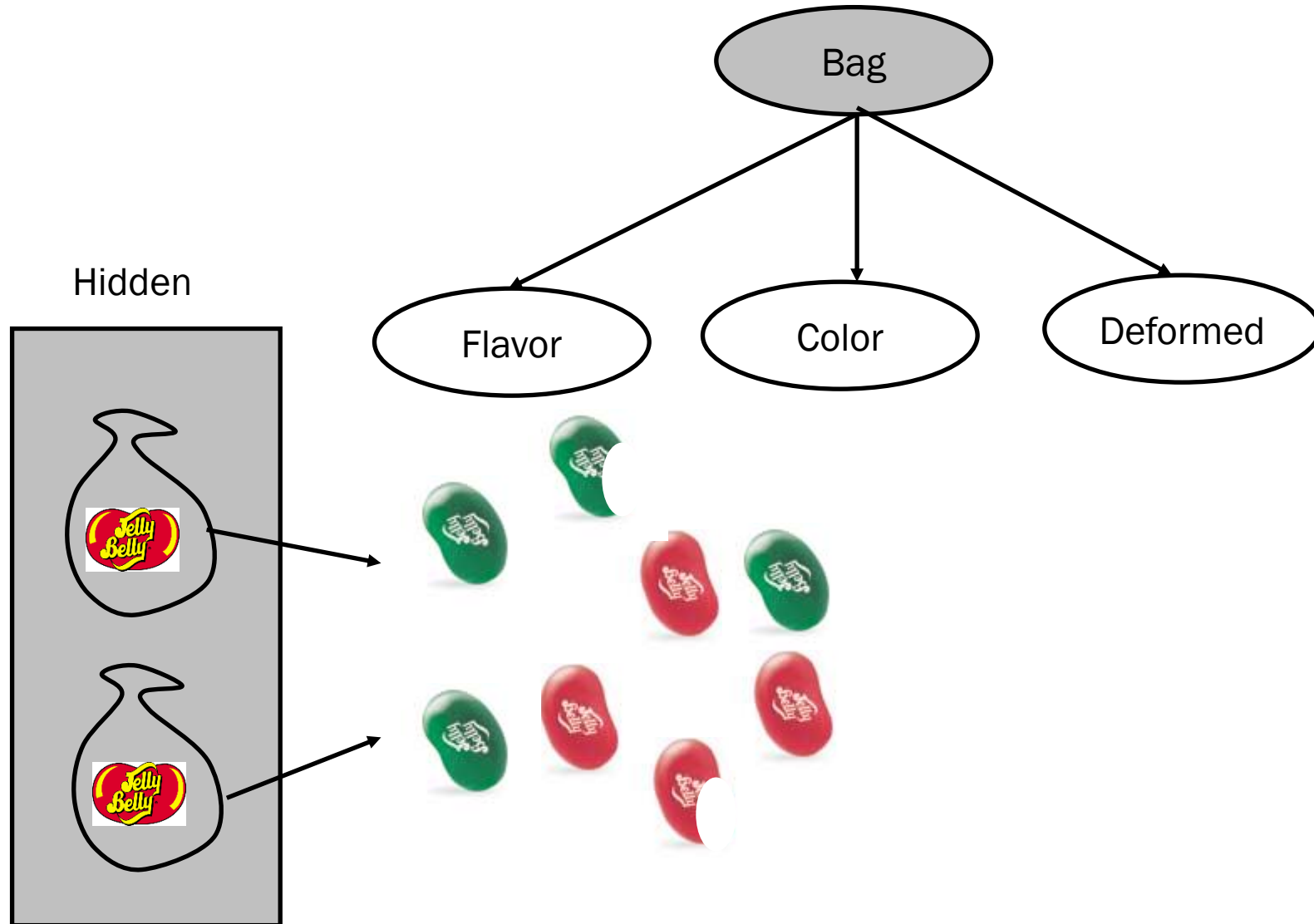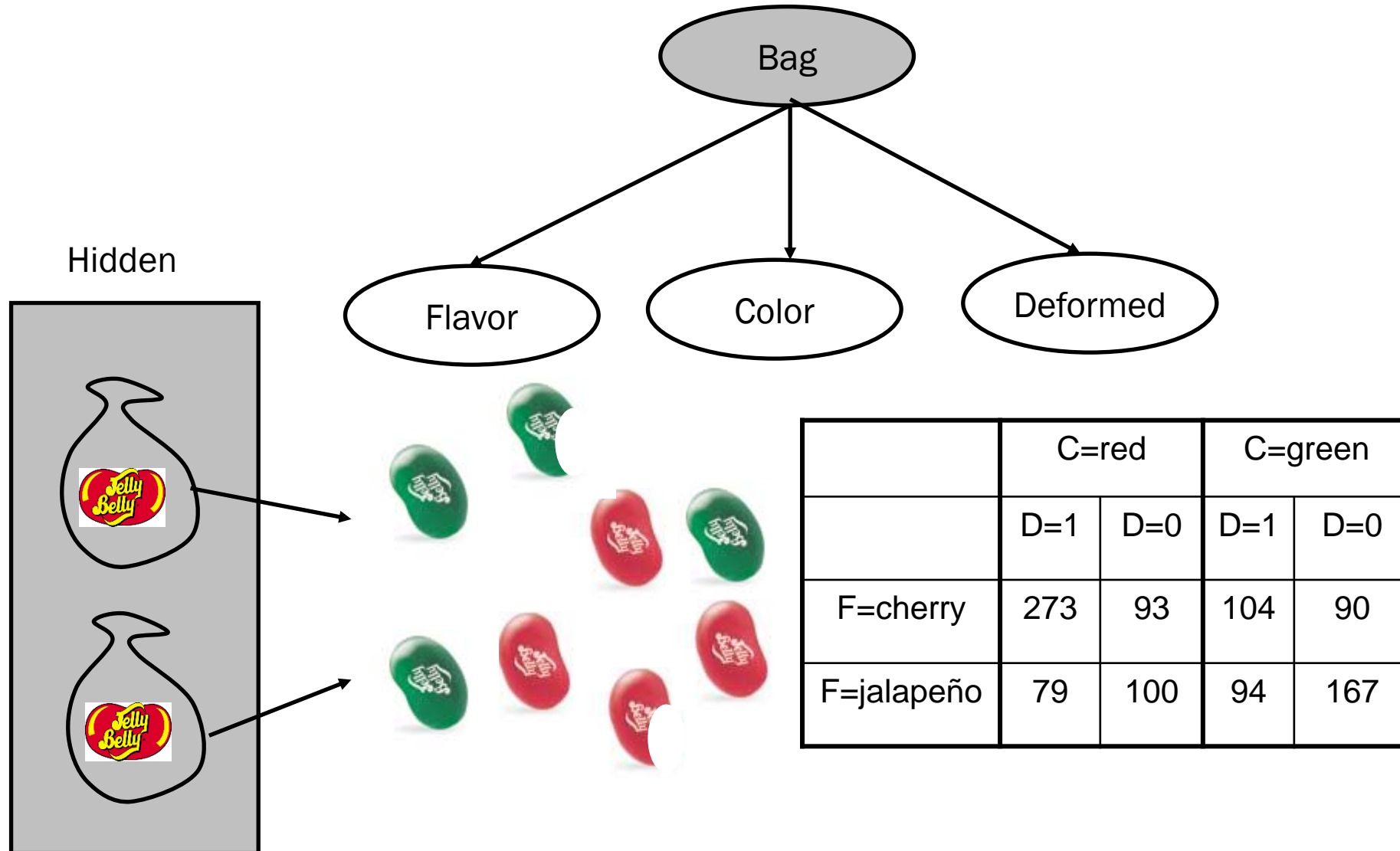# FUDAN SUMMER SCHOOL INTRODUCTION TO AI

Probabilistic Reasoning and Decision Making

Day 9 –Expectation Maximization (EM) Algorithm
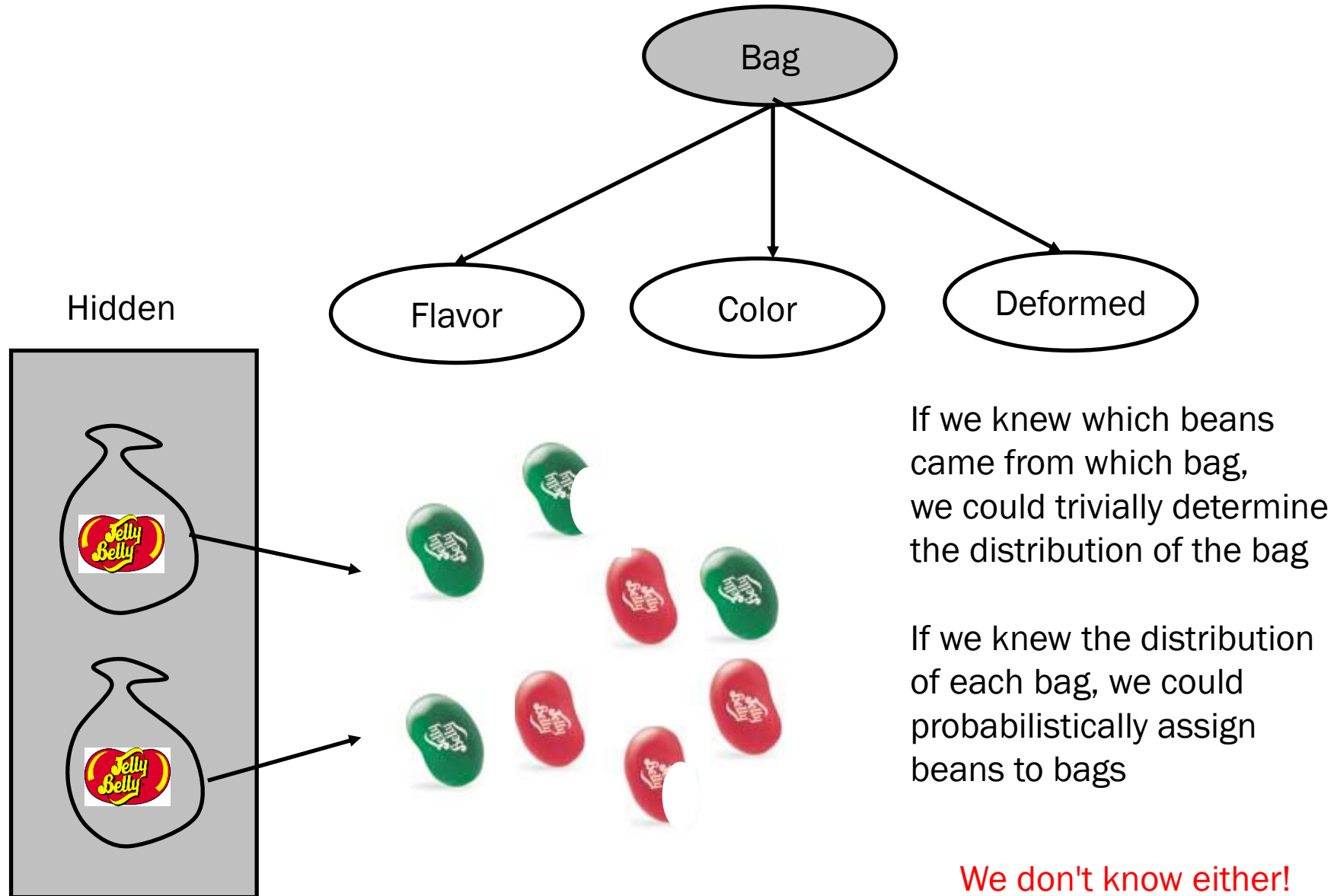
# Learning from Partially Observed Data, Example
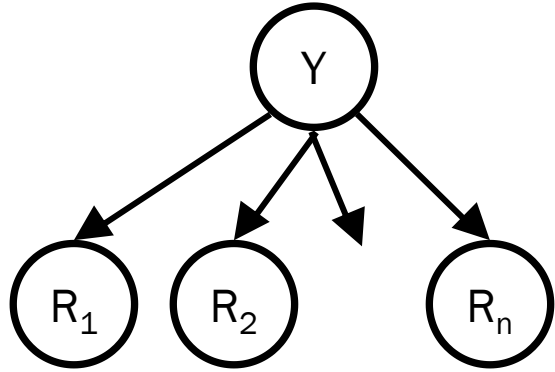
# Learning from Partially Observed Data, Example



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

# Learning from Partially Observed Data, Example



**Bag**

**Hidden**

**Flavor**　**Color**　**Deformed**

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

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

We don't know either!

# Learning with partial data

Another example: Movie Recommender System

# 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

# 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

$$\mathcal{L} = \log \prod_{t=1}^{T} P(V = v^{(t)})$$

$$= \sum_{t=1}^{T} \log P(V = v^{(t)})$$

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

# Learning with partial data, generally

Let $\{X_1, X_2, \ldots, X_n\}$ 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.

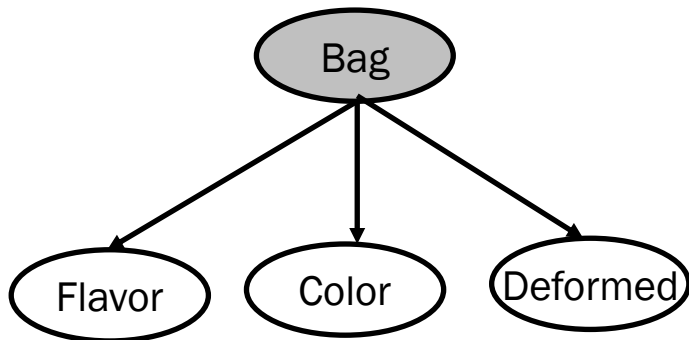$$V \cup H = \{X_1, X_2, \ldots, X_n\}$$

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

$$\mathcal{L} = \log \prod_{t=1}^{T} P\big(V = v^{(t)}\big)$$

$$= \sum_{t=1}^{T} \log P\big(V = v^{(t)}\big) = \sum_{t=1}^{T} \log \sum_{h} P\big(V = v^{(t)}, H = h\big)$$
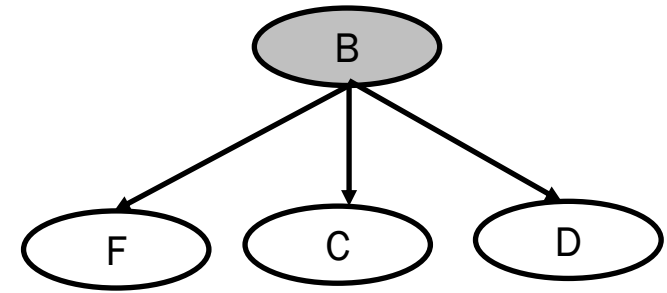
# Expectation Maximization: Central Idea

- 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
4. Lather, rinse, repeat!

# Specifying the problem

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

# EM: Example

Guess (initialize) the parameters of the model
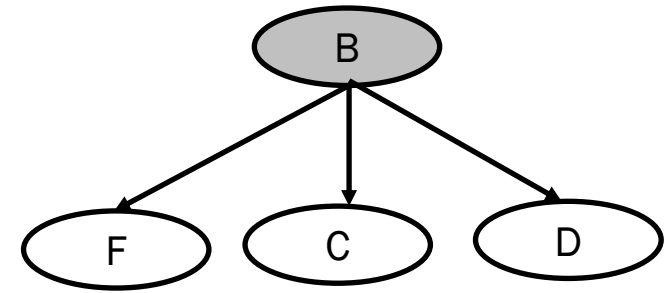Must be valid probabilities, don't have to be right
For simplicity we might choose:

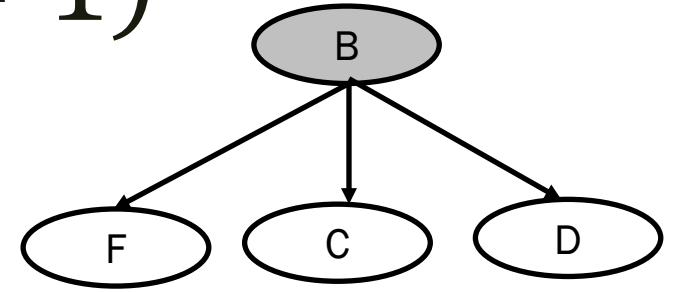$$P(B = 1) = 0.6$$
$$P(F = cherry|B = 1) = P(C = red|B = 1) = P(D = 1|B = 1) = 0.6$$
$$P(F = cherry|B = 2) = P(C = red|B = 2) = P(D = 1|B = 2) = 0.3$$

Initialization

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

What would be the ML estimate if we knew which beans came from which bags?
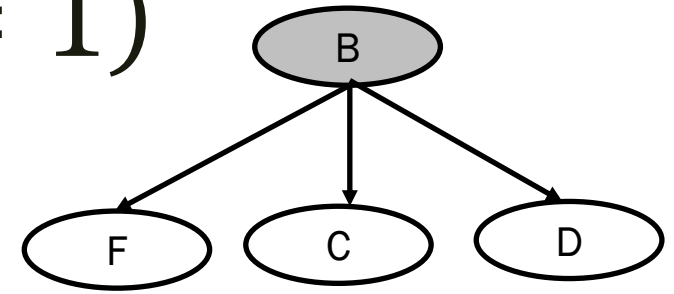
But we don't know this, so we need an estimate for this value.  Use the expectation!

E-step

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



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

E-step

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



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

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

$$= \sum_{t=1}^{T} \frac{P(B = 1)P\big(F = f^{(t)}\big|B = 1\big)P\big(C = c^{(t)}\big|B = 1\big)P\big(D = d^{(t)}\big|B = 1\big)}{\sum_b P(B = b)P\big(F = f^{(t)}\big|B = b\big)P\big(C = c^{(t)}\big|B = b\big)P\big(D = d^{(t)}\big|B = b\big)}$$
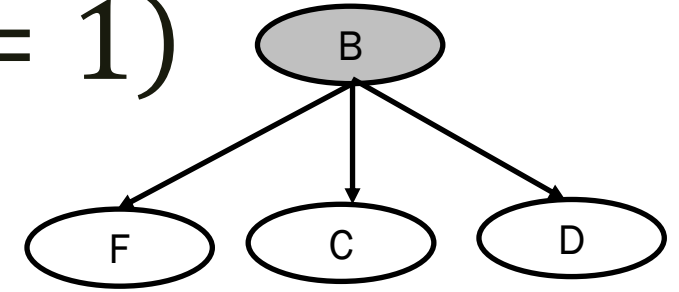
# EM: Example: Updating $P(Bag = 1)$



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

E-step

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

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

M-step

# EM: Initialization



| Param. | Iter 0 | Iter 1 | Iter 2 | ... |
|---|---|---|---|---|
| $P(B = 1)$ | | | | |
| $P(F = cherry|B = 1)$ | | | | |
| $P(F = cherry|B = 2)$ | | | | |
| $P(C = red|B = 1)$ | | | | |
| $P(C = red|B = 2)$ | | | | |
| $P(D = 1|B = 1)$ | | | | |
| $P(D = 1|B = 2)$ | | | | |

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 | | | |
| $P(F = cherry\|B = 1)$ | 0.6 | | | |
| $P(F = cherry\|B = 2)$ | 0.3 | | | |
| $P(C = red\|B = 1)$ | 0.6 | | | |
| $P(C = red\|B = 2)$ | 0.3 | | | |
| $P(D = 1\|B = 1)$ | 0.6 | | | |
| $P(D = 1\|B = 2)$ | 0.3 | | | |

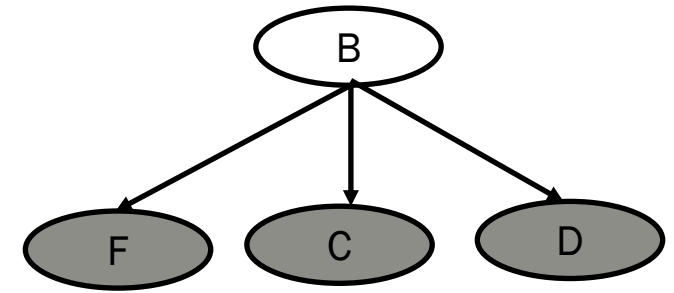EM algorithm (alternate expression):
Initialize the parameters of the model
While parameter values have not converged:
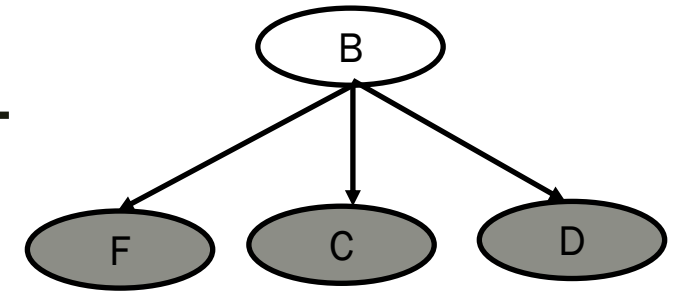   For each parameter, *p*:
   • E-step: Use parameter values from last iteration to calculate expected counts necessary to update *p*
   • M-step: Use those expected counts to update *p*

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

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

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



| Param. | Iter 0 | Iter 1 | Iter 2 | ... |
|---|---|---|---|---|
| $P(B = 1)$ | 0.6 | | | |
| $P(F = cherry \mid B = 1)$ | 0.6 | | | |
| $P(F = cherry \mid B = 2)$ | 0.3 | | | |
| $P(C = red \mid B = 1)$ | 0.6 | | | |
| $P(C = red \mid B = 2)$ | 0.3 | | | |
| $P(D = 1 \mid B = 1)$ | 0.6 | | | |
| $P(D = 1 \mid B = 2)$ | 0.3 | | | |

EM algorithm (alternate expression):
Initialize the parameters of the model
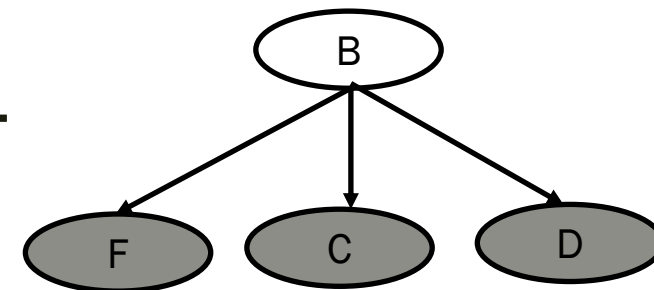While parameter values have not converged:
  For each parameter, *p*:
  - E-step: Use parameter values from last iteration to calculate expected counts necessary to update *p*
  - M-step: Use those expected counts to update *p*

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

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

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

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

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

Current parameter values:
$P(B=1) = 0.6$
$P(F=cherry|B=1) = P(C=red|B=1) = P(D=1|B=1) = 0.6$
$P(F=cherry|B=2) = P(C=red|B=2) = P(D=1|B=2) = 0.3$

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 |

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

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

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

| Param. | Iter 0 |
|---|---|
| $P(B = 1)$ | 0.6 |
| $P(F = cherry|B = 1)$ | 0.6 |
| $P(F = cherry|B = 2)$ | 0.3 |
| $P(C = red|B = 1)$ | 0.6 |
| $P(C = red|B = 2)$ | 0.3 |
| $P(D = 1|B = 1)$ | 0.6 |
| $P(D = 1|B = 2)$ | 0.3 |

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

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

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

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

| Param. | Iter 0 |
|---|---|
| $P(B = 1)$ | 0.6 |
| $P(F = cherry|B = 1)$ | 0.6 |
| $P(F = cherry|B = 2)$ | 0.3 |
| $P(C = red|B = 1)$ | 0.6 |
| $P(C = red|B = 2)$ | 0.3 |
| $P(D = 1|B = 1)$ | 0.6 |
| $P(D = 1|B = 2)$ | 0.3 |

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

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

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

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

| Param. | Iter 0 |
|---|---|
| $P(B = 1)$ | 0.6 |
| $P(F = cherry|B = 1)$ | 0.6 |
| $P(F = cherry|B = 2)$ | 0.3 |
| $P(C = red|B = 1)$ | 0.6 |
| $P(C = red|B = 2)$ | 0.3 |
| $P(D = 1|B = 1)$ | 0.6 |
| $P(D = 1|B = 2)$ | 0.3 |

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

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

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

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

| Param. | Iter 0 |
|---|---|
| $P(B = 1)$ | 0.6 |
| $P(F = cherry|B = 1)$ | 0.6 |
| $P(F = cherry|B = 2)$ | 0.3 |
| $P(C = red|B = 1)$ | 0.6 |
| $P(C = red|B = 2)$ | 0.3 |
| $P(D = 1|B = 1)$ | 0.6 |
| $P(D = 1|B = 2)$ | 0.3 |

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

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

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

| Param. | Iter 0 |
|---|---|
| $P(B = 1)$ | 0.6 |
| $P(F = cherry|B = 1)$ | 0.6 |
| $P(F = cherry|B = 2)$ | 0.3 |
| $P(C = red|B = 1)$ | 0.6 |
| $P(C = red|B = 2)$ | 0.3 |
| $P(D = 1|B = 1)$ | 0.6 |
| $P(D = 1|B = 2)$ | 0.3 |

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

| 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 | |

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

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

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

| Param. | Iter 0 |
|---|---|
| $P(B = 1)$ | 0.6 |
| $P(F = cherry|B = 1)$ | 0.6 |
| $P(F = cherry|B = 2)$ | 0.3 |
| $P(C = red|B = 1)$ | 0.6 |
| $P(C = red|B = 2)$ | 0.3 |
| $P(D = 1|B = 1)$ | 0.6 |
| $P(D = 1|B = 2)$ | 0.3 |

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

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

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

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

| Param. | Iter 0 |
|---|---|
| $P(B = 1)$ | 0.6 |
| $P(F = cherry|B = 1)$ | 0.6 |
| $P(F = cherry|B = 2)$ | 0.3 |
| $P(C = red|B = 1)$ | 0.6 |
| $P(C = red|B = 2)$ | 0.3 |
| $P(D = 1|B = 1)$ | 0.6 |
| $P(D = 1|B = 2)$ | 0.3 |

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

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

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

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

| Param. | Iter 0 |
|---|---|
| $P(B = 1)$ | 0.6 |
| $P(F = cherry|B = 1)$ | 0.6 |
| $P(F = cherry|B = 2)$ | 0.3 |
| $P(C = red|B = 1)$ | 0.6 |
| $P(C = red|B = 2)$ | 0.3 |
| $P(D = 1|B = 1)$ | 0.6 |
| $P(D = 1|B = 2)$ | 0.3 |

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

# What about the conditionals?

How do we calculate $P(F = cherry | B = 1)$

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

What is the estimate for count(*F=cherry, B=1*)?

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

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

E. I have no idea

# What about the conditionals?

How do we calculate $P(F = cherry | B = 1)$

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

# EM: General update formula for BN params

Nodes with parents:

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

Root nodes:

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

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

$X \rightarrow Y \rightarrow Z$

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

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

All variables are binary



Complete the formula below:

$$P(Y = 1 | x_1, x_2, \ldots, x_n) =$$

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

# EM for Learning Noisy-OR: Alternate model

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



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

# EM for Learning Noisy-OR: Alternate model

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

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

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

<u>What is $P(Y = 1 | x_1, x_2, \ldots, x_n)$ in this alternate model?</u>

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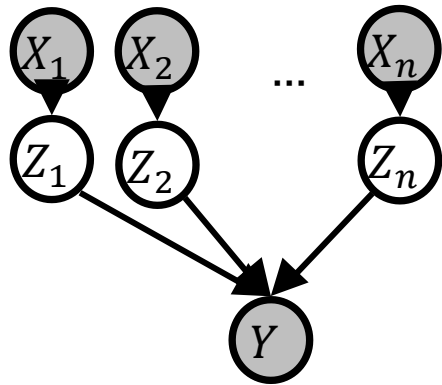


$P(Y | \vec{z}) = OR(\vec{z})$
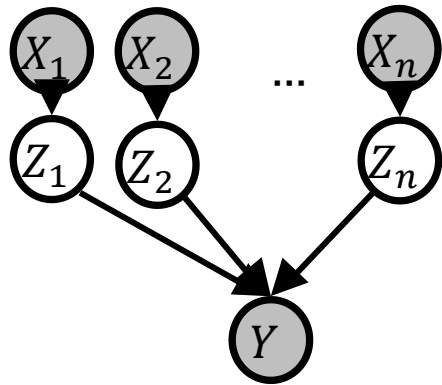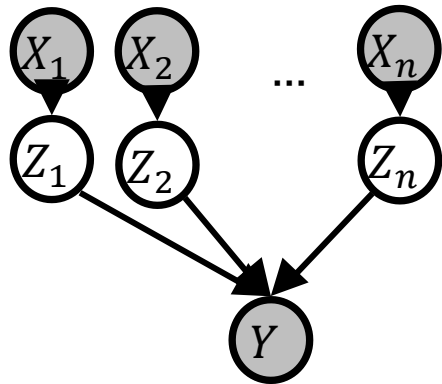
# EM for Learning Noisy-OR: Alternate model

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

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

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



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

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

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

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 $\left\{\left(x_1^{(t)}, x_n^{(t)}, \ldots, x_n^{(t)}, y^{(t)}\right)\right\}_{t=1}^{T}$, 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 $P(Y = 1 | x_1, x_2, \ldots, x_n)$ in this alternate model?

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



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

# EM for Learning Noisy-OR: Alternate model

Problem: From (complete) data $\left\{\left(x_1^{(t)}, x_n^{(t)}, \ldots, x_n^{(t)}, y^{(t)}\right)\right\}_{t=1}^{T}$, 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 $P(Y = 1 | x_1, x_2, \ldots, x_n)$ in this alternate model?

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



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

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
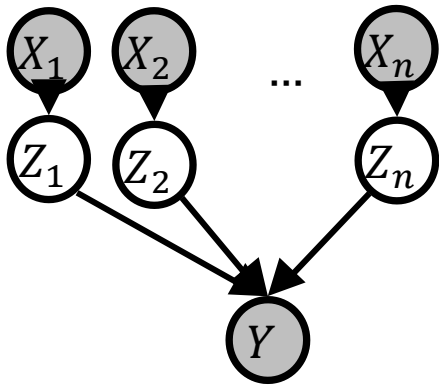
# EM for Learning Noisy-OR: Alternate model

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

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

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



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

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

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

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

- If we had fully observed data:
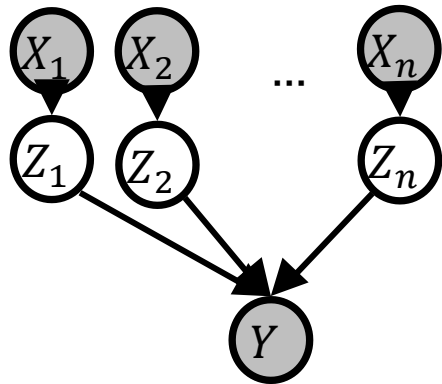
# Learning $P(Z_i = 1 | X_i = 1) = p_i$ using EM

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

# 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

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

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

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



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

$$P(Z_i = 1 | \vec{x}^{(t)}, y^{(t)}) = \frac{P(Y = y^{(t)} | Z_i = 1) P(Z_i = 1 | \vec{x}^{(t)})}{P(Y = y^{(t)} | \vec{x}^{(t)})}$$
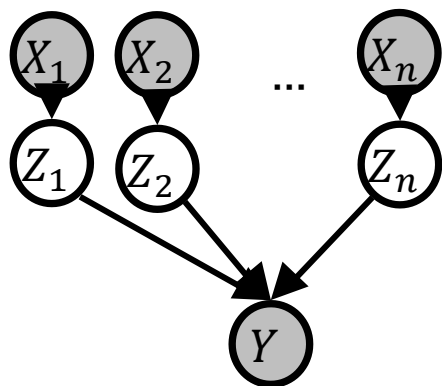
Which of the following is NOT a true statement?

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

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

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

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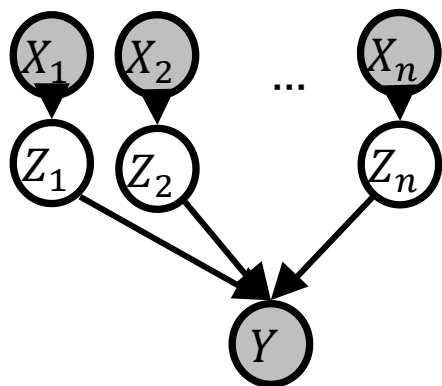


$P(Z_i = 1 | \vec{x}^{(t)}, y^{(t)}) = \dfrac{P(Y = y^{(t)} | Z_i = 1) P(Z_i = 1 | \vec{x}^{(t)})}{P(Y = y^{(t)} | \vec{x}^{(t)})}$

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

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

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

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



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

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

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