

estimate CPTs for B.N.
— all variables observed

FUDAN SUMMER SCHOOL

INTRODUCTION TO AI

Probabilistic Reasoning and Decision Making
Day 9 – Expectation Maximization (EM) Algorithm

MLE:

$$P(X=x_i) = \frac{\text{count}(X=x_i)}{T}$$

$$P(X=x_i | P_{ai}=r_u) = \frac{\text{count}(X=x_i, P_{ai}=r_u)}{\text{count}(P_{ai}=r_u)} = \frac{\sum_t I(x_i, x_i^{(t)}) I(P_{ai}, r_u^{(t)})}{\sum_t I(P_{ai}, r_u^{(t)})}$$

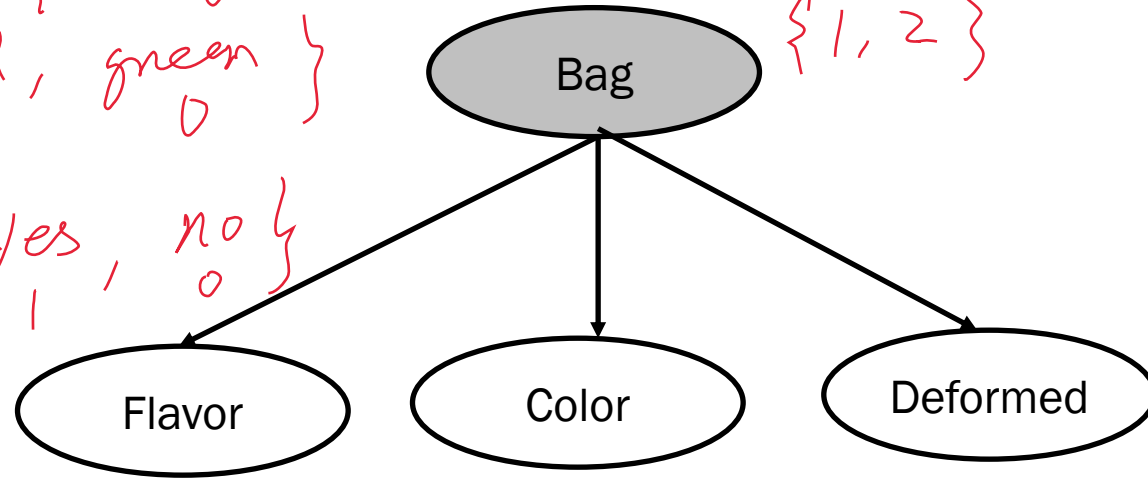
Learning from Partially Observed Data, Example

flavor: { jalapeno, cherry }

color: { red, green }

deformed: { yes, no }

Hidden



unobserved
{1, 2}

Goal: estimate
all the CPTs

Sample	flavor	color	Deformed
1	0	0	1
2	1	0	1
3	1	1	1
⋮	—	—	—

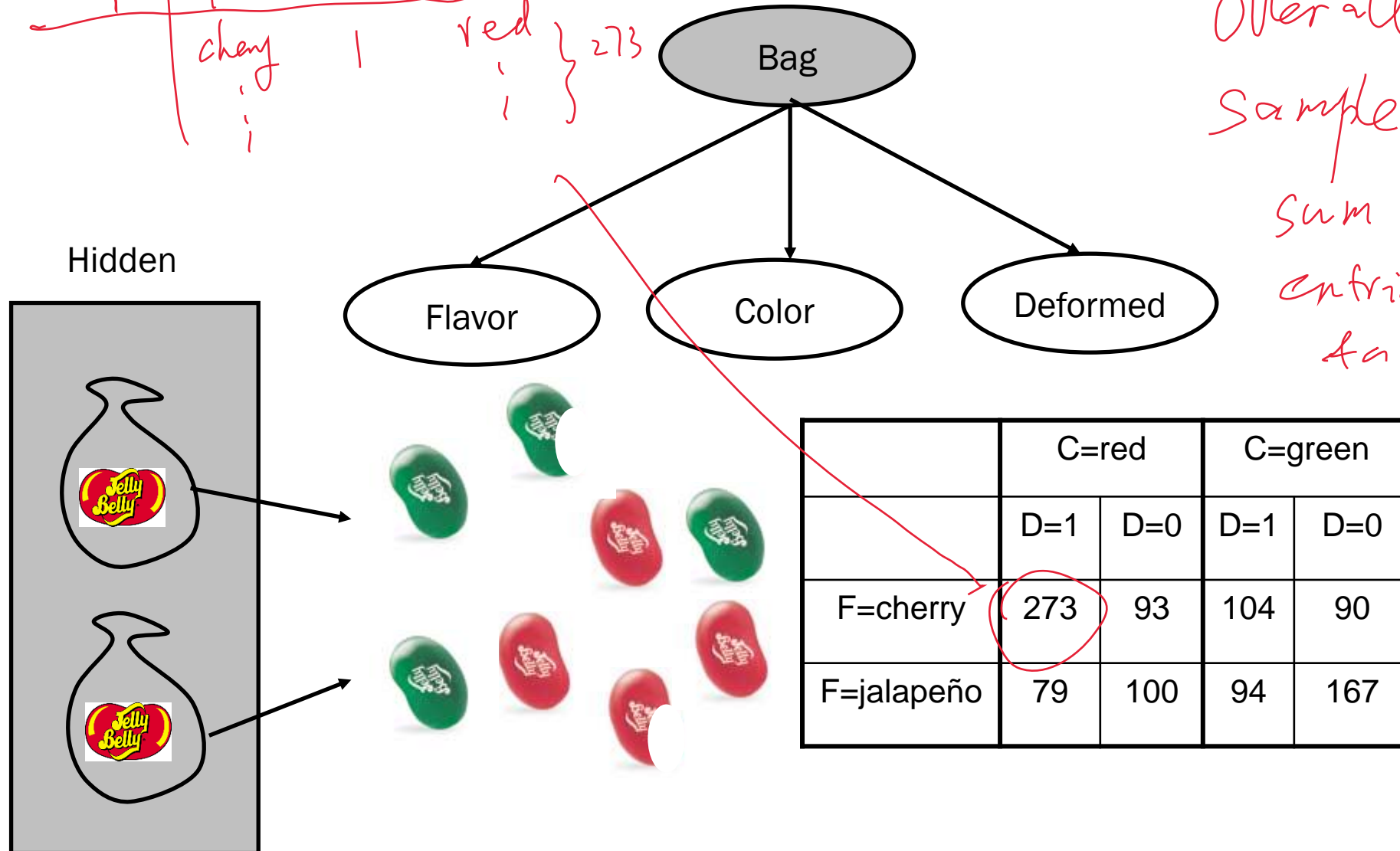
Learning from Partially Observed Data, Example

Handwritten table:

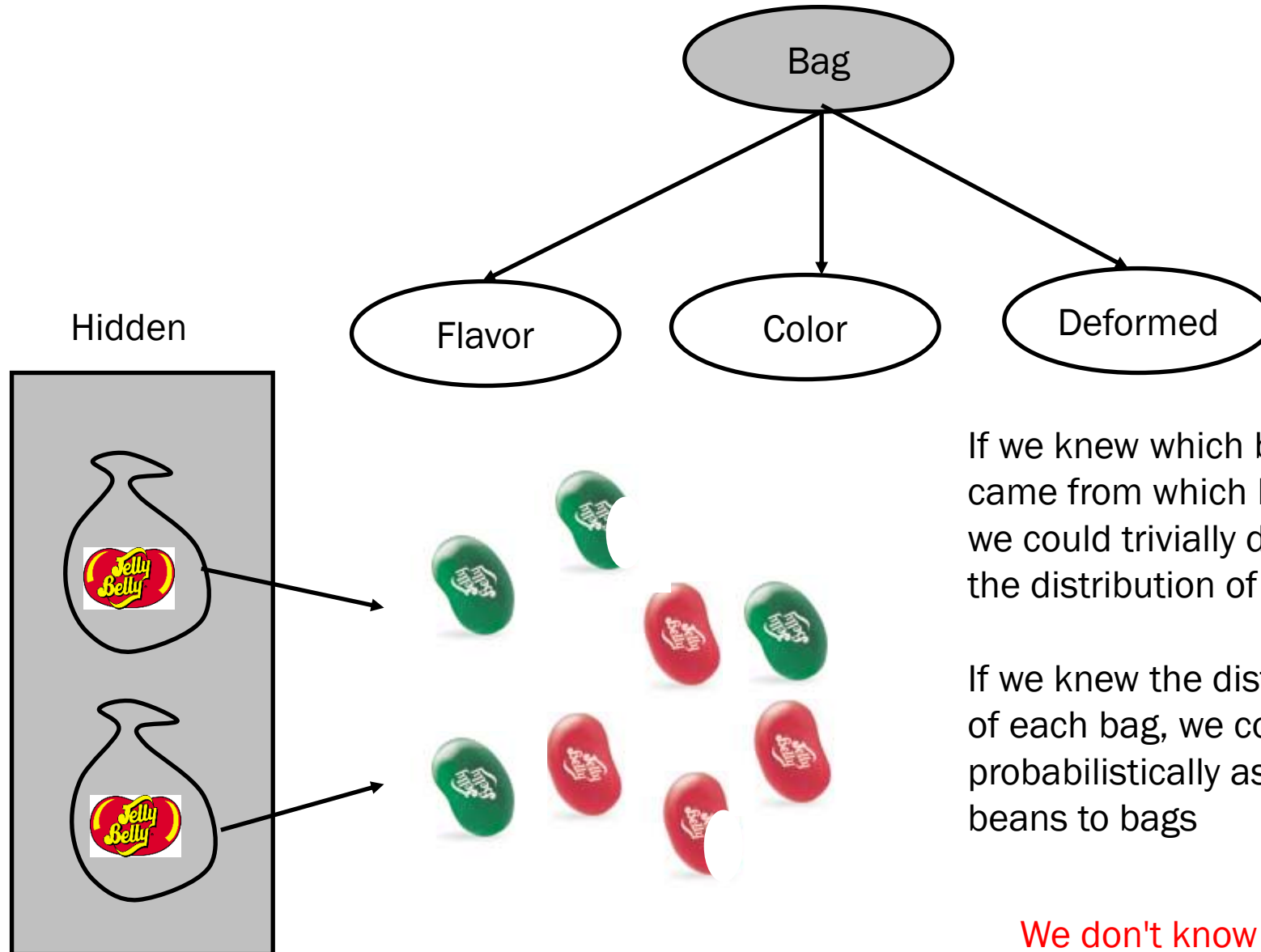
Samples	F	D	C
	cherry	1	red

} 273

Overall # of samples T ??
Sum of all entries in the table.



Learning from Partially Observed Data, Example



If we knew which beans came from which bag, we could trivially determine the distribution of the bag

MLE ✓

If we knew the distribution of each bag, we could probabilistically assign beans to bags

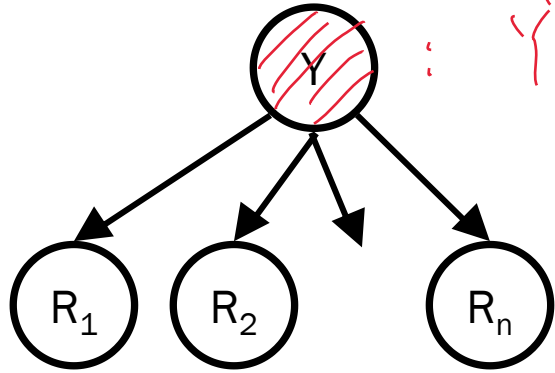
bag info is totally unknown

We don't know either!

Learning with partial data (basic naive bayes model)

$k-1$ probabilities $\sum_{y=1}^k P(Y=y) = 1$

Another example: Movie Recommender System



Y : type of movie goes $\{1, 2, \dots, k\}$ (unobserved)

R_i : rating for i th movie $\{0, 1\}$ (observed)

$P(R_i=1|Y=y)$

don't recommend / recommend

Goal: $P(Y=y | R_1, \dots, R_n)$

Y | $P(R_i=1|Y=y)$

k probabilities

total # of CPTs: $\frac{n \times k}{\text{children}} + \frac{k-1}{\text{parent}}$

Learning with partial data, generally

Let $\{X_1, X_2, \dots, X_n\}$ denote all the nodes in a Bayes Net.
Let H denote the subset of hidden (unobserved) nodes.
Let V denote the subset of visible (observed) nodes.
 $V \cup H = \{X_1, X_2, \dots, X_n\}$

Goal: Estimate the CPTs in the BN to maximize the probability of the partially observed data

$$V \cap H = \emptyset$$

$$\underset{\text{CPTs } t}{\operatorname{argmax}} \left(\prod_t P(V = V^{(t)}) \right)$$

$$\Downarrow$$
$$\underset{\text{CPTs}}{\operatorname{argmax}} P(\text{data}_{\text{vis}})$$

Learning with partial data, generally

Let $\{X_1, X_2, \dots, X_n\}$ denote all the nodes in a Bayes Net.
Let H denote the subset of hidden (unobserved) nodes.
Let V denote the subset of visible (observed) nodes.
 $V \cup H = \{X_1, X_2, \dots, X_n\}$

Goal: Estimate the CPTs in the BN to maximize the probability of the partially observed data

$$\begin{aligned}\mathcal{L} &= \log \prod_{t=1}^T P(V = v^{(t)}) \\ &= \sum_{t=1}^T \log P(V = v^{(t)})\end{aligned}$$

What should we do next?

- A. Use the product rule
- B. Express $P(V = v^{(t)})$ using the conditional independence in the BN
- C. Use marginalization
- D. Use Bayes Rule

MLE: $L = \sum_t \log P(\vec{x} = \vec{x}^{(t)})$
↑
full B.N.

full joint $P(\vec{V} = \vec{V}^{(t)}, \vec{H} = \vec{h})$
not the full B.N.

Learning with partial data, generally

Let $\{X_1, X_2, \dots, X_n\}$ denote all the nodes in a Bayes Net.
 Let H denote the subset of hidden (unobserved) nodes.
 Let V denote the subset of visible (observed) nodes.
 $V \cup H = \{X_1, X_2, \dots, X_n\}$

Goal: Estimate the CPTs in the BN to maximize the probability of the partially observed data

$$\begin{aligned} \mathcal{L} &= \log \prod_{t=1}^T P(V = v^{(t)}) \\ &= \sum_{t=1}^T \log P(V = v^{(t)}) = \sum_{t=1}^T \log \sum_h \underbrace{P(V = v^{(t)}, H = h)}_{\text{full joint}} \end{aligned}$$

$V, v^{(t)}$: vector of observed
 H, h : vector of hidden

$$= \sum_{t=1}^T \log \sum_h \left(\prod_{i=1}^n P(x = x_i | p_{ai} = \tau_i) \right)$$

\mathcal{L}

$\mathcal{L} P(x = x_i | p_{ai} = \tau_i)$

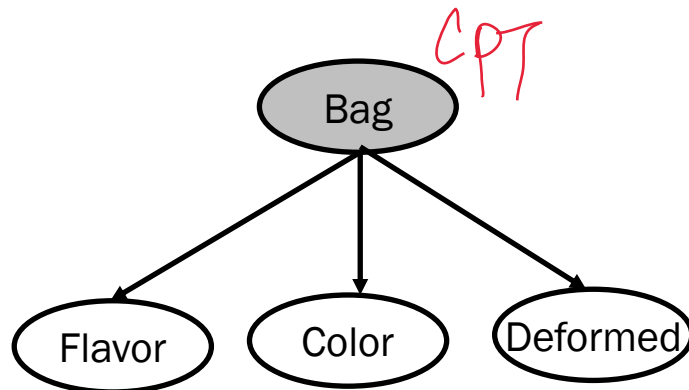
$\triangle \Rightarrow \bigcirc$

no closed form solution.

Iterative

Expectation Maximization: Central Idea

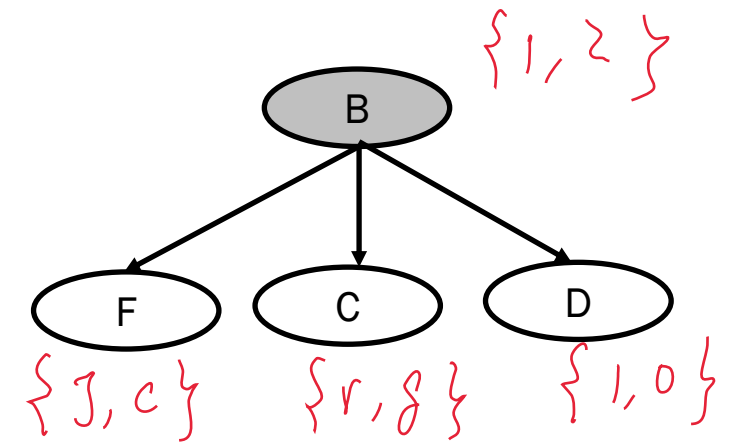
- Expectation Step (E-step): Compute *expected* values for missing data based on the observed data
- Maximization Step (M-step): Re-compute parameters using ML formula for fully observed data, using expected values as if they were observations



1. Guess the parameters of the model (really, just GUESS!)
2. Use those parameters to calculate expected counts of beans from each bag
3. Use those expected counts to update parameters (ML)
4. Lather, rinse, repeat!

Specifying the problem

What are the parameters of this model that we want to learn?



$$P(B=1)$$

✓

$$P(F=J|B=1)$$

$$P(F=J|B=2)$$

$$C: P(C=r|B=1)$$

$$P(C=g|B=2)$$

$$D: P(D=1|B=1)$$

$$P(D=1|B=2)$$

7 CPTs

EM: Example

Iteration 0

Guess (initialize) the parameters of the model

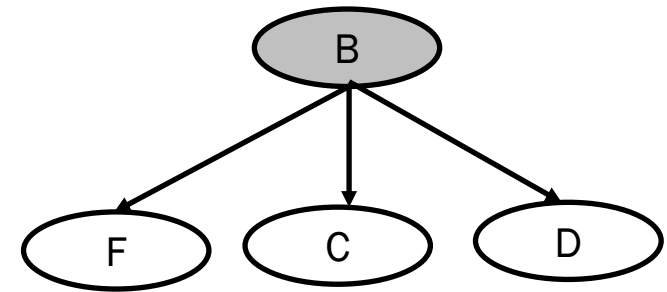
Must be valid probabilities, don't have to be right

For simplicity we might choose:

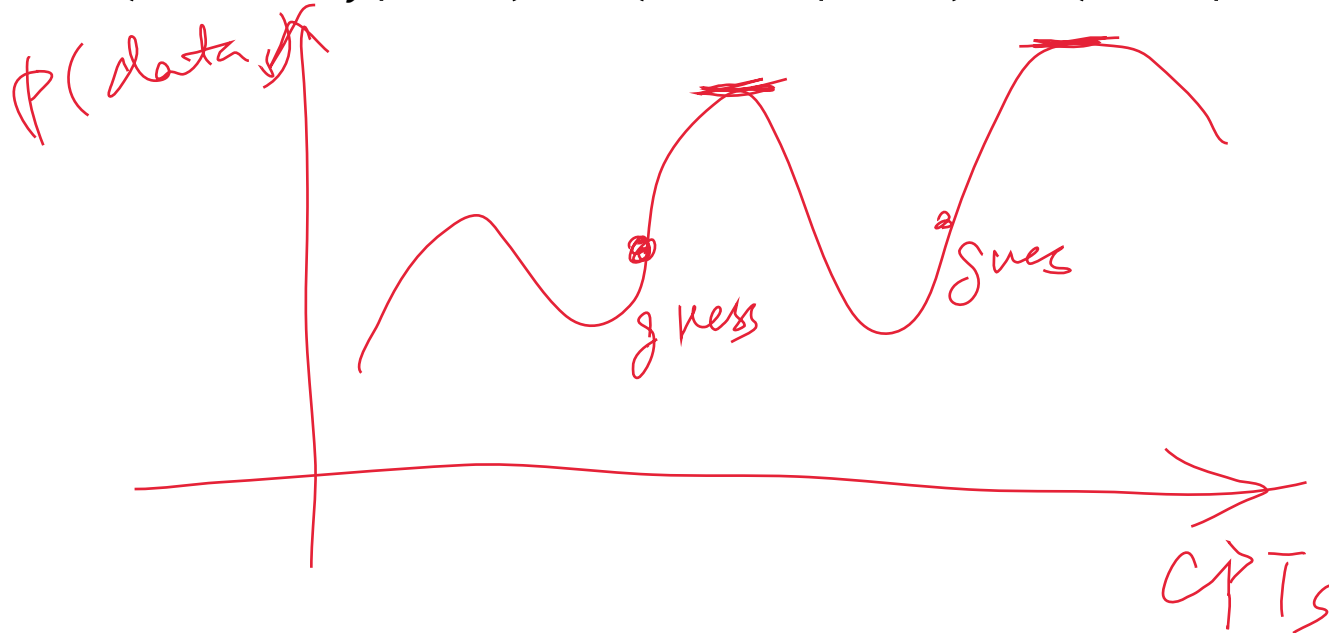
$$P(B = 1) = 0.6$$

$$P(F = \text{cherry} | B = 1) = P(C = \text{red} | B = 1) = P(D = 1 | B = 1) = 0.6$$

$$P(F = \text{cherry} | B = 2) = P(C = \text{red} | B = 2) = P(D = 1 | B = 2) = 0.3$$



Initialization



Guess

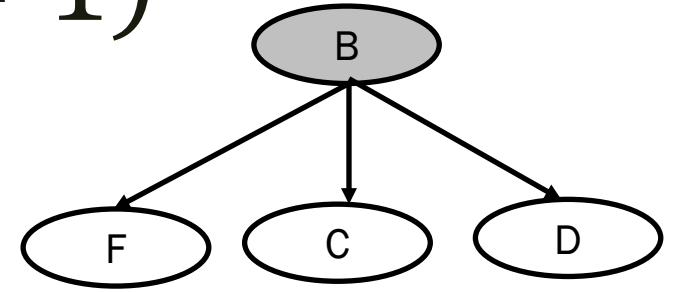
1) better random

2) must $[0, 1]$

Warning: EM may not always produce global max

EM: Example: Updating $P(Bag = 1)$

What would be the ML estimate if we knew which beans came from which bags?



$$P(B=1) = \frac{\text{count}(B=1)}{T} \quad \text{I don't have}$$

But we don't know this, so we need an estimate for this value. Use the expectation!

Expected count:

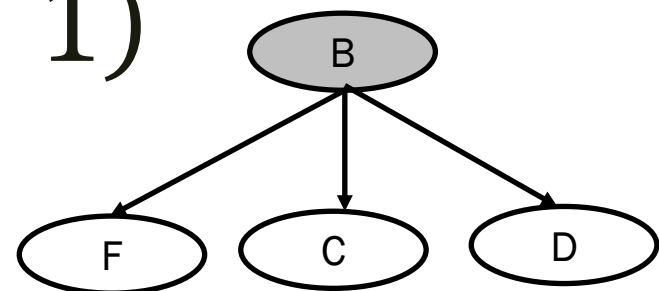
Expected count = $\sum_{t=1}^T P(B=1 | F=f^{(t)}, C=c^{(t)}, D=d^{(t)})$

the expected count for $B=1$ using B.N.
↓
may be decimals

E-step

EM: Example: Updating $P(Bag = 1)$

try to generate data for B using B.N
w/ guessed CPTs



$$\widehat{\text{count}}(B = 1) = \sum_{t=1}^T P(B = 1 | F = f^{(t)}, C = c^{(t)}, D = d^{(t)})$$

E-step

all known through the previous guess

product

$$P(B=1 | F=f^{(t)}, C=c^{(t)}, D=d^{(t)})$$

$$P(B=1, F=f^{(t)}, C=c^{(t)}, D=d^{(t)})$$

Marginalize

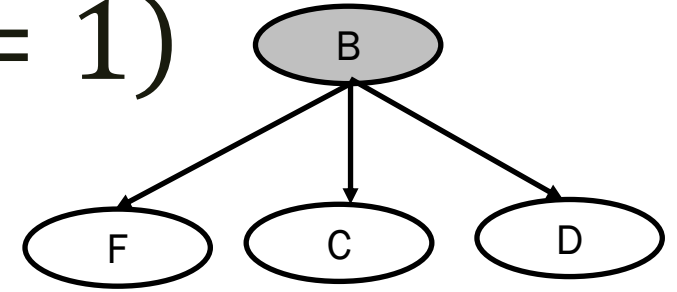
$$P(B=1, F=f^{(t)}, C=c^{(t)}, D=d^{(t)}) + P(B=2, F=f^{(t)}, C=c^{(t)}, D=d^{(t)})$$

$$P(B, F=f^{(t)}, C=c^{(t)}, D=d^{(t)}) = P(B) P(F=f^{(t)} | B) P(C=c^{(t)} | B) P(D=d^{(t)} | B)$$

EM: Example: Updating $P(Bag = 1)$

$$\widehat{\text{count}}(B = 1) = \sum_{t=1}^T P(B = 1 | F = f^{(t)}, C = c^{(t)}, D = d^{(t)})$$

$$= \sum_{t=1}^T \frac{P(B = 1)P(F = f^{(t)} | B = 1)P(C = c^{(t)} | B = 1)P(D = d^{(t)} | B = 1)}{\sum_b P(B = b)P(F = f^{(t)} | B = b)P(C = c^{(t)} | B = b)P(D = d^{(t)} | B = b)}$$



E-step

Where do we get these values?

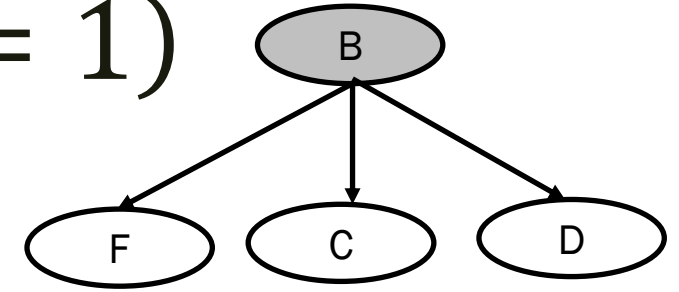
A. We can't, we're stuck

B. From the last iteration

C. Directly from the data

for iter 0, it
is pure guess

EM: Example: Updating $P(Bag = 1)$



E-step

*i+1*th iteration

$$\widehat{\text{count}}(B = 1) = \sum_{t=1}^T P(B = 1 | F = f^{(t)}, C = c^{(t)}, D = d^{(t)})$$

known

$$= \sum_{t=1}^T \frac{P(B = 1)P(F = f^{(t)} | B = 1)P(C = c^{(t)} | B = 1)P(D = d^{(t)} | B = 1)}{\sum_b P(B = b)P(F = f^{(t)} | B = b)P(C = c^{(t)} | B = b)P(D = d^{(t)} | B = b)}$$

new CPT

$$P(Bag = 1) \leftarrow \frac{\widehat{\text{count}}(B = 1)}{T}$$

M-step

For each iteration, you should update all the CPTs before moving on to the next step.

stop when converge.

EM: Initialization

Param.	Iter 0	Iter 1	Iter 2	...
$P(B = 1)$.6	?	}	
$P(F = \text{cherry} B = 1)$.6	}	}	
$P(F = \text{cherry} B = 2)$.3	}	}	
$P(C = \text{red} B = 1)$.6	}	}	
$P(C = \text{red} B = 2)$.3	}	}	
$P(D = 1 B = 1)$.6	}	}	
$P(D = 1 B = 2)$.3	✓	✓	

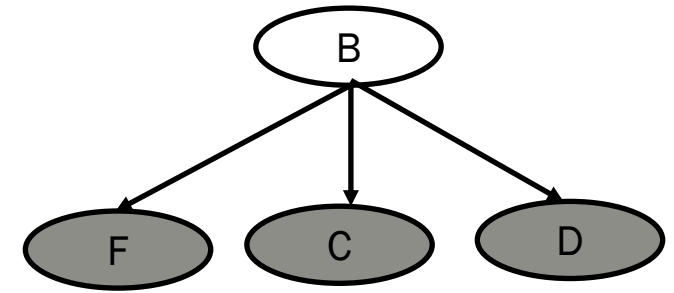
EM algorithm (alternate expression):

Initialize the parameters of the model

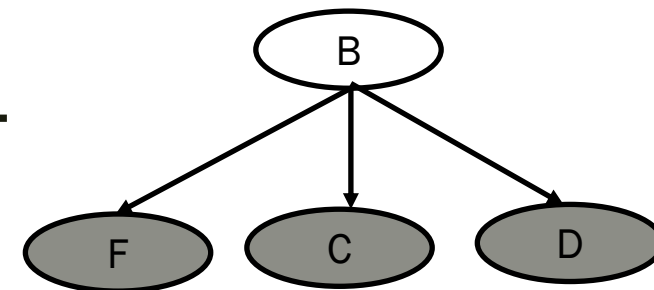
While parameter values have not converged:

For each parameter, p :

- E-step: Use parameter values from last iteration to calculate expected counts necessary to update p
- M-step: Use those expected counts to update p



EM: Calculate $P(B = 1)$ for iter 1



Param.	Iter 0	Iter 1	Iter 2	...
$P(B = 1)$	0.6			
$P(F = \text{cherry} B = 1)$	0.6			
$P(F = \text{cherry} B = 2)$	0.3			
$P(C = \text{red} B = 1)$	0.6			
$P(C = \text{red} B = 2)$	0.3			
$P(D = 1 B = 1)$	0.6			
$P(D = 1 B = 2)$	0.3			

EM algorithm (alternate expression):

Initialize the parameters of the model

While parameter values have not converged:

For each parameter, p :

- E-step: Use parameter values from last iteration to calculate expected counts necessary to update p
- M-step: Use those expected counts to update p

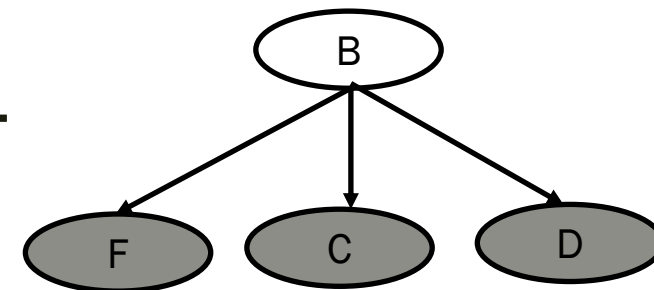
MLE

$$p = \begin{cases} \frac{\text{count}}{T} & \text{w/o parents} \\ \frac{\text{count}(x, \text{Par})}{\text{count}(\text{Par})} & \text{w/ parents} \end{cases}$$

expected $\widehat{\text{count}}(B = 1) = \sum_{t=1}^T P(B = 1 | F = f^{(t)}, C = c^{(t)}, D = d^{(t)})$

$$= \sum_{t=1}^T \frac{P(B = 1)P(F = f^{(t)}|B = 1)P(C = c^{(t)}|B = 1)P(D = d^{(t)}|B = 1)}{\sum_b P(B = b)P(F = f^{(t)}|B = b)P(C = c^{(t)}|B = b)P(D = d^{(t)}|B = b)}$$

EM: Calculate $P(B = 1)$ for iter 1



Param.	Iter 0	Iter 1	Iter 2	...
$P(B = 1)$	0.6			
$P(F = \text{cherry} B = 1)$	0.6			
$P(F = \text{cherry} B = 2)$	0.3			
$P(C = \text{red} B = 1)$	0.6			
$P(C = \text{red} B = 2)$	0.3			
$P(D = 1 B = 1)$	0.6			
$P(D = 1 B = 2)$	0.3			

EM algorithm (alternate expression):

Initialize the parameters of the model

While parameter values have not converged:

For each parameter, p :

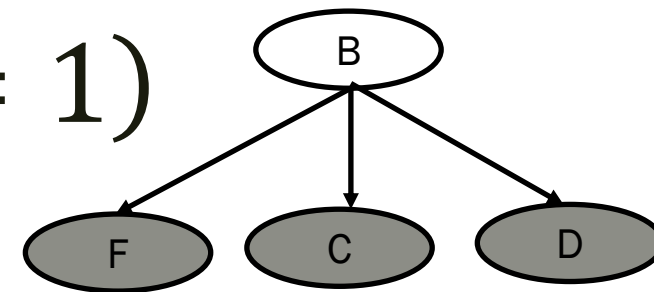
- E-step: Use parameter values from last iteration to calculate expected counts necessary to update p
- M-step: Use those expected counts to update p

$$\widehat{\text{count}}(B = 1) = \sum_{t=1}^T P(B = 1 | F = f^{(t)}, C = c^{(t)}, D = d^{(t)})$$

$$P(\text{Bag} = 1)_{\text{iter1}} \leftarrow \frac{\widehat{\text{count}}(B = 1)}{T}$$

$$= \sum_{t=1}^T \frac{P(B = 1)P(F = f^{(t)}|B = 1)P(C = c^{(t)}|B = 1)P(D = d^{(t)}|B = 1)}{\sum_b P(B = b)P(F = f^{(t)}|B = b)P(C = c^{(t)}|B = b)P(D = d^{(t)}|B = b)}$$

EM: Example: Updating $P(Bag = 1)$



Expected
current

E-step + M-step

$$P(Bag = 1) \leftarrow \frac{1}{T} \sum_{t=1}^T \frac{P(B = 1)P(F = f^{(t)}|B = 1)P(C = c^{(t)}|B = 1)P(D = d^{(t)}|B = 1)}{\sum_b P(B = b)P(F = f^{(t)}|B = b)P(C = c^{(t)}|B = b)P(D = d^{(t)}|B = b)}$$

Current parameter values:

$$P(B = 1) = 0.6$$

$$P(F = \text{cherry}|B = 1) = P(C = \text{red}|B = 1) = P(D = 1|B = 1) = 0.6$$

$$P(F = \text{cherry}|B = 2) = P(C = \text{red}|B = 2) = P(D = 1|B = 2) = 0.3$$

What is the contribution from the red, deformed, cherry candies to this total?

	C=red		C=green	
	D=1	D=0	D=1	D=0
F=cherry	273	93	104	90
F=jalapeño	79	100	94	167

each sample gives

$$\frac{.6 \times .6 \times .6 \times .6}{.6 \times .6 \times .6 \times .6 + .4 \times .3 \times .3 \times .3} = .9231$$

$$273 \times .9231 = \underline{252}$$

represent 273
samples
(cherry, red, damaged)

	C=red		C=green	
	D=1	D=0	D=1	D=0
F=cherry	273	93	104	90
F=jalapeño	79	100	94	167

$$P(Bag = 1)_{iter1} \leftarrow \frac{\widehat{\text{count}}(B = 1)}{T}$$

Param.	Iter 0
$P(B = 1)$	0.6
$P(F = \textit{cherry} B = 1)$	0.6
$P(F = \textit{cherry} B = 2)$	0.3
$P(C = \textit{red} B = 1)$	0.6
$P(C = \textit{red} B = 2)$	0.3
$P(D = 1 B = 1)$	0.6
$P(D = 1 B = 2)$	0.3

$$\begin{aligned} &\widehat{\text{count}}(B = 1) \\ &= \sum_{t=1}^T \frac{P(B = 1)P(F = f^{(t)}|B = 1)P(C = c^{(t)}|B = 1)P(D = d^{(t)}|B = 1)}{\sum_b P(B = b)P(F = f^{(t)}|B = b)P(C = c^{(t)}|B = b)P(D = d^{(t)}|B = b)} \end{aligned}$$

F	C	D	contribution
cherry	red	1	252
cherry	red	0	273 72
cherry	green	1	
cherry	green	0	
jalapeno	red	1	
jalapeno	red	0	
jalapeno	green	1	
jalapeno	green	0	

$$.6 \times .6 \times .6 \times .4$$

$$.6 \times .6 \times .6 \times .4 + .4 \times .3 \times .3 \times .7$$

$$= 1.7742 \times 93$$

$$A: 80$$

$$B: 72$$

$$C: 44$$

	C=red		C=green	
	D=1	D=0	D=1	D=0
F=cherry	273	93	104	90
F=jalapeño	79	100	94	167

$$P(\text{Bag} = 1)_{\text{iter}1} \leftarrow \frac{\widehat{\text{count}}(B = 1)}{T} \approx \frac{642}{1000} \approx 0.642$$

Sum

$$\widehat{\text{count}}(B = 1) = \sum_{t=1}^T \frac{P(B = 1)P(F = f^{(t)}|B = 1)P(C = c^{(t)}|B = 1)P(D = d^{(t)}|B = 1)}{\sum_b P(B = b)P(F = f^{(t)}|B = b)P(C = c^{(t)}|B = b)P(D = d^{(t)}|B = b)}$$

Param.	Iter 0
$P(B = 1)$	0.6
$P(F = \text{cherry} B = 1)$	0.6
$P(F = \text{cherry} B = 2)$	0.3
$P(C = \text{red} B = 1)$	0.6
$P(C = \text{red} B = 2)$	0.3
$P(D = 1 B = 1)$	0.6
$P(D = 1 B = 2)$	0.3

F	C	D	contribution
cherry	red	1	252
cherry	red	0	72
cherry	green	1	27 80.52
cherry	green	0	44.53
jalapeno	red	1	61.16
jalapeno	red	0	49.48
jalapeno	green	1	46.52
jalapeno	green	0	36.52

$$104 \times \frac{.6 \times .6 \times .4 \times .6}{.6 \times .6 \times .4 \times .6 + .4 \times .3 \times .7 \times .3} = 80.52$$

A: 61
B: 80
C: 46
D: 36

	C=red		C=green	
	D=1	D=0	D=1	D=0
F=cherry	273	93	104	90
F=jalapeño	79	100	94	167

$$P(Bag = 1)_{iter1} \leftarrow \frac{\widehat{\text{count}}(B = 1)}{T}$$

$$\begin{aligned} &\widehat{\text{count}}(B = 1) \\ &= \sum_{t=1}^T \frac{P(B = 1)P(F = f^{(t)}|B = 1)P(C = c^{(t)}|B = 1)P(D = d^{(t)}|B = 1)}{\sum_b P(B = b)P(F = f^{(t)}|B = b)P(C = c^{(t)}|B = b)P(D = d^{(t)}|B = b)} \end{aligned}$$

F	C	D	contribution
cherry	red	1	
cherry	red	0	
cherry	green	1	
cherry	green	0	
jalapeno	red	1	
jalapeno	red	0	
jalapeno	green	1	
jalapeno	green	0	

Param.	Iter 0
$P(B = 1)$	0.6
$P(F = \textit{cherry} B = 1)$	0.6
$P(F = \textit{cherry} B = 2)$	0.3
$P(C = \textit{red} B = 1)$	0.6
$P(C = \textit{red} B = 2)$	0.3
$P(D = 1 B = 1)$	0.6
$P(D = 1 B = 2)$	0.3

	C=red		C=green	
	D=1	D=0	D=1	D=0
F=cherry	273	93	104	90
F=jalapeño	79	100	94	167

$$P(Bag = 1)_{iter1} \leftarrow \frac{\widehat{\text{count}}(B = 1)}{T}$$

$$\begin{aligned} &\widehat{\text{count}}(B = 1) \\ &= \sum_{t=1}^T \frac{P(B = 1)P(F = f^{(t)}|B = 1)P(C = c^{(t)}|B = 1)P(D = d^{(t)}|B = 1)}{\sum_b P(B = b)P(F = f^{(t)}|B = b)P(C = c^{(t)}|B = b)P(D = d^{(t)}|B = b)} \end{aligned}$$

F	C	D	contribution
cherry	red	1	
cherry	red	0	
cherry	green	1	
cherry	green	0	
jalapeno	red	1	
jalapeno	red	0	
jalapeno	green	1	
jalapeno	green	0	

Param.	Iter 0
$P(B = 1)$	0.6
$P(F = \textit{cherry} B = 1)$	0.6
$P(F = \textit{cherry} B = 2)$	0.3
$P(C = \textit{red} B = 1)$	0.6
$P(C = \textit{red} B = 2)$	0.3
$P(D = 1 B = 1)$	0.6
$P(D = 1 B = 2)$	0.3

	C=red		C=green	
	D=1	D=0	D=1	D=0
F=cherry	273	93	104	90
F=jalapeño	79	100	94	167

$$P(Bag = 1)_{iter1} \leftarrow \frac{\widehat{\text{count}}(B = 1)}{T}$$

$$\begin{aligned} &\widehat{\text{count}}(B = 1) \\ &= \sum_{t=1}^T \frac{P(B = 1)P(F = f^{(t)}|B = 1)P(C = c^{(t)}|B = 1)P(D = d^{(t)}|B = 1)}{\sum_b P(B = b)P(F = f^{(t)}|B = b)P(C = c^{(t)}|B = b)P(D = d^{(t)}|B = b)} \end{aligned}$$

F	C	D	contribution
cherry	red	1	
cherry	red	0	
cherry	green	1	
cherry	green	0	
jalapeno	red	1	
jalapeno	red	0	
jalapeno	green	1	
jalapeno	green	0	

Param.	Iter 0
$P(B = 1)$	0.6
$P(F = \textit{cherry} B = 1)$	0.6
$P(F = \textit{cherry} B = 2)$	0.3
$P(C = \textit{red} B = 1)$	0.6
$P(C = \textit{red} B = 2)$	0.3
$P(D = 1 B = 1)$	0.6
$P(D = 1 B = 2)$	0.3

	C=red		C=green	
	D=1	D=0	D=1	D=0
F=cherry	273	93	104	90
F=jalapeño	79	100	94	167

$$P(Bag = 1)_{iter1} \leftarrow \frac{\widehat{\text{count}}(B = 1)}{T}$$

$$\begin{aligned} &\widehat{\text{count}}(B = 1) \\ &= \sum_{t=1}^T \frac{P(B = 1)P(F = f^{(t)}|B = 1)P(C = c^{(t)}|B = 1)P(D = d^{(t)}|B = 1)}{\sum_b P(B = b)P(F = f^{(t)}|B = b)P(C = c^{(t)}|B = b)P(D = d^{(t)}|B = b)} \end{aligned}$$

F	C	D	contribution
cherry	red	1	
cherry	red	0	
cherry	green	1	
cherry	green	0	
jalapeno	red	1	
jalapeno	red	0	
jalapeno	green	1	
jalapeno	green	0	

Param.	Iter 0
$P(B = 1)$	0.6
$P(F = \textit{cherry} B = 1)$	0.6
$P(F = \textit{cherry} B = 2)$	0.3
$P(C = \textit{red} B = 1)$	0.6
$P(C = \textit{red} B = 2)$	0.3
$P(D = 1 B = 1)$	0.6
$P(D = 1 B = 2)$	0.3

	C=red		C=green	
	D=1	D=0	D=1	D=0
F=cherry	273	93	104	90
F=jalapeño	79	100	94	167

$$P(Bag = 1)_{iter1} \leftarrow \frac{\widehat{\text{count}}(B = 1)}{T}$$

$$\begin{aligned} &\widehat{\text{count}}(B = 1) \\ &= \sum_{t=1}^T \frac{P(B = 1)P(F = f^{(t)}|B = 1)P(C = c^{(t)}|B = 1)P(D = d^{(t)}|B = 1)}{\sum_b P(B = b)P(F = f^{(t)}|B = b)P(C = c^{(t)}|B = b)P(D = d^{(t)}|B = b)} \end{aligned}$$

F	C	D	contribution
cherry	red	1	
cherry	red	0	
cherry	green	1	
cherry	green	0	
jalapeno	red	1	
jalapeno	red	0	
jalapeno	green	1	
jalapeno	green	0	

Param.	Iter 0
$P(B = 1)$	0.6
$P(F = \textit{cherry} B = 1)$	0.6
$P(F = \textit{cherry} B = 2)$	0.3
$P(C = \textit{red} B = 1)$	0.6
$P(C = \textit{red} B = 2)$	0.3
$P(D = 1 B = 1)$	0.6
$P(D = 1 B = 2)$	0.3

	C=red		C=green	
	D=1	D=0	D=1	D=0
F=cherry	273	93	104	90
F=jalapeño	79	100	94	167

$$P(Bag = 1)_{iter1} \leftarrow \frac{\widehat{\text{count}}(B = 1)}{T}$$

$$\begin{aligned} &\widehat{\text{count}}(B = 1) \\ &= \sum_{t=1}^T \frac{P(B = 1)P(F = f^{(t)}|B = 1)P(C = c^{(t)}|B = 1)P(D = d^{(t)}|B = 1)}{\sum_b P(B = b)P(F = f^{(t)}|B = b)P(C = c^{(t)}|B = b)P(D = d^{(t)}|B = b)} \end{aligned}$$

F	C	D	contribution
cherry	red	1	
cherry	red	0	
cherry	green	1	
cherry	green	0	
jalapeno	red	1	
jalapeno	red	0	
jalapeno	green	1	
jalapeno	green	0	

Param.	Iter 0
$P(B = 1)$	0.6
$P(F = \textit{cherry} B = 1)$	0.6
$P(F = \textit{cherry} B = 2)$	0.3
$P(C = \textit{red} B = 1)$	0.6
$P(C = \textit{red} B = 2)$	0.3
$P(D = 1 B = 1)$	0.6
$P(D = 1 B = 2)$	0.3

	iter 0	iter 1
$P(B=1)$.6	.642
	.6	→ .699
	.3	→
	.6	→
	.3	→
	.6	→
	.3	→

What about the conditionals?

How do we calculate $P(F = \text{cherry} | B = 1)$

What is the ML formula for this parameter, if we had fully observed data?

For MLE, $P(F = \text{cherry} | B = 1) = \frac{\text{count}(F = \text{cherry}, B = 1)}{\text{count}(B = 1)}$

Strange: F appears on both sides of the condition

??
 \Downarrow
 \uparrow known 642

What is the estimate for $\text{count}(F = \text{cherry}, B = 1)$?

A. $\sum_{t=1}^T P(B = 1, F = \text{cherry} | F = f^{(t)}, C = c^{(t)}, D = d^{(t)})$: expected count if $f^{(t)}, c^{(t)}, d^{(t)}$ are observed.

B. $\sum_{t=1}^T P(B = 1 | F = f^{(t)}, C = c^{(t)}, D = d^{(t)})$

C. $\sum_{t=1}^T P(F = f^{(t)} | B = 1)$

D. $\sum_{t=1}^T P(F = \text{cherry} | B = 1)$

E. I have no idea

A: change it to

$\sum_{t=1}^T \mathbb{I}(f^{(t)} = \text{cherry}) P(B = 1, F = \text{cherry} | F = f^{(t)}, C = c^{(t)}, D = d^{(t)})$

What about the conditionals?

How do we calculate $P(F = \text{cherry} | B = 1)$

$$\widehat{\text{count}}(B = 1, F = \text{cherry}) = \sum_{t=1}^T P(B = 1, F = \text{cherry} | F = f^{(t)}, C = c^{(t)}, D = d^{(t)})$$

$$= \frac{1}{T} \prod_{t=1}^T P(B = 1 | F = \text{cherry}, C = c^{(t)}, D = d^{(t)}) \mathbb{1}(f^{(t)} = \text{cherry})$$

252
72
80.52
44

1st four rows
where $F = \text{cherry}$

$$\frac{\text{Count}(B = 1, F = \text{cherry})}{\text{Count}(B = 1)} = \frac{449}{642} \approx .699$$

EM: General update formula for BN params

Nodes with parents:

$$P(X_i = x | Pa_i = \pi)_{iter_i+1} = \frac{\widehat{\text{count}}(X_i = x, Pa_i = \pi)}{\widehat{\text{count}}(Pa_i = \pi)}$$

Root nodes:

$$P(X_i = x)_{iter_i+1} = \frac{\widehat{\text{count}}(X_i = x)}{T}$$

EM algorithm (alternate expression):

Initialize the parameters of the model

While parameter values have not converged:

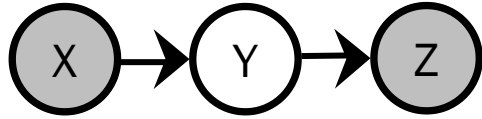
For each parameter, p :

- E-step: Use parameter values from last iteration to calculate expected counts necessary to update p
- M-step: Use those expected counts to update p

Properties of EM

- Monotonic convergence: Each iteration of EM increases (or does not change) log-likelihood of observed data.
- No tuning parameters, no learning rates, no backtracking
- Converges to a local or global maximum. Often depends on initialization values.

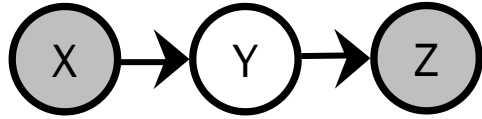
EM Practice on Second Simple Example



Which parameters of this network can you estimate directly from the data (in one step—no iteration required)?

- A. $P(X)$
- B. $P(Y|X)$
- C. $P(Z|Y)$
- D. Both A and C
- E. None of them

EM Practice on Second Simple Example



$$V = \{X, Z\} \quad H = \{Y\}$$

Your turn! Express $P(Z = z|Y = y)$ in terms of $I(x, x^{(t)})$, $I(z, z^{(t)})$ and $P(Y = y|X = x^{(t)}, Z = z^{(t)})$

Consider a model with the following structure:

$$A \leftarrow B \rightarrow C$$

Suppose all of the variables are binary, and suppose only A and C are observable in the dataset. If we have T observations, then the dataset is of the form $\{(a_t, c_t)\}_{t=1}^T$.

- (a) If we want to apply EM to this dataset, list all of the CPT entries that would need to be estimated.
(Hint: there are 5 CPT entries in this model.)

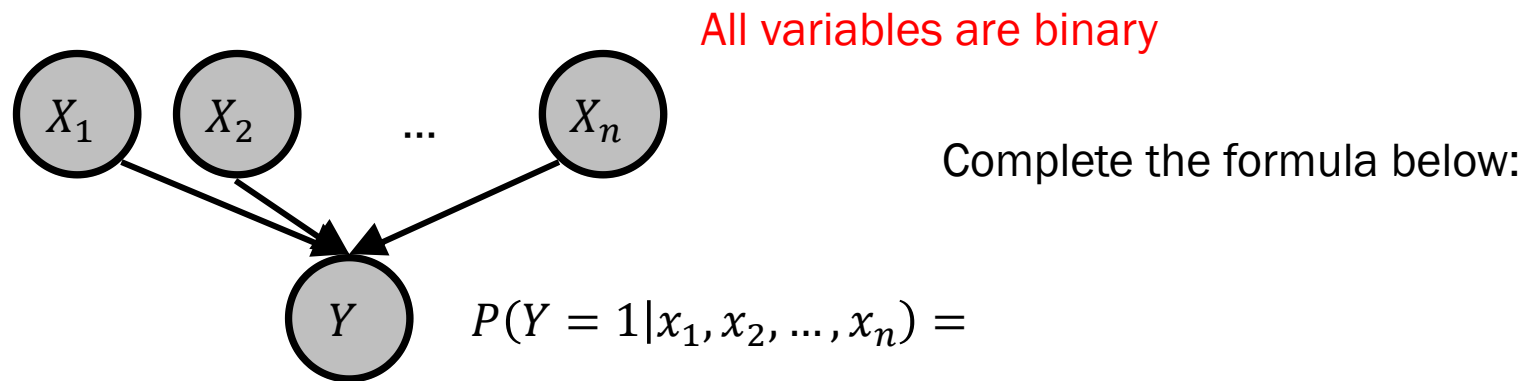
Consider a model with the following structure:

$$A \leftarrow B \rightarrow C$$

Suppose all of the variables are binary, and suppose only A and C are observable in the dataset. If we have T observations, then the dataset is of the form $\{(a_t, c_t)\}_{t=1}^T$.

- (b) Find formulas for the EM update rules for the CPT entries listed in part (a). (Hint: to simplify the formulas, you can first express the EM updates in terms of equality-testing functions and the probability $P(B|A, C)$. Then, you can find a formula for $P(B|A, C)$ in terms of the CPT entries.)

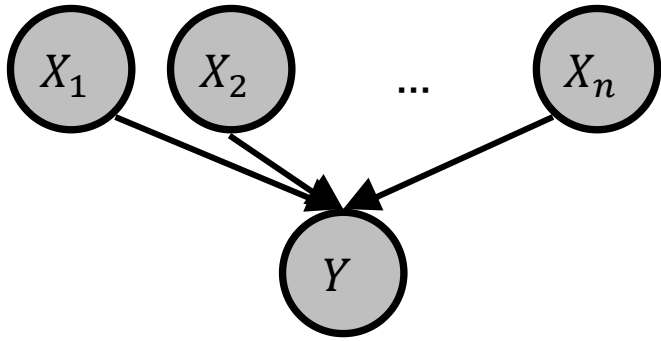
EM for Learning Noisy-OR model



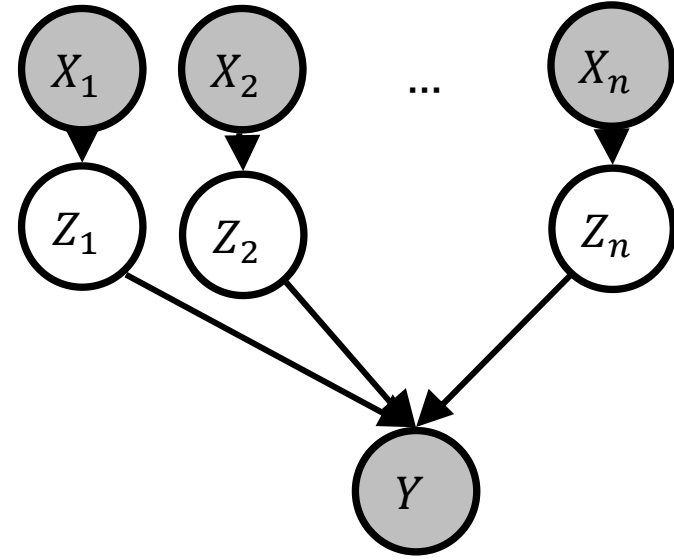
Problem: From *complete* data $\{(x_1^{(t)}, x_2^{(t)}, \dots, x_n^{(t)}, y^{(t)})\}_{t=1}^T$, estimate $p_i \in [0,1]$.

EM for Learning Noisy-OR: Alternate model

Problem: From (complete) data $\{(x_1^{(t)}, x_n^{(t)}, \dots, x_n^{(t)}, y^{(t)})\}_{t=1}^T$, estimate $p_i \in [0,1]$.



$$P(Y = 1 | x_1, x_2, \dots, x_n) = 1 - \prod_{i=1}^n (1 - p_i)^{x_i}$$



EM for Learning Noisy-OR: Alternate model

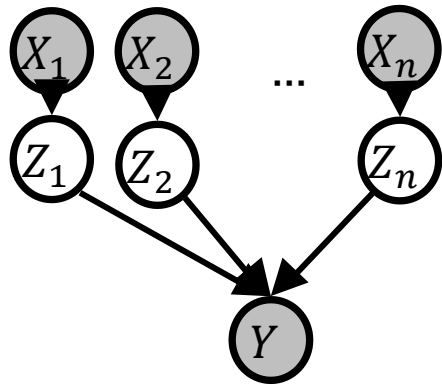
Problem: From (complete) data $\{(x_1^{(t)}, x_n^{(t)}, \dots, x_n^{(t)}, y^{(t)})\}_{t=1}^T$, estimate $p_i \in [0,1]$.

$$P(Z_i = 0 | X_i = 0) = 1$$

$$P(Z_i = 1 | X_i = 1) = p_i$$

What is $P(Y = 1 | x_1, x_2, \dots, x_n)$ in this alternate model?

First show that $P(Y = 1 | \vec{x}) = \sum_{\vec{z}} P(Y = 1 | \vec{z}) P(\vec{z} | \vec{x})$



$$P(Y | \vec{z}) = \text{OR}(\vec{z})$$

EM for Learning Noisy-OR: Alternate model

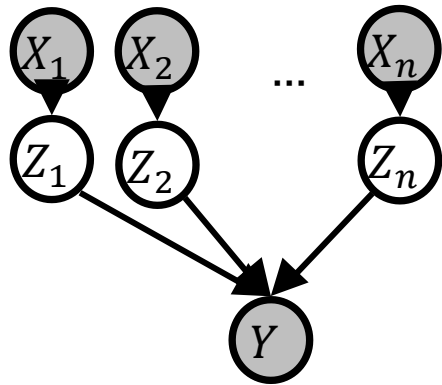
Problem: From (complete) data $\{(x_1^{(t)}, x_n^{(t)}, \dots, x_n^{(t)}, y^{(t)})\}_{t=1}^T$, estimate $p_i \in [0,1]$.

$$P(Z_i = 0 | X_i = 0) = 1$$

$$P(Z_i = 1 | X_i = 1) = p_i$$

What is $P(Y = 1 | x_1, x_2, \dots, x_n)$ in this alternate model?

$$P(Y = 1 | \vec{x}) = \sum_{\vec{z}} \boxed{P(Y = 1 | \vec{z}) P(\vec{z} | \vec{x})}$$



$$P(Y | \vec{z}) = \text{OR}(\vec{z})$$

When is the term in the red box 0?

- A. When at least one z_i is 0
- B. When all z_i s are 0
- C. You can't tell from the information given

EM for Learning Noisy-OR: Alternate model

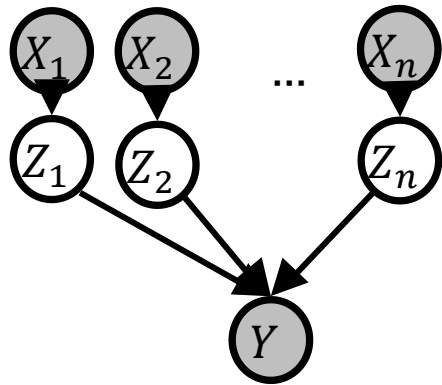
Problem: From (complete) data $\{(x_1^{(t)}, x_n^{(t)}, \dots, x_n^{(t)}, y^{(t)})\}_{t=1}^T$, estimate $p_i \in [0,1]$.

$$P(Z_i = 0 | X_i = 0) = 1$$

$$P(Z_i = 1 | X_i = 1) = p_i$$

What is $P(Y = 1 | x_1, x_2, \dots, x_n)$ in this alternate model?

$$P(Y = 1 | \vec{x}) = \sum_{\vec{z}} P(Y = 1 | \vec{z}) P(\vec{z} | \vec{x})$$



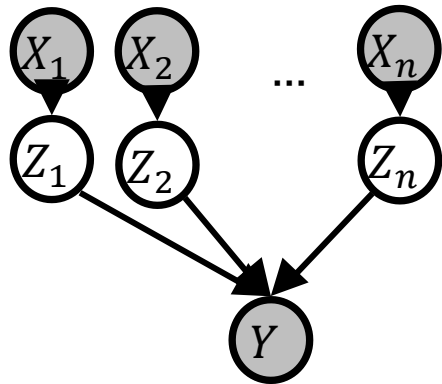
$$P(Y | \vec{z}) = \text{OR}(\vec{z})$$

EM for Learning Noisy-OR: Alternate model

Problem: From (complete) data $\{(x_1^{(t)}, x_n^{(t)}, \dots, x_n^{(t)}, y^{(t)})\}_{t=1}^T$, estimate $p_i \in [0,1]$.

$$P(Z_i = 0 | X_i = 0) = 1$$

$$P(Z_i = 1 | X_i = 1) = p_i$$



$$P(Y | \vec{Z}) = \text{OR}(\vec{Z})$$

What is $P(Y = 1 | x_1, x_2, \dots, x_n)$ in this alternate model?

$$P(Y = 1 | \vec{x}) = \sum_{\vec{z}} P(Y = 1 | \vec{z}) P(\vec{z} | \vec{x}) = \sum_{\vec{z} \neq 0} P(\vec{z} | \vec{x}) =$$

- A. $P(\vec{z} = 0 | \vec{x})$
- B. $1 - P(\vec{z} = 0 | \vec{x})$
- C. $P(\vec{z} | \vec{x} = 0)$
- D. $1 - P(\vec{z} | \vec{x} = 0)$
- E. None of these

EM for Learning Noisy-OR: Alternate model

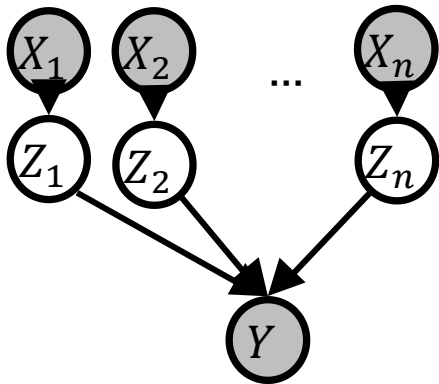
Problem: From (complete) data $\{(x_1^{(t)}, x_n^{(t)}, \dots, x_n^{(t)}, y^{(t)})\}_{t=1}^T$, estimate $p_i \in [0,1]$.

$$P(Z_i = 0 | X_i = 0) = 1$$

$$P(Z_i = 1 | X_i = 1) = p_i$$

What is $P(Y = 1 | x_1, x_2, \dots, x_n)$ in this alternate model?

$$P(Y = 1 | \vec{x}) = \sum_{\vec{z}} P(Y = 1 | \vec{z}) P(\vec{z} | \vec{x}) = \sum_{\vec{z} \neq 0} P(\vec{z} | \vec{x}) = 1 - P(\vec{z} = 0 | \vec{x})$$



$$P(Y | \vec{z}) = \text{OR}(\vec{z})$$

Learning $P(Z_i = 1|X_i = 1) = p_i$ using EM

- If we had fully observed data:

Learning $P(Z_i = 1|X_i = 1) = p_i$ using EM

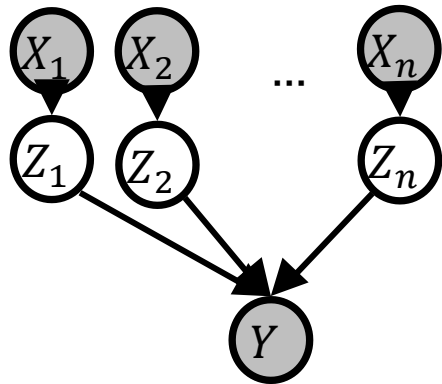
$$\widehat{count}(Z_i = 1, X_i = 1) = \sum_{t=1}^T P(Z_i = 1, X_i = 1 | \vec{x}^{(t)}, y^{(t)})$$

Learning $P(Z_i = 1 | X_i = 1) = p_i$ using EM:
 Calculate $P(Z_i = 1 | \vec{x}^{(t)}, y^{(t)})$ with current params

$$P(Z_i = 0 | X_i = 0) = 1$$

$$P(Z_i = 1 | X_i = 1) = p_i$$

$$\widehat{\text{count}}(Z_i = 1, X_i = 1) = \sum_{t=1}^T P(Z_i = 1, X_i = 1 | \vec{x}^{(t)}, y^{(t)}) = \sum_{t=1}^T x_i^{(t)} P(Z_i = 1 | \vec{x}^{(t)}, y^{(t)})$$



$$P(Y | \vec{Z}) = \text{OR}(\vec{Z})$$

$$P(Z_i = 1 | \vec{x}^{(t)}, y^{(t)}) = \frac{P(Y = y^{(t)} | Z_i = 1) P(Z_i = 1 | \vec{x}^{(t)})}{P(Y = y^{(t)} | \vec{x}^{(t)})}$$

Which of the following is NOT a true statement?

A. $P(Y = y^{(t)} | Z_i = 1) = 0$ when $y^{(t)} = 0$

B. $P(Y = y^{(t)} | Z_i = 1) = 1$ when $y^{(t)} = 1$

C. $P(Z_i = 1 | \vec{x}^{(t)}) = 1$ when $x_i^{(t)} = 1$

D. $P(Z_i = 1 | \vec{x}^{(t)}) = 0$ when $x_i^{(t)} = 0$

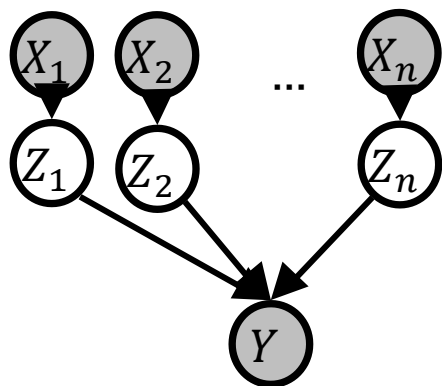
E. More than one of the above is NOT true

Learning $P(Z_i = 1 | X_i = 1) = p_i$ using EM:
 Calculate $P(Z_i = 1 | \vec{x}^{(t)}, y^{(t)})$ with current params

$$P(Z_i = 0 | X_i = 0) = 1$$

$$P(Z_i = 1 | X_i = 1) = p_i$$

$$\widehat{\text{count}}(Z_i = 1, X_i = 1) = \sum_{t=1}^T P(Z_i = 1, X_i = 1 | \vec{x}^{(t)}, y^{(t)}) = \sum_{t=1}^T x_i^{(t)} P(Z_i = 1 | \vec{x}^{(t)}, y^{(t)})$$



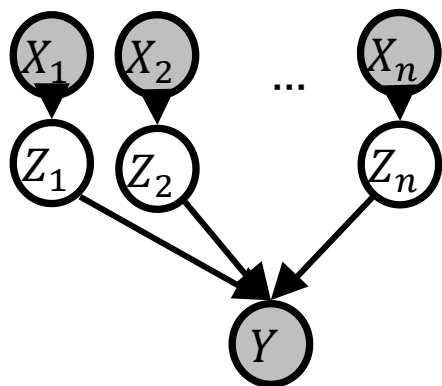
$$P(Y | \vec{Z}) = \text{OR}(\vec{Z})$$

$$P(Z_i = 1 | \vec{x}^{(t)}, y^{(t)}) = \frac{P(Y = y^{(t)} | Z_i = 1) P(Z_i = 1 | \vec{x}^{(t)})}{P(Y = y^{(t)} | \vec{x}^{(t)})}$$

Learning $P(Z_i = 1|X_i = 1) = p_i$ using EM: Update rule

$$P(Z_i = 0|X_i = 0) = 1$$

$$P(Z_i = 1|X_i = 1) = p_i$$



$$P(Y|\vec{Z}) = \text{OR}(\vec{Z})$$

$$\widehat{\text{count}}(Z_i = 1, X_i = 1) = \sum_{t=1}^T \left[\frac{p_i y^{(t)} x_i^{(t)}}{1 - \prod_{i=1}^n (1 - p_i) x_i^{(t)}} \right]$$

$$p_i \leftarrow \frac{\widehat{\text{count}}(Z_i = 1, X_i = 1)}{\widehat{\text{count}}(X_i = 1)} = \frac{p_i}{T_i} \sum_{t=1}^T \left[\frac{y^{(t)} x_i^{(t)}}{1 - \prod_{i=1}^n (1 - p_i) x_i^{(t)}} \right]$$