

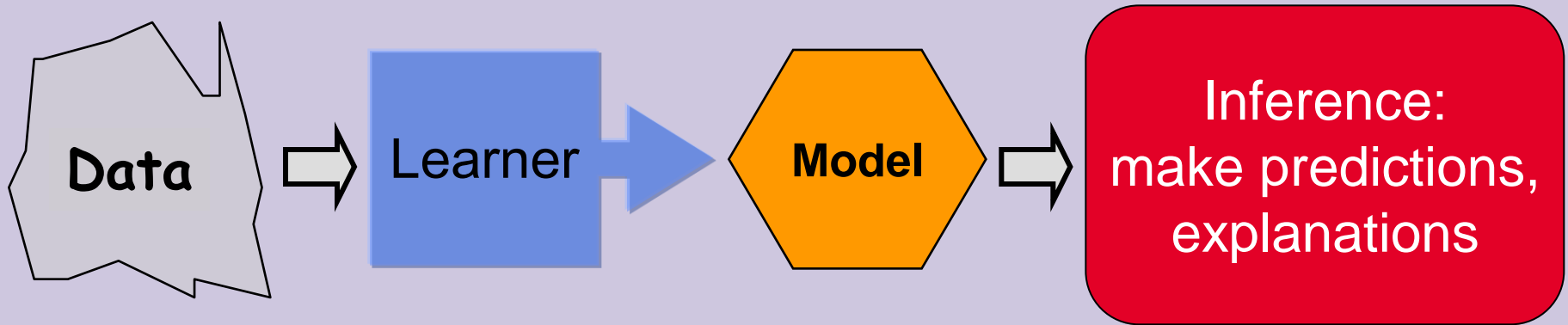
Learning Graphical Models

Gal Elidan

Hebrew University

(few slides are thanks to Nir Friedman and Daphne Koller)

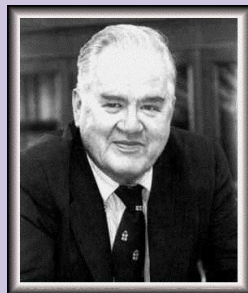
Why Machine Learning



Statistics and ML have many common goals:

models that fit data, prediction, probabilistic explanation, ways to cope with uncertainty, discovering truths about the data...

John Tukey
("bit", "software")



"Exploratory data analysis is an attitude, a flexibility..." (1980)

- Learning = hypothesis exploration + estimation
- Algorithms cope with high-dimensional domains

Bayesian networks are everywhere

Describe the child

in the drop-down boxes at the right. Relevant information will appear below.

Age: Sex:

Complaint:

Localized pain: Can the child localize, or point to, the site of the pain?

- No, unable to localize
- Below the navel to the child's left
- Above the child's navel
- Either of the child's sides
- Below the navel to the child's right
- Above the navel to the child's right

Results so far

Disorder	Rel
Viral gastroenteritis	<input checked="" type="checkbox"/>
Psychosomatic pain	<input checked="" type="checkbox"/>
Urinary tract infection	<input checked="" type="checkbox"/>

FIXIT - FL

File Window Help

3200L 9:37:04am

What To ASK

Have customer make a copy to the fax machine producing black lines?

QUESTIONS To ASK

- TX-Make Cpy-BkLns?
- Exposure Glass
- TX White Lines
- TX Black Copie
- Red Indicators
- Memory Full?
- TX Distorted In
- Noise (general)
- Receive File In

Angiosperms, Gymnosperms, Neurospora, Green Algae, Brown Algae, Red Algae, Diatoms, Yeasts, Sponges, Nematode, Arthropods, Ascobolus, Vertebrates, Urochordates, Echinoderms, Mollusks, Slime Molds, Amoeba, Heliozoans, Cilates, Dinoflagellates, Archaeobacteria, Myxobacteria, Gram Positive Bacteria, Cyanobacteria, Purple Bacteria, Chloroplasts, Mitochondria.

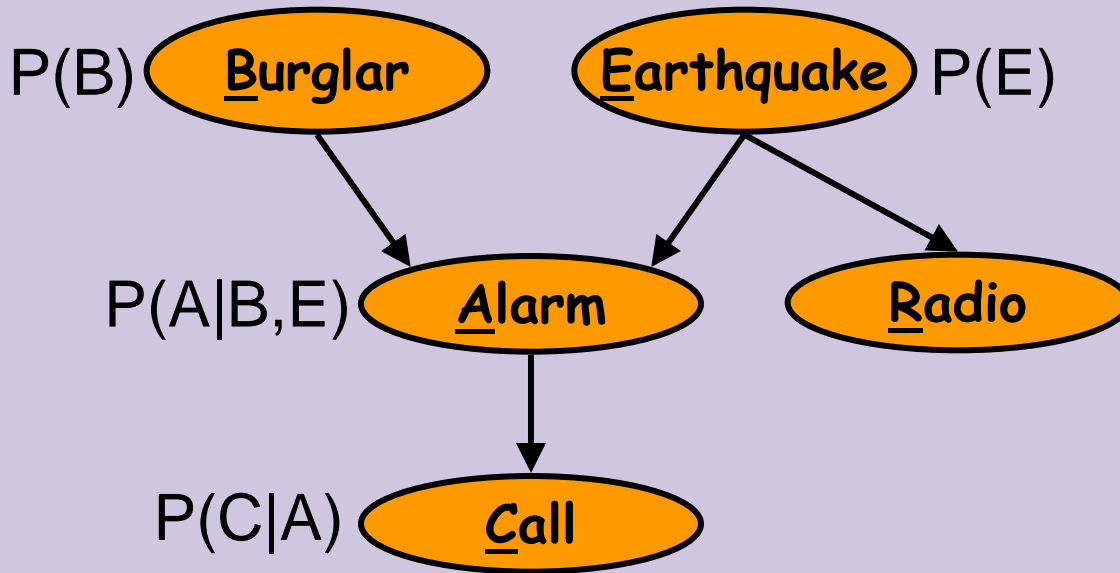
NASA Mission Control

ASCENT / ENTRY PROP		RR1104 CH	
LEFT	RIGHT	FORMER	ABORT
246	246	246	246
244	244	244	244
244	244	244	244
244	244	244	244
244	244	244	244
58.6	58.8	51.8	51.8
58.8	58.8	51.8	51.8
1118	1837	1113	1119
2378	2378	2378	2378
188.8	188.8	188.8	188.8
265	266	265	266
265	266	265	266
54.8	76		
67	67		
2188	2188	2188	2188
327	327		
0:00:00:00	0:00:00:00		
5.9	5.9	5.9	5.5
5.8	5.4	5.8	5.4
474	474	474	474
467	464	467	464
184	163		
499	298	499	298
268:04:55:17.31	268:04:55:17.31		
268:04:55:16.78	268:04:55:16.78		
268:04:54:11.20	268:04:54:11.20		
268:04:29:13.57	268:04:29:13.57		
268:04:56:26.97	268:04:56:26.97		

Overview

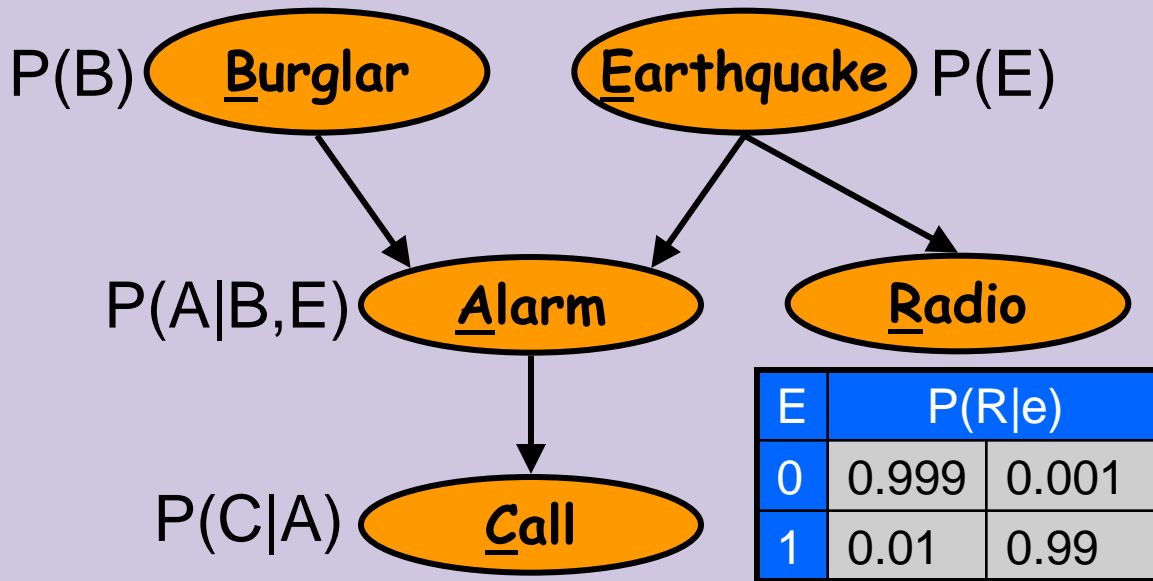
- Introduction
- Inference
- Parameter Estimation
- Model Selection

Bayesian Networks



Independence assumptions:
 $C \perp E, B, R \mid A$

Bayesian Networks



Independence assumptions:

$$X_i \perp \text{NonDesc}_i \mid \text{Par}_i$$

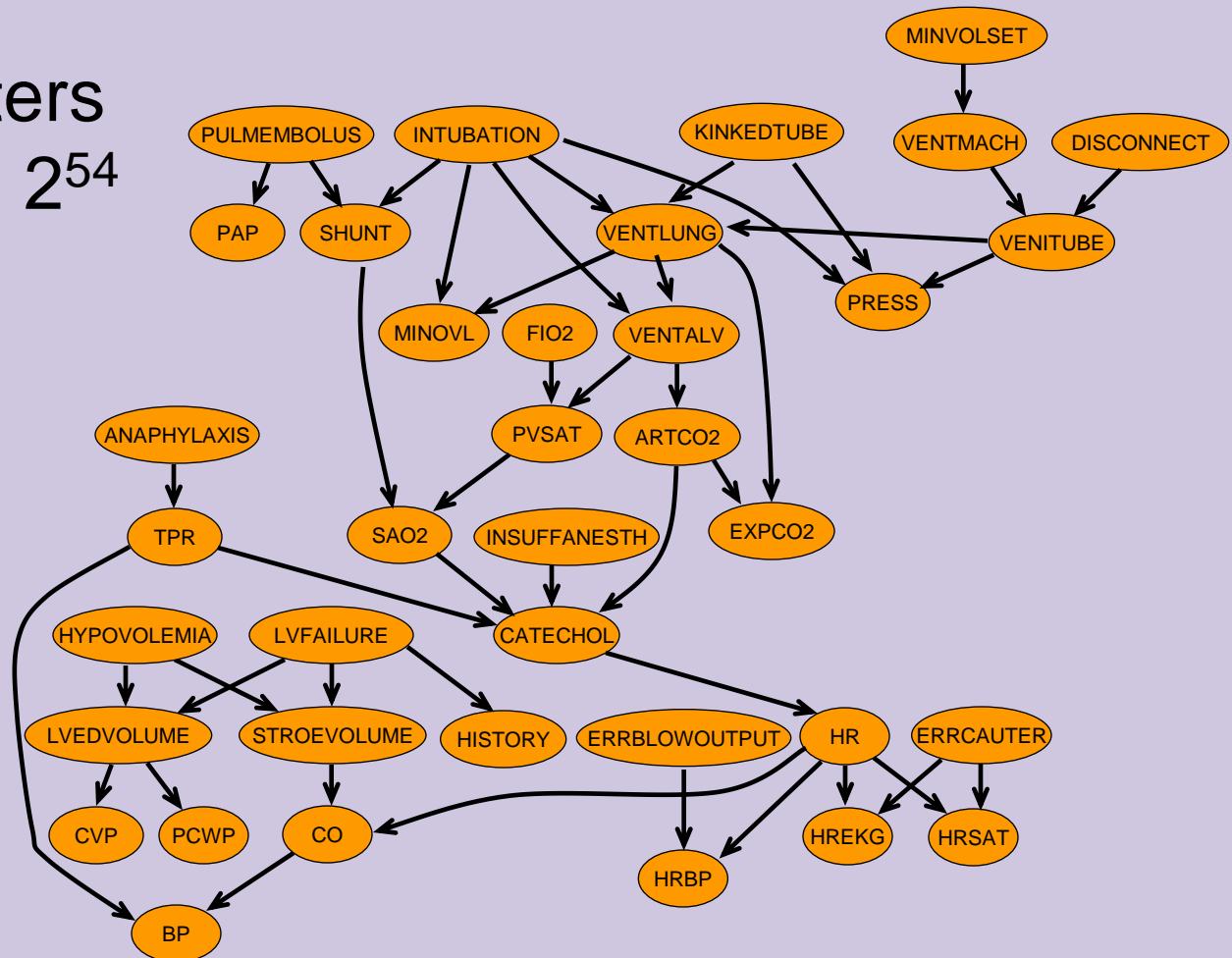
➔
$$P(\cdot) = \prod_i P(X_i \mid \text{Par}_i) = P(B)P(E)P(A \mid B, E)P(R \mid E)P(C \mid A)$$

What are the implications of this?

Example: “ICU Alarm” network

Domain: Monitoring Intensive-Care Patients

- 37 variables
- 509 parameters
...instead of 2^{54}



Proof

- w.l.o.g. let X_1, \dots, X_n be an order in which a parent appears before a child (topological ordering)
- assume $X_i \perp \text{NonDesc}_i \mid \text{Par}_i$

$$P(\cdot) = \prod_i P(X_i \mid X_1, \dots, X_{i-1})$$

chain rule

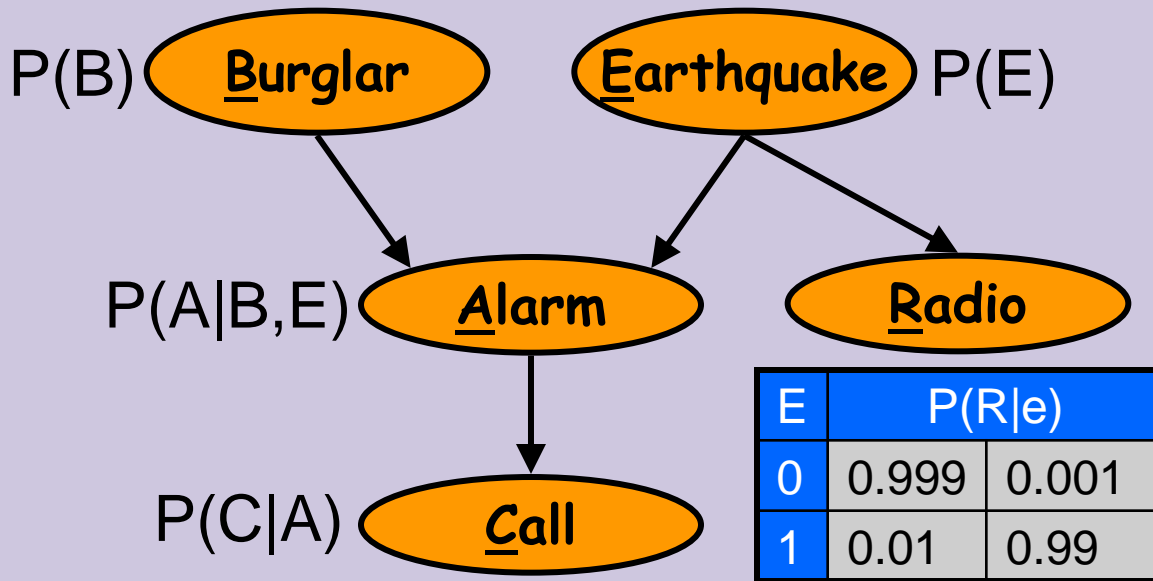
Topological ordering
 $S_i \subseteq \text{ND}_i$

$$= \prod_i P(X_i \mid \text{Par}_i \cup S_i)$$

$$= \prod_i P(X_i \mid \text{Par}_i)$$

Independence
assumptions

Bayesian Networks



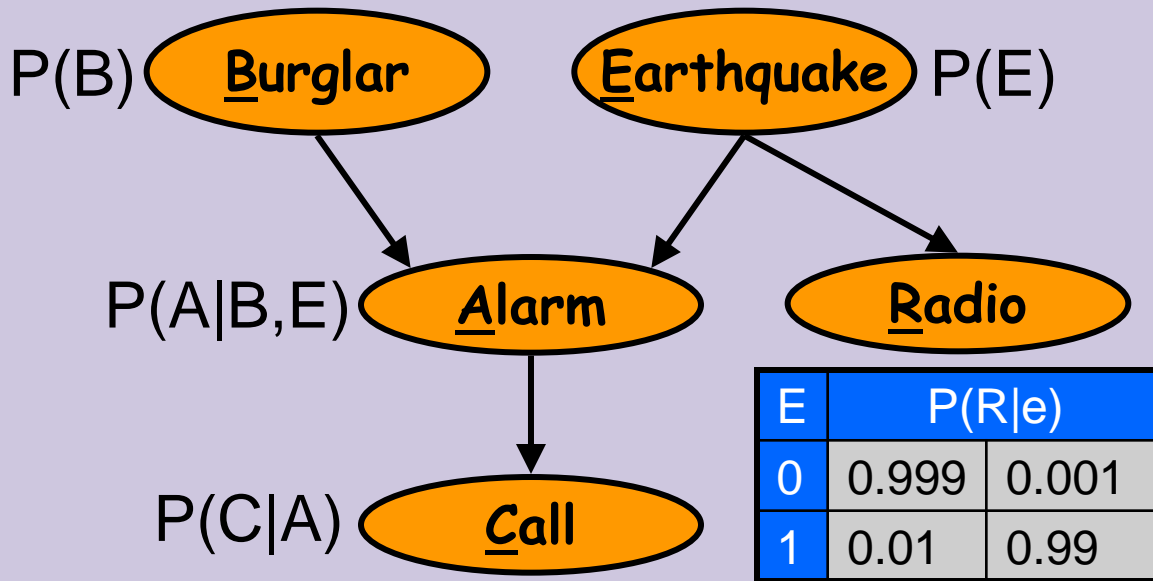
Independence assumptions:

$$X_i \perp \text{NonDesc}_i \mid \text{Par}_i$$

➔
$$P(\cdot) = \prod_i P(X_i \mid \text{Par}_i) = P(B)P(E)P(A \mid B, E)P(R \mid E)P(C \mid A)$$

- ✓ Compact representation of uncertainty
- ✓ Intuitive and interpretable representation
- ✓ Bidirectional inferences (prediction, explanation)
- ✓ Amenable to inference and learning algorithms

Bayesian Networks



Independence assumptions:

$$X_i \perp \text{NonDesc}_i \mid \text{Par}_i$$

➔
$$P(\cdot) = \prod_i P(X_i \mid \text{Par}_i) = P(B)P(E)P(A \mid B, E)P(R \mid E)P(C \mid A)$$

Formalism captures many common models:
 mixture/clustering, hierarchical Bayes,
 logistic regression, HMMs, factor analysis...

Take Home Problems

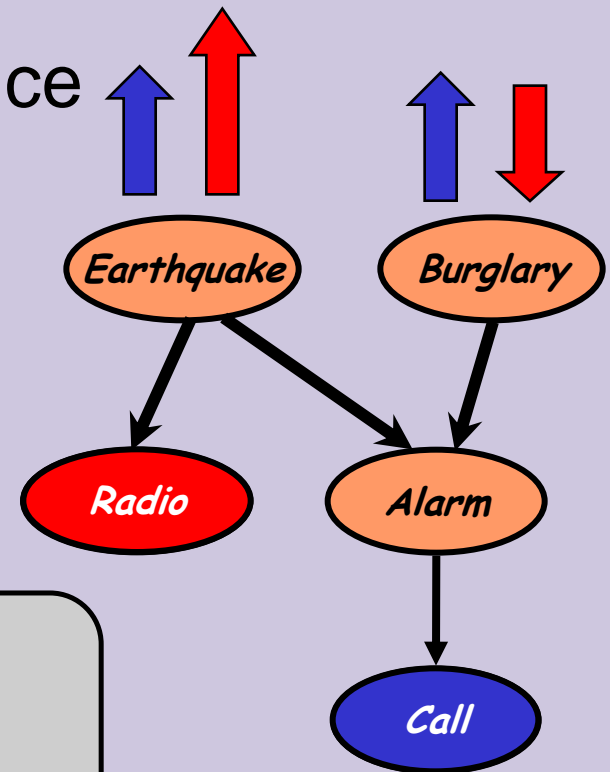
- 1) Assume we are given valid $\{P(X_i | Par_i)\}$
Prove that $P_B(\cdot) = \prod_i P(X_i | Par_i)$ is a distribution

Overview

- Introduction
- **Inference**
- Parameter Estimation
- Model Selection

Inference

- **Posterior probabilities**
 - Probability of any event given any evidence
- **Most likely explanation**
 - Scenario that explains evidence
- **Rational decision making**
 - Maximize expected utility
 - Value of Information
- **Effect of intervention**

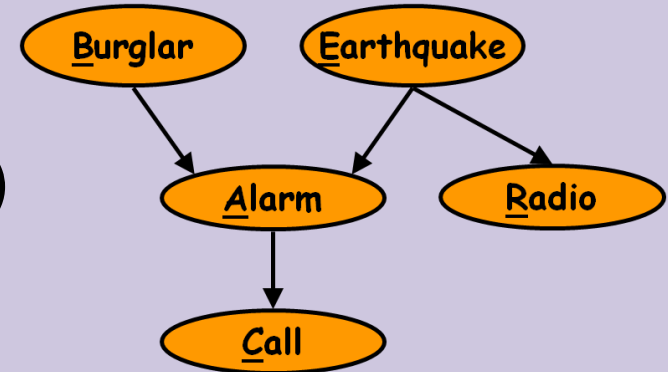


why is this difficult?

$$P(X_i) = \sum_{X_1, \dots, X_{i-1}, X_{i+1}, X_n} P(X_1, \dots, X_n)$$

Does Decomposition Help?

Let's say we are interested in $P(C)$



$$\begin{aligned} P(C) &= \sum_{b,e,a,r} P(B)P(E)P(A|B,E)P(R|E)P(C|A) \\ &= P(C|A) \sum_b P(B) \sum_e P(E) \sum_a P(A|B,E) \sum_r P(R|E) \end{aligned}$$

Still difficult in general...

Take Home Problems

1) Assume we are given valid $\{P(X_i | Par_i)\}$

Prove that $P_B(\cdot) = \prod_i P(X_i | Par_i)$ is a distribution

2) Let $P_B(X_1, \dots, X_n) = \prod_i P(X_i | X_{i-1})$ be a distribution represented by a chain network. How many operations (+, x) are required to compute $P_B(X_n)$ naively? By taking advantage of decomposition?

Overview

- Introduction
- Inference
- **Parameter Estimation**
- Model Selection

Why learn from data?

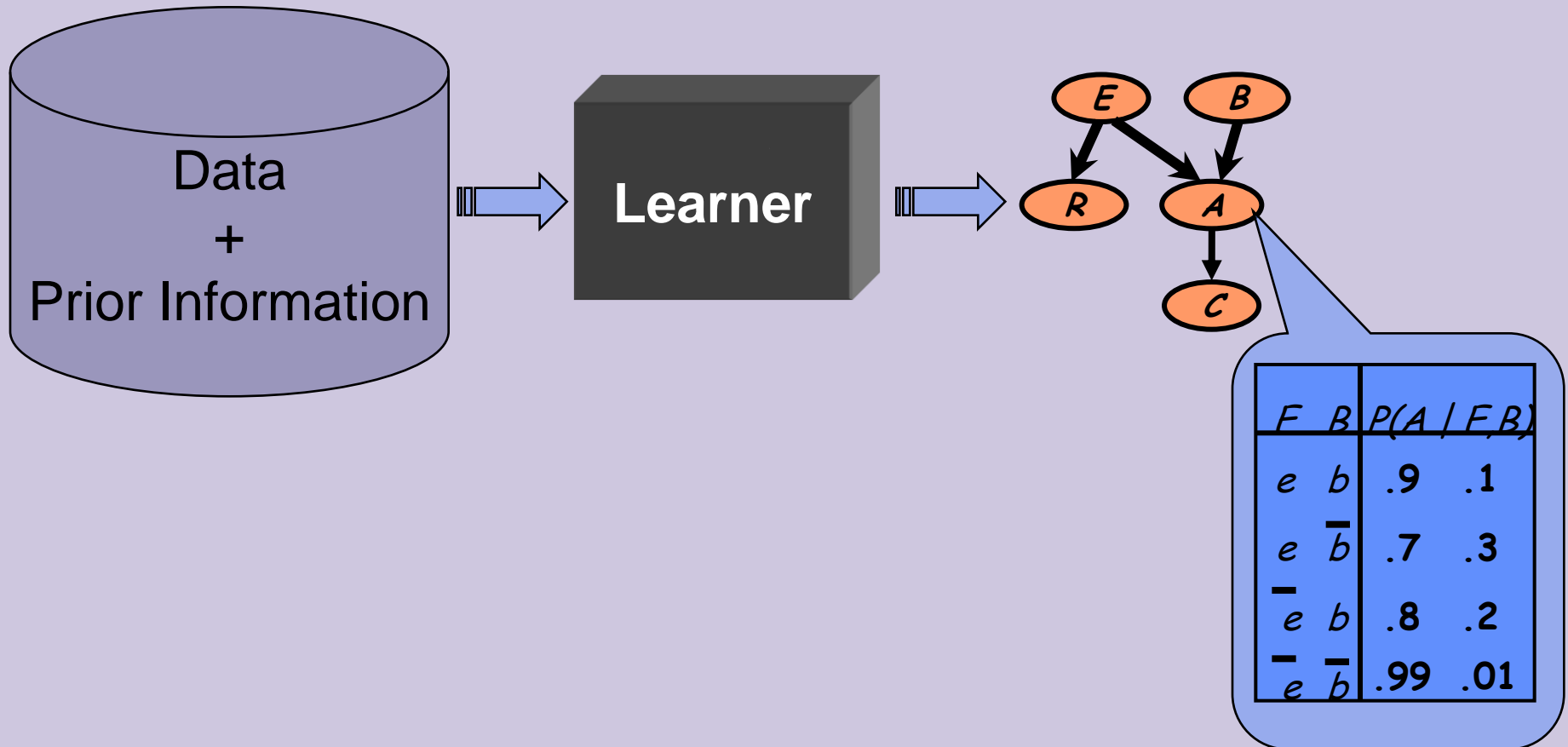
Knowledge acquisition bottleneck

- Knowledge acquisition is an expensive process
- Often we don't have an expert
- Robust encoding is often quite challenging (hard for humans to estimate global effects)

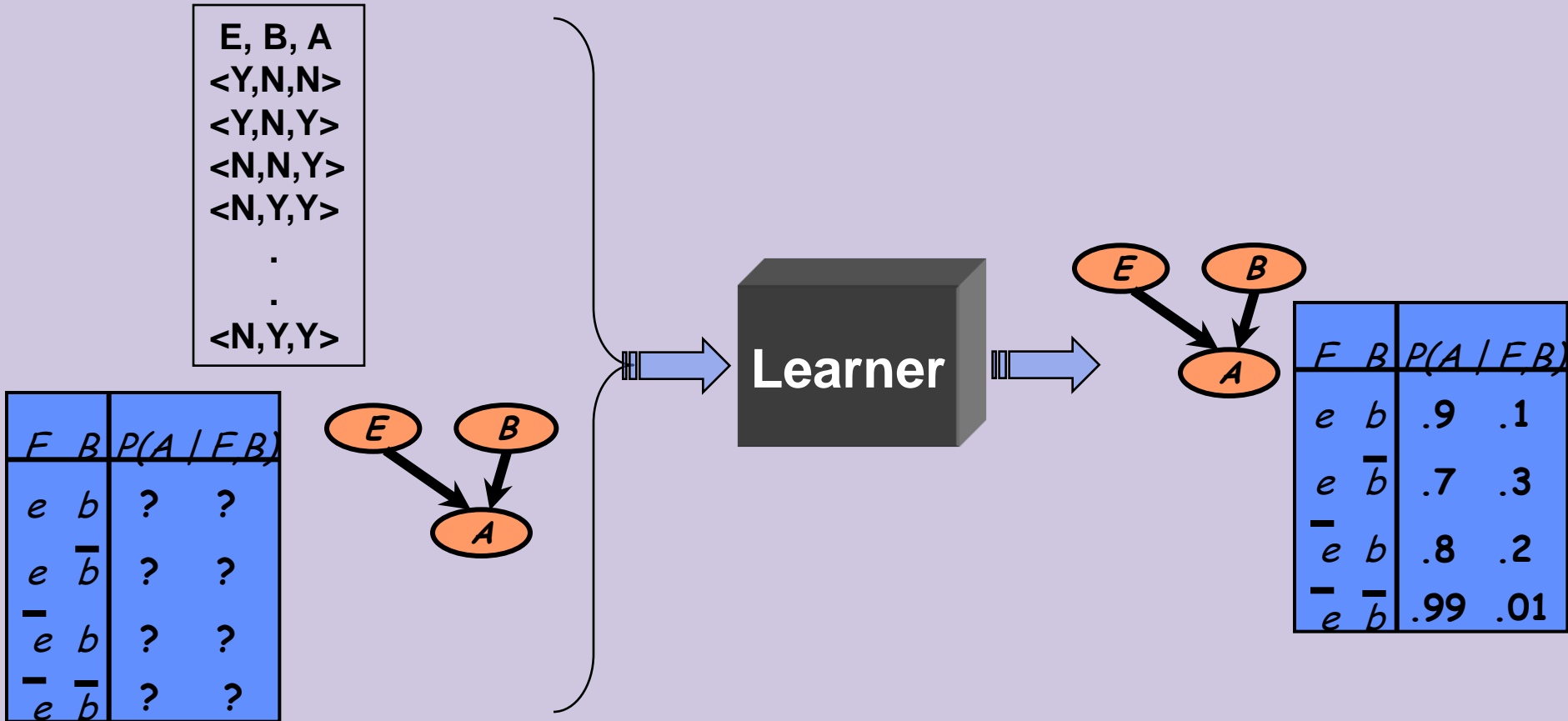
Data is cheap

- Amount of available information growing rapidly
- Learning allows us to construct models from raw data

Learning Bayesian networks

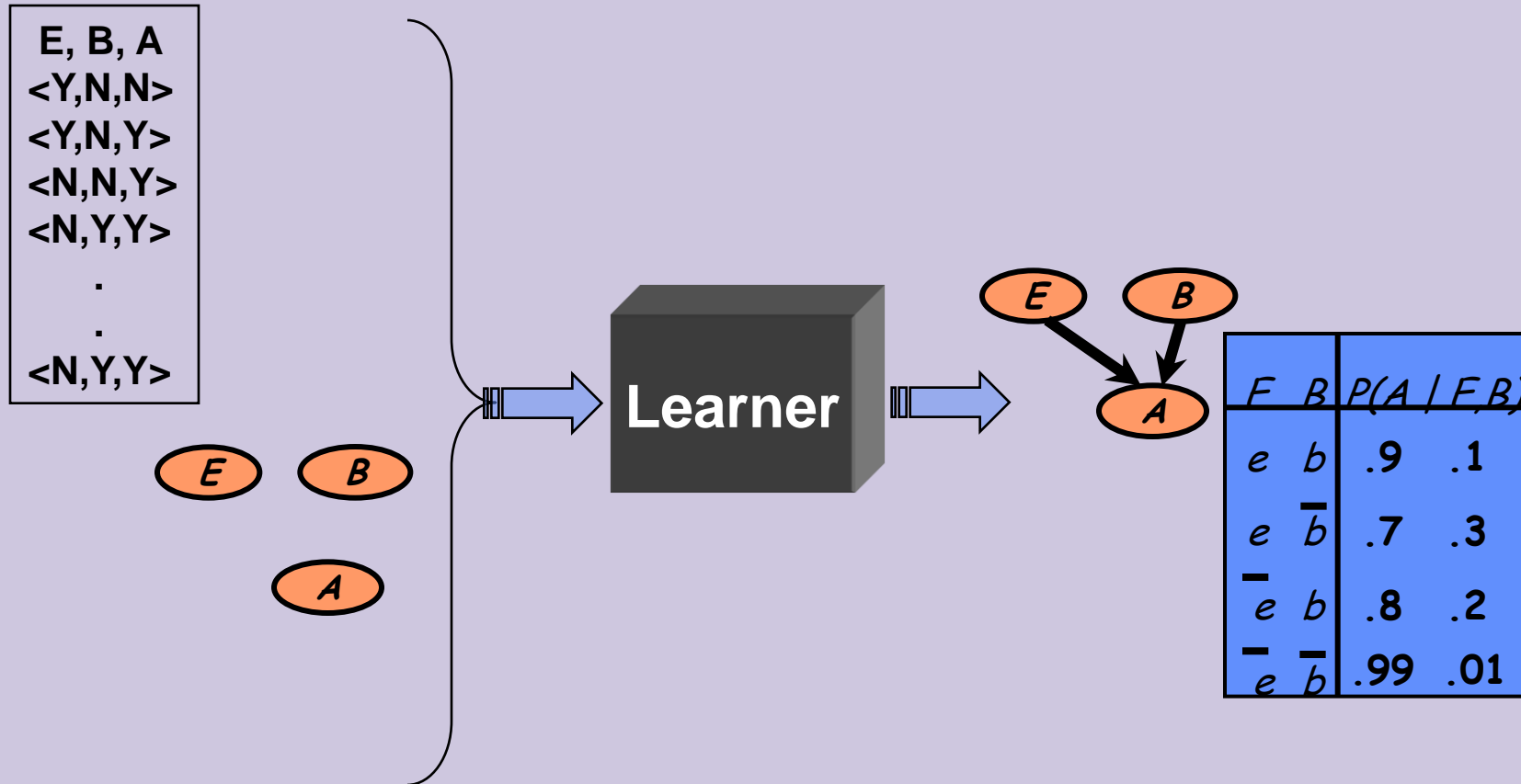


Known Structure, Complete Data



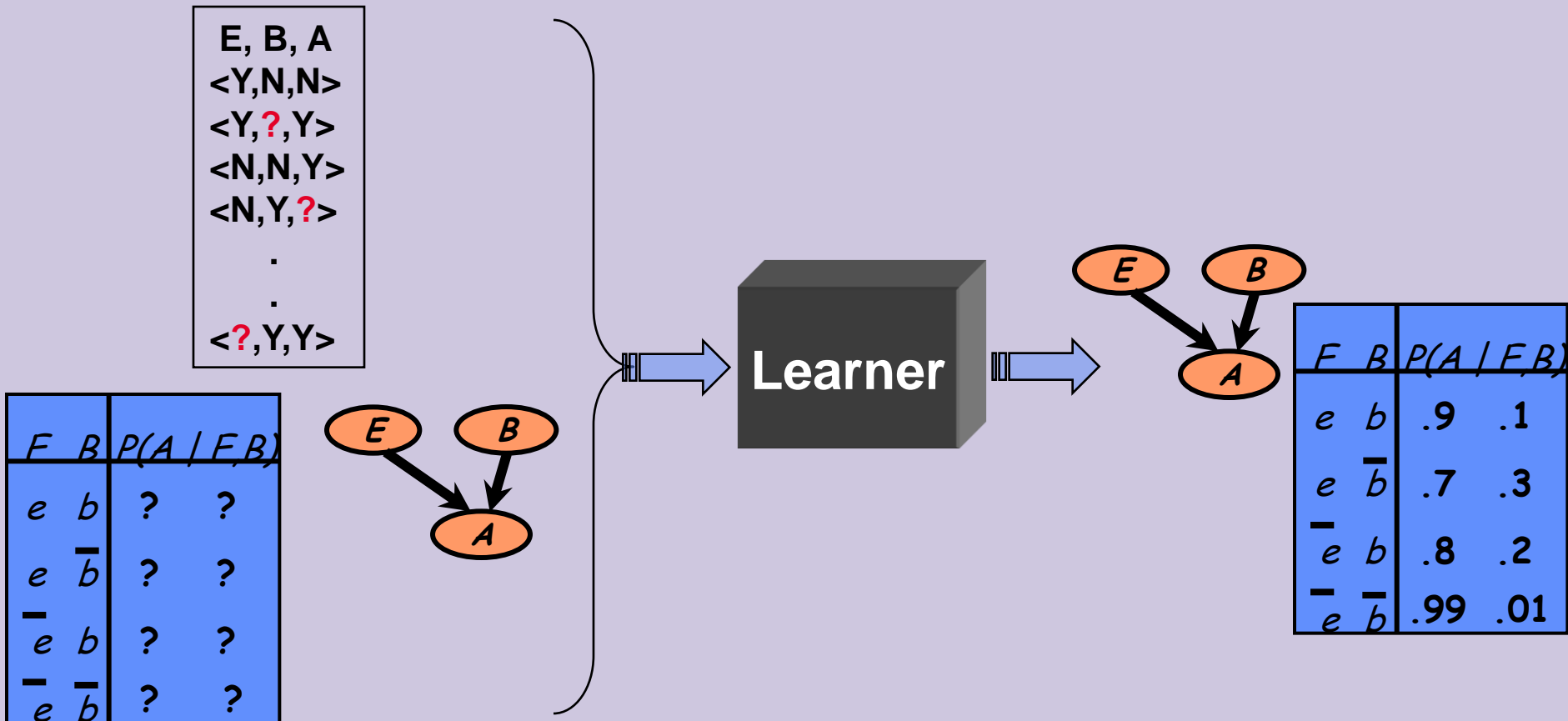
- Network structure is specified
 - Inducer needs to estimate parameters
- Data does not contain missing values

Unknown Structure, Complete Data



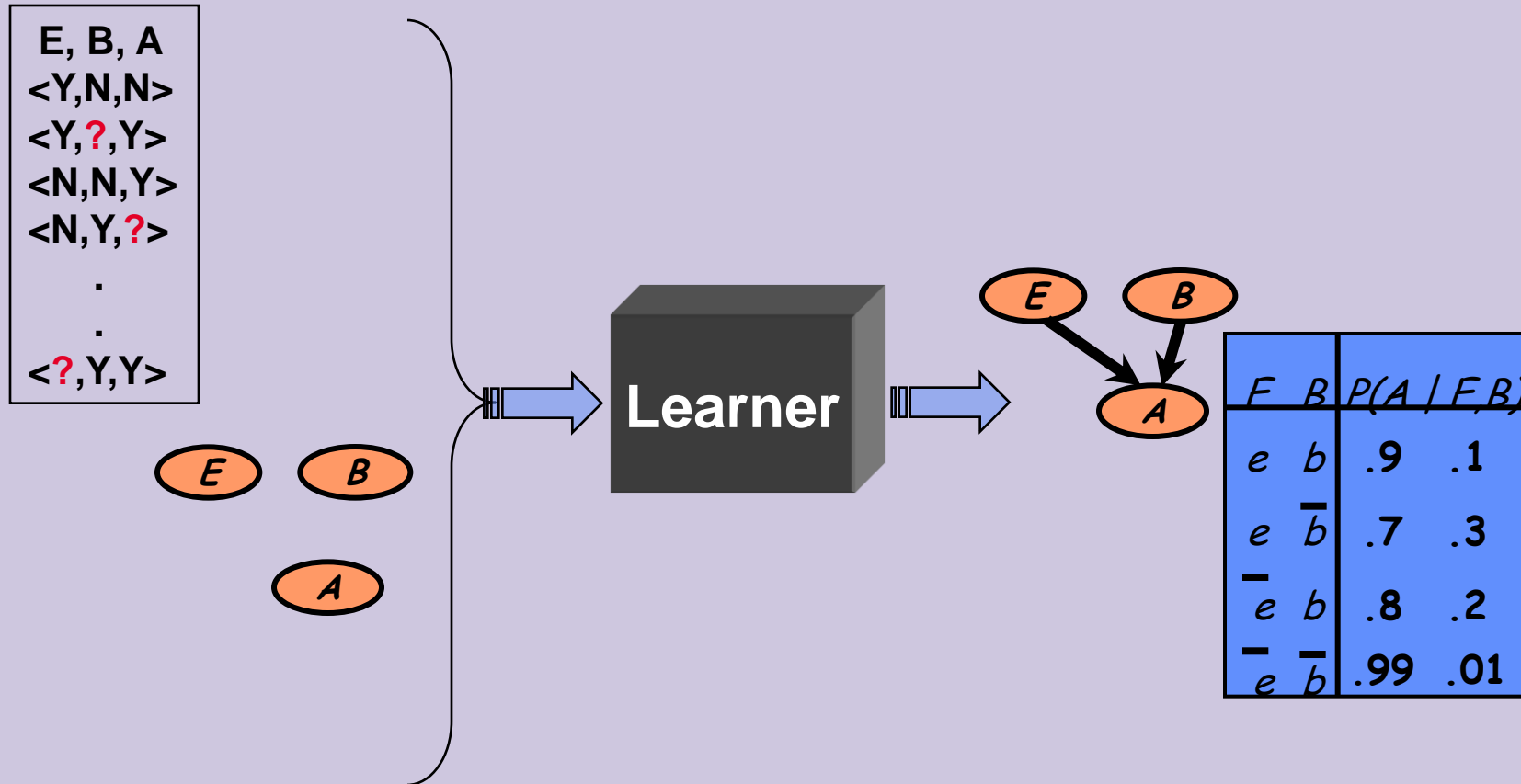
- Network structure is not specified
 - Inducer needs to select arcs & estimate parameters
- Data does not contain missing values

Known Structure, Incomplete Data



- Network structure is specified
- Data contains missing values
 - Need to consider assignments to missing values

Unknown Structure, Incomplete Data

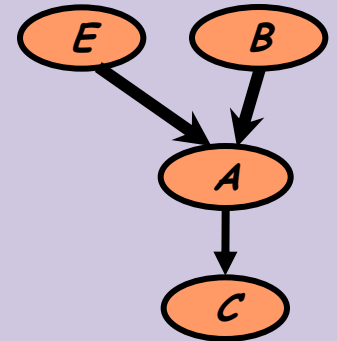


- Network structure is not specified
- Data contains missing values
 - Need to consider assignments to missing values

Learning Parameters

Training data has the form:

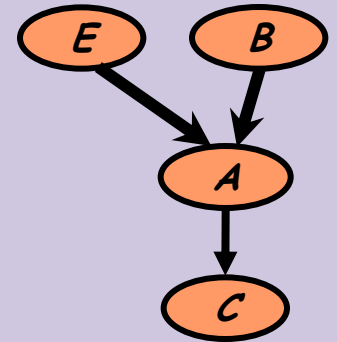
$$D = \begin{bmatrix} E[1] & B[1] & A[1] & C[1] \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ E[M] & B[M] & A[M] & C[M] \end{bmatrix}$$



Likelihood Function

- Assume i.i.d. samples
- Likelihood function is

$$L(\Theta : D) = \prod_m P(E[m], B[m], A[m], C[m] : \Theta)$$

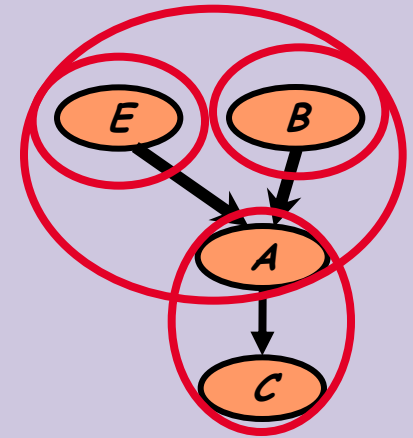


Likelihood Function

By definition of network, we get

$$L(\Theta : D) = \prod_m P(E[m], B[m], A[m], C[m] : \Theta)$$

$$= \prod_m \begin{pmatrix} P(E[m] : \Theta) \\ P(B[m] : \Theta) \\ P(A[m] | B[m], E[m] : \Theta) \\ P(C[m] | A[m] : \Theta) \end{pmatrix}$$



$E[1]$	$B[1]$	$A[1]$	$C[1]$
.	.	.	.
.	.	.	.
$E[M]$	$B[M]$	$A[M]$	$C[M]$

Likelihood Function

Rewriting terms, we get

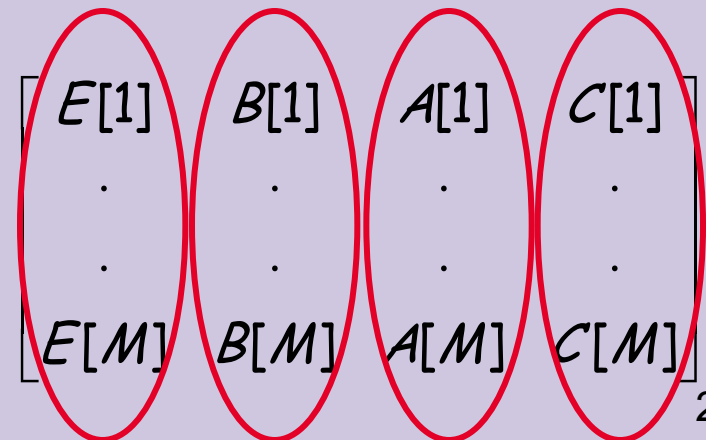
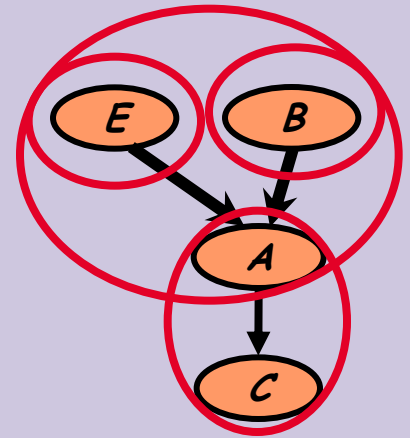
$$L(\Theta : D) = \prod_m P(E[m], B[m], A[m], C[m] : \Theta)$$

$$\prod_m P(E[m] : \Theta)$$

$$\prod_m P(B[m] : \Theta)$$

$$= \prod_m P(A[m] | B[m], E[m] : \Theta)$$

$$\prod_m P(C[m] | A[m] : \Theta)$$



General Bayesian Networks

Generalizing for any Bayesian network:

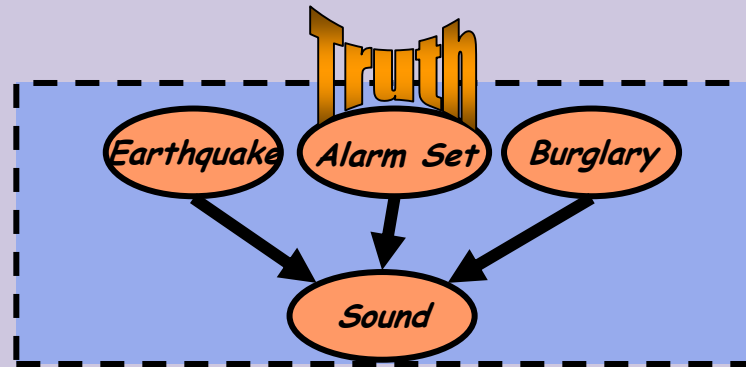
$$\begin{aligned} \mathcal{L}(\Theta : D) &= \prod_m P(x_1[m], \dots, x_n[m] : \Theta) \\ &= \prod_i \prod_m P(x_i[m] | Pa_i[m] : \Theta_i) \\ &= \prod_i \mathcal{L}_i(\Theta_i : D) \end{aligned}$$

Now you can use all that you have learned...

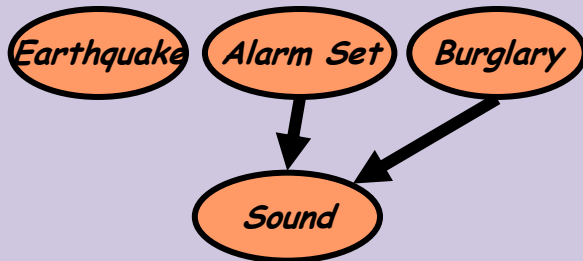
Overview

- Introduction
- Inference
- Parameter Learning
- **Model Selection**
 - Scoring function
 - Structure search

Why Struggle for Accurate Structure?

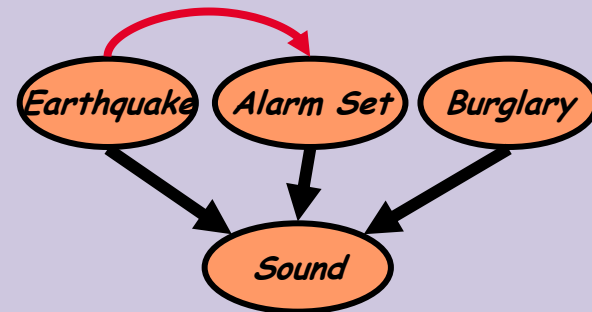


Missing an arc



- Cannot be compensated for by fitting parameters
- Wrong assumptions about domain structure

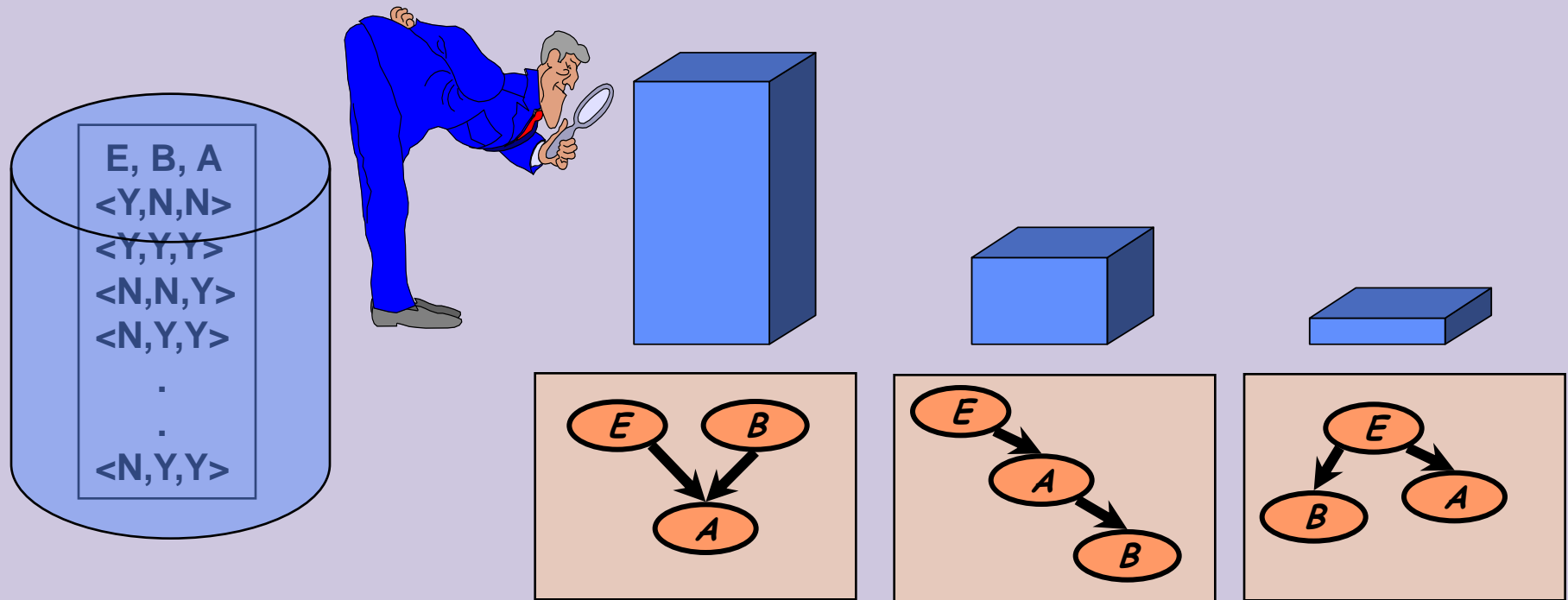
Adding an arc



- Increases the number of parameters to be estimated
- Wrong assumptions about domain structure

Score-based Learning

Define scoring function that evaluates how well a structure matches the data



Search for a structure that maximizes the score

Likelihood Score for Structure

$$\ell(G : D) = \log L(G : D) = M \sum_i (I(X_i; Pa_i^G) - H(X_i))$$

Mutual information between
 X_i and its parents

- Larger dependence of X_i on $Pa_i \Rightarrow$ higher score
- Adding arcs always helps
 - $I(X; Y) \leq I(X; \{Y, Z\})$
 - Max score attained by fully connected network
 - Overfitting: A bad idea...

Bayesian Score

Likelihood score: $L(G : D) = P(D | G, \hat{\theta}_G)$

Max likelihood params

Bayesian approach:

Deal with uncertainty by assigning probability to all possibilities

$$P(D | G) = \int P(D | G, \theta) P(\theta | G) d\theta$$

Marginal Likelihood

Likelihood

Prior over parameters

Bayesian Score

Likelihood score: $L(G : D) = P(D | G, \hat{\theta}_G)$

Max likelihood params

Bayesian approach:

Deal with uncertainty by assigning probability to all possibilities

$$P(D | G) = \int P(D | G, \theta) P(\theta | G) d\theta$$

Fortunately, in many cases integral has closed form.

Asymptotically we get:

$$\log P(D | G) = \ell(G : D) - \frac{\log M}{2} \dim(G) + O(1)$$

Fit empirical distribution

Complexity penalty

Structure Search as Optimization

Input:

- Training data
- Scoring function
- Set of possible structures

Output:

- A network that maximizes the score

Key Computational Property: Decomposability:

$$\text{score}(G) = \sum \text{score} (\text{family of } X \text{ in } G)$$

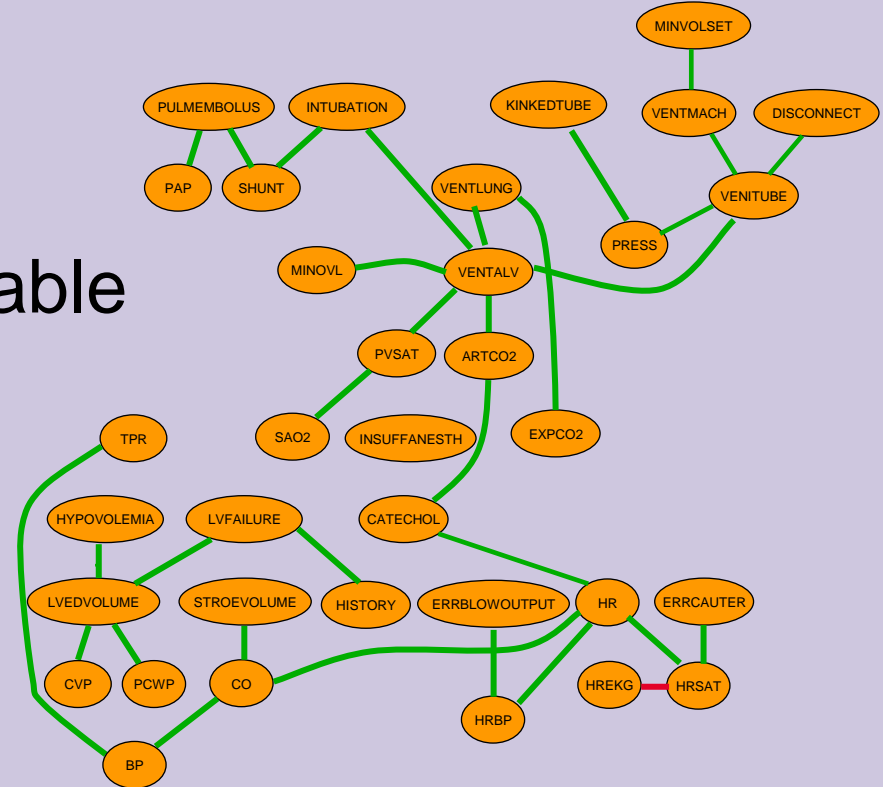
Tree-Structured Networks

Trees:

At most one parent per variable

Why trees?

- Elegant math
 - ⇒ we can solve the optimization problem
- Sparse parameterization
 - ⇒ avoid over-fitting



Learning Trees

- Let $p(i)$ denote parent of X_i
- We can write the Bayesian score as

$$Score(G : D) = \sum_i \log P(X_i : Pa_i) - Pen_i$$

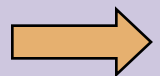
$$= \sum_i Score(X_i : Pa_i)$$

$$= \sum_i (Score(X_i : X_{p(i)}) - Score(X_i)) + \sum_i Score(X_i)$$

Improvement over
"empty" network

Score of "empty" network

Score = sum of edge scores + constant



Can find the optimal tree using
max-spanning tree algorithm

Beyond Trees

Essentially everything else is computationally difficult:

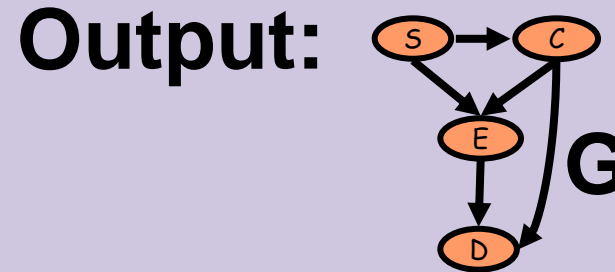
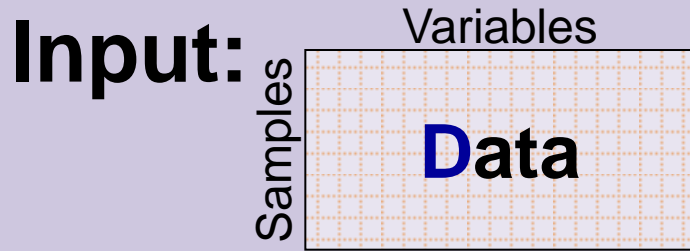
- Learning the optimal chain is NP-hard (exponential in the number of variables)
- Learning the optimal poly-tree is NP-hard
- Learning the optimal Bayesian network with at most k parents per node is NP-hard for $k > 1$

 This is where computer science comes in...

Heuristic Search

- Define a search space:
 - search states are possible structures
 - operators make small changes to structure
- Traverse space looking for high-scoring structures
- Search techniques:
 - Greedy hill-climbing
 - Best first search
 - Annealing
 - ...

Hypothesis Exploration (Local Search)

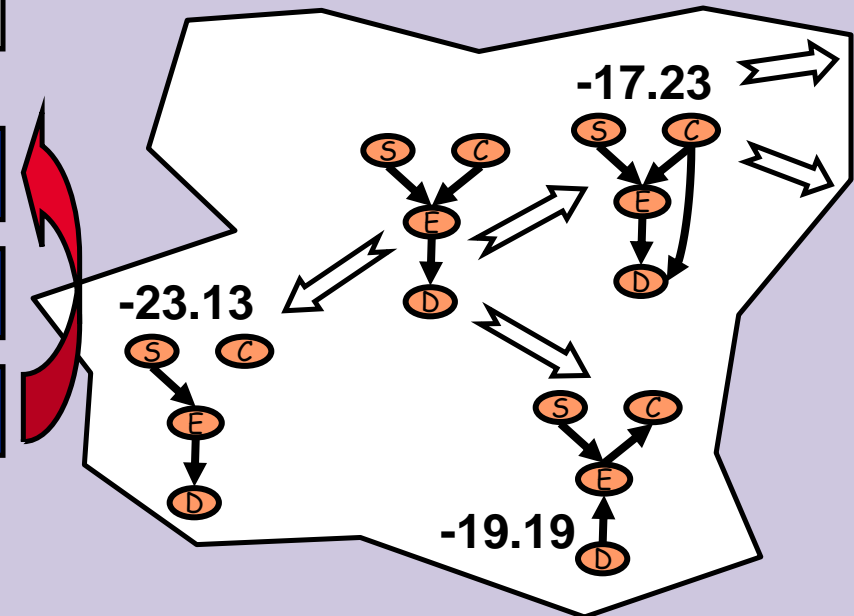


Init: Start with initial structure

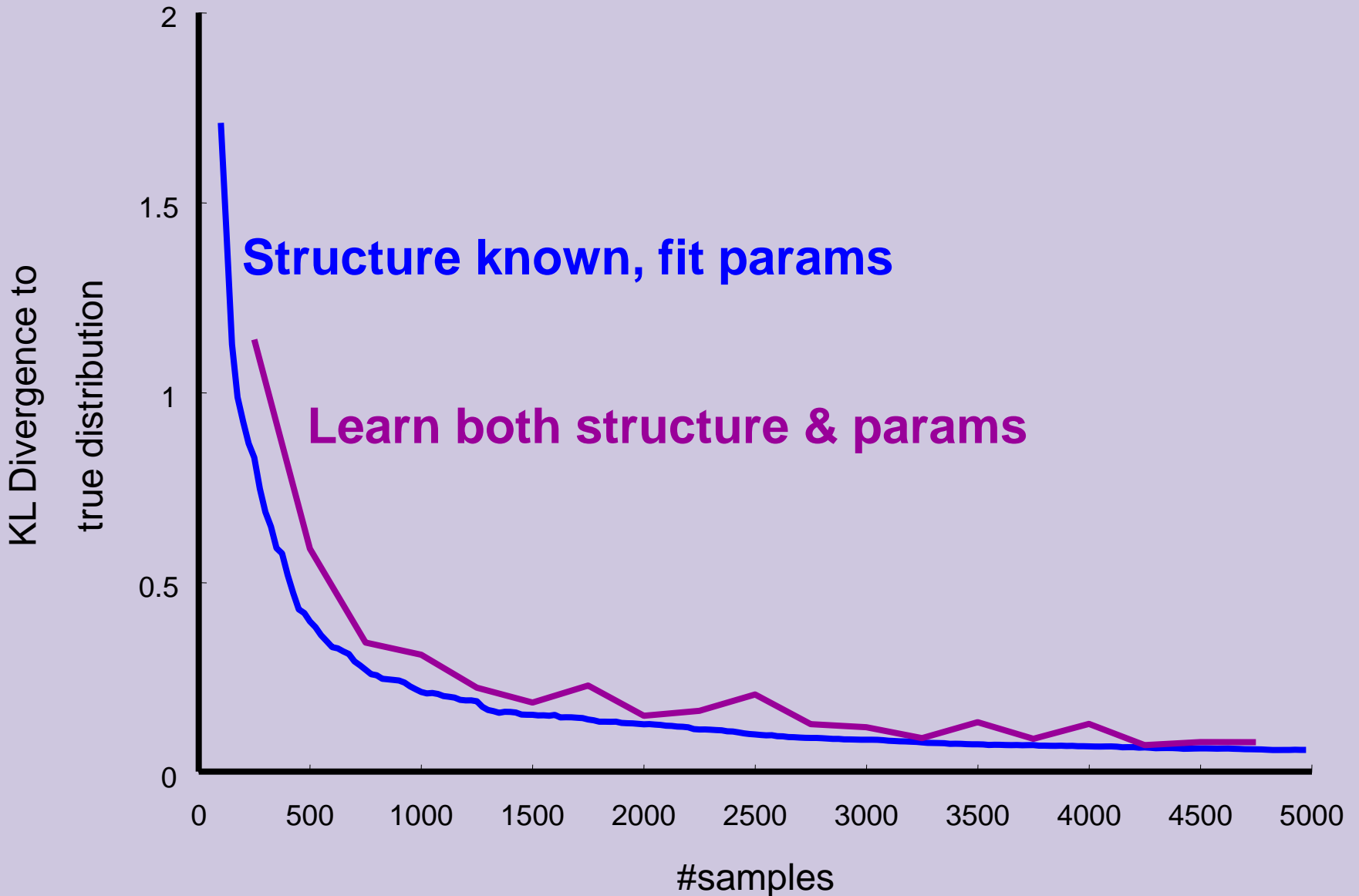
1 Consider local changes

2 Score each candidate

3 Apply best modification



Learning in Practice: Alarm domain



Local Search: Possible Pitfalls

- Local search can get stuck in:
 - **Local Maxima:**
 - All one-edge changes reduce the score
 - **Plateau:**
 - Some one-edge changes leave the score unchanged
- Standard heuristics can escape both
 - Random restarts
 - TABU search
 - Simulated annealing

Structure Search: Summary

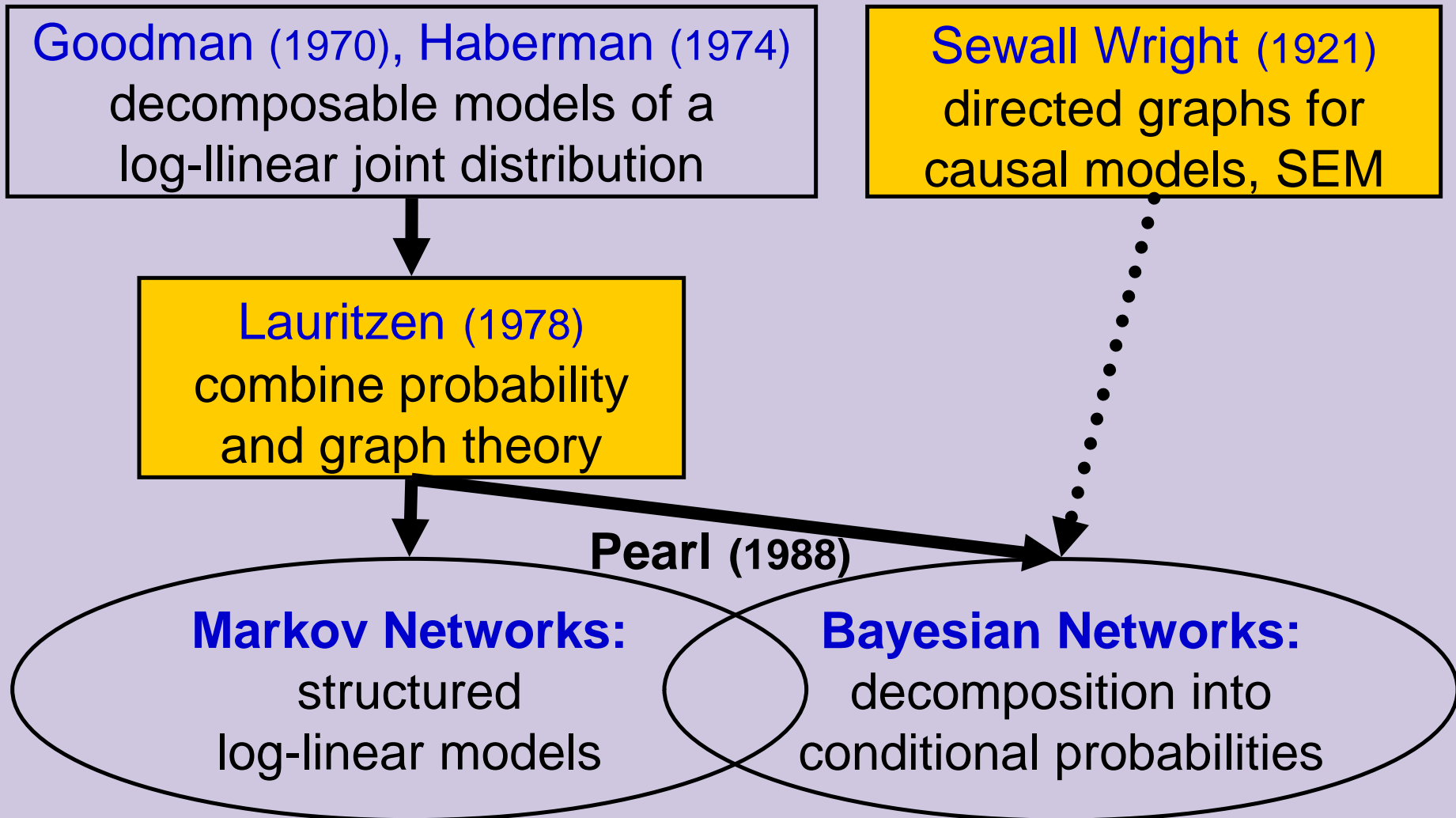
- Discrete optimization problem
- In some cases, optimization problem is easy
 - Example: learning trees
- In general, NP-Hard
 - Need to resort to heuristic search
 - In practice, search is relatively fast (~100 vars in ~2-5 min):
 - Decomposability
 - Sufficient statistics
 - Adding randomness to search is critical

Take Home Problems

- 1) Assume we are given valid $\{P(X_i | Par_i)\}$
Prove that $P_B(\cdot) = \prod_i P(X_i | Par_i)$ is a distribution
- 2) Let $P_B(X_1, \dots, X_n) = \prod_i P(X_i | X_{i-1})$ be a distribution represented by a chain network. How many operations (+, x) are required to compute $P_B(X_n)$ naively? By taking advantage of decomposition?
- 3) When adding/deleting an edge in the search we need to compute the score of the resulting graph. Explain precisely how does decomposability helps in this computation? What if we reverse an edge?

Graphical Models

Goal: make a joint distribution more amenable



Conclusion

- Many distributions have a dependency structure
- Utilizing this structure is good
- Discovering this structure has implications:
 - To density estimation
 - To knowledge discovery
 - To marginal computations
- Many applications
 - Medicine
 - Biology
 - Web
 - ...

Conclusion

- Many distributions have a dependency structure
- Utilizing this structure is good
- Discovering this structure has implications:
 - To density estimation
 - To knowledge discovery

Statistics and Learning

on the continuum of data analysis techniques

- **Common goals and similar challenges**
- **Both rely on probability theory**
- **Different methodologies (a plus!)**