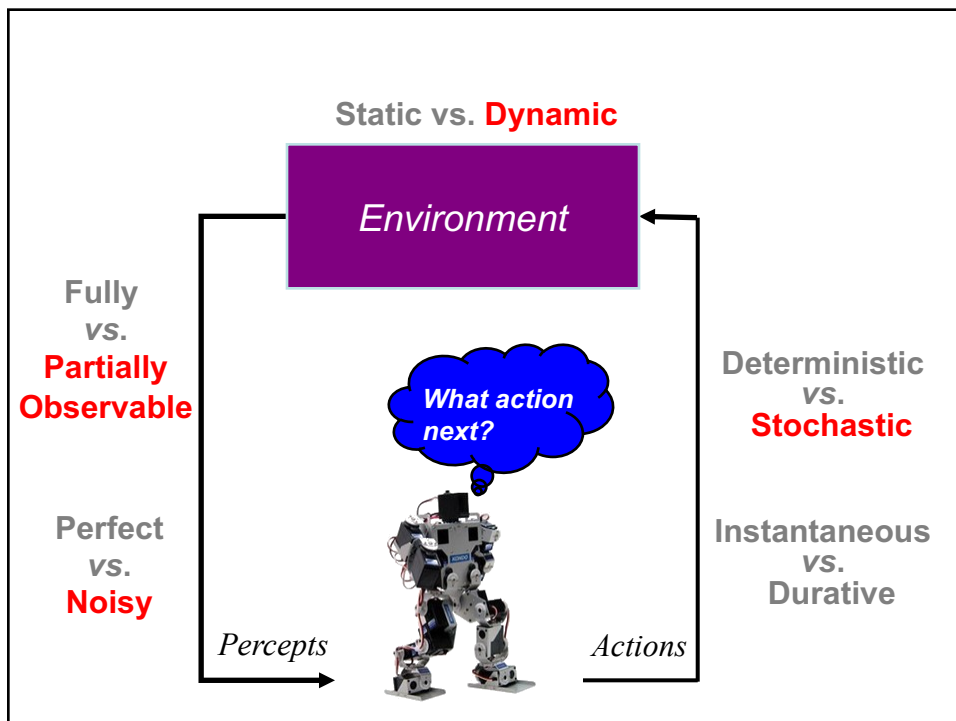


# CSE 473: Artificial Intelligence

## Bayesian Networks - Learning

Dan Weld

Slides adapted from Jack Breese, Dan Klein, Daphne Koller, Stuart Russell, Andrew Moore & Luke Zettlemoyer



## AI Topics

---

- **Search**
  - Problem Spaces
  - BFS, DFS, UCS, A\* (tree and graph)
  - Completeness and Optimality
  - Heuristics: admissibility and consistency
- **CSPs**
  - Constraint graphs, backtracking search
  - Forward checking, AC3 constraint propagation, ordering heuristics
- **Games**
  - Minimax, Alpha-beta pruning, Expectimax, Evaluation Functions
- **MDPs**
  - Bellman equations
  - Value iteration & policy iteration
  - RTDP, LAO\* & UCT
  - POMDPs
- **Reinforcement Learning**
  - Exploration vs. Exploitation
  - Model-based vs. model-free
  - Q-learning
  - Linear value function approx.
- **Hidden Markov Models**
  - Markov chains
  - Forward algorithm
  - Particle Filter
- **Bayesian Networks**
  - Basic definition, independence (d-sep)
  - Variable elimination
  - Gibbs sampling
- **Learning**
  - BN parameters with data complete & incomplete (Expectation Maximization)
  - Search thru space of BN structures

## Search thru a Problem Space / State Space

Ex. Proving a trig identity, e.g.  $\sin^2(x) = \frac{1}{2} - \frac{1}{2} \cos(2x)$

### • Input:

- Set of states
- Operators [and costs]
- Start state
- Goal state [test]

### • Output:

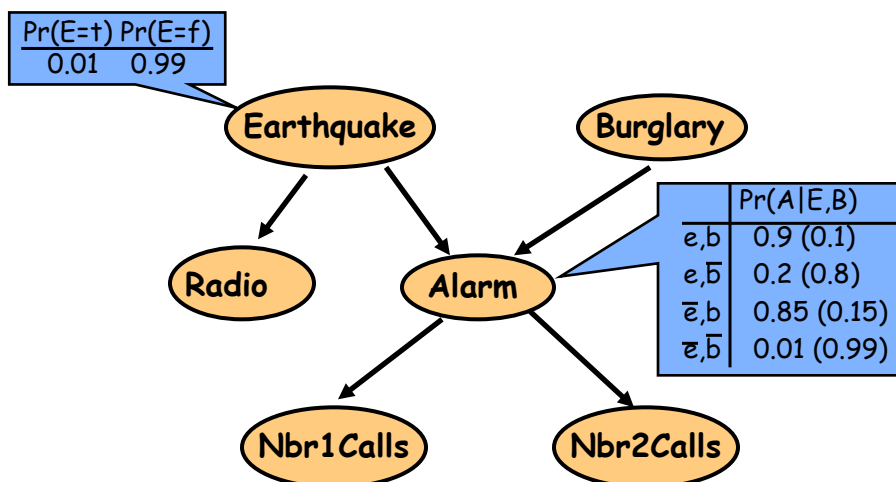
- Path: start  $\Rightarrow$  a state satisfying goal test
- [May require shortest path]
- [Sometimes just need state passing test]

## Today

- Bonus Topic – Hybrid Bayes Nets
- Learning
  - Parameter Learning & Priors
  - Expectation Maximization
  - Structure Learning

5

## Bayes Nets



© Daniel S. Weld

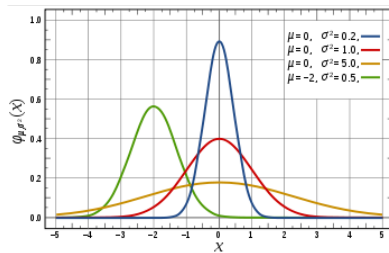
6

## Continuous Variables

Pr(E=t)	Pr(E=f)
0.01	0.99

Earthquake

So far: assuming variables have discrete values  
 Could also allow continuous values,  $E \in \mathbb{R}$ ,  
 How specify probabilities? (explicit CPT would be infinitely large)



$$P(x | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

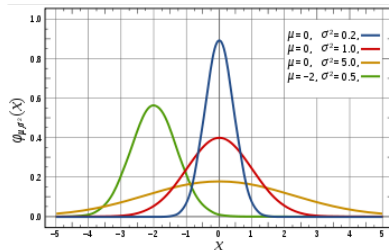
© Daniel S. Weld

## Continuous Variables

Pr(E=t)	Pr(E=f)
0.01	0.99

Earthquake

So far: assuming variables have discrete values  
 Could also allow continuous values,  $E \in \mathbb{R}$ ,  
 And specify probabilities using a continuous distribution, such as a Gaussian



$$P(x | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

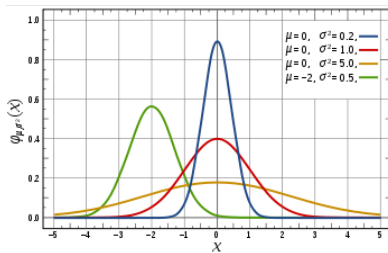
© Daniel S. Weld

# Continuous Variables

Earthquake

$\Pr(E=x)$   
 mean:  $\mu = 6$   
 variance:  $\sigma = 2$

So far: assuming variables have discrete values  
 Could also allow continuous values,  $E \in \mathbb{R}$ ,  
 And specify probabilities using a continuous distribution, such as a Gaussian



$$P(x | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

© Daniel S. Weld

# Continuous Variables

$\Pr(A=t) \Pr(A=f)$   
 0.01 0.99

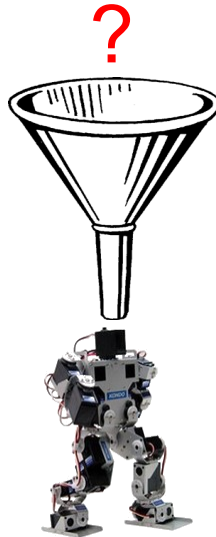
Aliens

Earthquake

	$\Pr(E A)$
$a$	$\mu = 6$ $\sigma = 2$
$\bar{a}$	$\mu = 1$ $\sigma = 3$

© Daniel S. Weld

# Learning



## Supremacy of Machine Learning

- **Machine learning is preferred approach to**
  - Speech recognition, Natural language processing
  - Web search – result ranking
  - Computer vision
  - Medical outcomes analysis
  - Robot control
  - Computational biology
  - Sensor networks
  - ...
- **This trend is accelerating**
  - Improved machine learning algorithms
  - Improved data capture, networking, faster computers
  - Software too complex to write by hand
  - New sensors / IO devices
  - Demand for self-customization to user, environment

## What is Machine Learning ?



## Machine Learning

Study of algorithms that

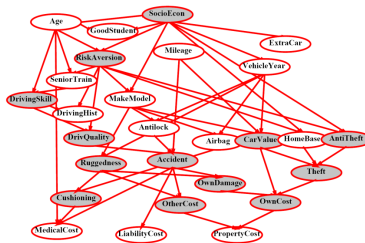
- improve their performance *Ability to accumulate reward*
- at some task *Executing actions*
- with experience *Executing actions*

*?? Reinforcement Learning ??*

# Machine Learning

Study of algorithms that

- improve their performance *Prediction accuracy*
- at some task *Answering probabilistic queries*
- with experience *Seeing labeled data*



id	version	id	name	url	description	price	year	make	model	color	mpg
1	10001	1	Ford Taurus	10001	10001	12742	1995	Ford	Taurus	Blue	24
2	10002	1	Ford Taurus	10002	10002	12742	1995	Ford	Taurus	Blue	24
3	10003	1	Ford Taurus	10003	10003	12742	1995	Ford	Taurus	Blue	24
4	10004	1	Ford Taurus	10004	10004	12742	1995	Ford	Taurus	Blue	24
5	10005	1	Ford Taurus	10005	10005	12742	1995	Ford	Taurus	Blue	24
6	10006	1	Ford Taurus	10006	10006	12742	1995	Ford	Taurus	Blue	24
7	10007	1	Ford Taurus	10007	10007	12742	1995	Ford	Taurus	Blue	24
8	10008	1	Ford Taurus	10008	10008	12742	1995	Ford	Taurus	Blue	24
9	10009	1	Ford Taurus	10009	10009	12742	1995	Ford	Taurus	Blue	24
10	10010	1	Ford Taurus	10010	10010	12742	1995	Ford	Taurus	Blue	24
11	10011	1	Ford Taurus	10011	10011	12742	1995	Ford	Taurus	Blue	24
12	10012	1	Ford Taurus	10012	10012	12742	1995	Ford	Taurus	Blue	24
13	10013	1	Ford Taurus	10013	10013	12742	1995	Ford	Taurus	Blue	24
14	10014	1	Ford Taurus	10014	10014	12742	1995	Ford	Taurus	Blue	24
15	10015	1	Ford Taurus	10015	10015	12742	1995	Ford	Taurus	Blue	24

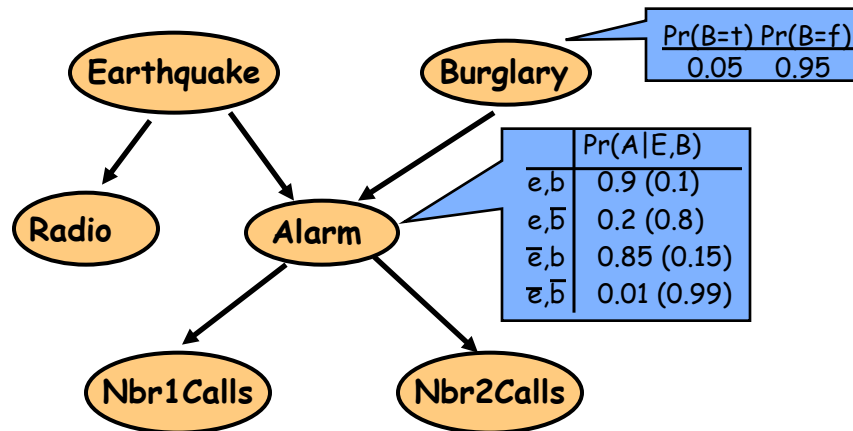
20

# Learning Bayes Networks

- Learning Parameters for a Bayesian Network
  - Fully observable variables
    - Maximum Likelihood (ML), MAP & Bayesian estimation
    - Example: Naïve Bayes for text classification
  - Hidden variables
    - Expectation Maximization (EM)
- Learning Structure of Bayesian Networks



## The Origin of Bayes Nets



© Daniel S. Weld

22

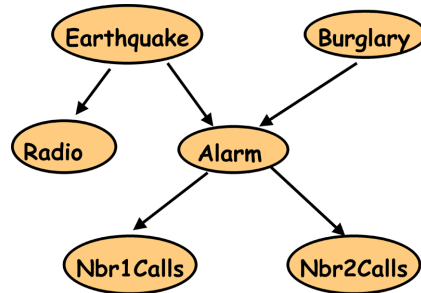
## Learning Bayes Nets

Suppose ...

1. Know structure & get complete observations of every var
2. Know structure & get observations only of **some** vars  
Others are hidden (learn with EM)
3. Don't even know structure!

© Daniel S. Weld

## Parameter Estimation and Bayesian Networks

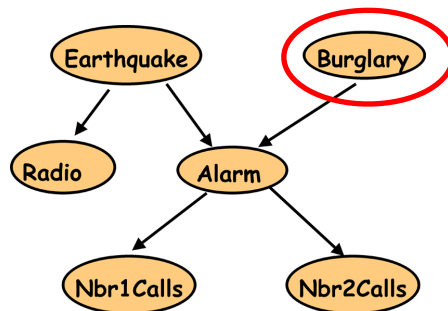


E	B	R	A	J	M
T	F	T	T	F	T
F	F	F	F	F	T
F	T	F	T	T	T
F	F	F	T	T	T
F	T	F	F	F	F
...					

We have:

- Bayes Net **structure** and **observations**
- We need: Bayes Net **parameters**

## Parameter Estimation and Bayesian Networks

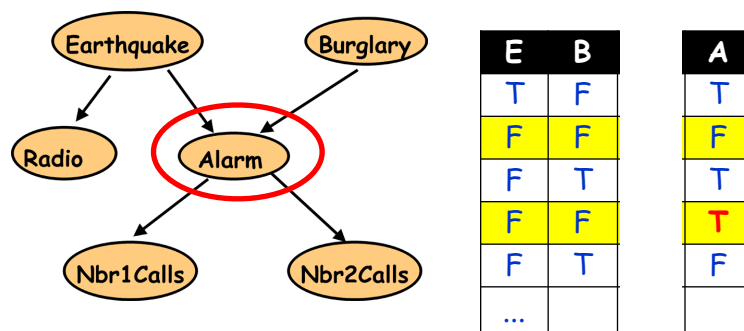


B
F
F
T
F
T

$$P(B) = ? \quad = 0.4$$

$$P(\neg B) = 1 - P(B) = 0.6$$

## Parameter Estimation and Bayesian Networks



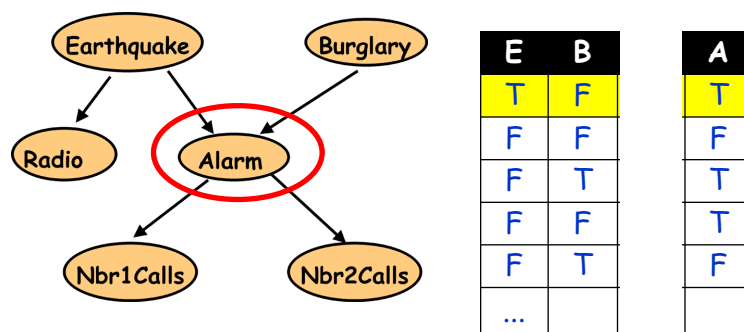
$$P(A|E,B) = ?$$

$$P(A|E,\neg B) = ?$$

$$P(A|\neg E,B) = ?$$

$$P(A|\neg E,\neg B) = 0.5$$

## Parameter Estimation and Bayesian Networks



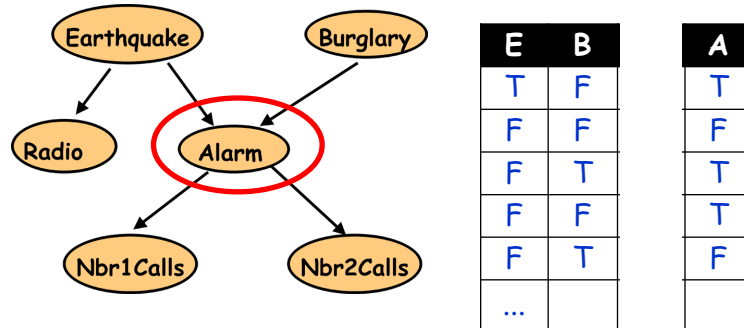
$$P(A|E,B) = ?$$

$$P(A|E,\neg B) = 1.0 ?$$

$$P(A|\neg E,B) = ?$$

$$P(A|\neg E,\neg B) = ?$$

## Parameter Estimation and Bayesian Networks



$$P(A|E,B) = ?$$

$$P(A|E,\neg B) = ?$$

$$P(A|\neg E,B) = ?$$


$$P(A|\neg E,\neg B) = ?$$

## Parameter Estimation and Bayesian Networks

Coin


**Coin Flip**

$C_1$




$P(H|C_1) = 0.1$

$C_2$



$P(H|C_2) = 0.5$

$C_3$



$P(H|C_3) = 0.9$

**Which coin will I use?**

$P(C_1) = 1/3$


$P(C_2) = 1/3$

$P(C_3) = 1/3$

**Prior:** Probability of a hypothesis before we make any observations


**Coin Flip**

$C_1$




$P(H|C_1) = 0.1$

$C_2$



$P(H|C_2) = 0.5$

$C_3$



$P(H|C_3) = 0.9$

**Which coin will I use?**

$P(C_1) = 1/3$

$P(C_2) = 1/3$

$P(C_3) = 1/3$

**Uniform Prior:** All hypothesis are equally likely before we make any observations

## Experiment 1: Heads

Which coin did I use?

$$P(C_1|H) = ? \quad P(C_2|H) = ? \quad P(C_3|H) = ?$$

$$P(C_1|H) = \frac{P(H|C_1)P(C_1)}{P(H)} \quad P(H) = \sum_{i=1}^3 P(H|C_i)P(C_i)$$

$C_1$



$$P(H|C_1) = 0.1$$

$$P(C_1) = 1/3$$

$C_2$



$$P(H|C_2) = 0.5$$

$$P(C_2) = 1/3$$

$C_3$



$$P(H|C_3) = 0.9$$

$$P(C_3) = 1/3$$

## Experiment 1: Heads

Which coin did I use?

$$P(C_1|H) = 0.066 \quad P(C_2|H) = 0.333 \quad P(C_3|H) = 0.6$$

**Posterior:** Probability of a hypothesis given data

$C_1$



$$P(H|C_1) = 0.1$$

$$P(C_1) = 1/3$$

$C_2$



$$P(H|C_2) = 0.5$$

$$P(C_2) = 1/3$$

$C_3$



$$P(H|C_3) = 0.9$$

$$P(C_3) = 1/3$$

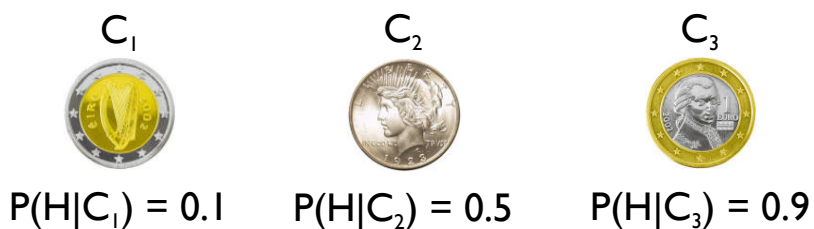
## Using Prior Knowledge

- Should we always use a **Uniform Prior** ?
- Background knowledge:

Heads => we have to buy Dan chocolate

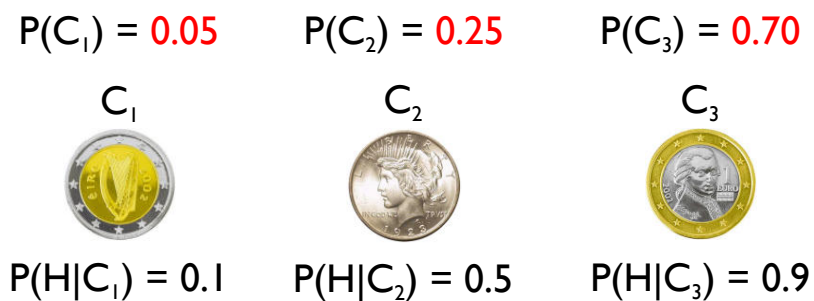
Dan **likes** chocolate...

=> Dan is more likely to use a coin biased in his favor



## Using Background Knowledge

We can encode it in the **prior**:



## Experiment 1: Heads

### Which coin did I use?

$$P(C_1|H) = 0.006 \quad P(C_2|H) = 0.165 \quad P(C_3|H) = 0.829$$

Compare with ML posterior after Exp 1:

$$P(C_1|H) = 0.066 \quad P(C_2|H) = 0.333 \quad P(C_3|H) = 0.600$$



$$P(H|C_1) = 0.1$$

$$P(C_1) = 0.05$$



$$P(H|C_2) = 0.5$$

$$P(C_2) = 0.25$$



$$P(H|C_3) = 0.9$$

$$P(C_3) = 0.70$$

## Probabilistic Estimation

Easy to compute

Maximum Likelihood Estimate (MLE)

Maximum A Posteriori Estimate (MAP)

Bayesian Estimate

Prior

Hypothesis

Prior	Hypothesis
Uniform	The most likely
Any	The most likely
Any	Weighted combination

Still easy to compute  
Incorporates prior knowledge

Minimizes error  
Great when data is scarce  
Potentially much harder to compute





## Bayesian Learning

Use Bayes rule:

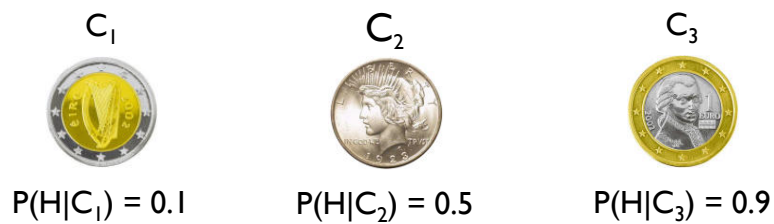
$$P(Y | X) = \frac{P(X | Y) P(Y)}{P(X)}$$

Diagram illustrating the components of Bayes' rule:

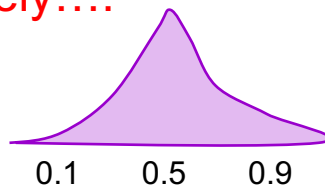
- Data Likelihood** points to  $P(X | Y)$ .
- Prior** points to  $P(Y)$ .
- Normalization** points to  $P(X)$ .
- Posterior** points to  $P(Y | X)$ .

Or equivalently:  $P(Y | X) \propto P(X | Y) P(Y)$

## Really? Only 3 Coins?



More Likely....

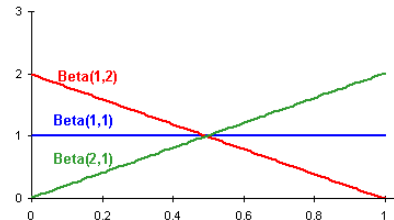


## What Prior to Use?

- Two common priors for *continuous variables*

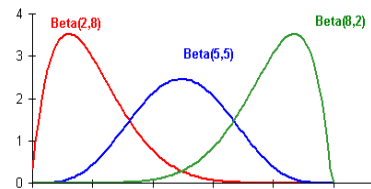
- Binary variable Beta**

- Posterior distribution is binomial
- Easy to compute posterior
- Easy to compute MAP estimate
  - MAP  $E[\text{Beta}(a, b)] = a/(a+b)$



- Discrete variable Dirichlet**

- Posterior distribution is multinomial
- Easy to compute posterior



© Danjel S. Weld  
41

## Estimation: Laplace Smoothing

- Laplace's estimate:

pretend you saw every outcome  
once more than you actually did



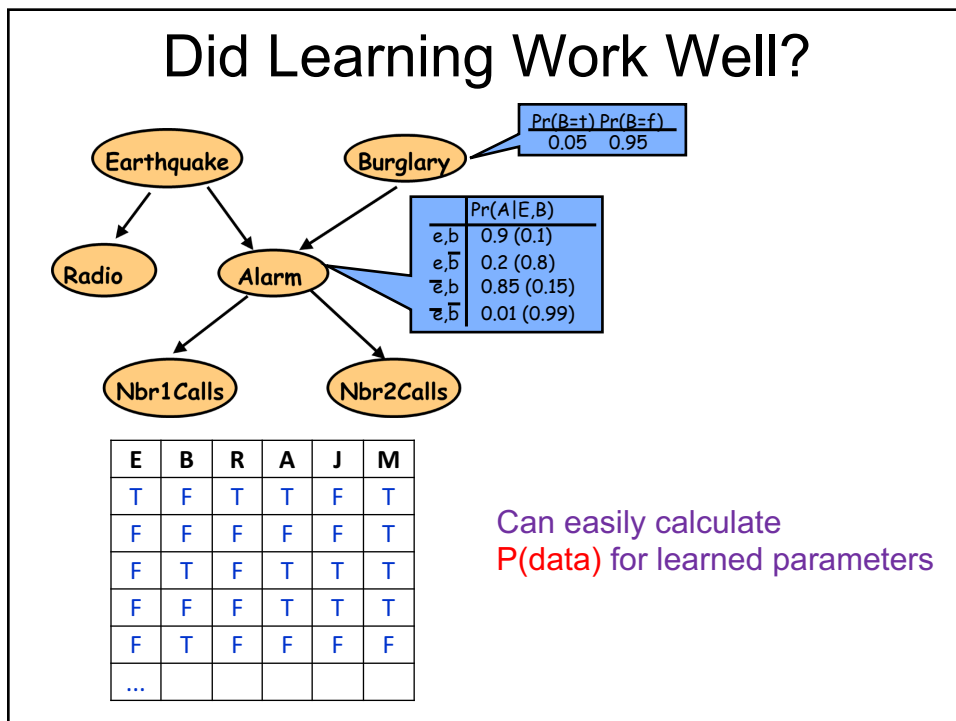
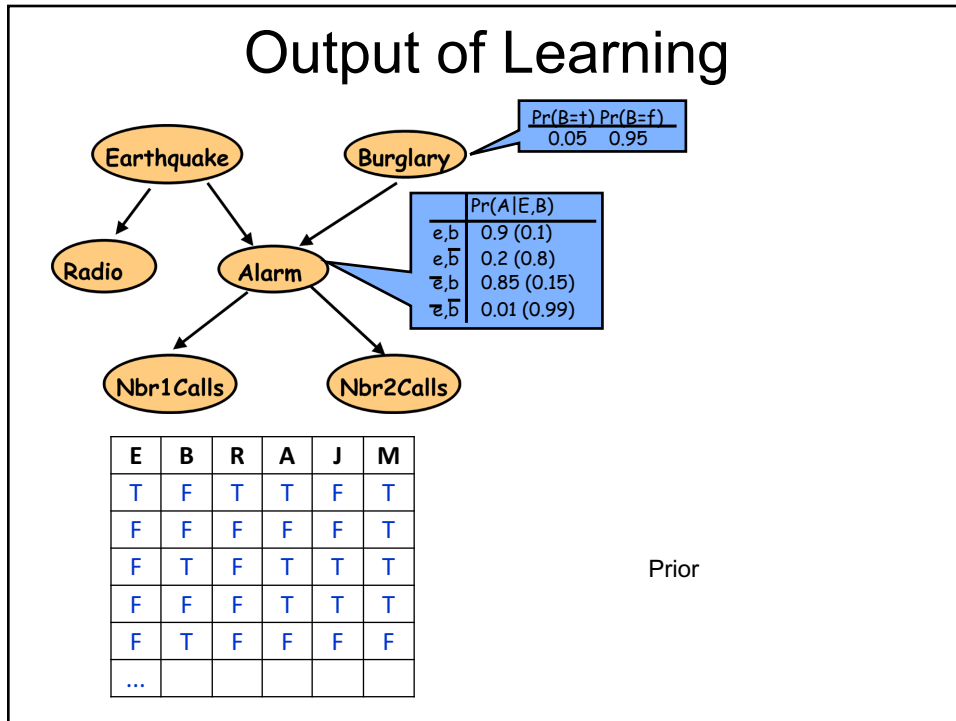
$$P_{LAP}(x) = \frac{c(x) + 1}{\sum_x [c(x) + 1]}$$

$$= \frac{c(x) + 1}{N + |X|}$$

$$P_{LAP}(H) = (2+1) / (3+2)$$

$$= 3/5$$

Another name for computing the MAP estimate with *Dirichlet priors*  
(Bayesian justification)

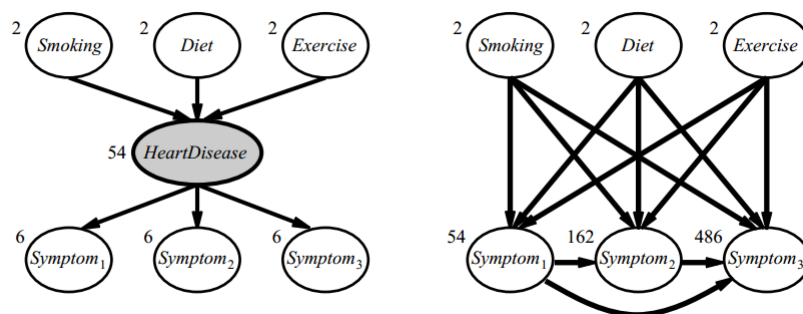


## Topics

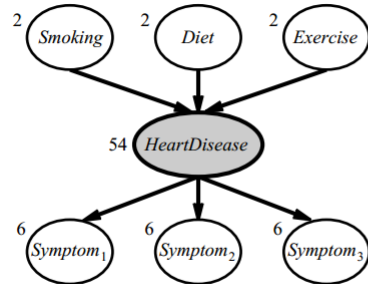
- Another Useful Bayes Net
  - Hybrid Discrete / Continuous
- Learning Parameters for a Bayesian Network
  - Fully observable
  - Hidden variables (EM algorithm)
- Learning Structure of Bayesian Networks

© Daniel S. Weld

## Why Learn Hidden Variables?



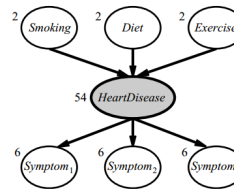
## How Learn Hidden Variables?



## Chicken & Egg Problem

- If we knew whether patient had disease

- It would be easy to learn CPTs
- But we can't observe states, so we don't!



- If we knew CPTs

- It would be easy to predict if patient had disease
- But we don't, so we can't!

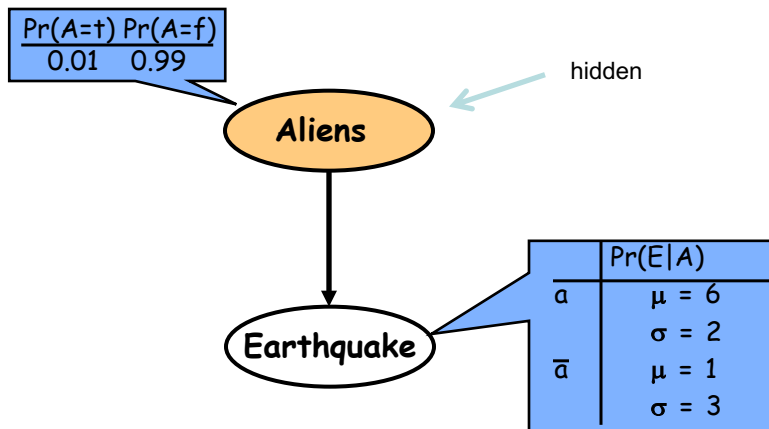
Face It...



58

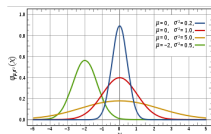


## Continuous Variables



© Daniel S. Weld

## Learning with Continuous Variables



Earthquake

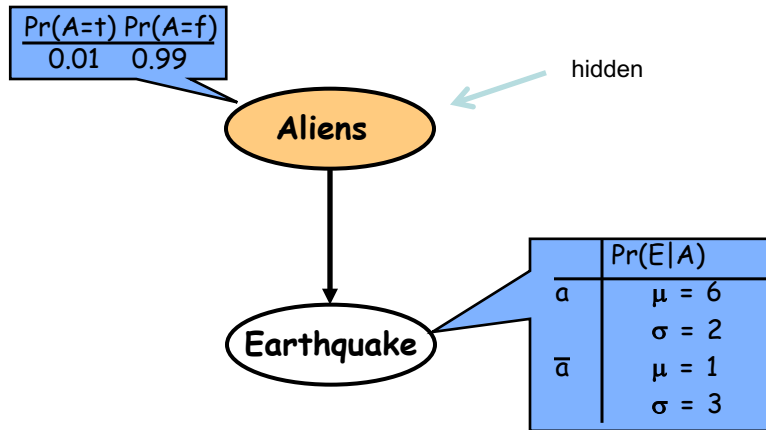
Pr(E=x)
mean: $\mu = ?$
variance: $\sigma = ?$

$$\hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2$$

© Daniel S. Weld

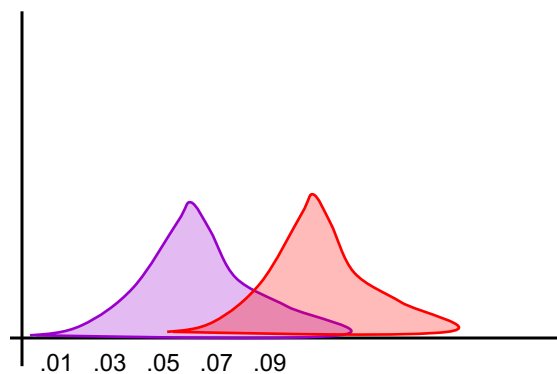
## Continuous Variables



© Daniel S. Weld

## Simplest Version

- Mixture of two distributions



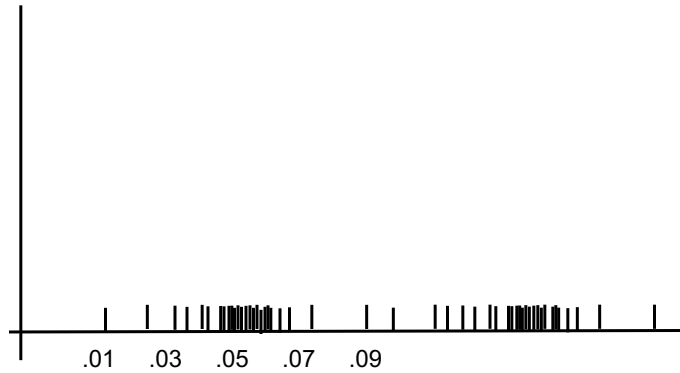
- Know: form of distribution & variance,  $\sigma = .5$
- Just need *mean* of each distribution

Slide by Daniel S. Weld

63



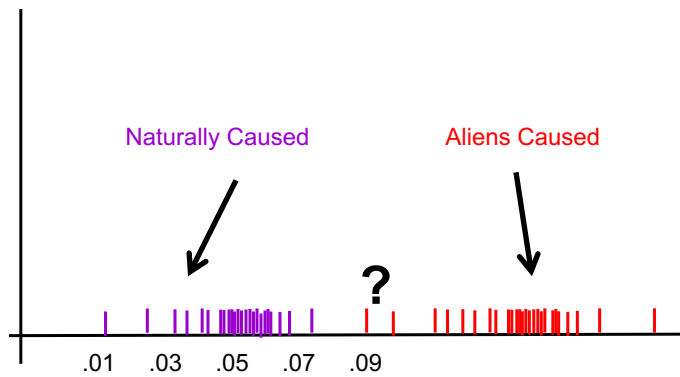
# Input Looks Like



Slide by Daniel S. Weld

64

# We Want to Predict

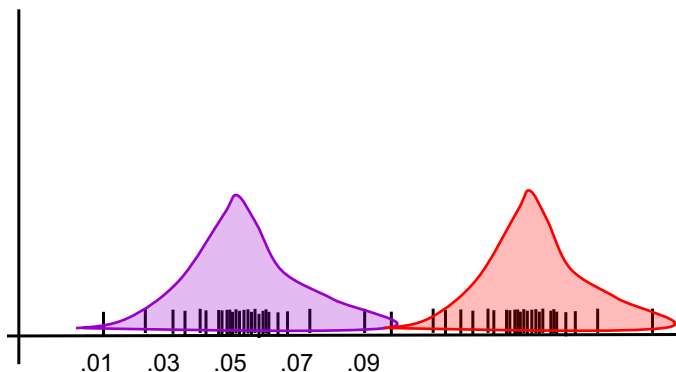


Slide by Daniel S. Weld

65

## Chicken & Egg

Note that coloring instances would be easy  
*if* we knew Gaussians....

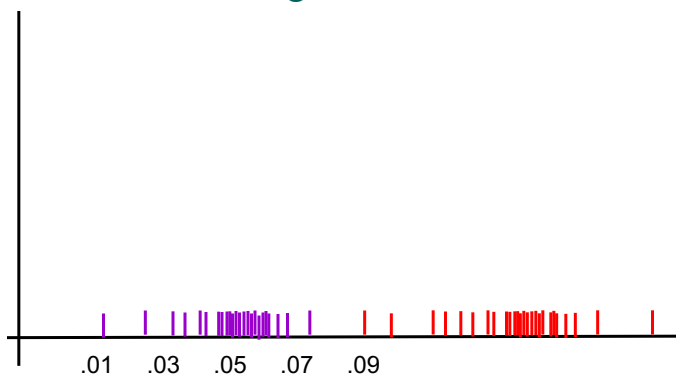


Slide by Daniel S. Weld

66

## Chicken & Egg

And finding Gaussian parameters would be easy  
*if* we knew the coloring

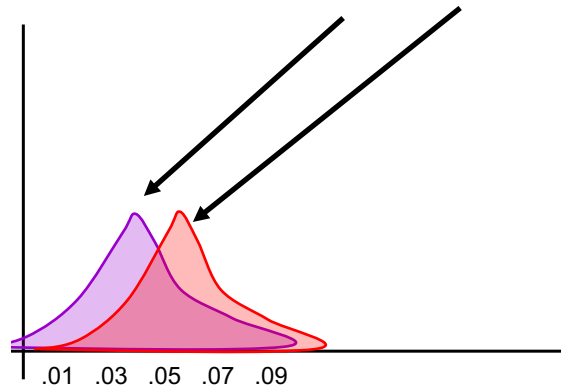


Slide by Daniel S. Weld

67

## Expectation Maximization (EM)

- Pretend we *do* know the parameters
  - Initialize randomly: set  $\theta_1=?$ ;  $\theta_2=?$

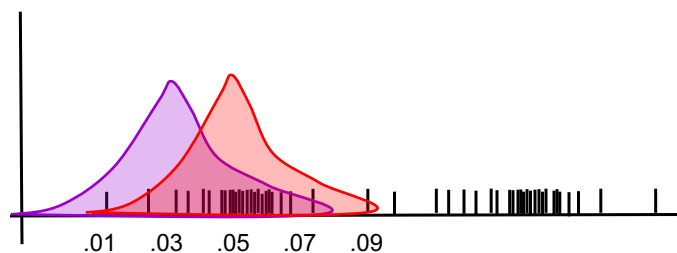


Slide by Daniel S. Weld

68

## Expectation Maximization (EM)

- Pretend we *do* know the parameters
  - Initialize randomly
- [E step] Compute probability of instance having each possible value of the hidden variable

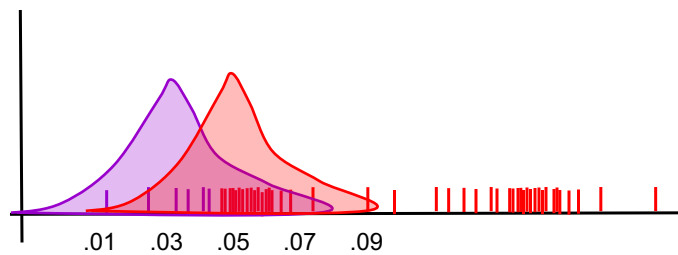


Slide by Daniel S. Weld

69

## Expectation Maximization (EM)

- Pretend we *do* know the parameters
  - Initialize randomly
- [E step] Compute probability of instance having each possible value of the hidden variable



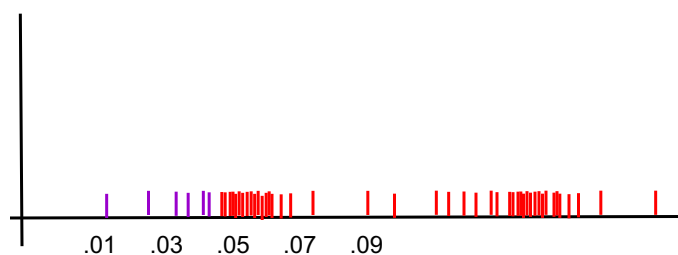
Slide by Daniel S. Weld

70

## Expectation Maximization (EM)

- Pretend we *do* know the parameters
  - Initialize randomly
- [E step] Compute probability of instance having each possible value of the hidden variable

[M step] Treating each instance as *fractionally* having **both** values compute the new parameter values

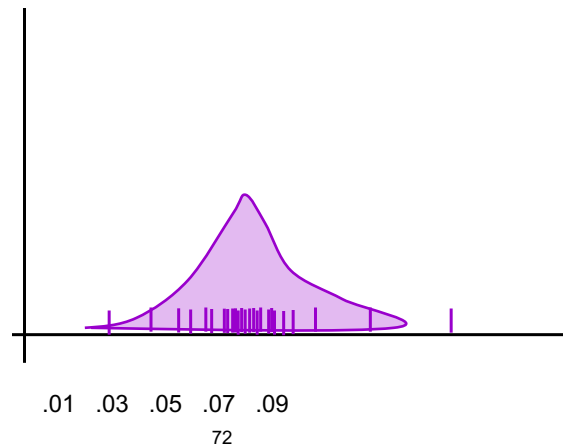


Slide by Daniel S. Weld

71

## ML Mean of Single Gaussian

$$U_{ml} = \operatorname{argmin}_u \sum_i (x_i - u)^2$$

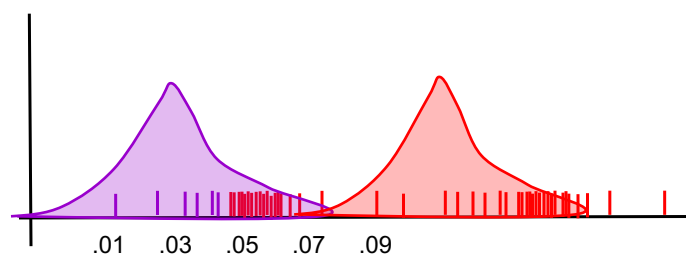


Slide by Daniel S. Weld

72

## Expectation Maximization (EM)

■  
**[M step]** Treating each instance as fractionally having **both** values compute the new parameter values

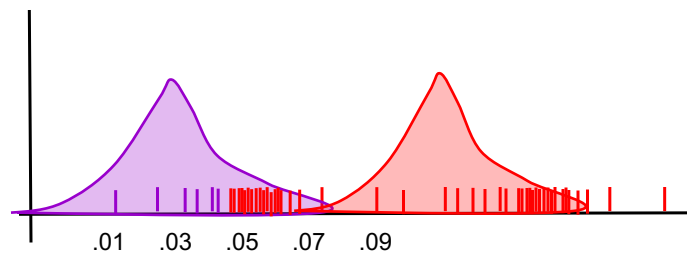


Slide by Daniel S. Weld

73

## Expectation Maximization (EM)

- **[E step]** Compute probability of instance having each possible value of the hidden variable



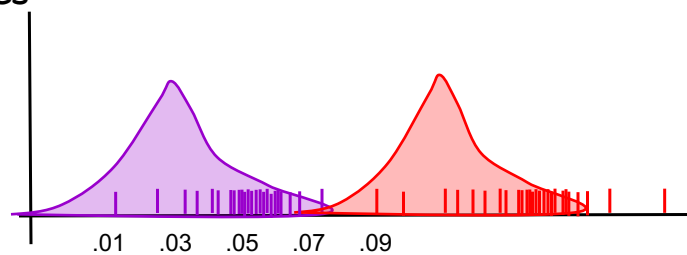
Slide by Daniel S. Weld

74

## Expectation Maximization (EM)

- **[E step]** Compute probability of instance having each possible value of the hidden variable

**[M step]** Treating each instance as fractionally having both values compute the new parameter values



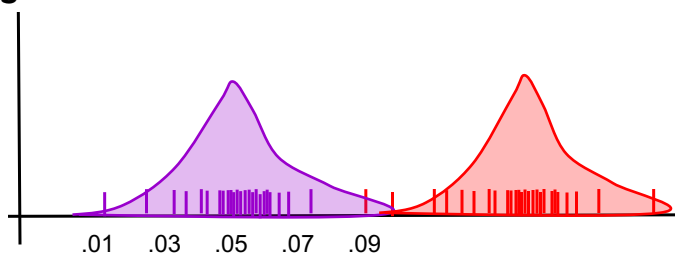
Slide by Daniel S. Weld

75

## Expectation Maximization (EM)

- **[E step]** Compute probability of instance having each possible value of the hidden variable

**[M step]** Treating each instance as fractionally having both values compute the new parameter values



Slide by Daniel S. Weld

76

## Topics

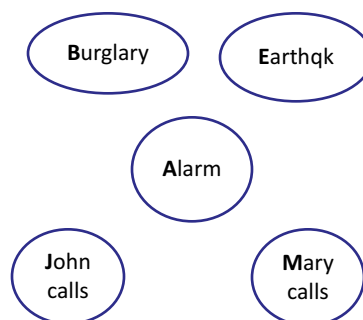
- **Another Useful Bayes Net**
  - Hybrid Discrete / Continuous
- **Learning Parameters for a Bayesian Network**
  - Fully observable
    - Maximum Likelihood (ML),
    - Maximum A Posteriori (MAP)
  - Hidden variables (EM algorithm)
- **Learning Structure of Bayesian Networks**

© Daniel S. Weld

What if we *don't* know structure?

## Learning The Structure of Bayesian Networks

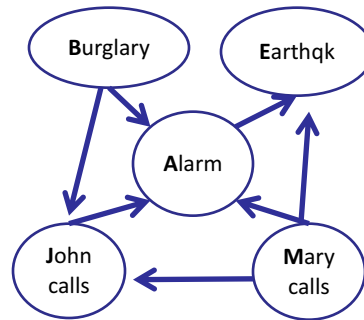
E	B	R	A	J	M
T	F	T	T	F	T
F	F	F	F	F	T
F	T	F	T	T	T
F	F	F	T	T	T
F	T	F	F	F	F
...					





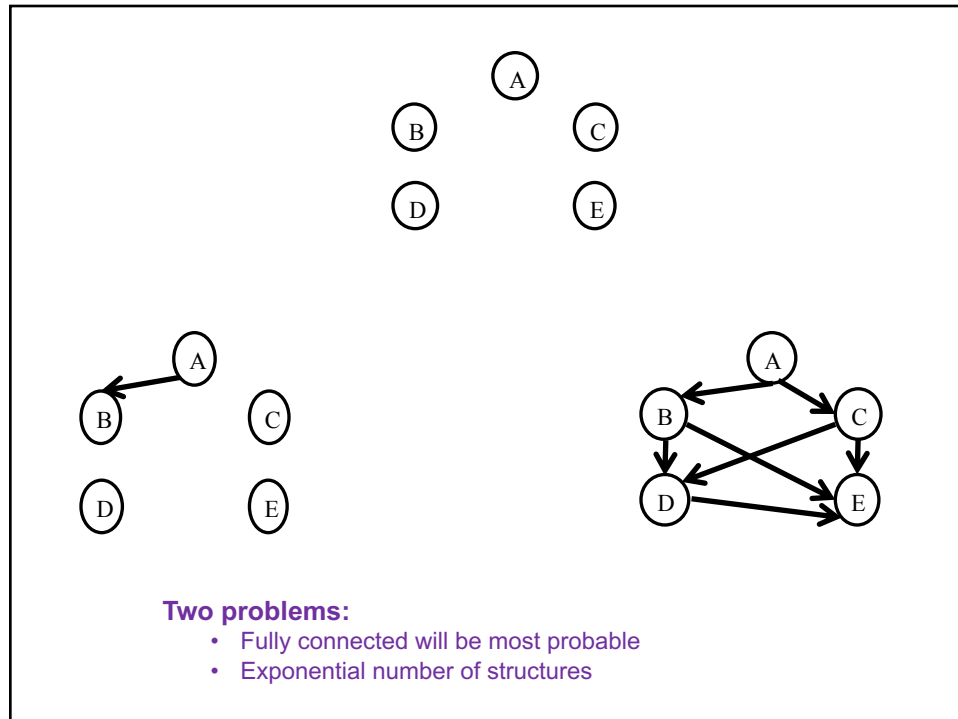
## Learning The Structure of Bayesian Networks

E	B	R	A	J	M
T	F	T	T	F	T
F	F	F	F	F	T
F	T	F	T	T	T
F	F	F	T	T	T
F	T	F	F	F	F
...					



## Learning The Structure of Bayesian Networks

- Search thru the space...
  - of possible network structures!
- For each structure, learn parameters
  - As just shown...
- Pick the one that fits observed data best
  - Calculate  $P(\text{data})$



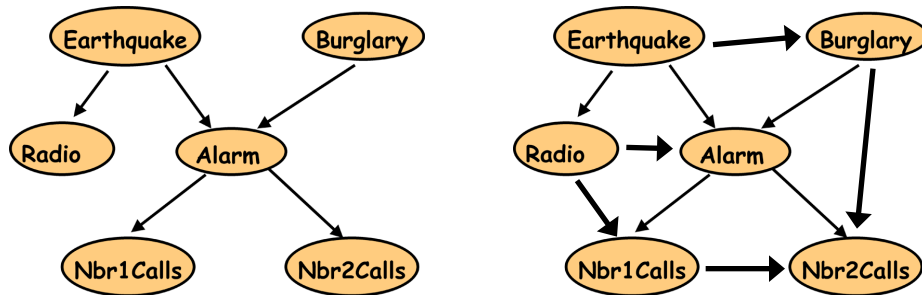
## Learning The Structure of Bayesian Networks

- Search thru the space...
  - of possible network structures!
- For each structure, learn parameters
  - As just shown...
- Pick the one that fits observed data best
  - Calculate  $P(\text{data})$

### Two problems:

- Fully connected will be most probable
  - Add penalty term (regularization)  $\propto$  model complexity
- Exponential number of structures
  - Local search

## Overfitting



Can represent strictly more P distributions

Can represent NOISE in training data

Often preforms WORSE on test data

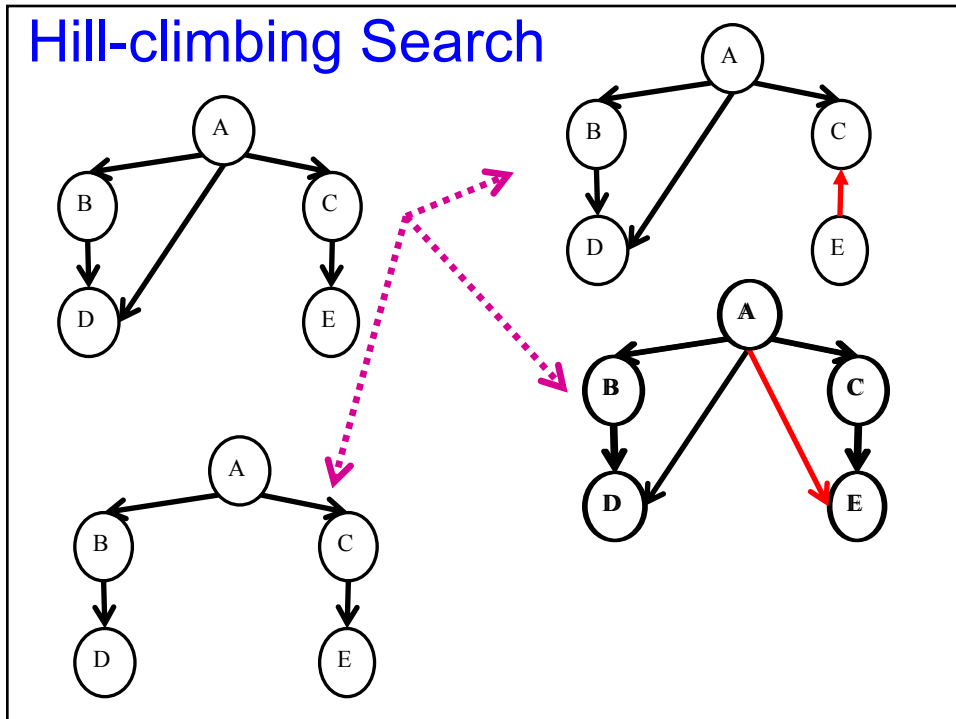
## Augment Score Function

- Bayesian Information Criterion (BIC)
  - $P(D | BN)$  – penalty
  - Penalty =  $\alpha$  complexity
  - $= \alpha \left[ \frac{1}{2} (\# \text{ parameters}) \right] \text{Log} (\# \text{ data points})$

Instance of “*regularization*”

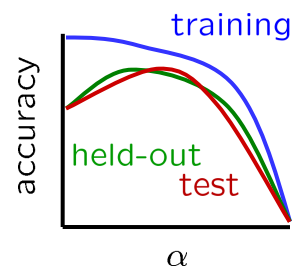
Solves problem of “*overfitting*”

## Hill-climbing Search



## Tuning on Held-Out Data

- Now we've got two kinds of unknowns
  - Parameters: the probabilities  $P(Y|X)$ ,  $P(Y)$
  - Hyperparameters, like
    - the amount of smoothing to do:  $k$ , or
    - regularization penalty,  $\alpha$
- Where to learn?
  - Learn parameters from training data
  - Must tune hyperparameters on different data
    - Why?
  - For each value of the hyperparameters, train and test on the held-out data
  - Choose the best value and do a final test on the test data



# Baselines

---

- **First step: get a baseline**
  - Baselines are very simple “straw man” procedures
  - Help determine how hard the task is
  - Help know what a “good” accuracy is
- **Weak baseline: most frequent label classifier**
  - Gives all test instances whatever label was most common in the training set
  - E.g. for spam filtering, might label everything as ham
  - Accuracy might be very high if the problem is skewed
  - E.g. calling everything “spam” gets 86%, so a classifier that gets 90% isn’t very good...
- **For real research, usually use previous work as a (strong) baseline**