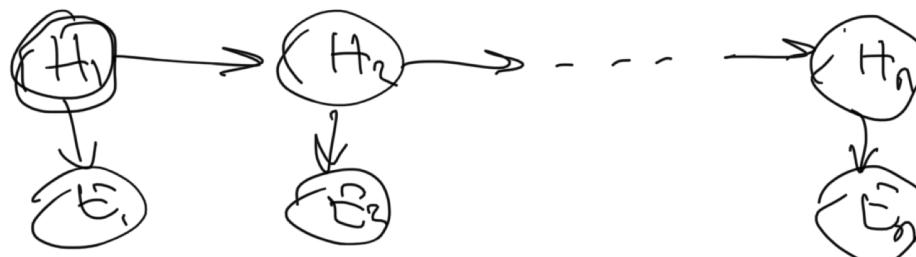
Parameters: collection of distributions $\theta = \{p_d : d \in D\}$
(e.g., $D = \{\text{start}, \text{trans}, \text{emit}\}$)

Each variable X_i is generated from distribution p_{d_i} :

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

Parameter sharing: d_i could be same for multiple i

$$d_{E_1} = \text{emit}$$
$$\vdots$$
$$d_{E_n} = \text{emit}$$



$$d_{H_1} = \text{start}$$
$$d_{H_2} = \text{trans}$$
$$d_{H_n} = \text{trans}$$

Maximum likelihood

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

$$\max_{p_G(\cdot), p_R(\cdot|c), p_R(\cdot|d)} (p_G(d)p_R(4|d)p_G(d)p_R(5|d)p_G(c)p_R(5|c))$$

subject to $\sum_g p_G(g) = 1$

$$\sum_r p_R(r|G=c) = 1$$

$$\sum_r p_R(r|G=d) = 1$$

Maximum likelihood

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

$$\max_{p_G(\cdot)} (p_G(d)p_G(d)p_G(c)) \quad \max_{p_R(\cdot|c)} p_R(5|c) \quad \max_{p_R(\cdot|d)} (p_R(4|d)p_R(5|d))$$
$$\sum_g p_G(g) = 1 \quad \sum_r p_R(r|G=c) = 1 \quad \sum_r p_R(r|G=d) = 1$$

Regularization: Laplace smoothing



Key idea: Laplace smoothing

For each distribution d and partial assignment $(x_{\text{Parents}(i)}, x_i)$, add λ to $\text{count}_d(x_{\text{Parents}(i)}, x_i)$.
Then normalize to get probability estimates.

Interpretation: hallucinate λ occurrences of each local assignment

Larger $\lambda \Rightarrow$ more smoothing \Rightarrow probabilities closer to uniform.

for each $(x_{\text{parents}}, y_i)$ we add ^{x imaginary}
data points where $(x_{\text{parents}}, x_i)$ was
"observed".

Maximum likelihood

$$\mathcal{D}_{\text{train}} = \{(d, 4), (d, 5), (c, 5)\} \Rightarrow p_G(d)^2 p_G(c)^1$$
$$\max_{p_G(\cdot)} (p_G(d)p_G(d)p_G(c)) \quad \max_{p_R(\cdot|c)} p_R(5|c) \quad \max_{p_R(\cdot|d)} (p_R(4|d)p_R(5|d))$$
$$p_G(d) + p_G(c) \leq 1 \quad p_R(1|c) + p_R(2|c) + \dots + p_R(5|c) \leq 1$$

Solution:

$$p_G(d) = \frac{2}{3}, p_G(c) = \frac{1}{3}, p_R(5|c) = 1, p_R(4|d) = \frac{1}{2}, p_R(5|d) = \frac{1}{2}$$

- Key: decomposes into subproblems, one for each distribution d and assignment x_{Parents}
- For each subproblem, solve in closed form

$$\max x^5 y^2 z^3 \text{ subject to } x + y + z = 1$$

Maximum likelihood

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

$$\max_{p_G(\cdot)} (p_G(d)p_G(d)p_G(c)) \max_{p_R(\cdot|c)} p_R(5 | c) \max_{p_R(\cdot|d)} (p_R(4 | d)p_R(5 | d))$$

Solution:

$$p_G(d) = \frac{2}{3}, p_G(c) = \frac{1}{3}, p_R(5 | c) = 1, p_R(4 | d) = \frac{1}{2}, p_R(5 | d) = \frac{1}{2}$$

- Key: decomposes into subproblems, one for each distribution d and assignment x_{Parents}
- For each subproblem, solve in closed form

$$\max x^5 y^2 z^3 \text{ subject to } x + y + z = 1$$

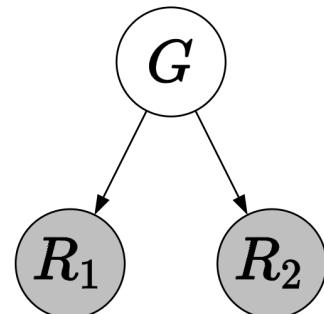
$x = \frac{5}{5+2+3}$ $y = \frac{2}{5+2+3}$ $z = \frac{3}{5+2+3}$

is achieved when $x \propto 5, y \propto 2, z \propto 3$

Maximum marginal likelihood

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

Example:



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

$$P(E=e) = \sum_h P(E=e \wedge H=h)$$

Maximum marginal likelihood objective:

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

Expectation Maximization (EM)

$$\alpha P(H=h_1 | E=e_1) \\ \alpha P(H=h_2 | E=e_2) \\ \vdots \\ \alpha P(H=h_s | E=e_s)$$

Intuition: generalization of the K-means algorithm

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



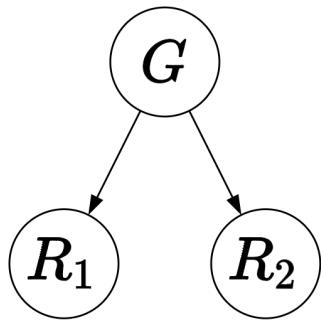
Algorithm: Expectation Maximization (EM)

Initialize θ

E-step:

- Compute $q(h) = \mathbb{P}(H = h | E = e; \theta)$ for each h
(use any probabilistic inference algorithm)
- Create weighted points: (h, e) with weight $q(h)$

Example: one iteration of EM



$$\mathcal{D}_{\text{train}} = \{(\textcolor{red}{?}, 2, 2), (\textcolor{red}{?}, 1, 2)\}$$

$$\propto P(L|G=g, R_1=r_1, R_2=r_2)$$

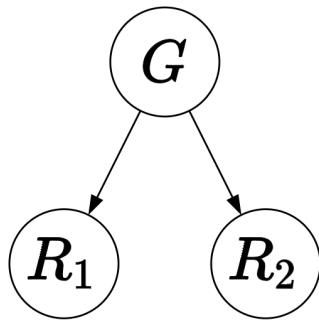
$\theta:$	g	$p_G(g)$
c	0.5	0.5
c	2	0.6
d	1	0.6
d	2	0.4

g	r	$p_R(r g)$
c	1	0.4
c	2	0.6
d	1	0.6
d	2	0.4

E-step

(r_1, r_2)	g	$\mathbb{P}(G=g, R_1=r_1, R_2=r_2)$	$q(g)$
(2, 2)	c	$0.5 \cdot 0.6 \cdot 0.6 = 0.18$	$\frac{0.18}{0.18+0.08} = 0.69$
(2, 2)	d	$0.5 \cdot 0.4 \cdot 0.4 = 0.08$	$\frac{0.08}{0.18+0.08} = 0.31$
(1, 2)	c	$0.5 \cdot 0.4 \cdot 0.6 = 0.12$	$\frac{0.12}{0.12+0.12} = 0.5$
(1, 2)	d	$0.5 \cdot 0.6 \cdot 0.4 = 0.12$	$\frac{0.12}{0.12+0.12} = 0.5$

Example: one iteration of EM



$$\mathcal{D}_{\text{train}} = \{(\textcolor{red}{?}, 2, 2), (\textcolor{red}{?}, 1, 2)\}$$

$\theta:$	g	$p_G(g)$
	g	0.5
	c	0.5
	d	0.5

g	r	$p_R(r g)$
g	r	$p_R(r g)$
c	1	0.4
c	2	0.6
d	1	0.6
d	2	0.4

E-step

(r_1, r_2)	g	$\mathbb{P}(G = g, R_1 = r_1, R_2 = r_2)$	$q(g)$
(2, 2)	c	$0.5 \cdot 0.6 \cdot 0.6 = 0.18$	$\frac{0.18}{0.18+0.08} = 0.69$
(2, 2)	d	$0.5 \cdot 0.4 \cdot 0.4 = 0.08$	$\frac{0.08}{0.18+0.08} = 0.31$
(1, 2)	c	$0.5 \cdot 0.4 \cdot 0.6 = 0.12$	$\frac{0.12}{0.12+0.12} = 0.5$
(1, 2)	d	$0.5 \cdot 0.6 \cdot 0.4 = 0.12$	$\frac{0.12}{0.12+0.12} = 0.5$

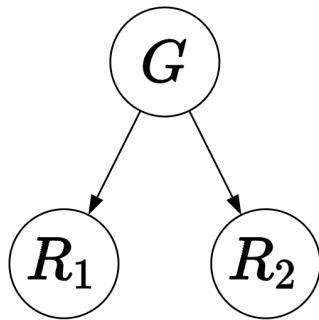
M-step

$\rightarrow \theta:$

g	count	$p_G(g)$
c	1	0.5
c	$0.69 + 0.5$	0.59
d	$0.31 + 0.5$	0.41

g	r	count	$p_R(r g)$
c	1	0.5	0.21
c	2	$0.5 + 0.69 + 0.69$	0.79
d	1	0.5	0.31
d	2	$0.5 + 0.31 + 0.31$	0.69

Example: one iteration of EM



$$\mathcal{D}_{\text{train}} = \{(\textcolor{red}{?}, 2, 2), (\textcolor{red}{?}, 1, 2)\}$$

$\theta:$	g	$p_G(g)$
	g	0.5
	c	0.4
	c	0.5
	d	0.6
	d	0.6
	d	0.4

g	r	$p_R(r g)$
c	1	0.4
c	2	0.6
d	1	0.6
d	2	0.4

E-step

(r_1, r_2)	g	$\mathbb{P}(G = g, R_1 = r_1, R_2 = r_2)$	$q(g)$
(2, 2)	c	$0.5 \cdot 0.6 \cdot 0.6 = 0.18$	$\frac{0.18}{0.18+0.08} = 0.69$
(2, 2)	d	$0.5 \cdot 0.4 \cdot 0.4 = 0.08$	$\frac{0.08}{0.18+0.08} = 0.31$
(1, 2)	c	$0.5 \cdot 0.4 \cdot 0.6 = 0.12$	$\frac{0.12}{0.12+0.12} = 0.5$
(1, 2)	d	$0.5 \cdot 0.6 \cdot 0.4 = 0.12$	$\frac{0.12}{0.12+0.12} = 0.5$

M-step

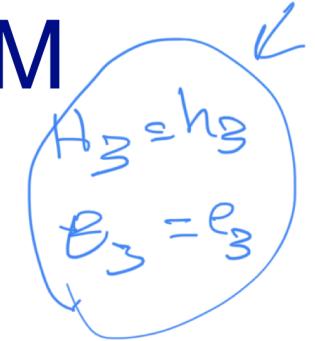
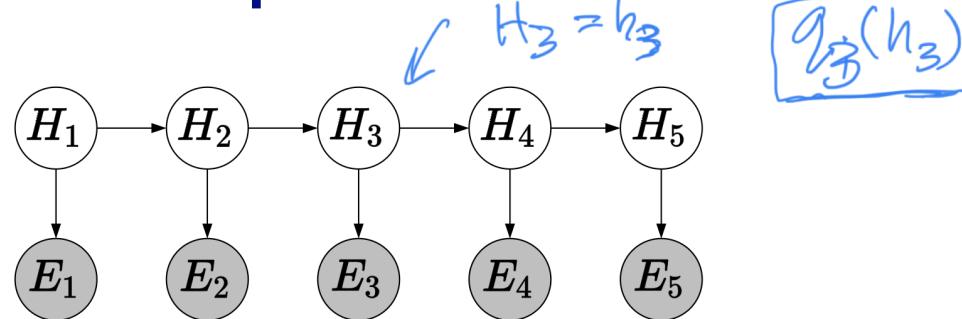
$\rightarrow \theta:$

g	count	$p_G(g)$
c	1	0.59
c	0.69 + 0.5	0.59
d	0.31 + 0.5	0.41

g	r	count	$p_R(r g)$
c	1	0.5	0.21
c	2	0.5 + 0.69 + 0.69	0.79
d	1	0.5	0.31
d	2	0.5 + 0.31 + 0.31	0.69

Increase marginal likelihood

Application: decipherment as an HMM



E-step: forward-backward algorithm computes

$$q_i(h) \stackrel{\text{def}}{=} \mathbb{P}(H_i = h \mid E_1 = e_1, \dots, E_n = e_n)$$

M-step: count (fractional) and normalize

$$\text{count}_{\text{emit}}(h, e) = \sum_{i=1}^n q_i(h) \cdot [e_i = e]$$

$$p_{\text{emit}}(e \mid h) \propto \text{count}_{\text{emit}}(h, e)$$

