

CSC242: Introduction to Artificial Intelligence

Lecture 4.5

Announcements

- Unit 4 Exam: Thu 27 Apr 940AM
- Project 4 due Thu 27 Apr 1159PM
- Final Exam: Thu 11 May 1600
Douglass Ballroom (not Dewey 1-101)
 - **BRING ID**




Bags


Agent, process, disease, ...




Candies

Actions, effects,
symptoms, results of
tests, ...

$D_1 =$ 

$D_2 =$ 

Observations

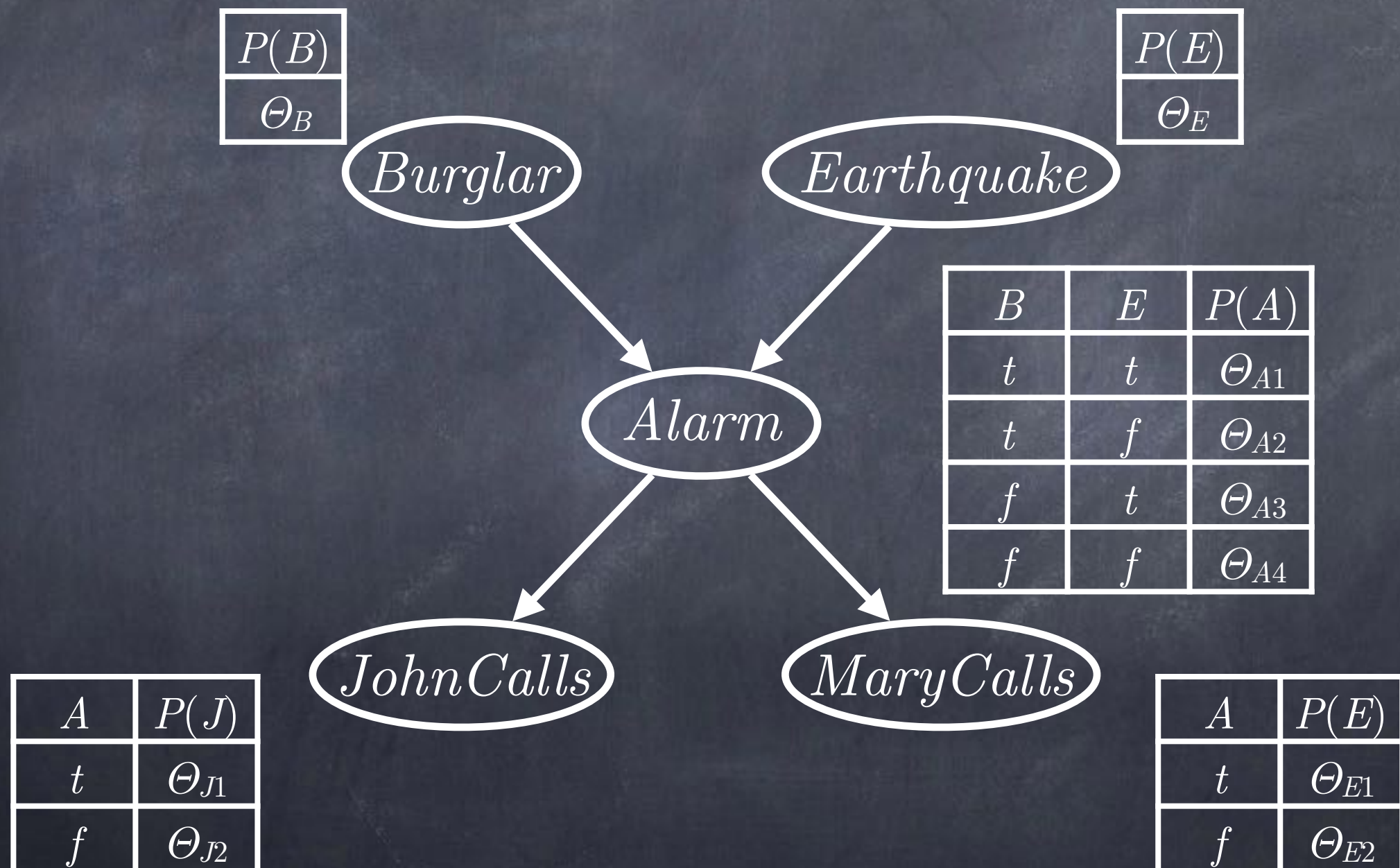
$D_3 =$ 

Goal

Predict next
candy

Predict agent's next move
Predict next output of process
Predict disease given symptoms
and tests

Parameter Learning (in Bayesian Networks)



Independent Identically Distributed (i.i.d.)

- Probability of a sample is independent of any previous samples

$$\mathbf{P}(D_i | D_{i-1}, D_{i-2}, \dots) = \mathbf{P}(D_i)$$

- Probability distribution doesn't change among samples

$$\mathbf{P}(D_i) = \mathbf{P}(D_{i-1}) = \mathbf{P}(D_{i-2}) = \dots$$

Maximum Likelihood Hypothesis

$$\operatorname{argmax}_{\Theta} P(\mathbf{d} \mid h_{\Theta})$$



h_{Θ}

$P(F=cherry)$	$P(F=lime)$
Θ	$1-\Theta$

Flavor

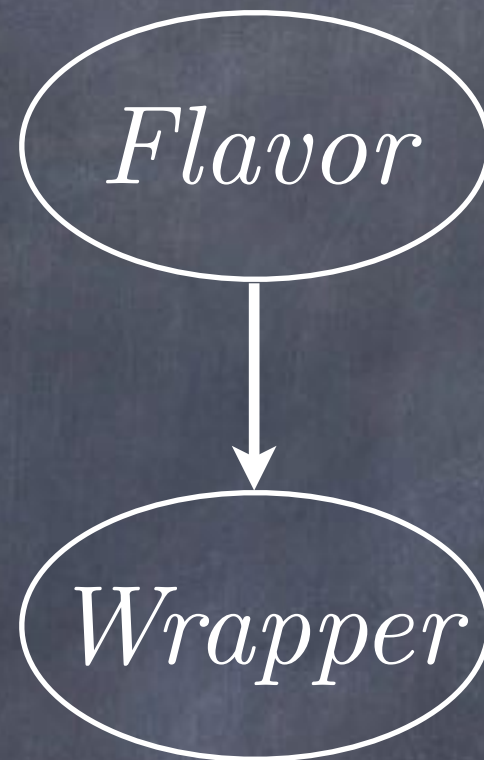
Maximum Likelihood Hypothesis

$$L(\mathbf{d} \mid h_{\Theta}) = c \log \Theta + l \log(1 - \Theta)$$

$$\operatorname{argmax}_{\Theta} L(\mathbf{d} \mid h_{\Theta}) = \frac{c}{c + l} = \frac{c}{N}$$

$h_{\Theta, \Theta_1, \Theta_2}$

$P(F=c)$	$P(F=l)$
Θ	$1-\Theta$



F	$P(W=r F)$	$P(W=g F)$
c	Θ_1	$1-\Theta_1$
l	Θ_2	$1-\Theta_2$

Maximum Likelihood Hypothesis

$$\Theta = \frac{c}{c+l} = \frac{c}{N}$$

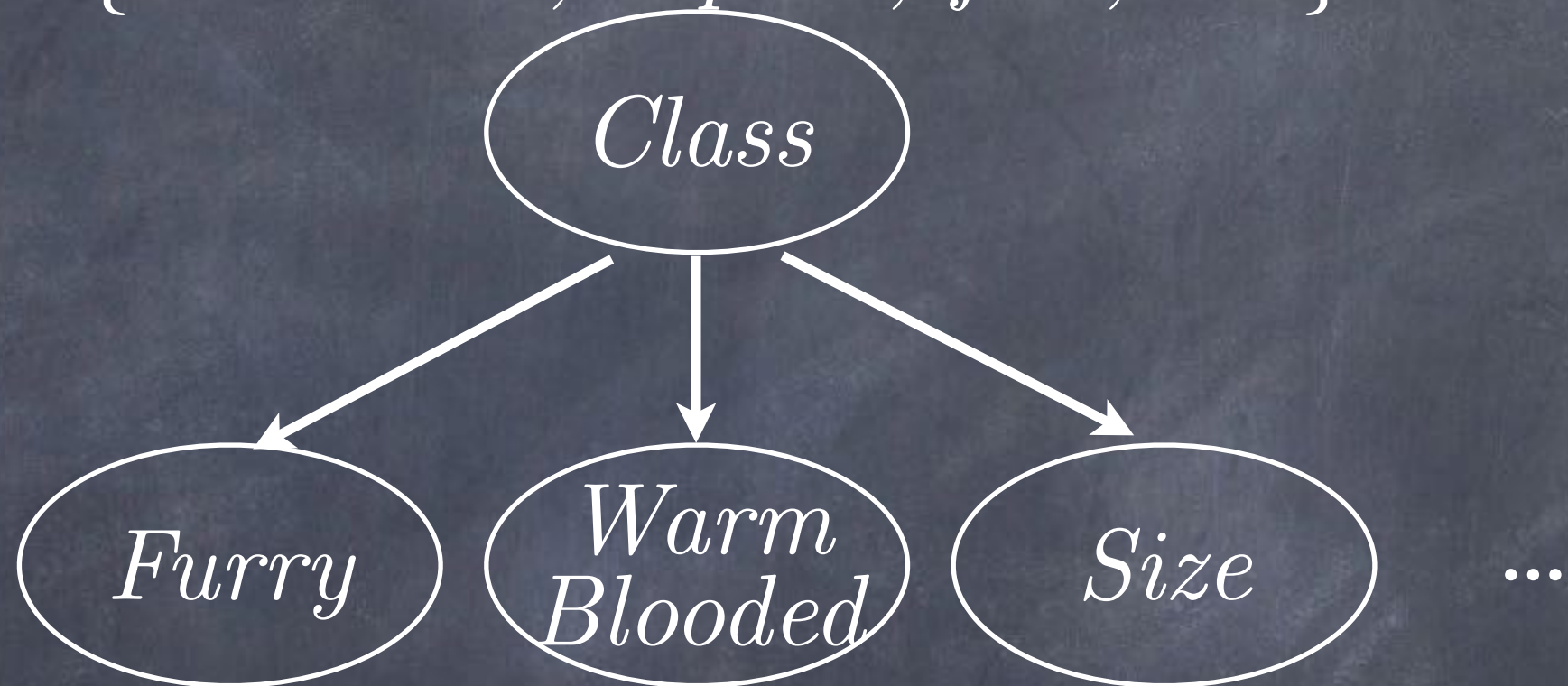
$$\Theta_1 = \frac{r_c}{r_c + g_c} = \frac{r_c}{c}$$

$$\Theta_2 = \frac{r_l}{r_l + g_l} = \frac{r_l}{l}$$

Observed frequencies are the BEST hypothesis, in terms of maximizing the likelihood of the data.

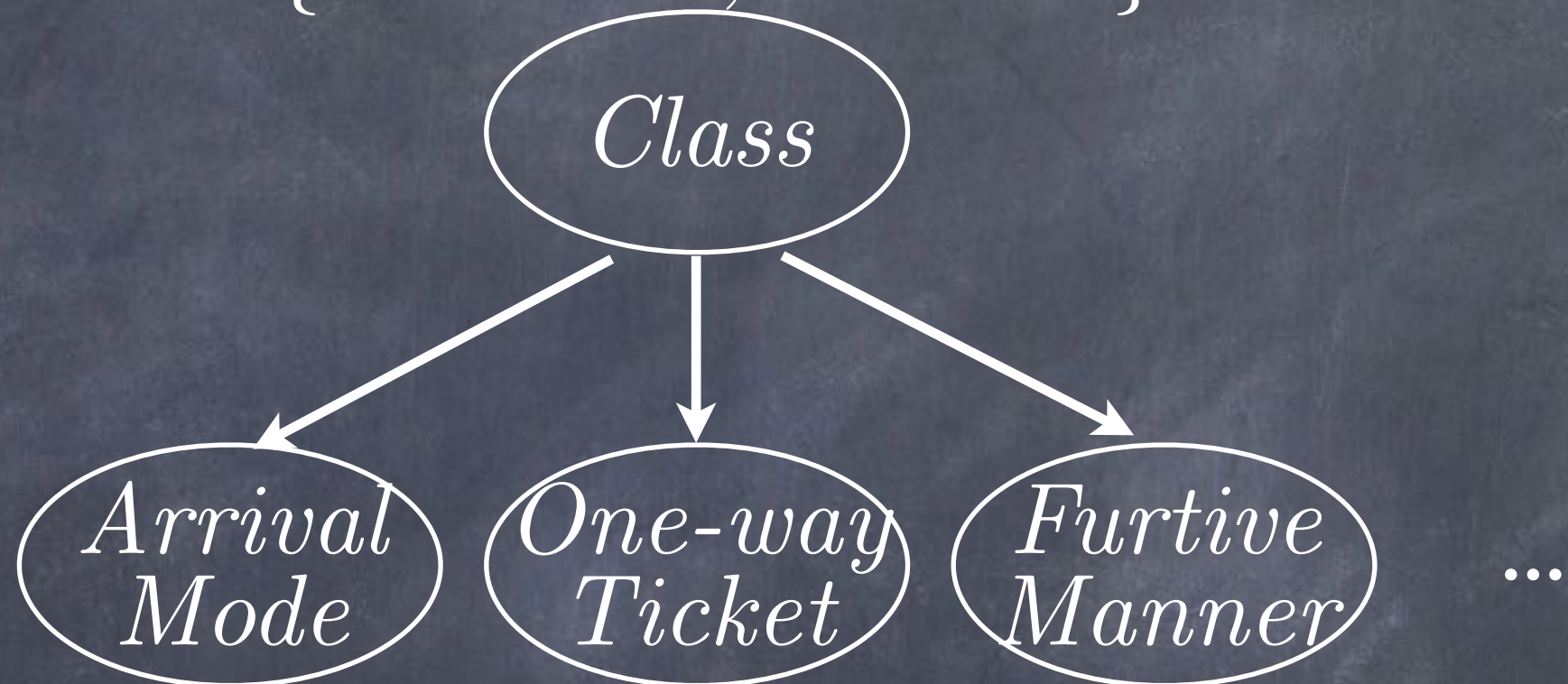
Naive Bayes Models

$\{ \textit{mammal}, \textit{reptile}, \textit{fish}, \dots \}$

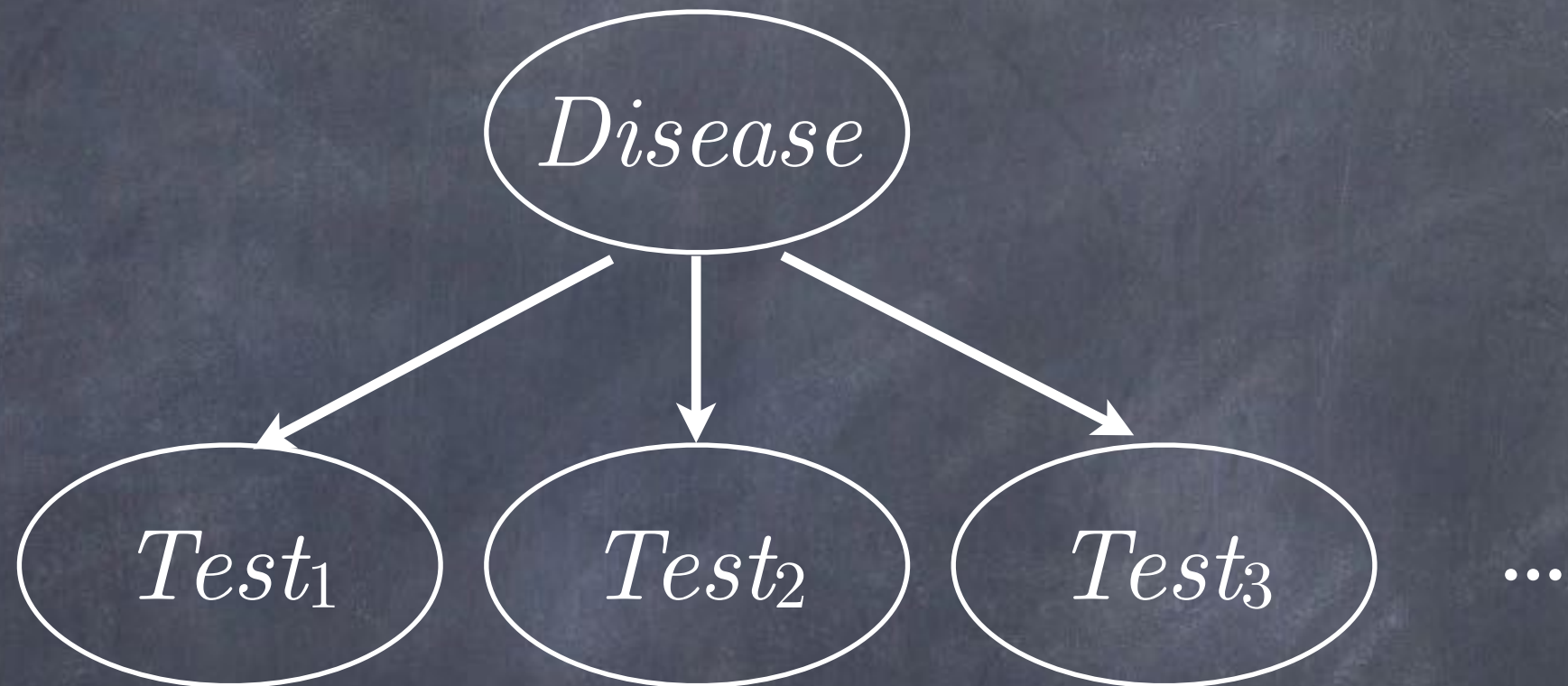


Naive Bayes Models

$\{ \textit{terrorist}, \textit{tourist} \}$



Naive Bayes Models



Learning Naive Bayes Models

- Naive Bayes model with n Boolean attributes requires $2n+1$ parameters
- Maximum likelihood hypothesis can be found with no search
 - Probabilities are observed frequencies
 - Scales to large problems
- Robust to noisy or missing data

Naive Bayes Classifier

New input: x_1, \dots, x_n

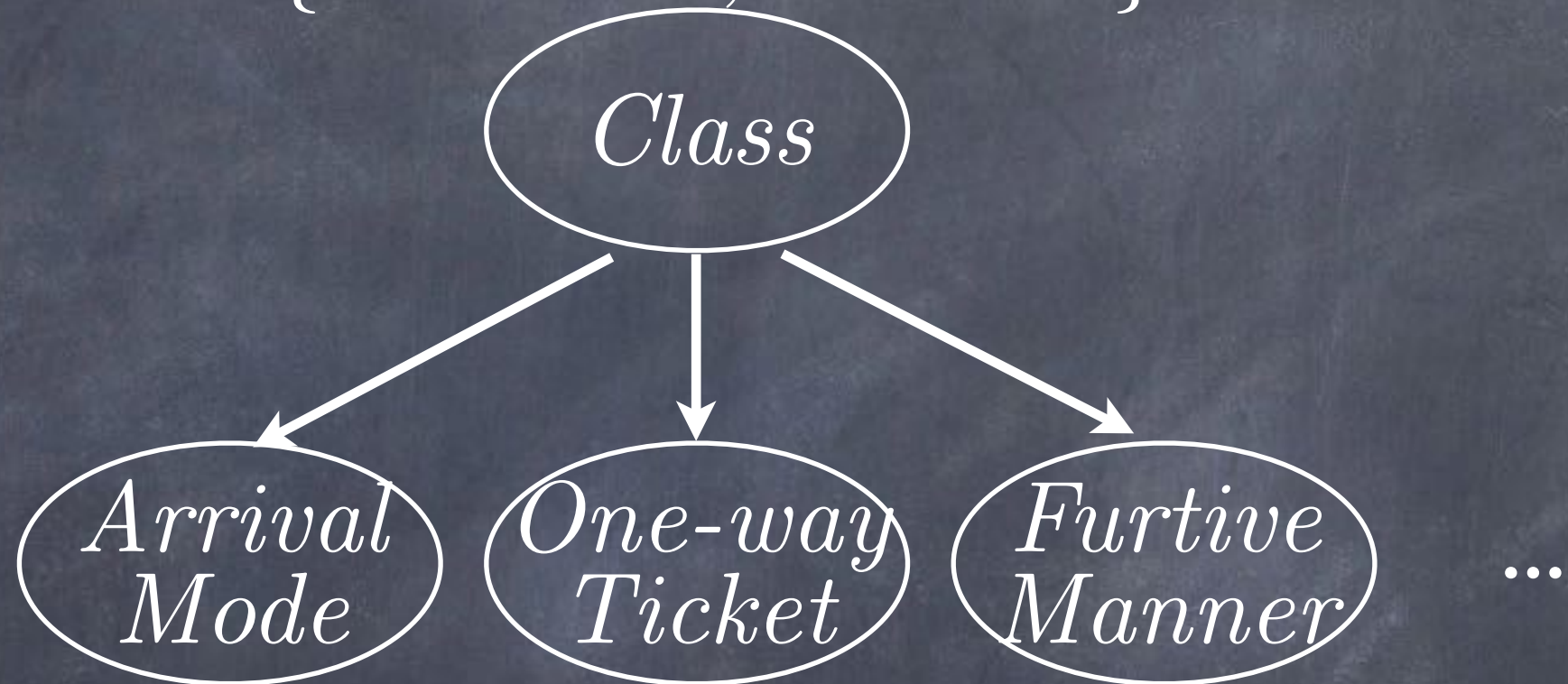
$$\begin{aligned}\mathbf{P}(C \mid x_1, \dots, x_n) &= \alpha \mathbf{P}(x_1, \dots, x_n \mid C) \mathbf{P}(C) \\ &= \alpha \mathbf{P}(C) \prod_i \mathbf{P}(x_i \mid C)\end{aligned}$$

Parameter Learning in Bayesian Networks

- Can learn the CPTs for a Bayes Net from observations that include values for all variables
- Finding maximum likelihood parameters decomposes into separate problems, one for each parameter
- Parameter values for a variable given its parents are the observed frequencies

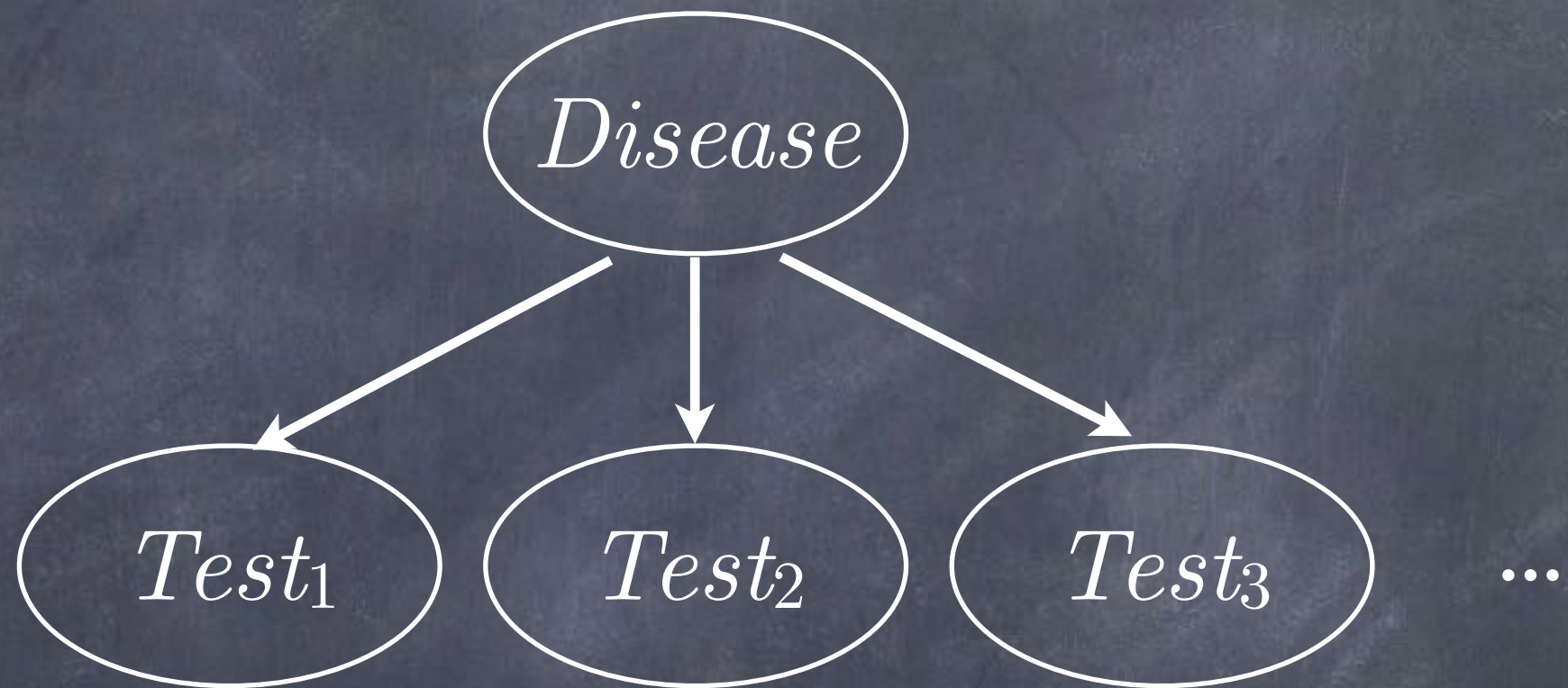


$\{ \textit{terrorist}, \textit{tourist} \}$

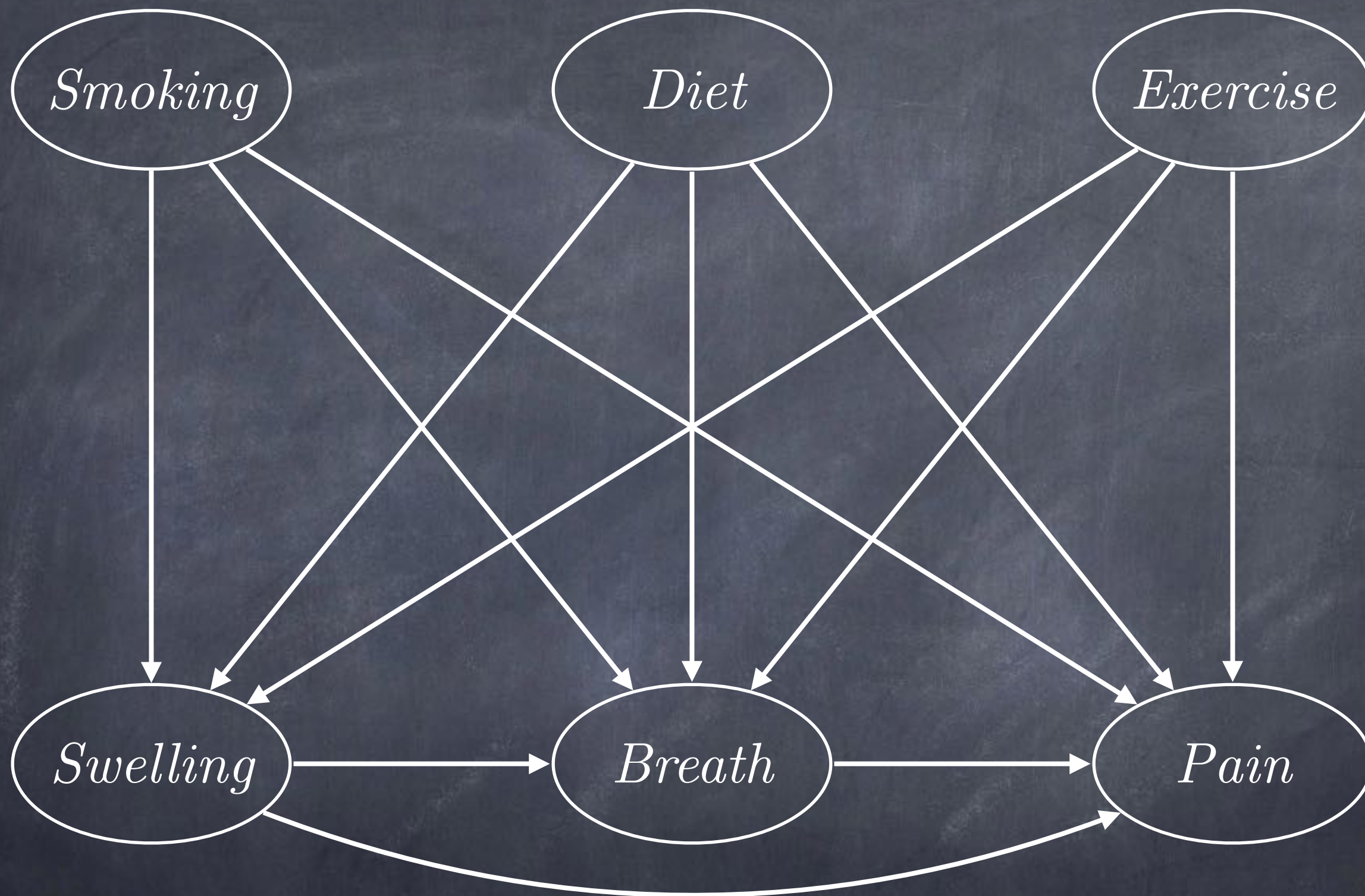


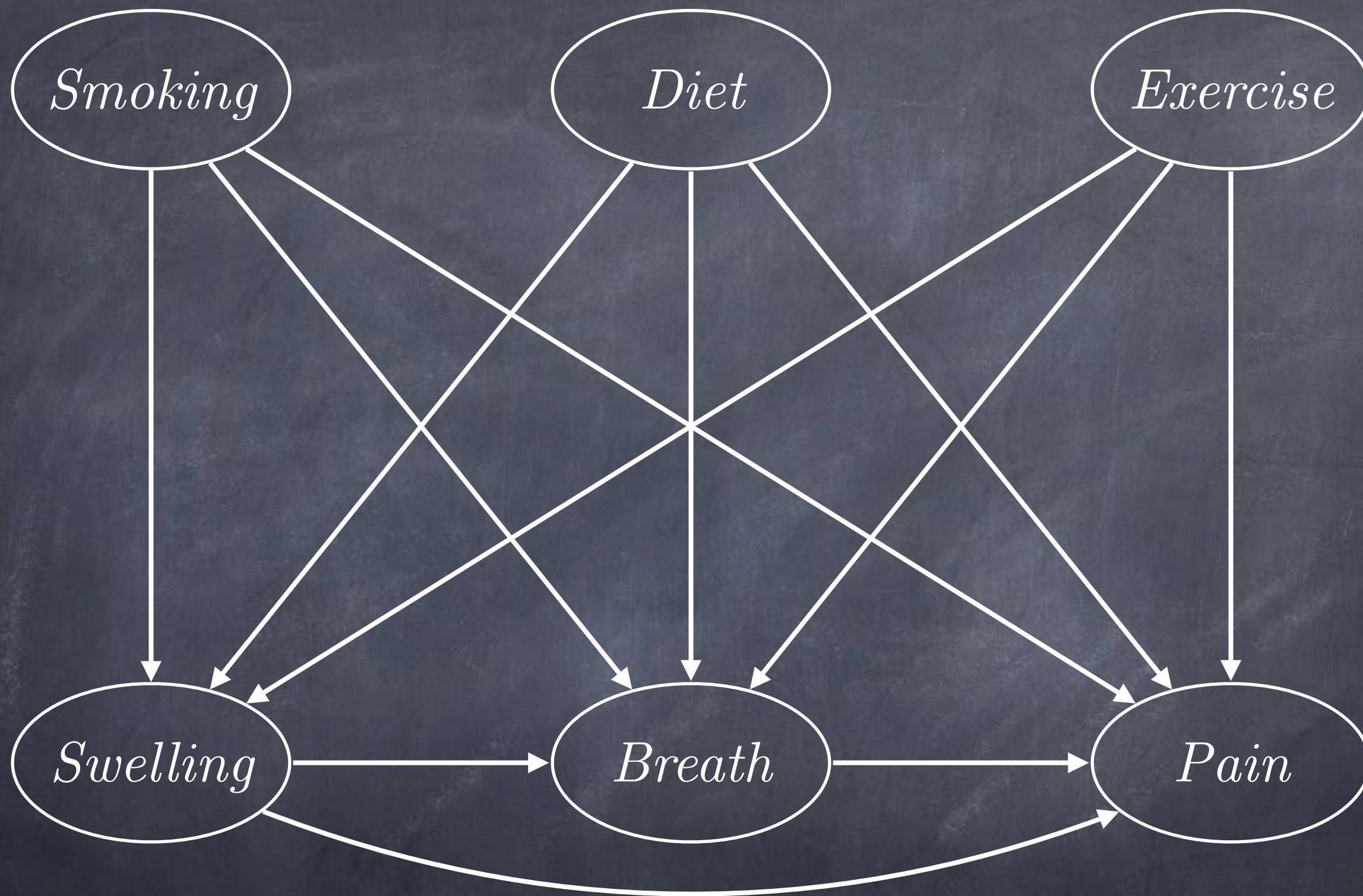
<i>Arrival</i>	<i>One-Way</i>	<i>Furtive</i>	<i>...</i>	<i>Class</i>
<i>taxi</i>	<i>yes</i>	<i>very</i>	<i>...</i>	<i>terrorist</i>
<i>car</i>	<i>no</i>	<i>none</i>	<i>...</i>	<i>tourist</i>
<i>car</i>	<i>yes</i>	<i>very</i>	<i>...</i>	<i>terrorist</i>
<i>car</i>	<i>yes</i>	<i>some</i>	<i>...</i>	<i>tourist</i>
<i>walk</i>	<i>yes</i>	<i>none</i>	<i>...</i>	<i>student</i>
<i>bus</i>	<i>no</i>	<i>some</i>	<i>...</i>	<i>tourist</i>

<i>Arrival</i>	<i>One-Way</i>	<i>Furtive</i>	<i>...</i>	<i>Class</i>
<i>taxi</i>	<i>yes</i>	<i>very</i>	<i>...</i>	<i>terrorist?</i>
<i>car</i>	<i>no</i>	<i>none</i>	<i>...</i>	<i>tourist?</i>
<i>car</i>	<i>yes</i>	<i>very</i>	<i>...</i>	<i>terrorist?</i>
<i>car</i>	<i>yes</i>	<i>some</i>	<i>...</i>	<i>tourist?</i>
<i>walk</i>	<i>yes</i>	<i>none</i>	<i>...</i>	<i>student?</i>
<i>bus</i>	<i>no</i>	<i>some</i>	<i>...</i>	<i>tourist?</i>

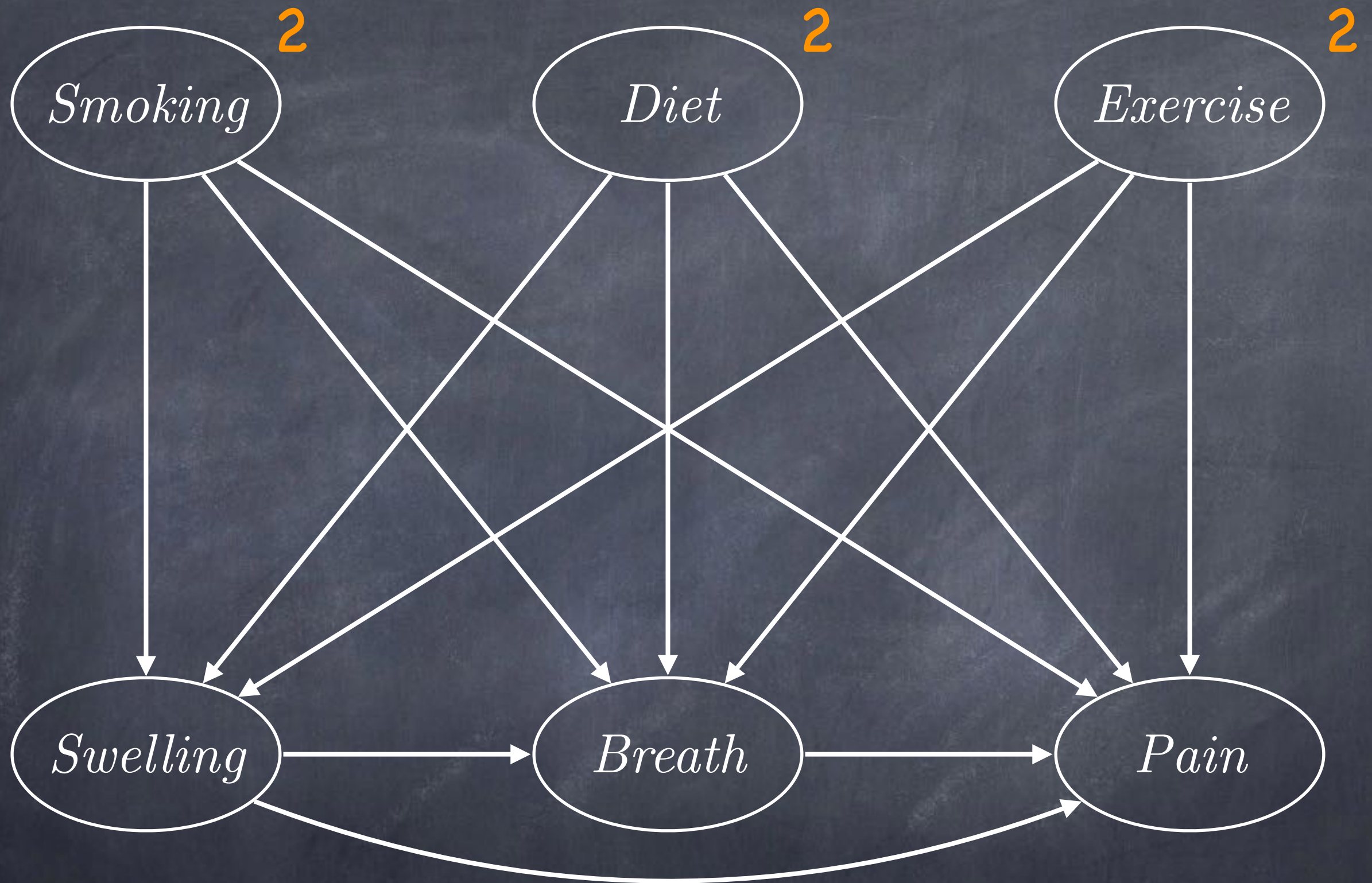


<i>Test</i>	<i>Test2</i>	<i>Test3</i>	<i>...</i>	<i>Disease</i>
<i>T</i>	<i>F</i>	<i>T</i>	<i>...</i>	<i>?</i>
<i>T</i>	<i>F</i>	<i>F</i>	<i>...</i>	<i>?</i>
<i>F</i>	<i>F</i>	<i>T</i>	<i>...</i>	<i>?</i>
<i>T</i>	<i>T</i>	<i>T</i>	<i>...</i>	<i>?</i>
<i>F</i>	<i>T</i>	<i>F</i>	<i>...</i>	<i>?</i>
<i>T</i>	<i>F</i>	<i>T</i>	<i>...</i>	<i>?</i>

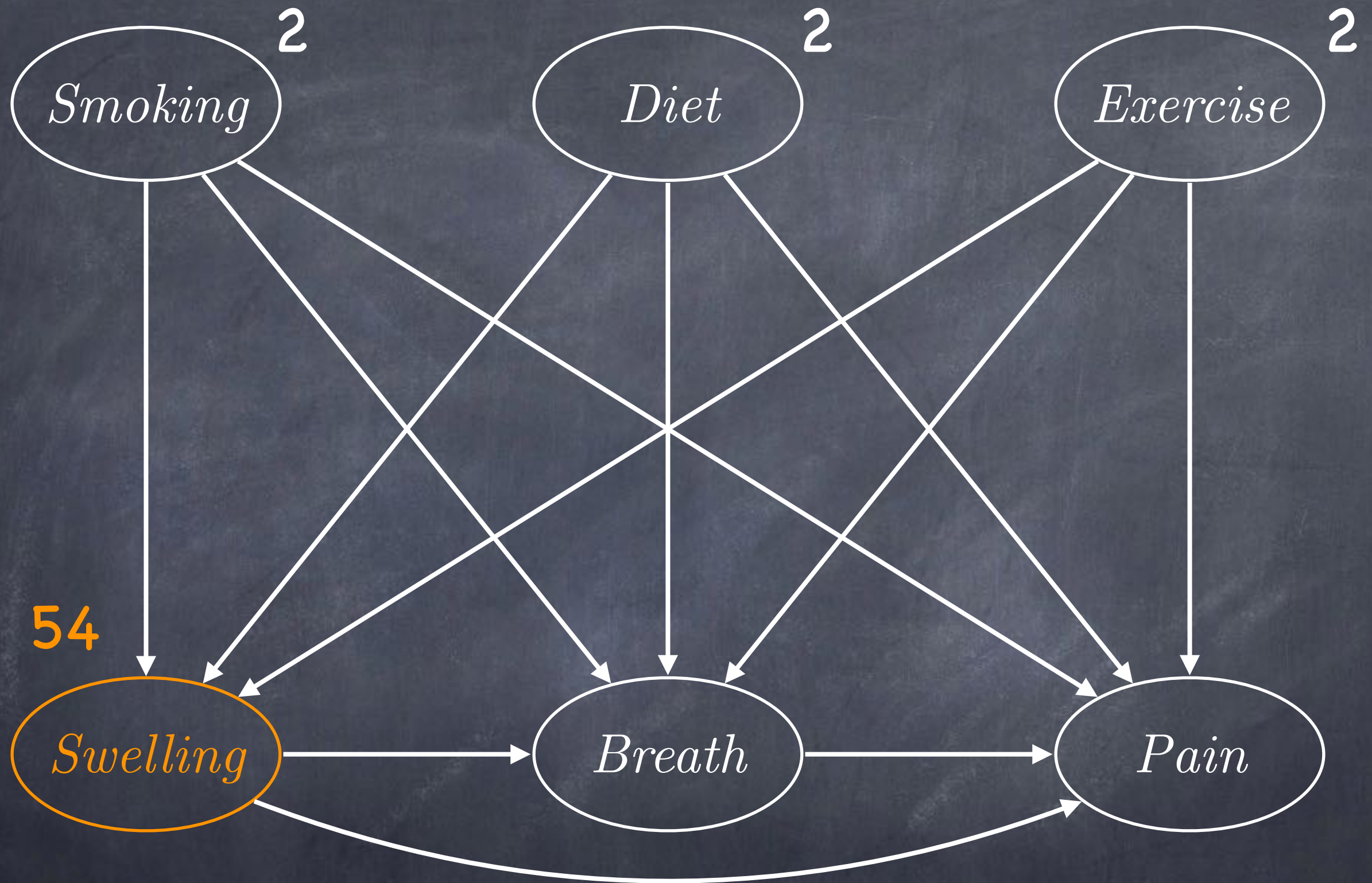




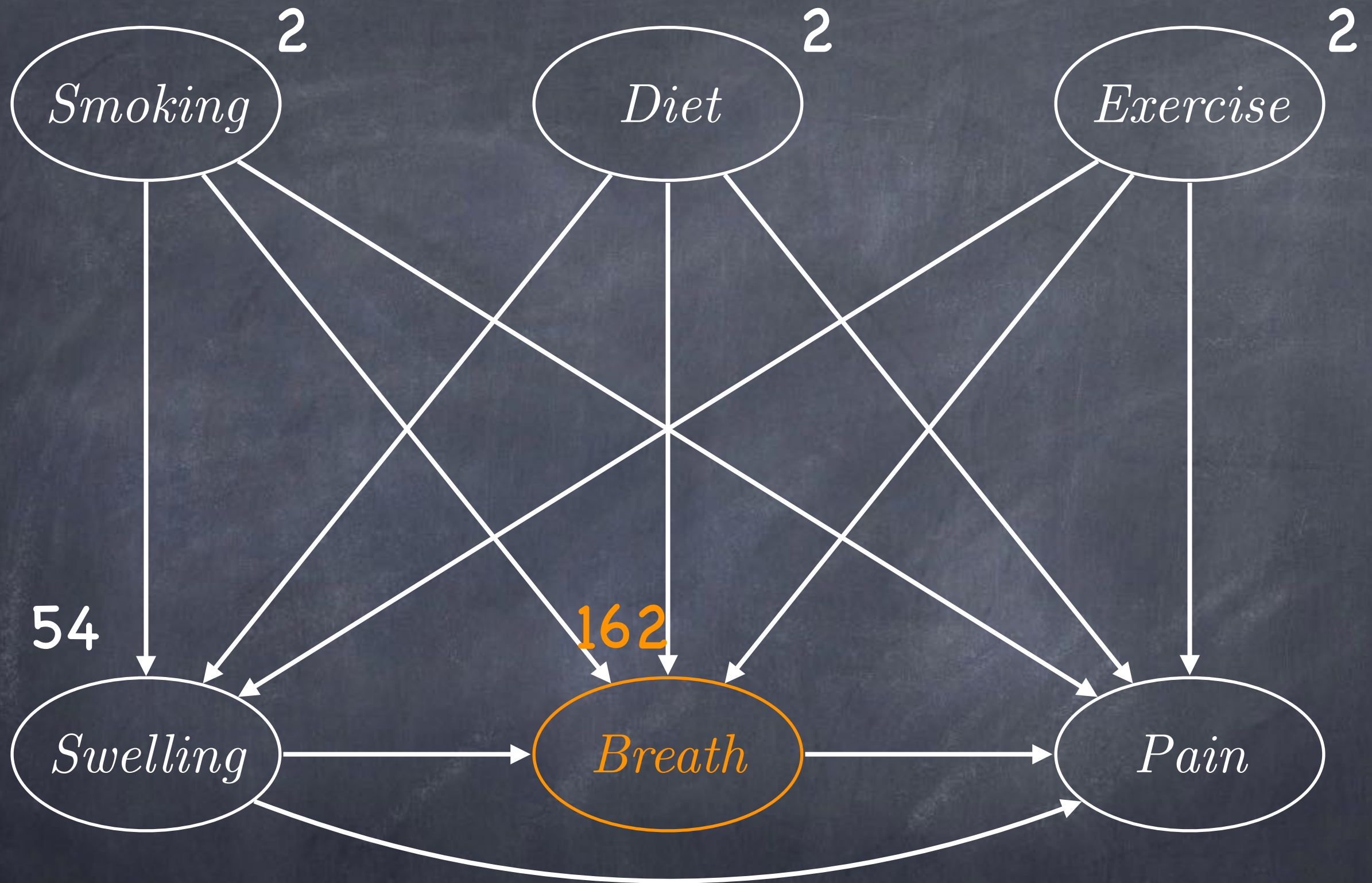
3 values/variable



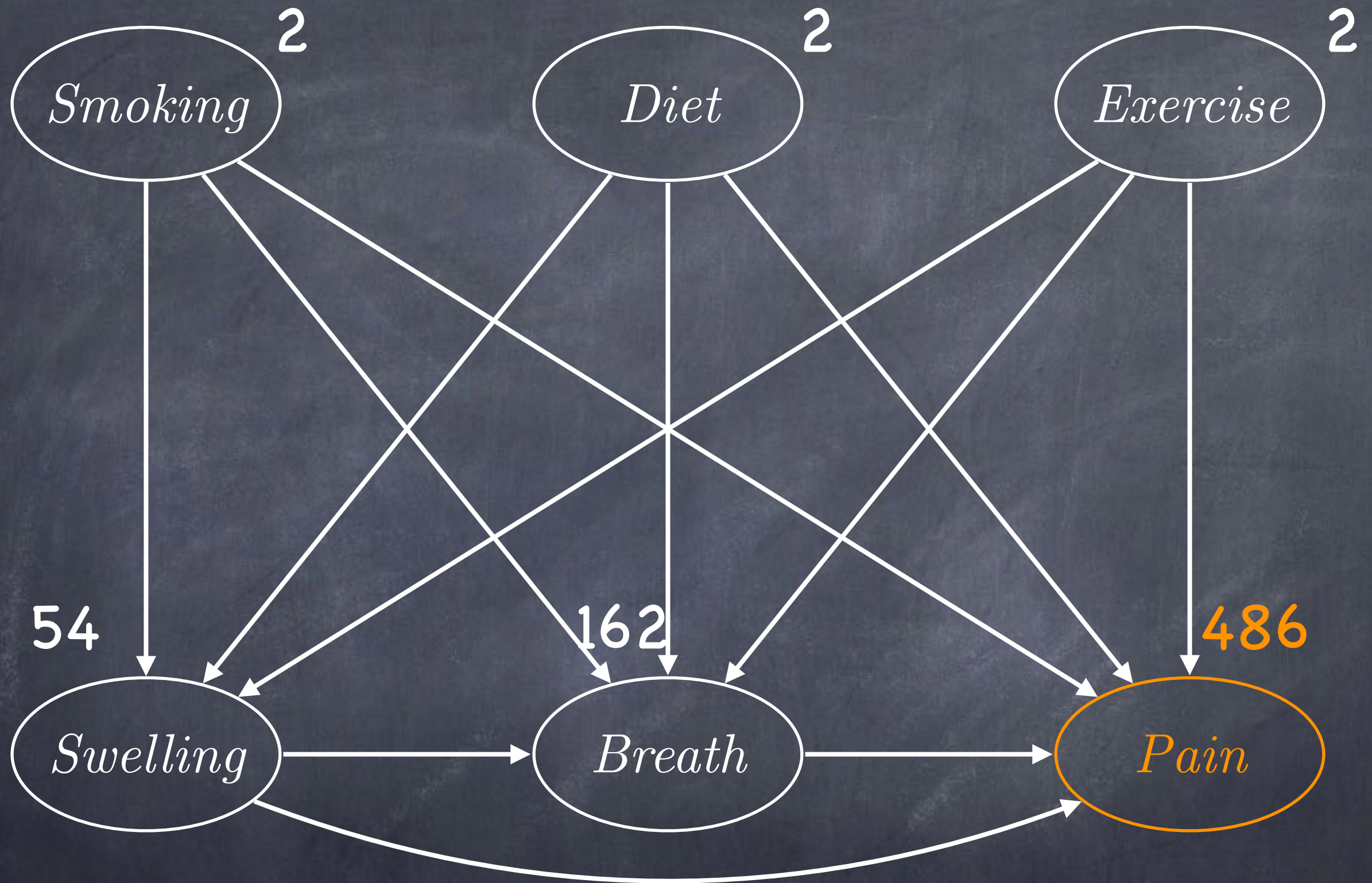
3 values/variable



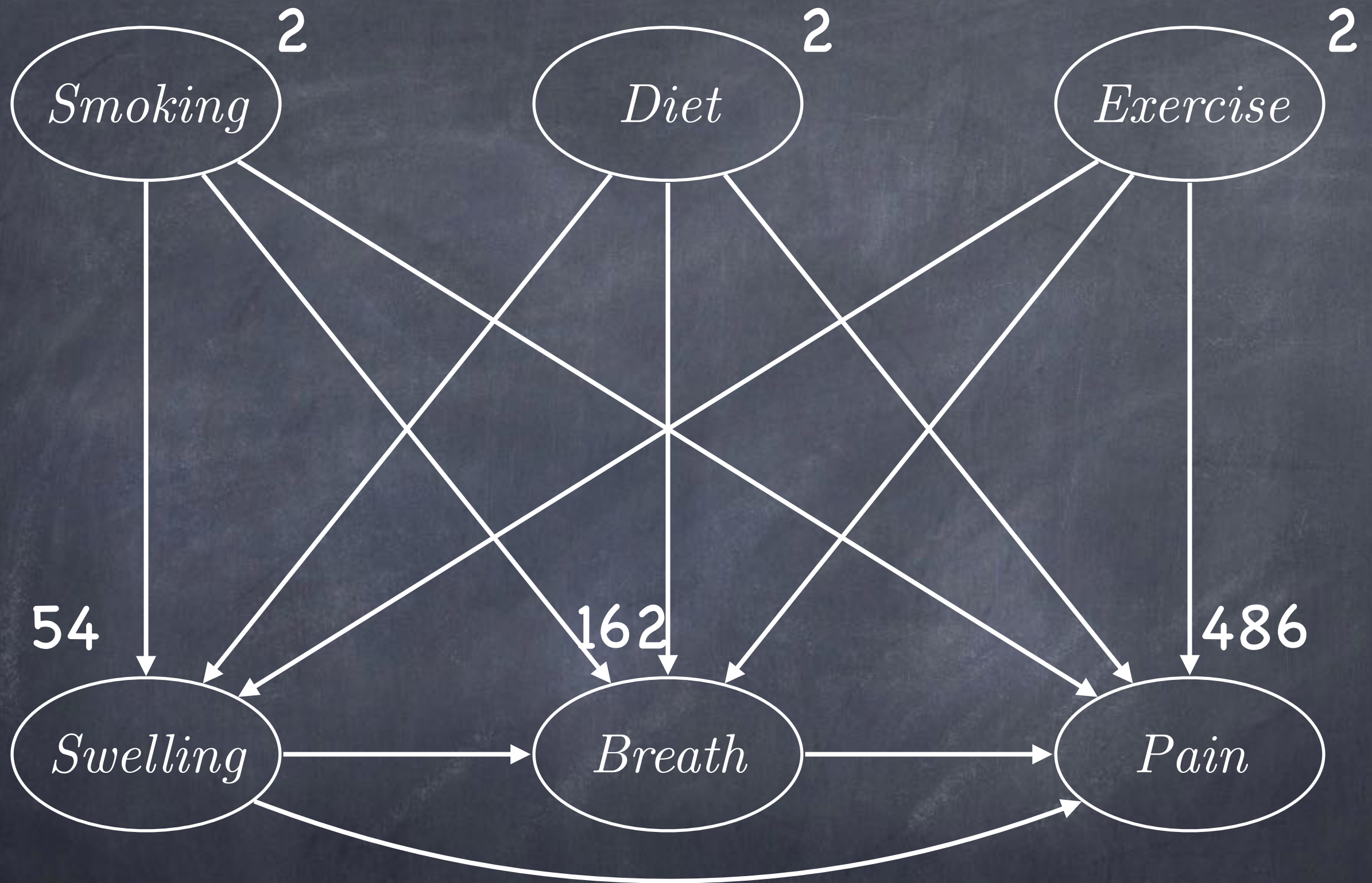
3 values/variable



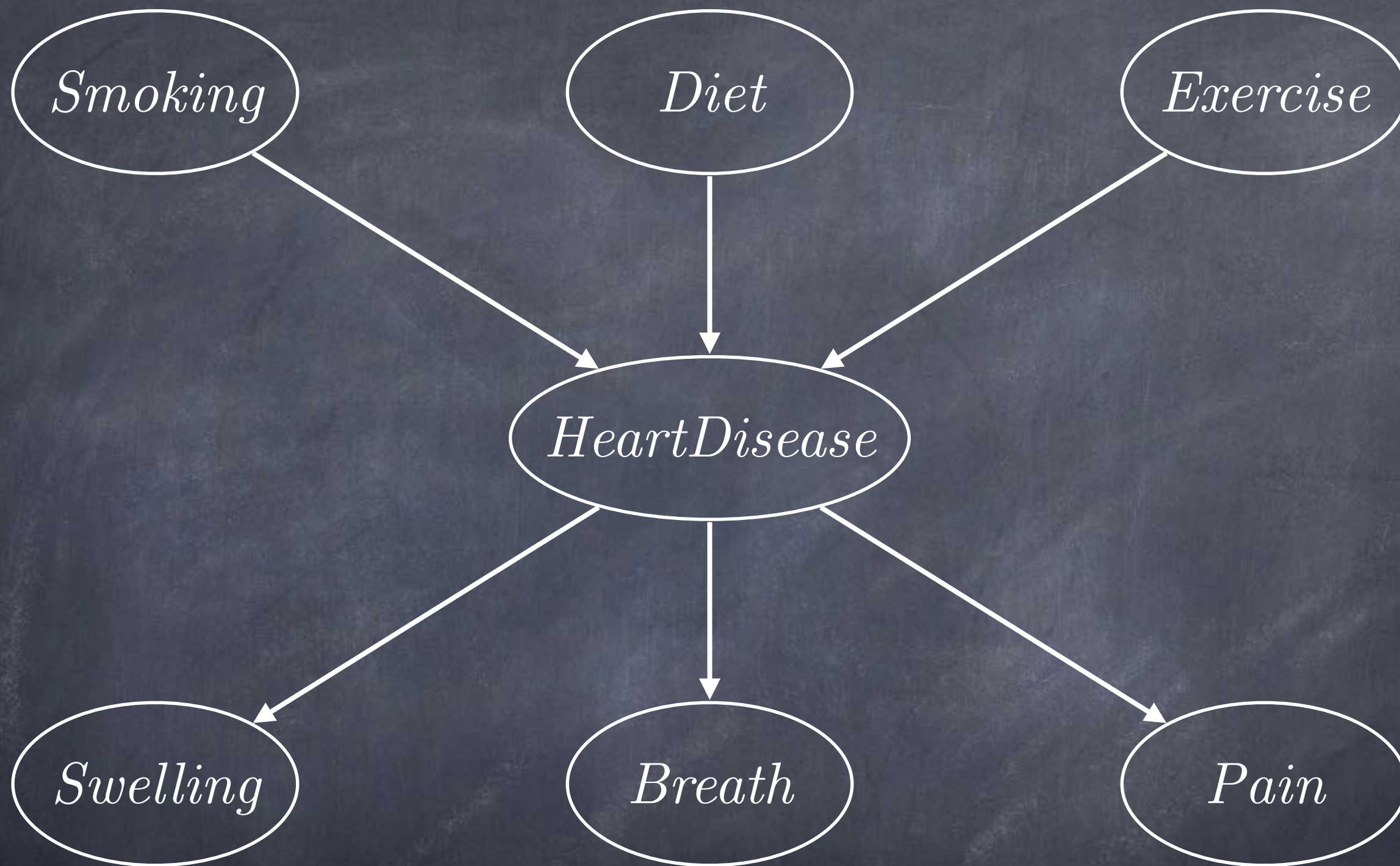
3 values/variable

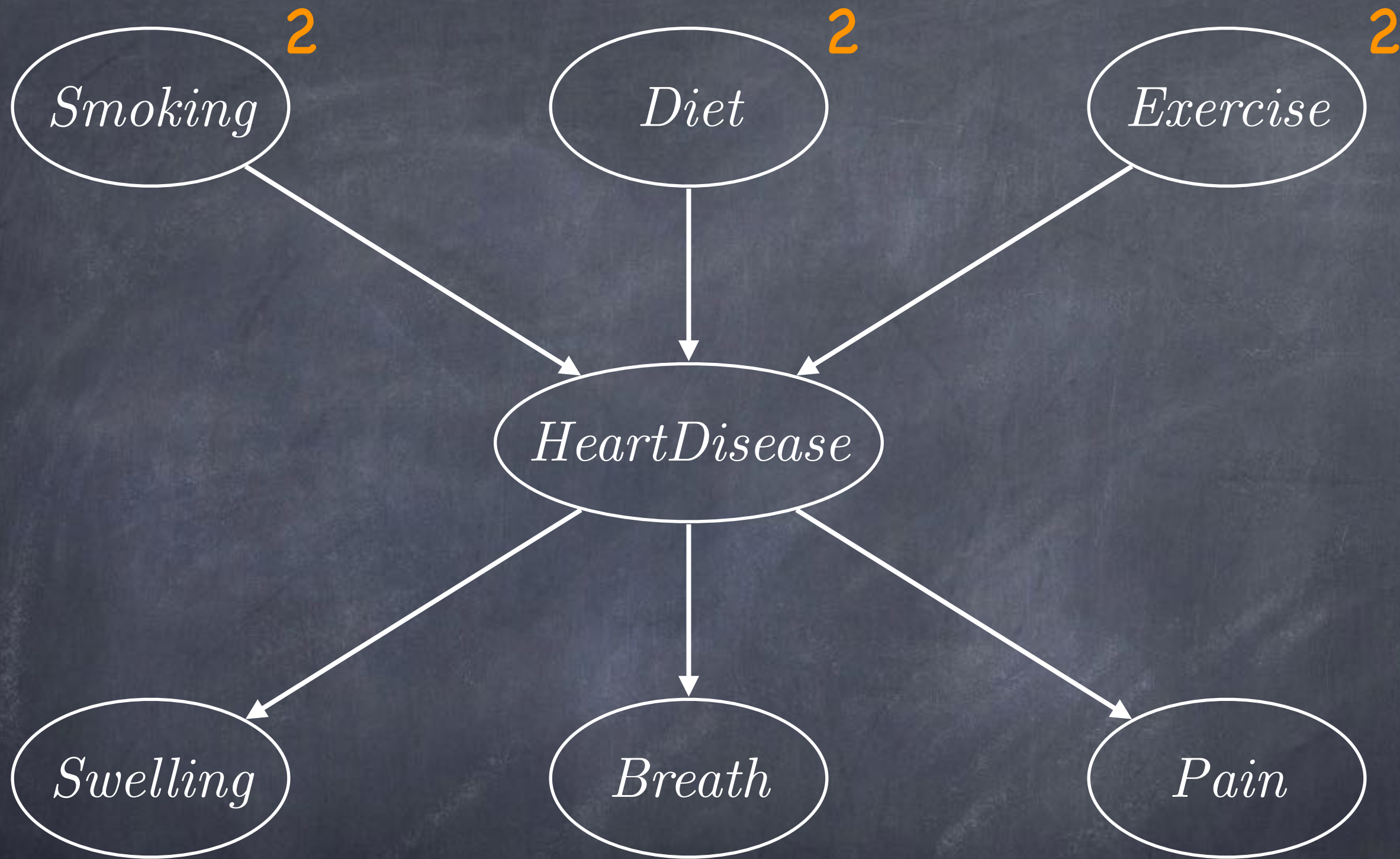


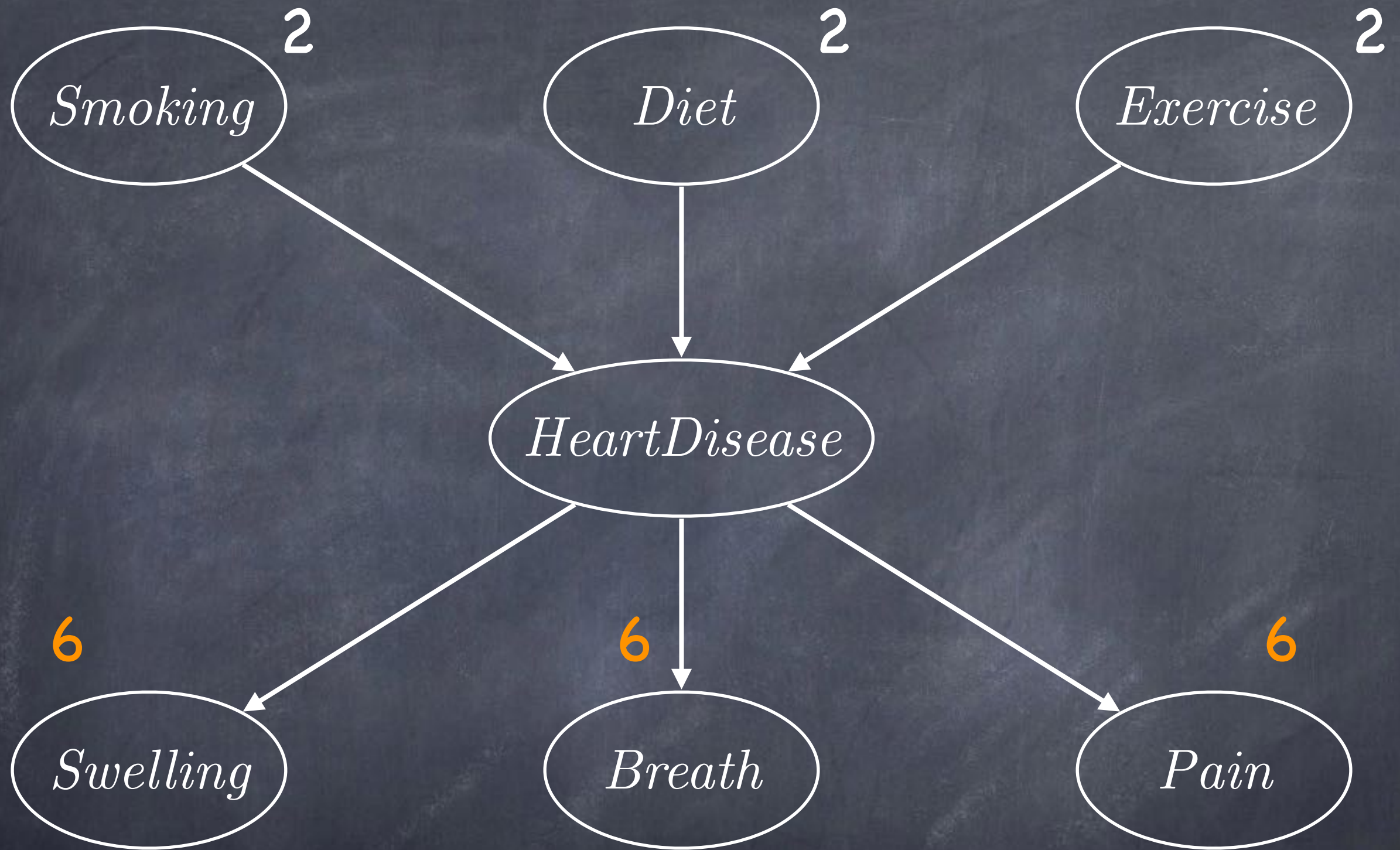
3 values/variable

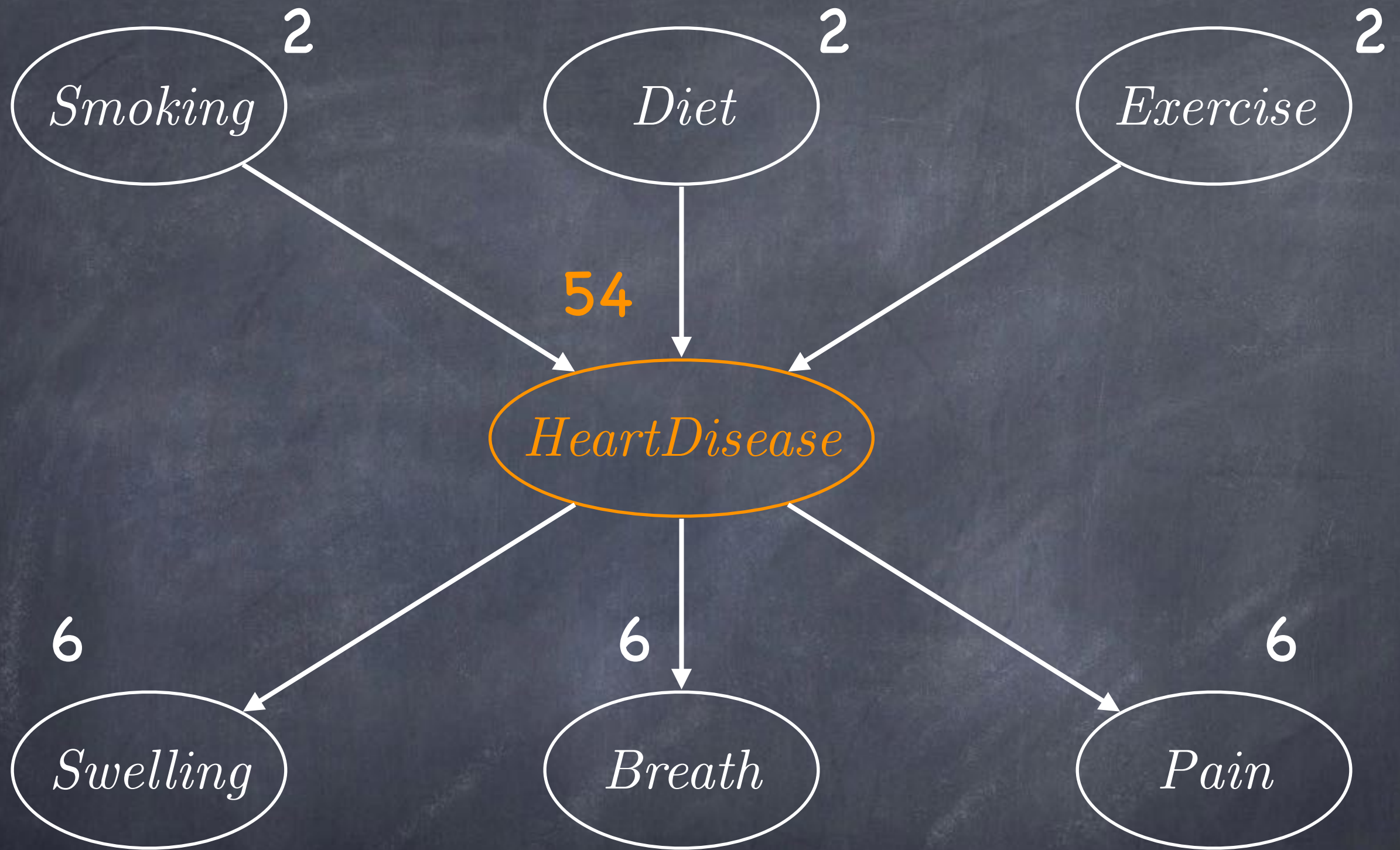


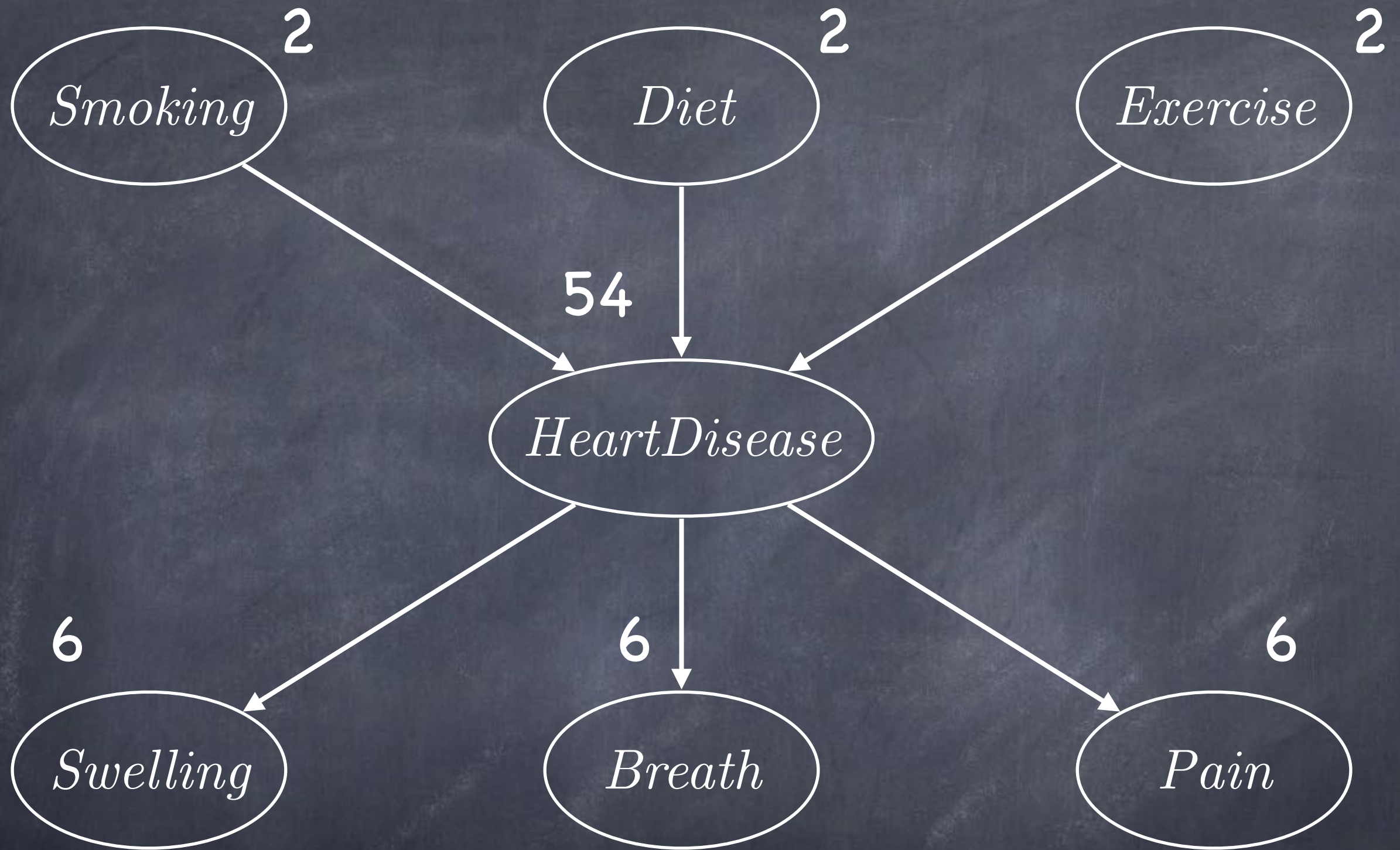
708 parameters!











78 parameters!

<i>Smoking</i>	<i>Diet</i>	<i>Exercise</i>	<i>Swelling</i>	<i>Breath</i>	<i>Pain</i>	<i>Disease</i>
<i>low</i>	<i>low</i>	<i>high</i>	<i>low</i>	<i>low</i>	<i>low</i>	<i>?</i>
<i>high</i>	<i>med</i>	<i>high</i>	<i>med</i>	<i>high</i>	<i>med</i>	<i>?</i>
<i>low</i>	<i>low</i>	<i>med</i>	<i>low</i>	<i>low</i>	<i>med</i>	<i>?</i>
<i>...</i>	<i>...</i>	<i>...</i>	<i>...</i>	<i>...</i>	<i>...</i>	<i>?</i>

Hidden!

Hidden (Latent) Variables

- Can dramatically reduce the number of parameters required to specify a Bayes net
 - Reduces amount of data required to learn the parameters
- Values of hidden variables not present in training data (observations)
 - “Complicates” the learning problem

EM

Expectation–Maximization

- Repeat
 - Expectation: “Pretend” we know the parameters and compute (or estimate) likelihood of data given model
 - Maximization: Recompute parameters using expected values as if they were observed values
- Until convergence

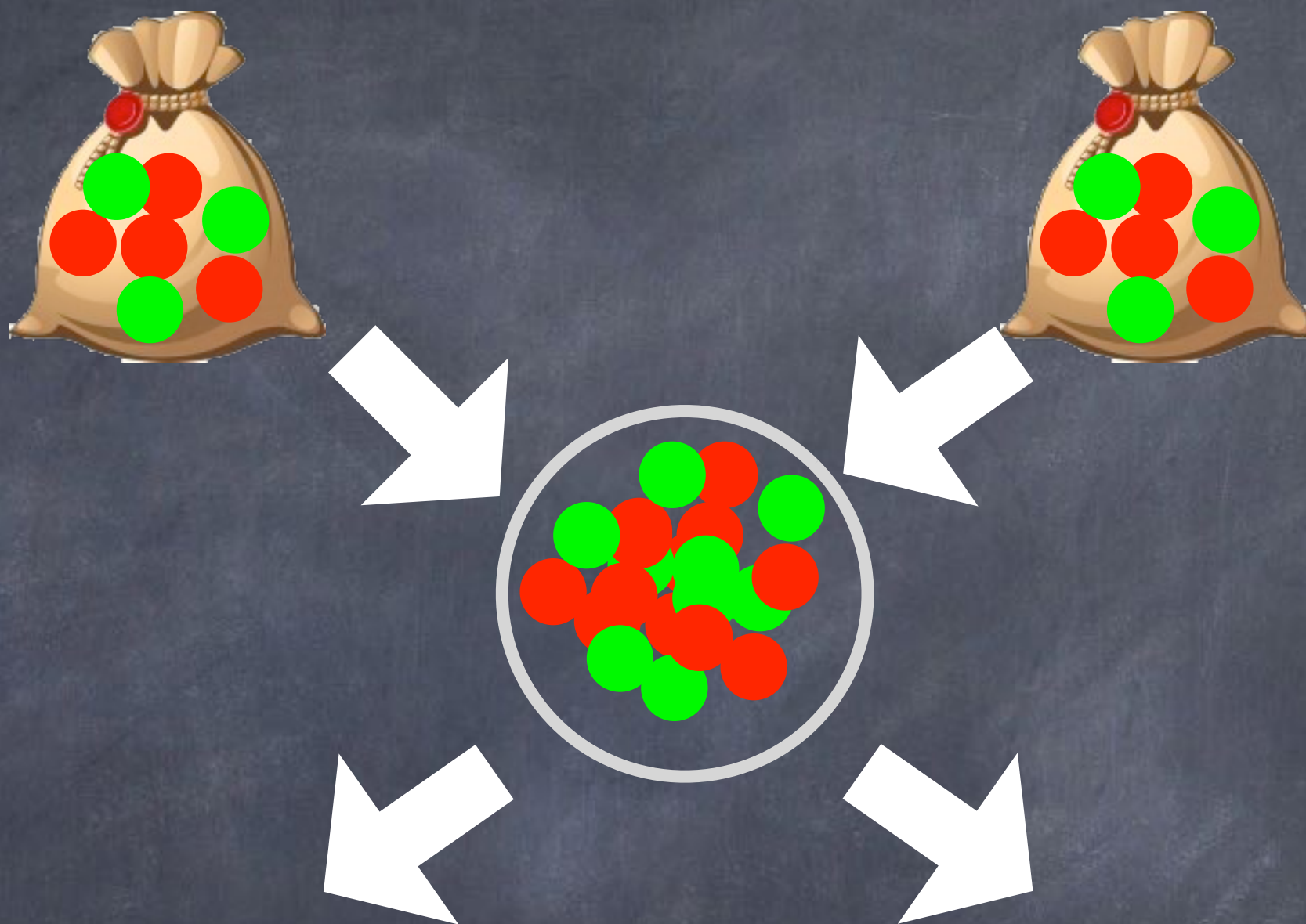


<i>Flavor</i>	<i>cherry</i>	<i>lime</i>
<i>Wrapper</i>	<i>red</i>	<i>green</i>
<i>Hole</i>	<i>true</i>	<i>false</i>



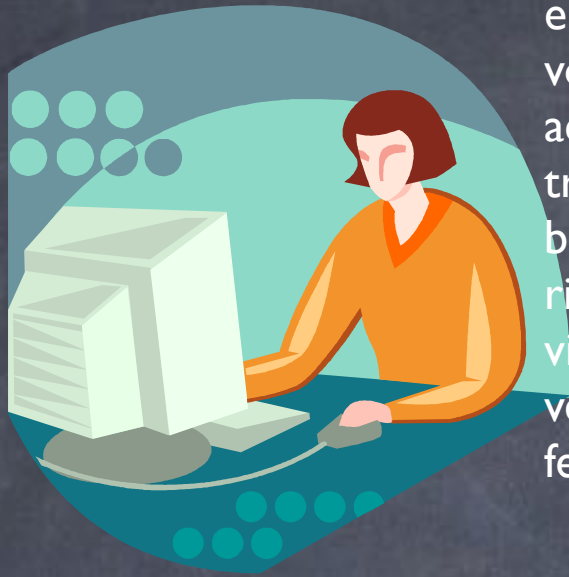
$\mathbf{P}(F, W, H)$

<i>Flavor</i>	<i>Wrapper</i>	<i>Hole</i>	$P(f, w, h)$
<i>cherry</i>	<i>red</i>	<i>t</i>	$p_{c,r,t}$
		<i>f</i>	$p_{c,r,f}$
	<i>green</i>	<i>t</i>	$p_{c,g,t}$
		<i>f</i>	$p_{c,g,f}$
<i>lime</i>	<i>red</i>	<i>t</i>	$p_{l,r,t}$
		<i>f</i>	$p_{l,r,f}$
	<i>green</i>	<i>t</i>	$p_{l,g,t}$
		<i>f</i>	$p_{l,g,f}$



$$P_1(F, W, H)$$

$$P_2(F, W, H)$$



Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam euismod euismod facilisis. Aliquam erat volutpat. Maecenas nisl ligula, dignissim et volutpat ac, pharetra blandit augue. Maecenas id ligula in leo tristique viverra. Curabitur lacinia nulla in nibh bibendum laoreet. Morbi a est mi, mattis imperdiet risus. Quisque quam felis, facilisis ac semper vel, viverra vitae nulla. Donec nisl lectus, faucibus vehicula tincidunt nec, ultrices nec eros. Proin non felis nec urna pellentesque tempor at sit amet est.



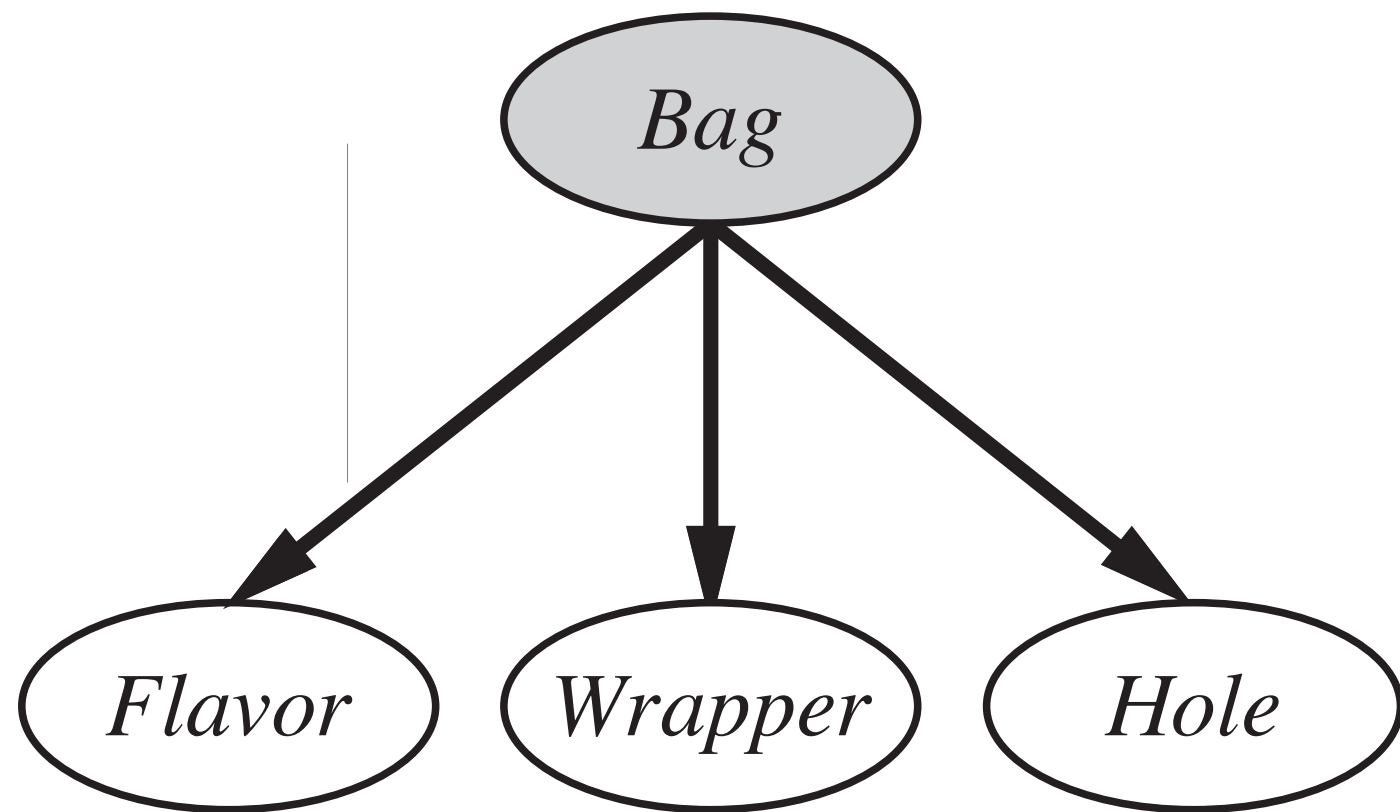
$$P_1(X_1, X_2, X_3)$$

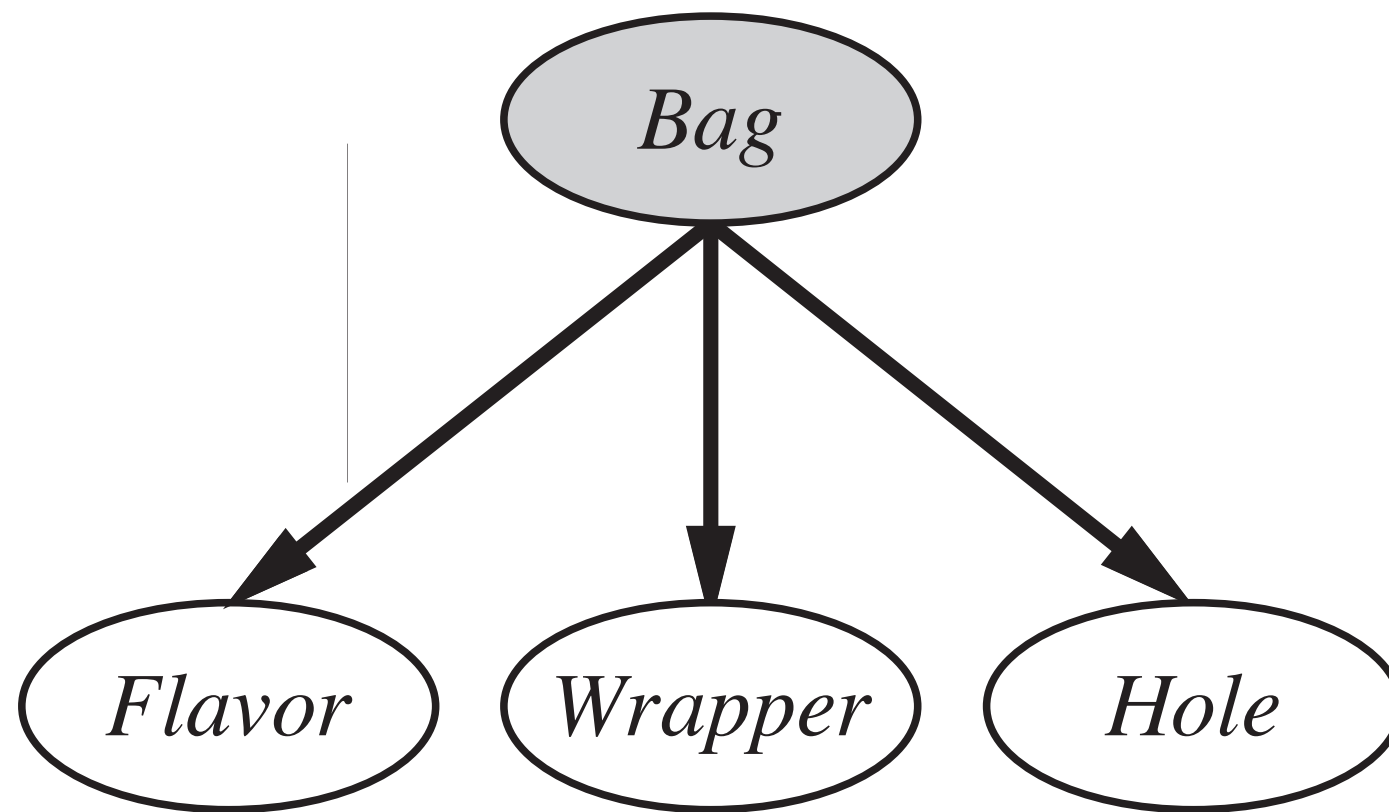
$$P_2(X_1, X_2, X_3)$$



$$\mathbf{P}(F, W, H)$$

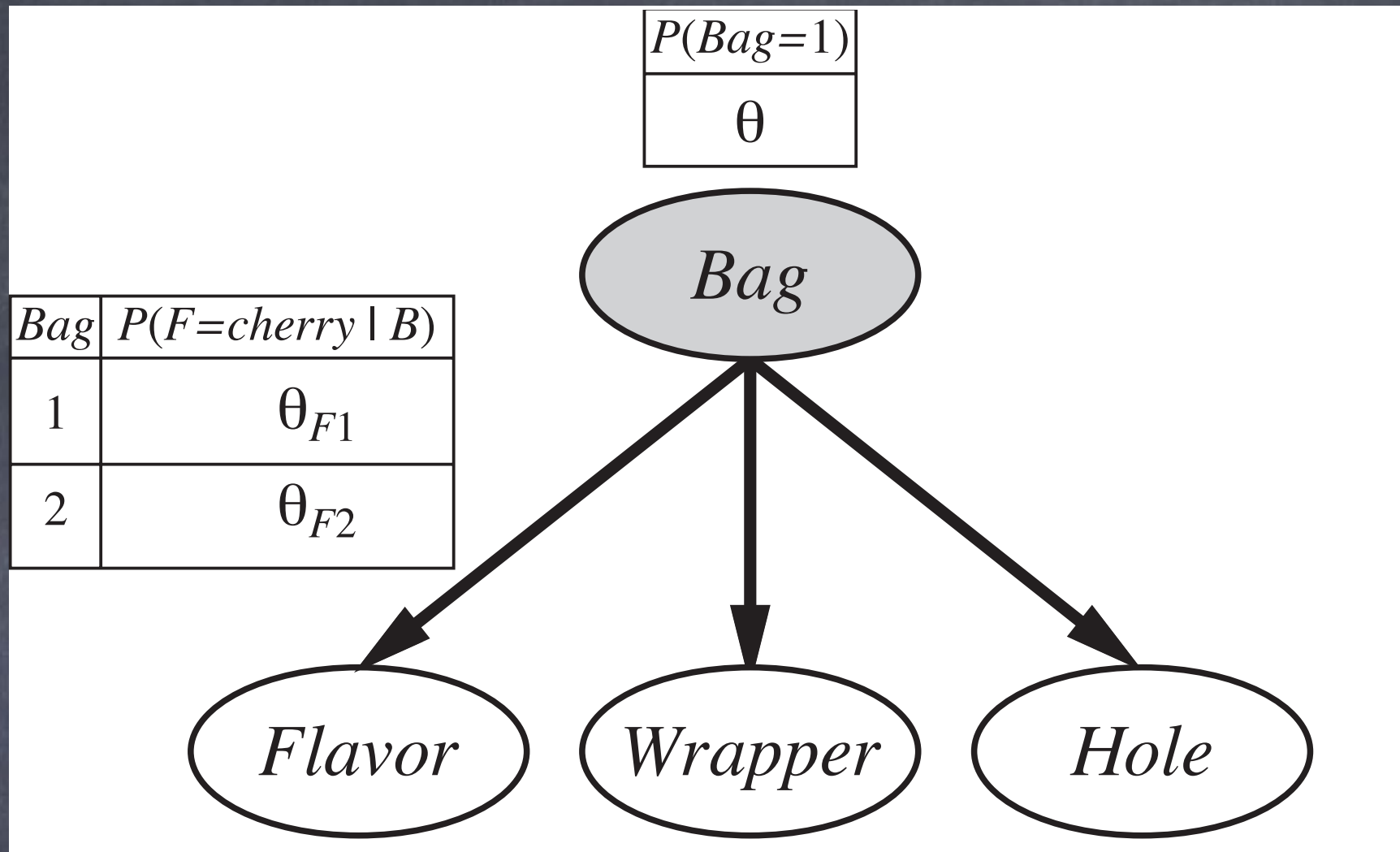
<i>Flavor</i>	<i>Wrapper</i>	<i>Hole</i>	$P(f, w, h)$
<i>cherry</i>	<i>red</i>	<i>t</i>	$p_{c,r,t}$
		<i>f</i>	$p_{c,r,f}$
	<i>green</i>	<i>t</i>	$p_{c,g,t}$
		<i>f</i>	$p_{c,g,f}$
<i>lime</i>	<i>red</i>	<i>t</i>	$p_{l,r,t}$
		<i>f</i>	$p_{l,r,f}$
	<i>green</i>	<i>t</i>	$p_{l,g,t}$
		<i>f</i>	$p_{l,g,f}$





<i>Bag</i>	$P(W=red B)$
1	Θ_{W1}
2	Θ_{W2}

<i>Bag</i>	$P(H=true B)$
1	Θ_{H1}
2	Θ_{H2}



Bag	$P(W=red B)$
1	Θ_{W1}
2	Θ_{W2}

Bag	$P(H=true B)$
1	Θ_{H1}
2	Θ_{H2}

7 parameters

<i>Flavor</i>	<i>Wrapper</i>	<i>Hole</i>	<i>Bag</i>
<i>cherry</i>	<i>red</i>	<i>true</i>	<i>?</i>
<i>cherry</i>	<i>red</i>	<i>true</i>	<i>?</i>
<i>lime</i>	<i>green</i>	<i>false</i>	<i>?</i>
<i>cherry</i>	<i>green</i>	<i>true</i>	<i>?</i>
<i>lime</i>	<i>green</i>	<i>true</i>	<i>?</i>
<i>cherry</i>	<i>red</i>	<i>false</i>	<i>?</i>
<i>lime</i>	<i>red</i>	<i>true</i>	<i>?</i>

Hidden!

$N=1000$	$W=red$		$W=green$	
	$H=yes$	$H=no$	$H=yes$	$H=no$
$F=cherry$	273	93	104	90
$F=lime$	79	100	94	167

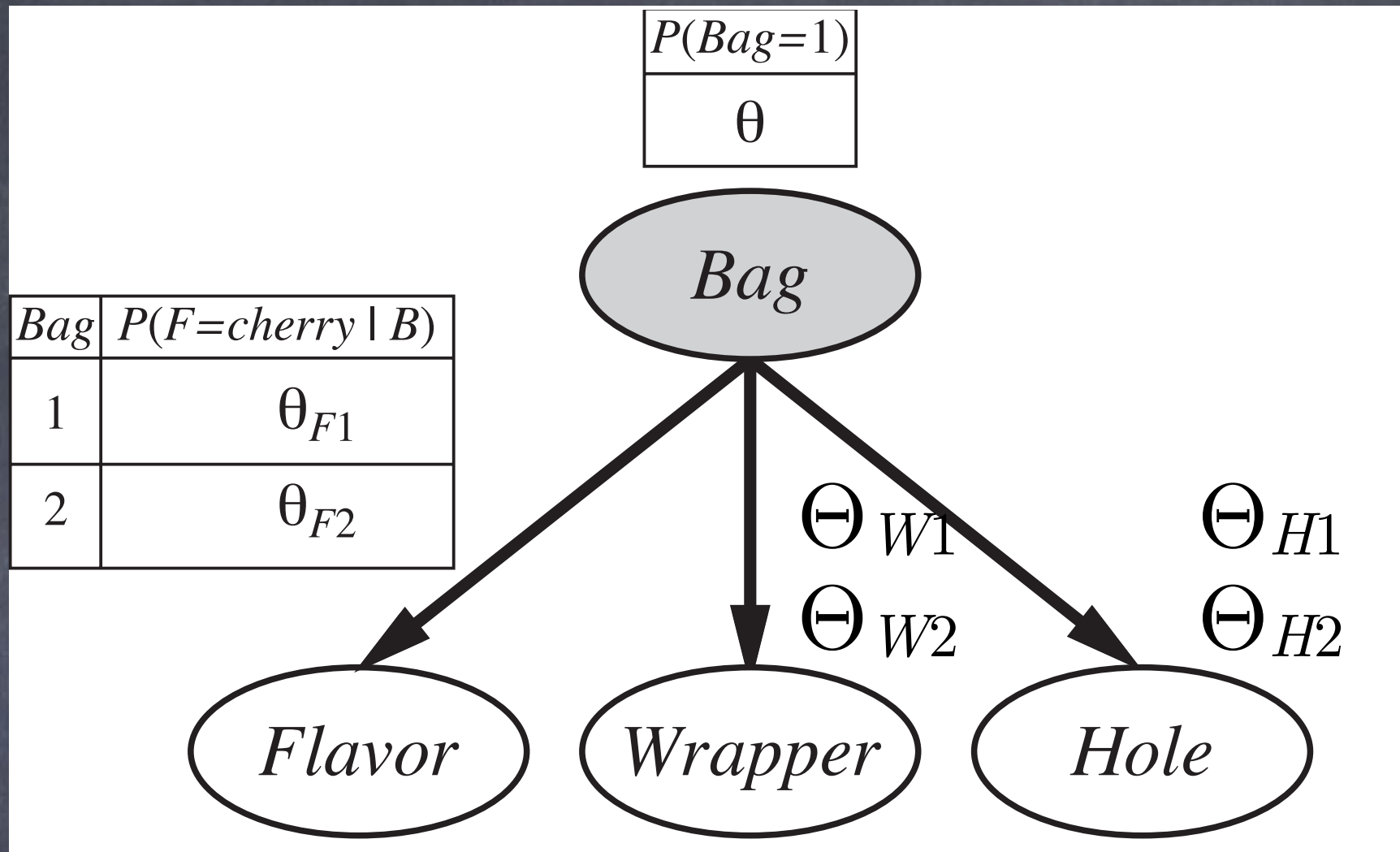
$N=1000$	$W=red$		$W=green$	
	$H=yes$	$H=no$	$H=yes$	$H=no$
$F=cherry$	273	93	104	90
$F=lime$	79	100	94	167

No values for *Bag*

EM

Expectation–Maximization

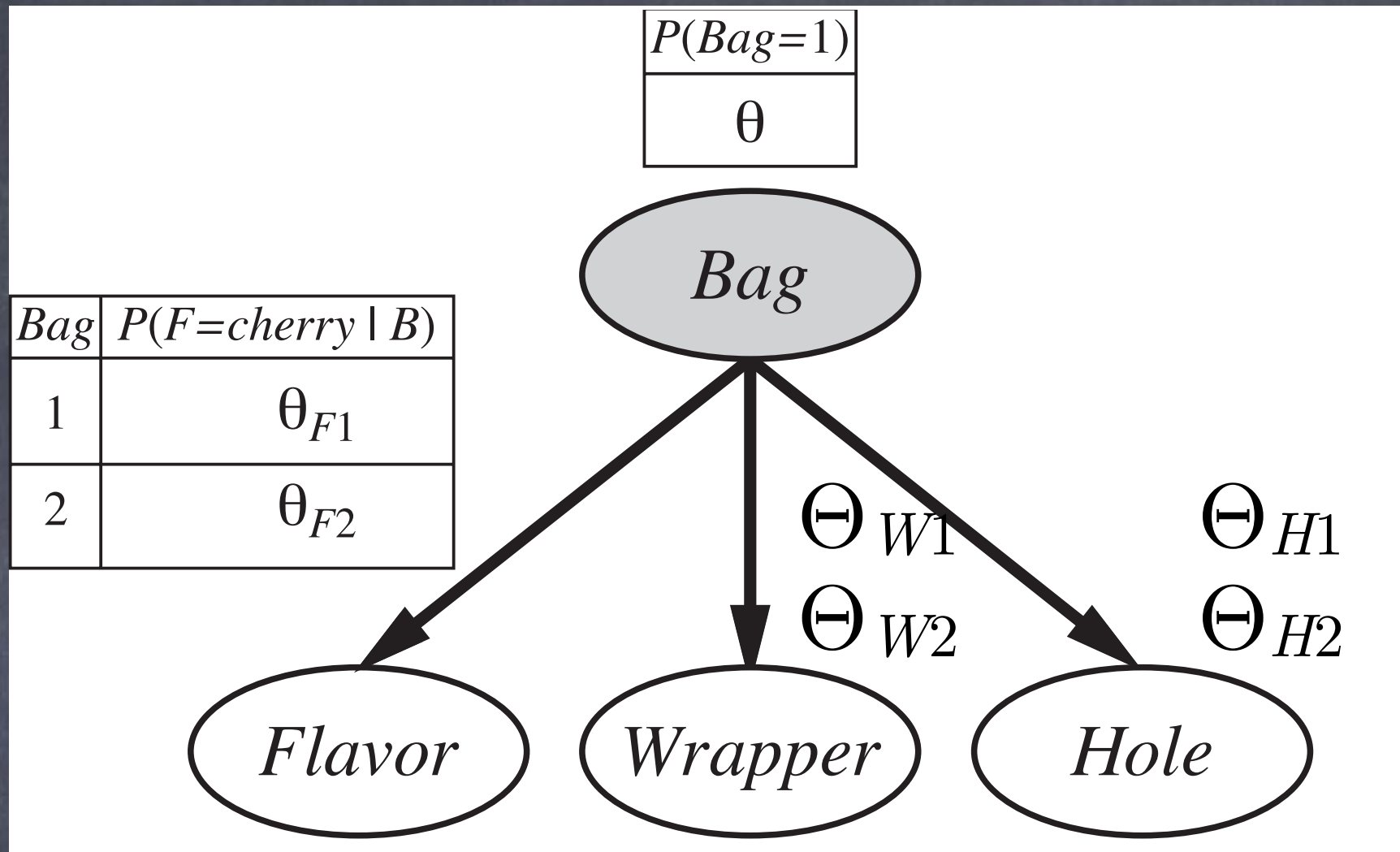
- Repeat
 - E: Use the current values of the parameters to compute the expected values of the hidden variables
 - M: Recompute the parameters to maximize the likelihood of the data given the values of (all) the variables
- Until convergence



$$\Theta = N(B=1) / N$$

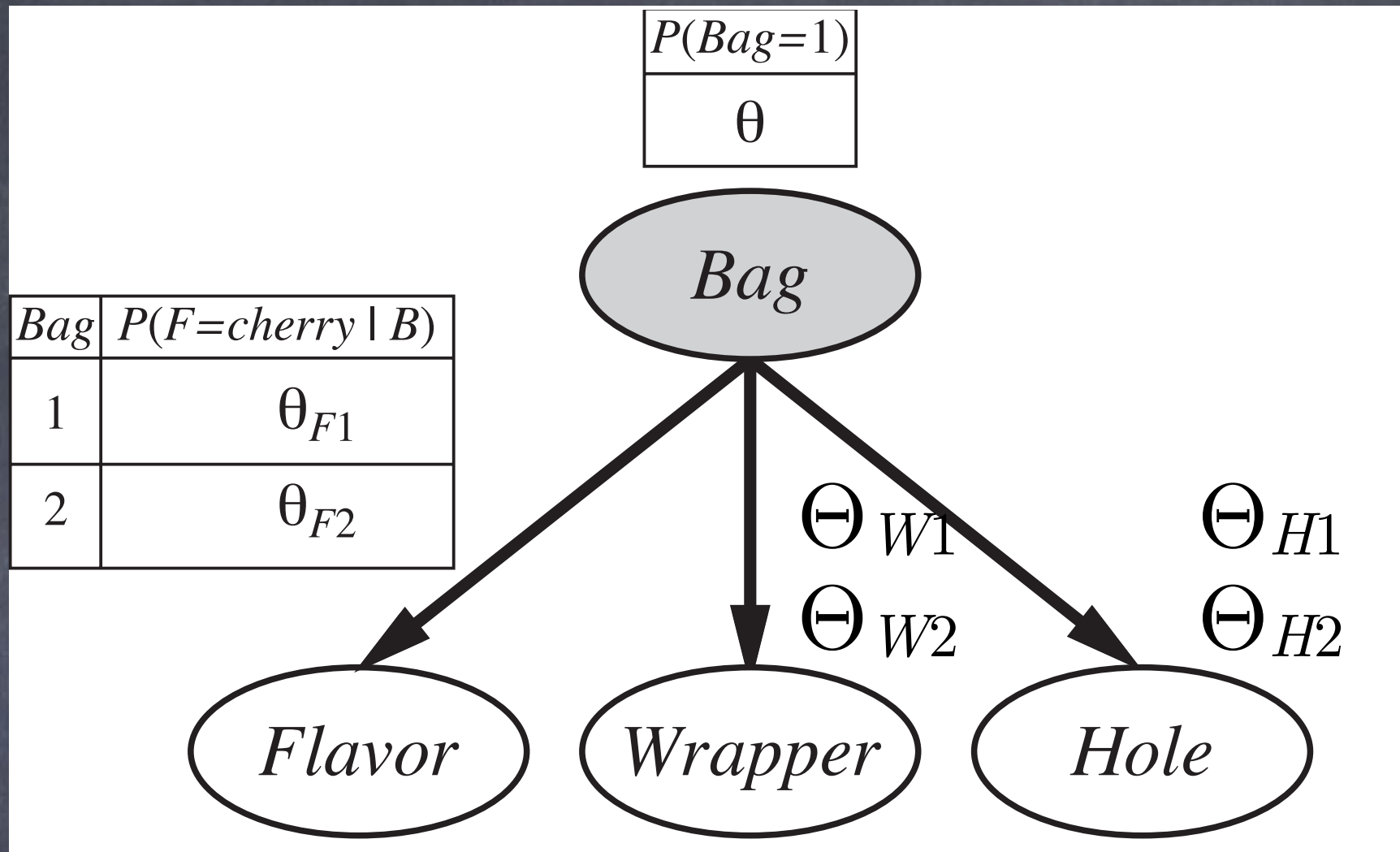
$$\Theta_{F1}, \Theta_{W1}, \Theta_{H1}$$

$$\Theta_{F2}, \Theta_{W2}, \Theta_{H2}$$

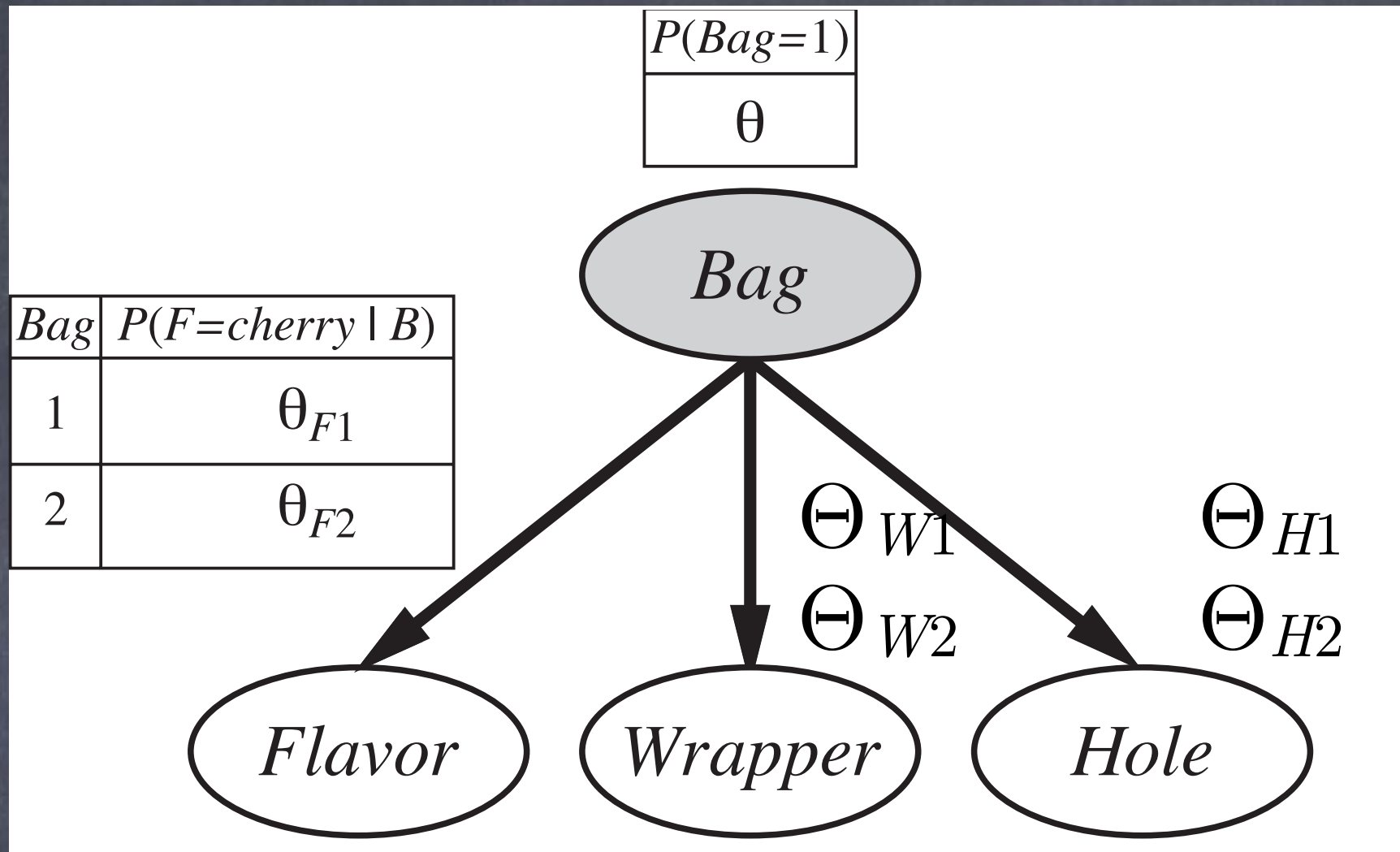


We don't know $N(B=1)$

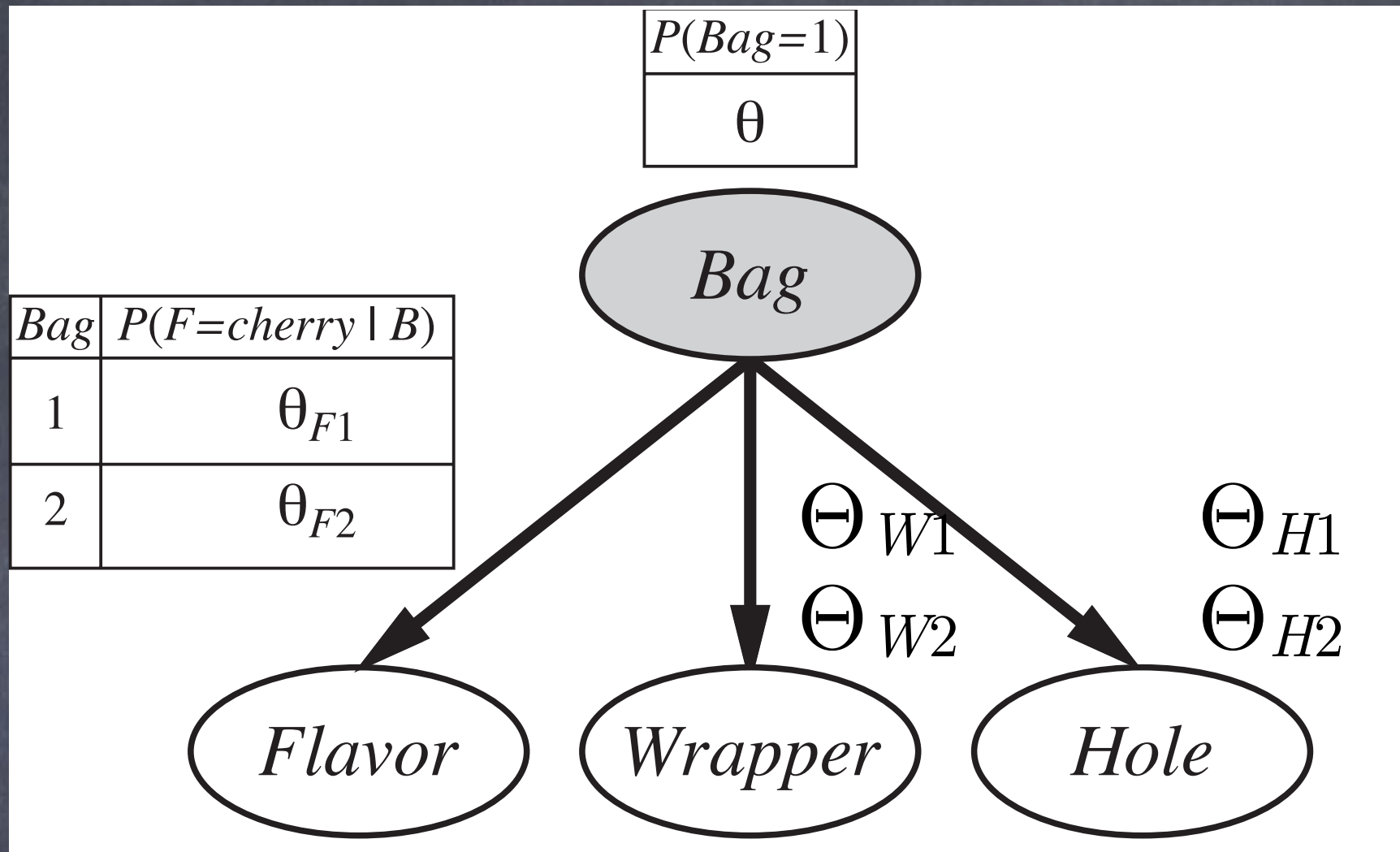
Estimate: $\hat{N}(B=1)$



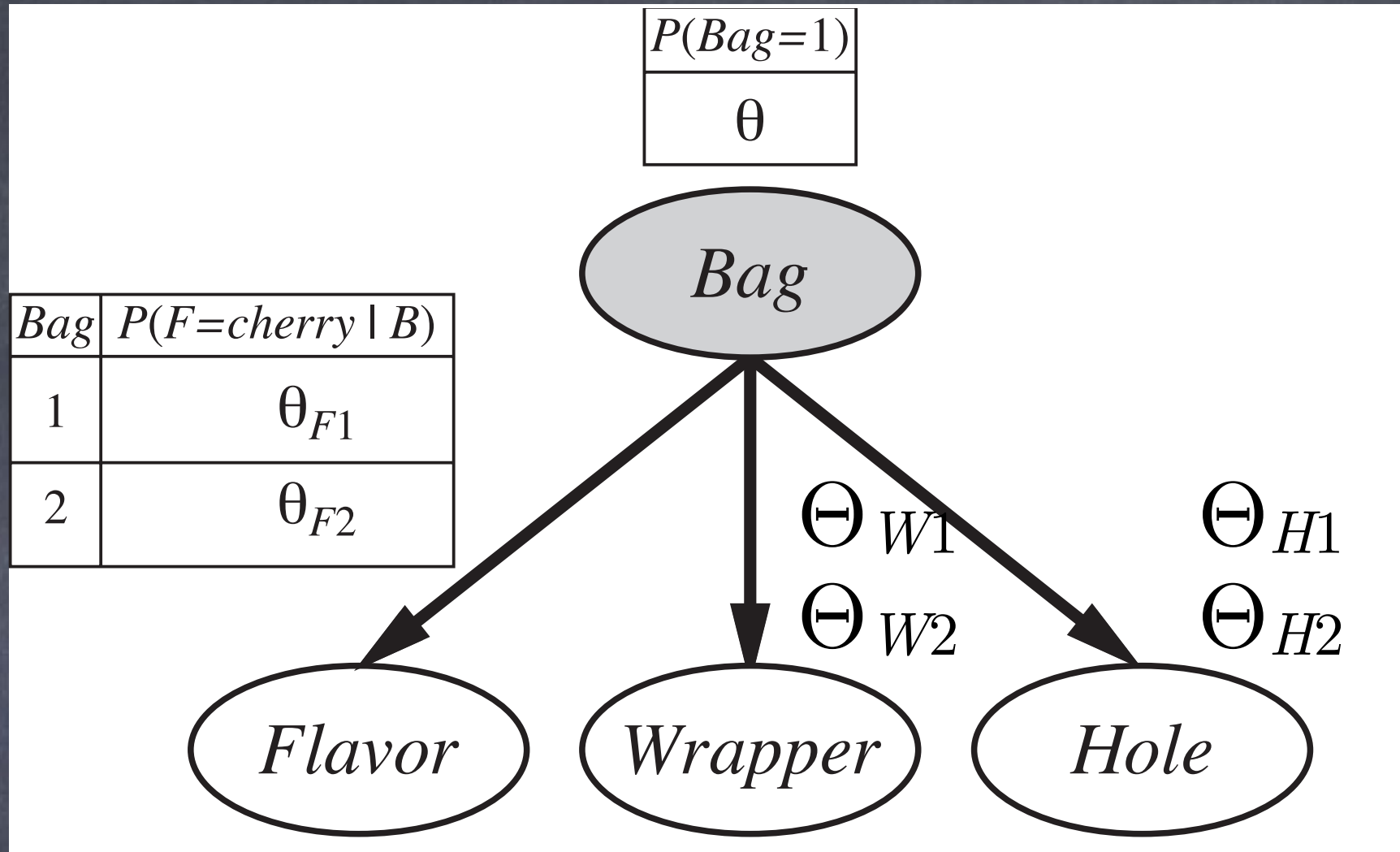
Flavor=cherry, Wrapper=red, Hole=true



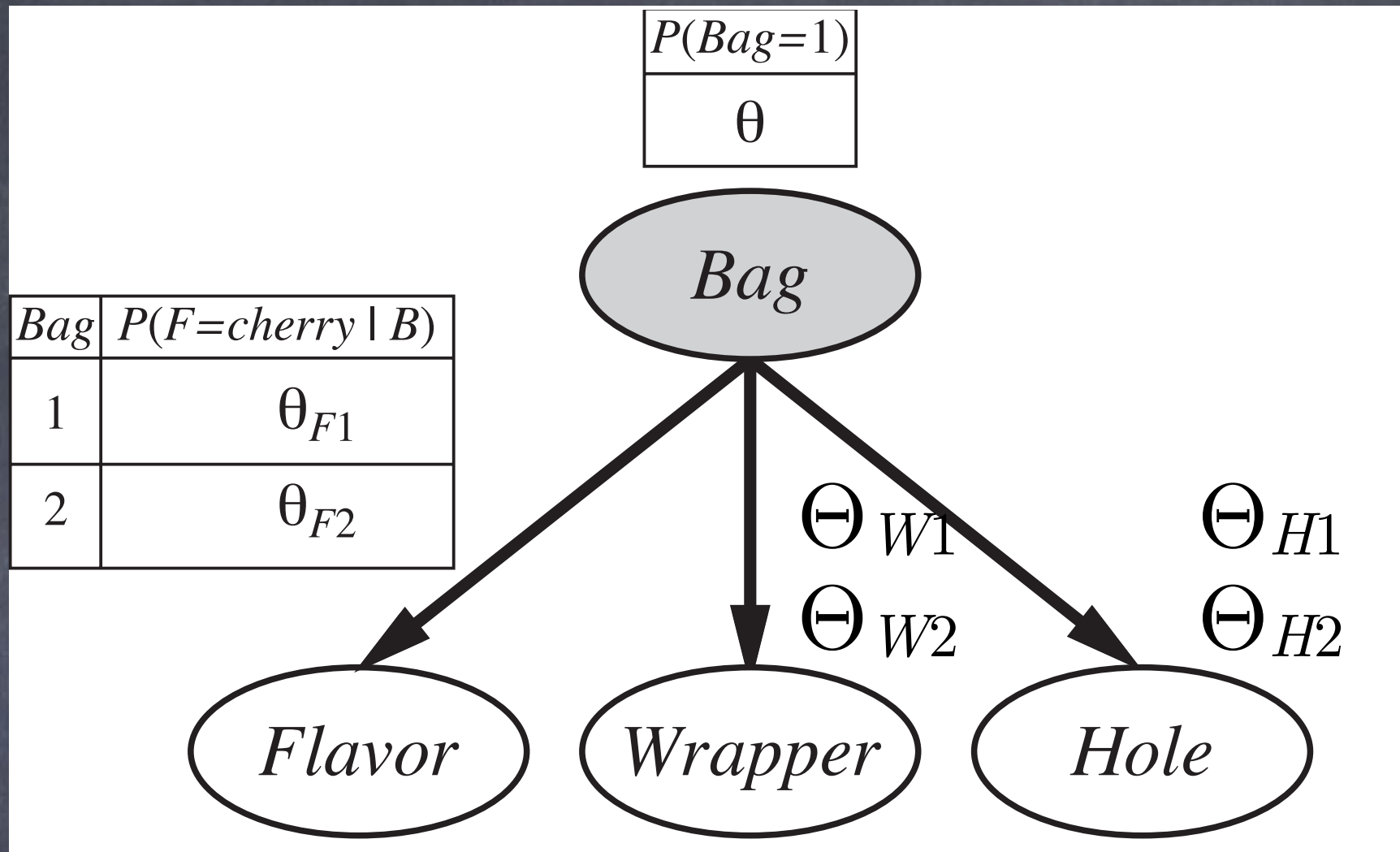
$P(Bag=1 \mid Flavor=cherry, Wrapper=red, Hole=true)?$



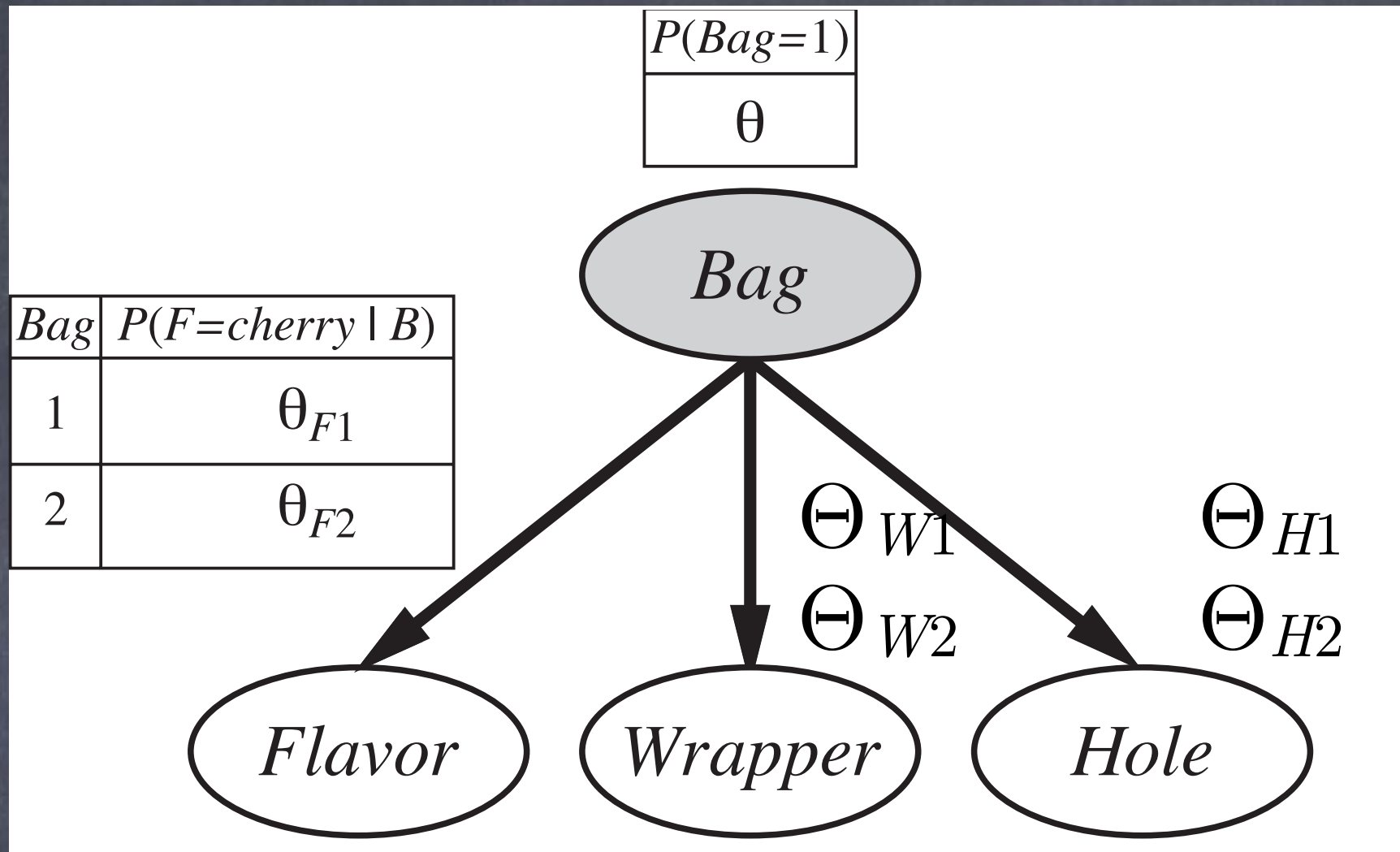
$$\begin{aligned}
 &P(B = 1 \mid F = c, W = r, H = t) \\
 &= \alpha P(B = 1, F = c, W = r, H = t) \\
 &= \alpha P(B = 1) P(F = c \mid B = 1) \\
 &\quad P(W = r \mid B = 1) P(H = t \mid B = 1) \\
 &= \alpha \Theta \Theta_{F1} \Theta_{W1} \Theta_{H1}
 \end{aligned}$$



$$\hat{N}(B = 1) = \sum_{j=1}^N \alpha P(B = 1 \mid F = f_j, W = w_j, H = h_j)$$



$$\Theta \leftarrow \hat{N}(B=1)/N$$



$$\Theta \leftarrow \hat{N}(B=1)/N$$

$$\Theta_{F1}, \Theta_{W1}, \Theta_{H1}$$

$$\Theta_{F2}, \Theta_{W2}, \Theta_{H2}$$

$$\Theta^{(0)} = 0.6$$

$$\Theta_{F1}^{(0)} = \Theta_{W1}^{(0)} = \Theta_{H1}^{(0)} = 0.6$$

$$\Theta_{F2}^{(0)} = \Theta_{W2}^{(0)} = \Theta_{H2}^{(0)} = 0.4$$

$$L(\mathbf{d}|h) \approx -2044$$

Actual:

$$\Theta^{(1)} = 0.6124$$

$$\Theta_{F1}^{(1)} = 0.6684$$

$$\Theta_{W1}^{(1)} = 0.6483$$

$$\Theta_{H1}^{(1)} = 0.6558$$

$$\Theta_{F2}^{(1)} = 0.3887$$

$$\Theta_{W2}^{(1)} = 0.3817$$

$$\Theta_{H2}^{(1)} = 0.3827$$

$$\Theta = 0.5$$

$$\Theta_{F1} = \Theta_{W1} = \Theta_{H1} = 0.8$$

$$\Theta_{F2} = \Theta_{W2} = \Theta_{H2} = 0.3$$

$$L(\mathbf{d}|h) \approx -2021$$

EM

Expectation–Maximization

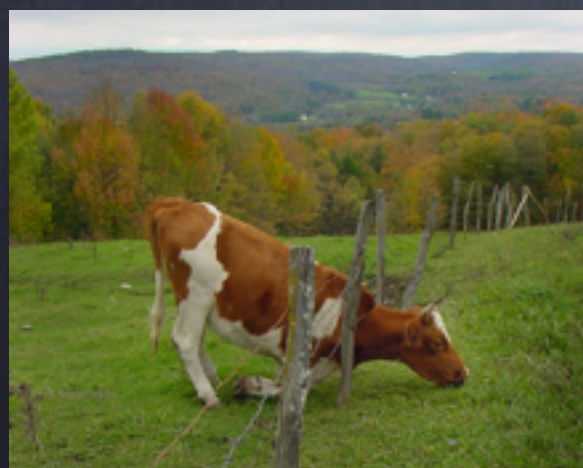
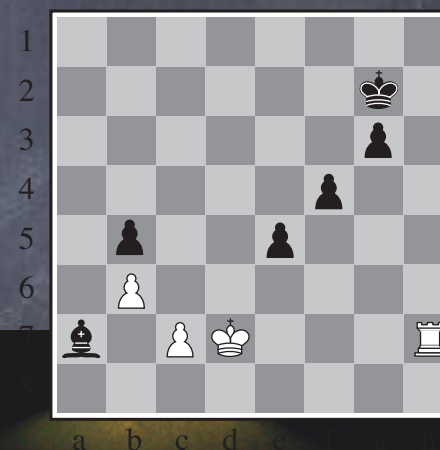
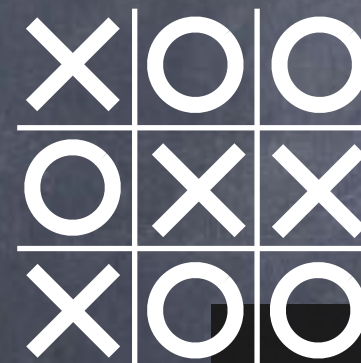
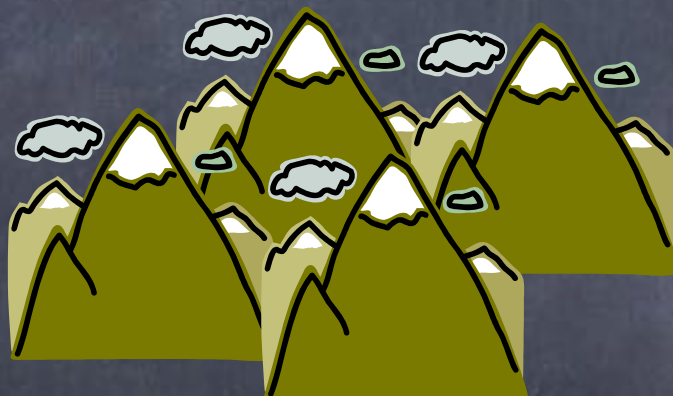
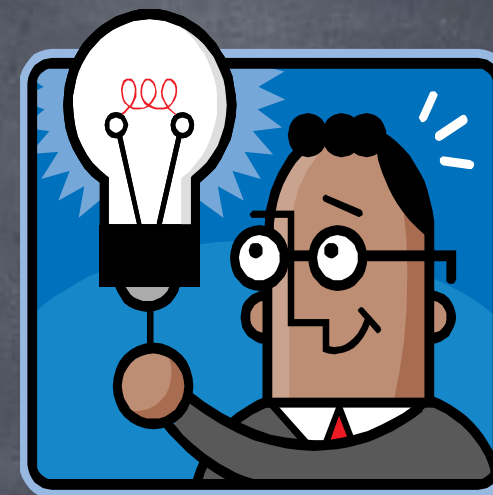
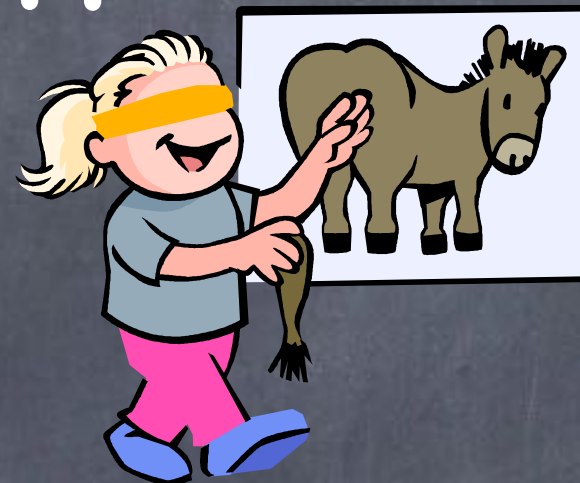
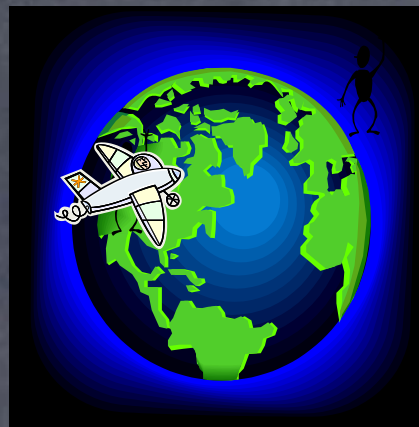
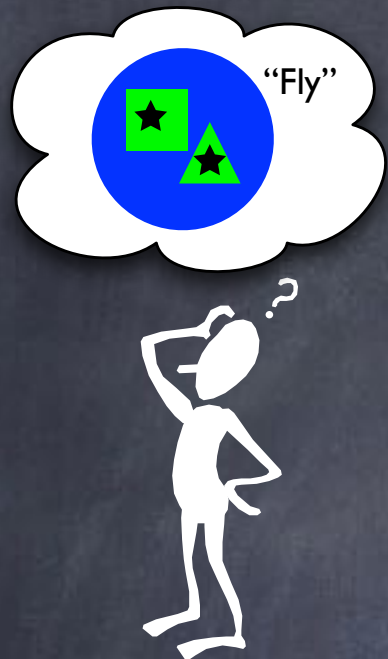
- Repeat
 - E: Pretend we know the values of the parameters and infer the expected values of the hidden variables
 - M: Update the parameters to maximize the likelihood of the data given the values of (all) the variables
- Until convergence

Learning Probabilistic Models

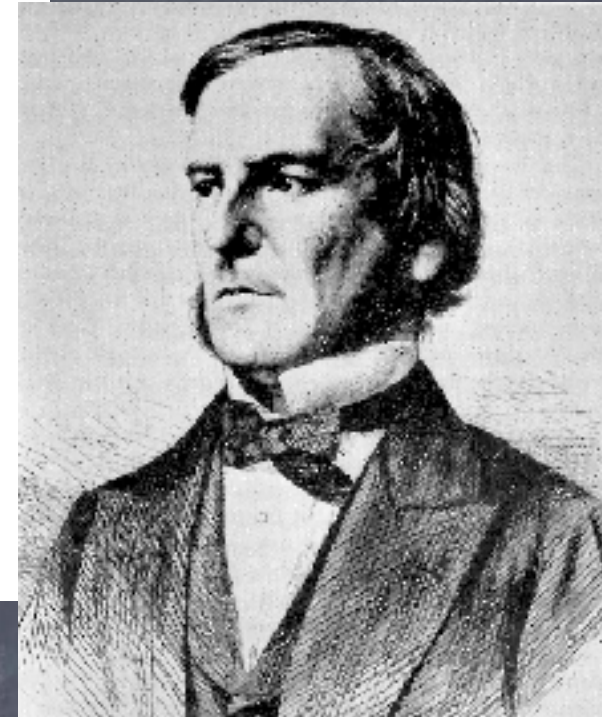
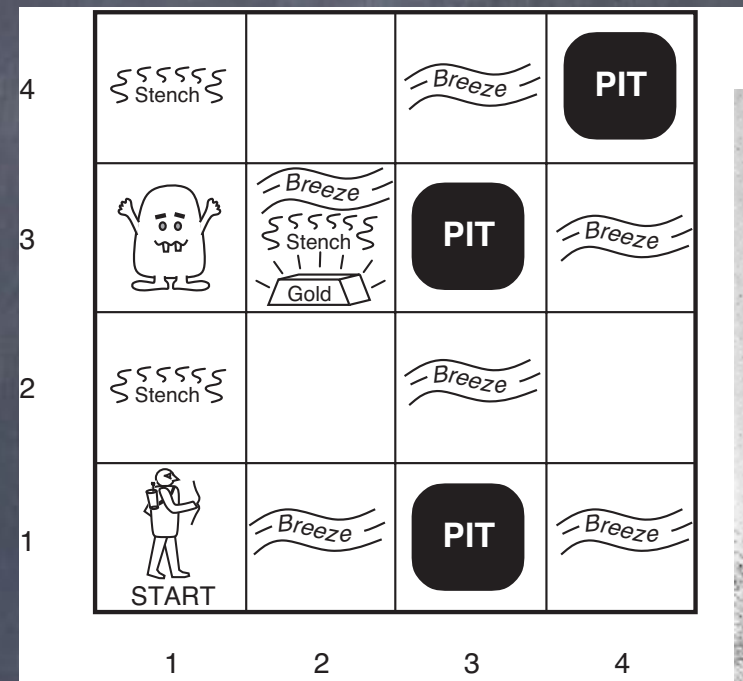
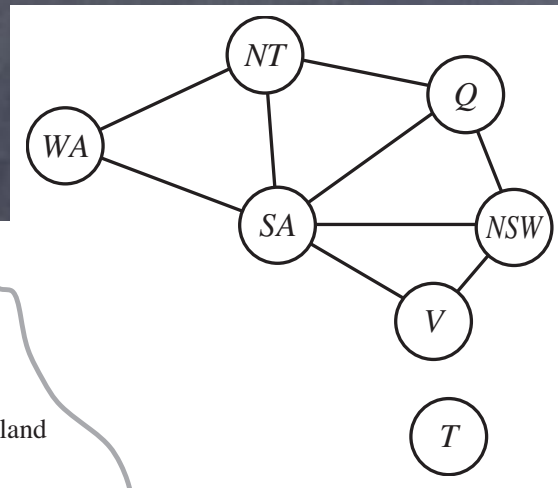
- Learning parameters of probability distributions
- Maximum Likelihood Hypothesis: maximizes the likelihood of the data
- Learning parameters of Bayes Nets
- Naive Bayes classifiers
- EM: Learning with incomplete data

Intro to AI

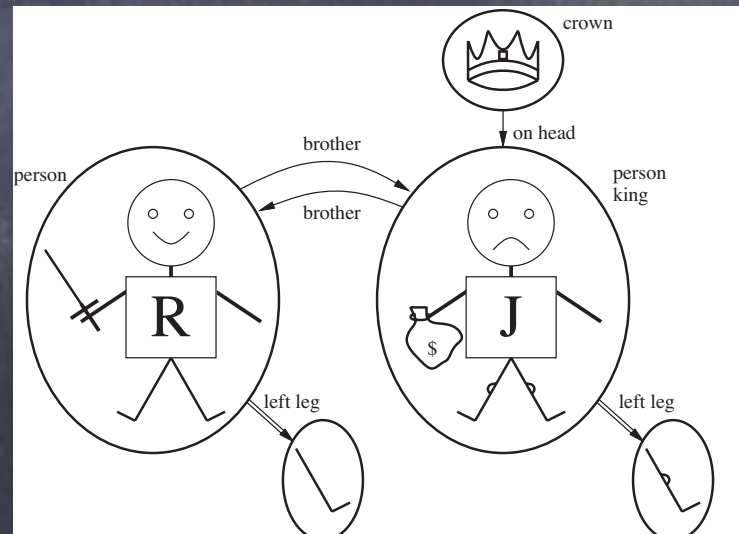
Search



Representation



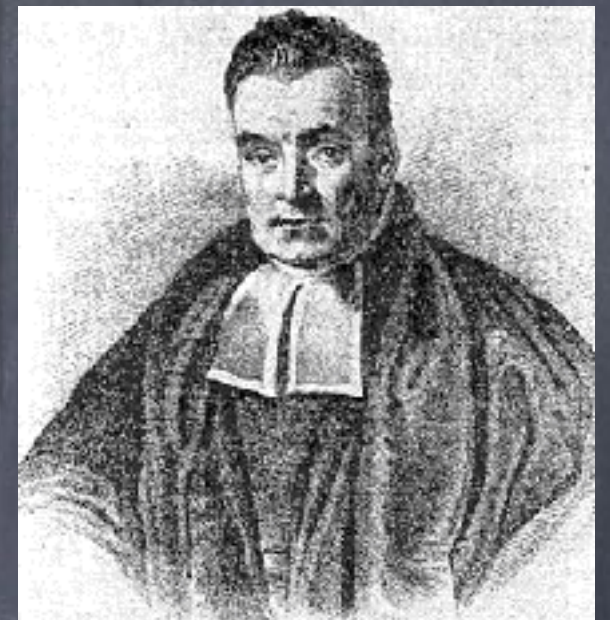
Hungry \vee *Cranky*
 \neg *Hungry*
Cranky



if $\alpha \vdash \beta$ then $\alpha \models \beta$
 if $\alpha \models \beta$ then $\alpha \vdash \beta$

Uncertainty

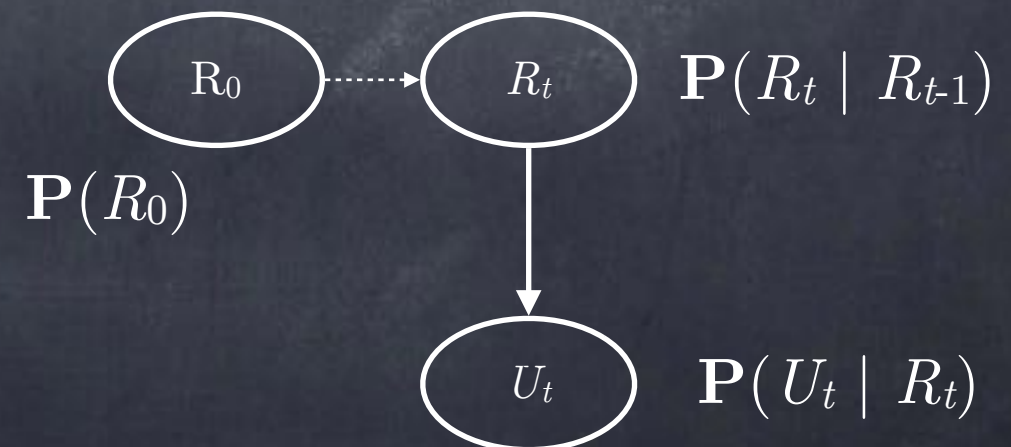
1,4 ?	2,4 ?	3,4 ?	4,4 ?
1,3 ?	2,3 ?	3,3 ?	4,3 ?
1,2 ?	2,2 ?	3,2 ?	4,2 ?
1,1 ?	2,1 ?	3,1 ?	4,1 ?



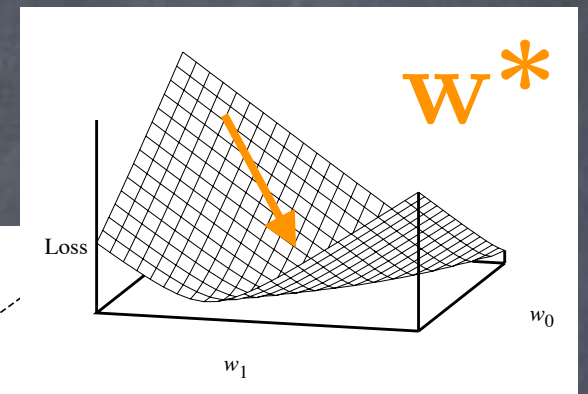
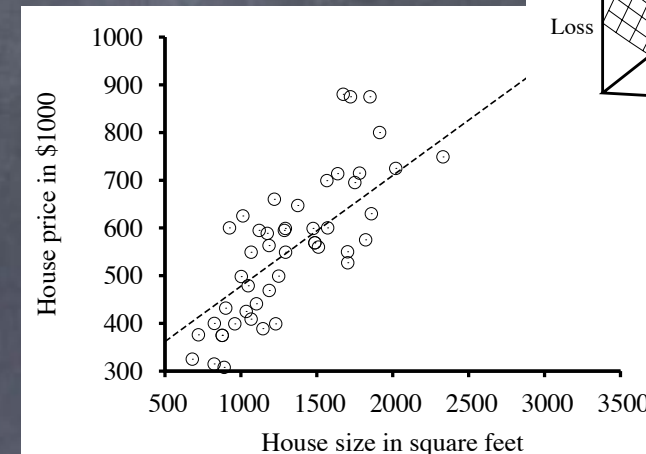
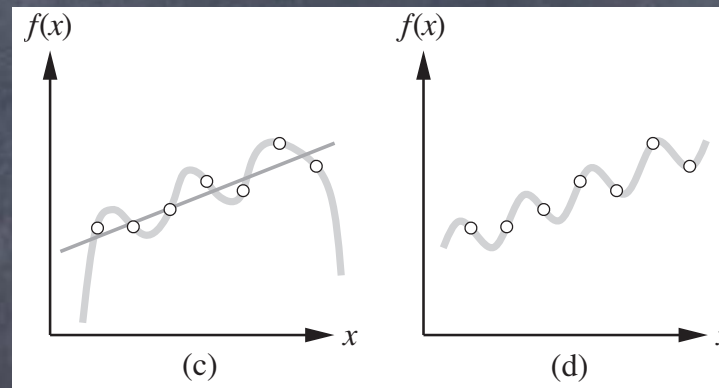
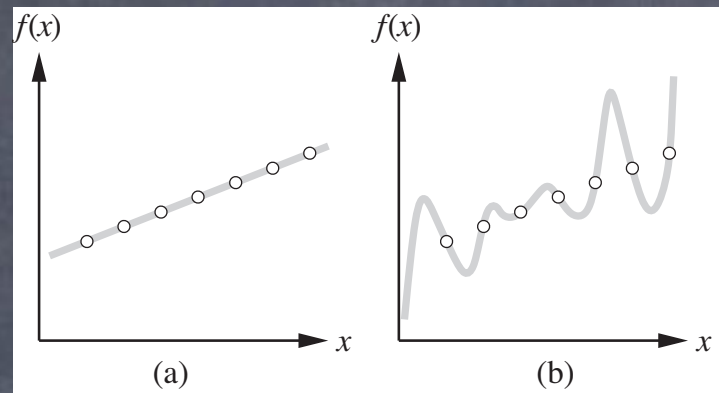
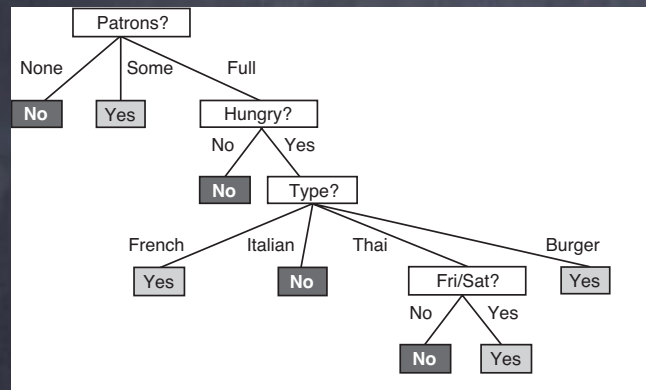
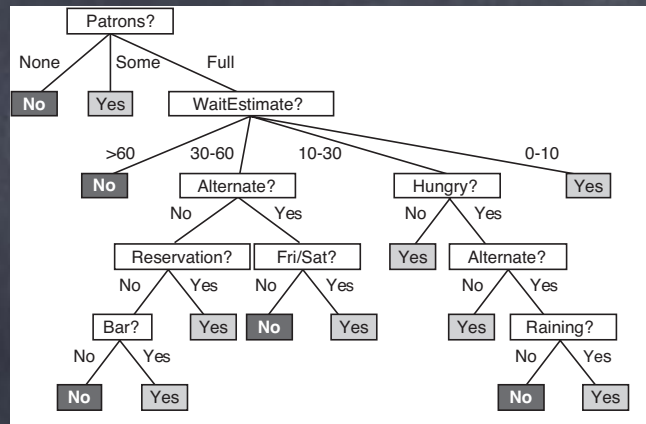
$$\mathbf{P}(X \mid \mathbf{e}) = \alpha \mathbf{P}(X, \mathbf{e}) = \alpha \sum_{\mathbf{y}} \mathbf{P}(X, \mathbf{e}, \mathbf{y})$$



$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i \mid \text{parents}(X_i))$$



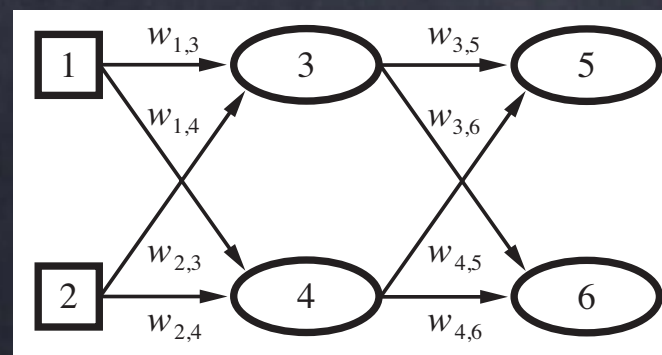
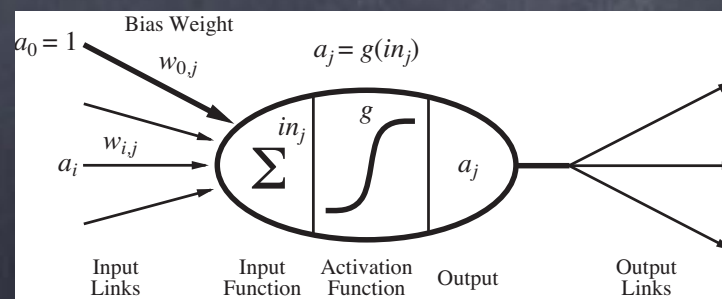
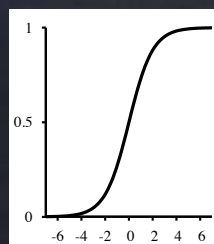
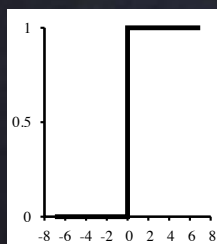
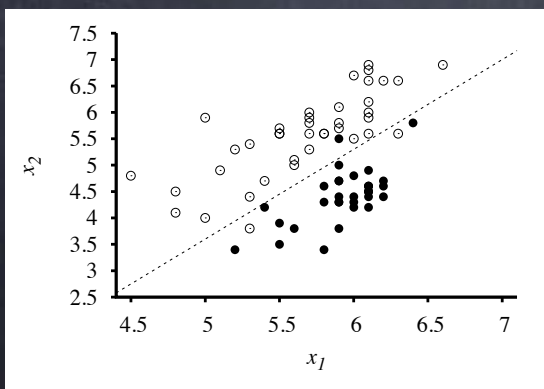
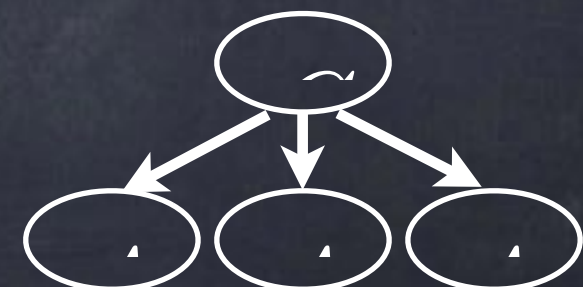
Learning



$$\underset{\Theta}{\operatorname{argmax}} P(\mathbf{d} \mid h_{\Theta})$$



$$\Theta = \frac{c}{c+l} = \frac{c}{N}$$



Intro to AI

- Searching HUGE spaces for solutions to problems
- Representation matters!
 - Search through space of proofs
- Uncertainty: Can be done but..
 - Independence assumptions are the key
- Learning: Search for parameters

For Next Time

Unit 4 Exam

No class Tue 2 May

Final Exam: Thu 11 May 1600

Douglass Ballroom — BRING ID