

Bayesian Parameter Estimation in Bayesian Networks

Sargur Srihari

srihari@cedar.buffalo.edu

Topics

- Bayesian network where parameters are variables
- Global parameter independence
 - Leads to global decomposition
- How to choose priors for Bayesian learning
 - K2 Prior
 - BDe Prior
- Comparison of Bayesian and MLE in ICU example

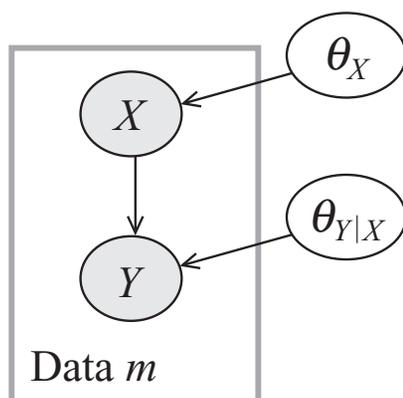
Fully Bayesian Approach

- In the full Bayesian approach to BN learning:
 - Parameters are considered to be random variables
- Need a joint distribution over unknown parameters θ and data instances D
- This joint distribution itself can be represented as a Bayesian network
 - instances and parameters of variables

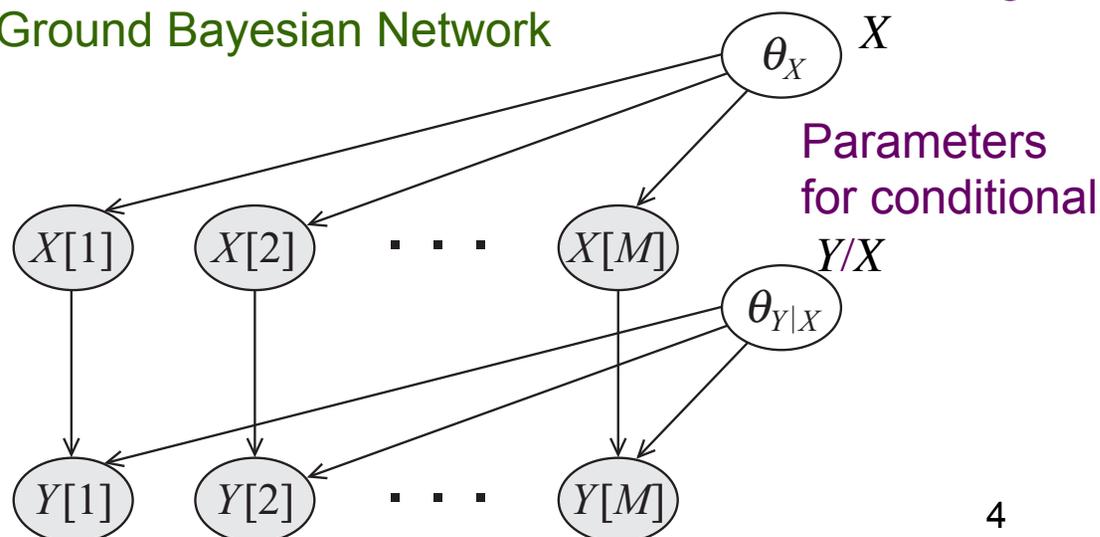
A simple example with 2 variables

- Consider the network $X \rightarrow Y$
- Training data $D = \{X[m], Y[m]\}$ for $m = 1, \dots, M$
- Unknown Parameter vectors θ_X and $\theta_{Y|X}$
- Meta network (BN of BN) for describing learning set-up

Plate Model



Ground Bayesian Network



Global Parameter Independence

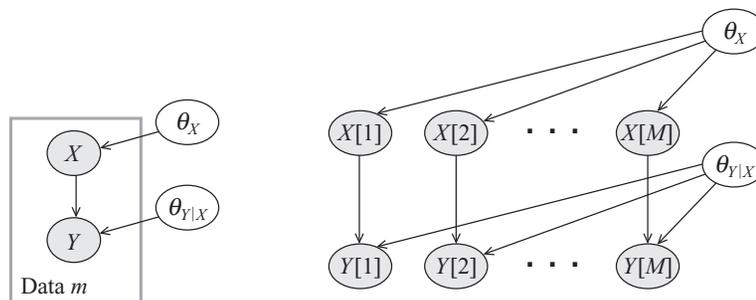
- If G is a Bayesian network with parameters

$$\theta = (\theta_{X_1|Pa_{X_1}}, \dots, \theta_{X_n|Pa_{X_n}})$$

- Global Parameter independence assumption:
 - Priors for individual parameters tells us nothing about another
 - A prior $P(\theta)$ satisfies global independence if

$$P(\theta) = \prod_i P(\theta_{X_i|Pa_{X_i}})$$

- Network with global parameter independence:



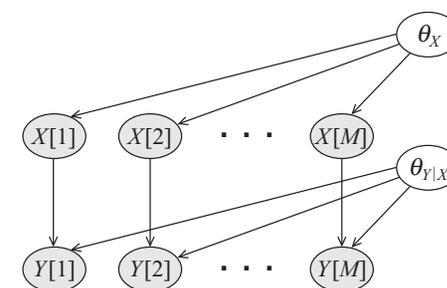
Caution on Parameter Independence

- Although we assume global parameter independence it should be considered with care
- Extension of student example
 - Student takes multiple classes
 - We want to learn distribution of grade given intelligence and difficulty
 - Same instructor for two courses may make parameters dependent

Global Decomposition

- Important conclusion if we assume global parameter independence
- E.g., if $x[m]$ and $y[m]$ are observed for all m , then θ_X and $\theta_{Y|X}$ are d-separated

$$P(\theta_X, \theta_{Y|X} | \mathcal{D}) = P(\theta_X | \mathcal{D})P(\theta_{Y|X} | \mathcal{D})$$



- Practical ramification

- Given data set \mathcal{D} , we can determine the posterior over θ_X independently of the posterior over $\theta_{Y|X}$

General Networks

- We considered simple network $X \rightarrow Y$
 - Concluded that each problem (of determining posterior) can be handled separately and combined
- Now we generalize conclusion
 - For general networks
- \mathcal{G} is a Bayesian network with parameters θ
- Need to specify prior $P(\theta)$ over all possible parameterizations of network

Prediction Problem

- Decomposition of posterior allows simplification of prediction
 - In $X \rightarrow Y$ example we want to compute probability of instance $x[M+1], y[M+1]$ given samples $x[1], y[1], \dots, x[M], y[M]$
 - We need to integrate over the parameter

$$P(x[M+1], y[M+1] | D) = \int P(x[M+1], y[M+1] | D, \theta) P(\theta | D) d\theta$$
 - Since posterior probability decomposes into a product we can show that

$$P(x[M+1], y[M+1] | D) = \left(\int P(x[M+1] | \theta_x) P(\theta_x | D) d\theta_x \right) \left(\int P(y[M+1] | \theta_{y|x}) P(\theta_{y|x} | D) d\theta_{y|x} \right)$$

- We have that the prediction problems for X and Y can be solved separately

Priors for BN Learning

- Each node X in a BN has a set of multinomial distributions $\theta_{X_i|pa_{X_i}}$, one for each instantiation pa_{X_i} of X_i 's parents Pa_{X_i}

- Each of these parameters has a separate Dirichlet prior governed by hyperparameters

$$\alpha_{X_i|pa_{X_i}} = \left(\alpha_{x_i^1|pa_{X_i}}, \dots, \alpha_{x_i^{K_i}|pa_{X_i}} \right) \quad \alpha \text{ are pseudocounts}$$

- Where K_i is the number of values of X_i

$$\theta \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k) \text{ if } P(\theta) \propto \prod_k \theta_k^{\alpha_k - 1}$$

$$\alpha = \sum_j \alpha_j \quad \text{Note that } \alpha \text{ is the total no of imaginary samples which does not sum to } 1$$

- How to determine α and α_j s?

K2 Prior

- In a software system called K2
- Use a fixed prior

$$\alpha_{x_i^j | pa_{X_i}} = 1$$

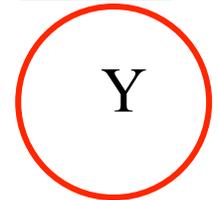
Add 1 to every cell in CPD

- For all hyperparameters in the network
- Has consequences that are conceptually unsatisfying

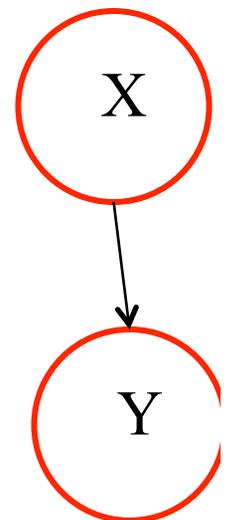
Inconsistency of K2 Prior

- For every multinomial parameter
 - a fixed *Dirichlet* $D(\alpha, \dots, \alpha)$
 - Typically $\alpha=1$, called K2 prior
 - K2 prior is inconsistent
- Ex: a binary variable Y with no parents
 - If we take *Dirichlet*(1,1) for θ_Y our imaginary sample size is 2
- Consider $X \rightarrow Y$ where X has four values
 - If *Dirichlet*(1,1) is prior for all parameters $\theta_{Y|x^i}$
 - Effectively two imaginary samples for each context x^i , for a total of eight

y1	y2
+1	+1



	x1	x2	x3	x4
y1	+1	+1	+1	+1
y2	+1	+1	+1	+1



Another intuitive approach

- Hyperparameter α_{xk} is an imaginary count of prior experience
- Have a data set D' of prior examples

- We set

$$\alpha_{x_i|pa_{X_i}} = \alpha \left[x_i, pa_{X_i} \right]$$

– where $\alpha [x_i, pa_{X_i}]$ is the no. of times $X_i=x_i$ and $Pa_{X_i}=pa_{X_i}$ in D'

- Count the parent-child combinations in imaginary data
- Equivalent to MLE with D and D'
- Requires storing a large set of pseudoinstances

BDe (Bayesian Dirichlet equivalent) prior

- Store size of imaginary data set α and $P'[X_1, \dots, X_n]$ of frequencies of events
 - probability of every possible ξ

- Set the parameters as follows:

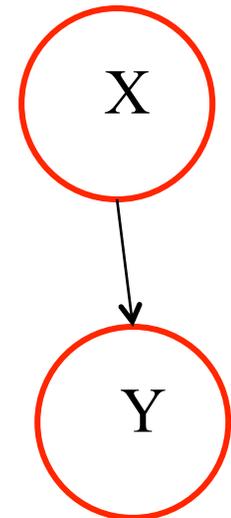
$$\alpha_{x_i | pa_{X_i}} = \alpha \cdot P'(x_i, pa_{X_i})$$

- This will divide data set proportionately

- Apply this to $X \rightarrow Y$ problem

$$\alpha_y = \alpha \cdot P'(y) = \sum_{x^i} \alpha \cdot P'(y, x^i) = \sum_{x^i} \alpha_{y|x^i}$$

- No of imaginary samples for different choices of parents will be identical



How do we represent P' for BDe?

- Use a Bayesian network so we can efficiently infer the quantities $P'(x_i, pa_{X_i})$
- When we have no prior knowledge we can set P' to be uniform
 - Implies an empty BN with uniform distributions for each variable
- Network structure is used only to provide parameters' priors not to guide structure search

BDe satisfies Score Equivalence

- BDe has a property useful in structure learning

- Bayesian score

$$P(G | D) = \frac{P(D | G)P(G)}{P(D)}$$

$$\text{score}_B(G : D) = \log P(D | G) + \log P(G)$$

- Score of a graph with Dirichlet Prior

$$\text{score}_{BIC}(G : D) = \ell(\theta_G : D) - \frac{\log M}{2} \text{Dim}[G]$$

- Score decomposability

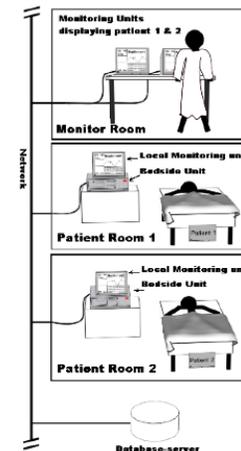
$$\text{score}(G) = \sum_{i=1}^n \text{score}(X_i, \text{Pa}_G(X_i)).$$

- Score equivalence

- Networks with same Bde score are I-equivalent

Application to ICU monitoring

- Decision support system (medical diagnostics)
- Monitoring severely ill patients in ICU
- Goal is to calculate probabilities for a differential diagnosis based on available evidence

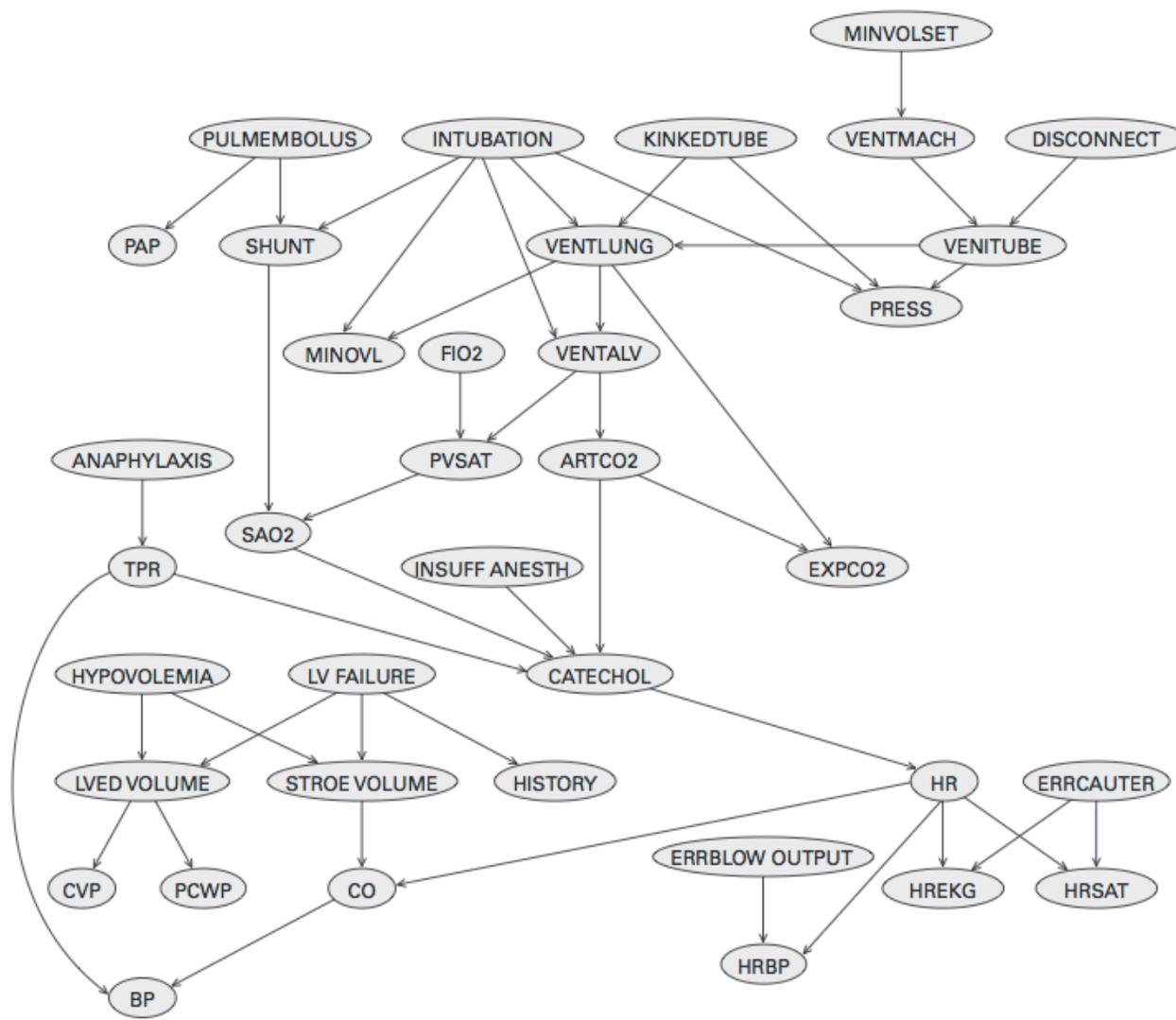


In differential diagnosis there are multiple alternatives which are systematically eliminated

Variables in ICU Alarm

- Medical knowledge encoded as 37 variables:
 - 8 diagnoses (at top level, no parents)
 - 16 findings
 - 13 intermediate variables
- Some top-level variables:
 - Anaphylaxis (severe allergic reaction)
 - Hypovolemia (shock of decreased blood volume)
 - Kinked Tube
 - Disconnect
 - Left Ventricle Failure

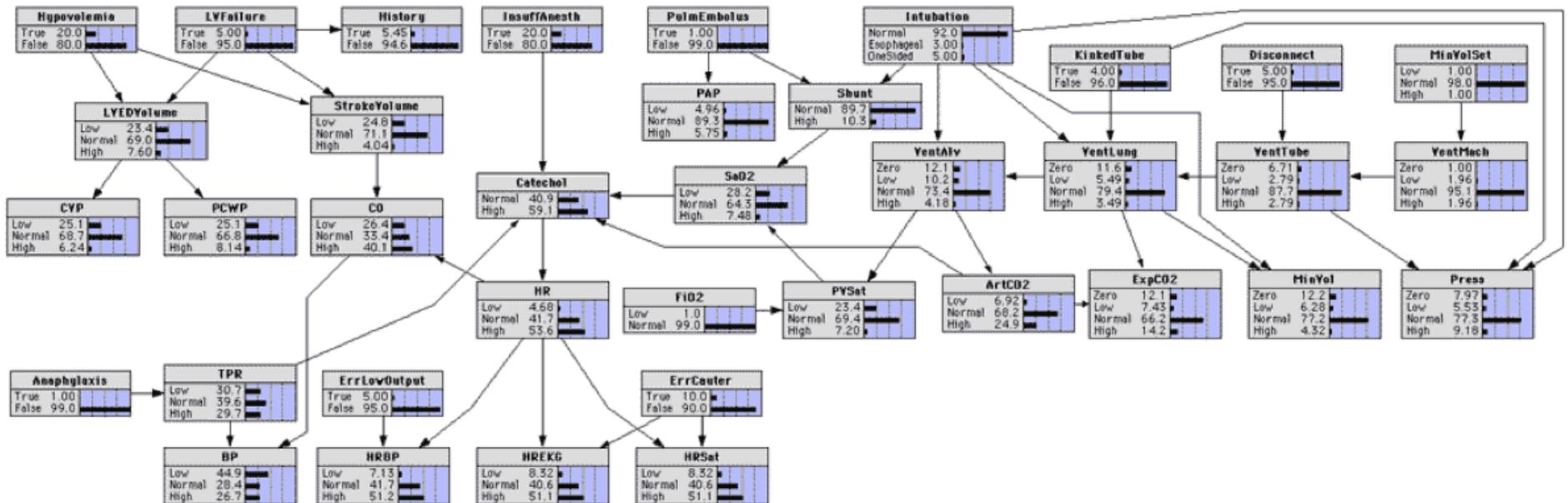
Bayesian Network for ICU-Alarm



Hand-constructed
by experts
for monitoring patients
in an
Intensive Care Unit

Benchmark for
Bayesian Network
Learning
(Parameter
Estimation)

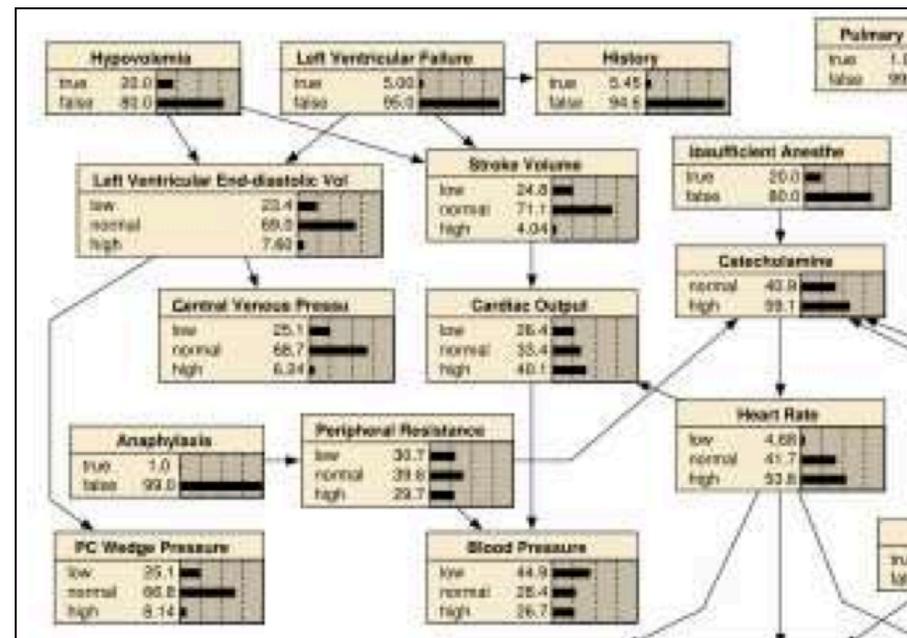
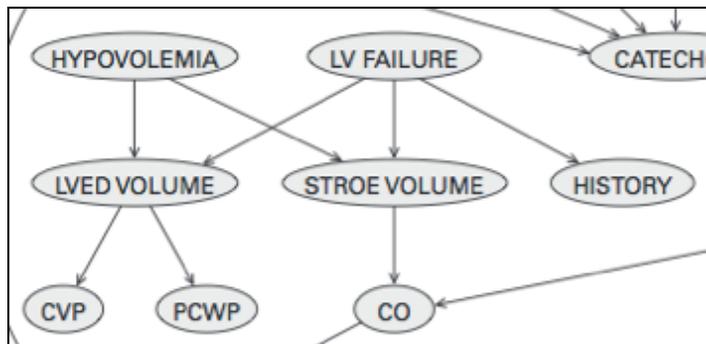
ICU alarm network



ALARM network: A Logical Alarm Reduction Mechanism

Marginal probabilities of each variable are shown

Contingency Tables between Dependent Variables



A contingency tables usually shows dependency between a pair of variables
 These show marginal probabilities of each variable

Parameters for ICU Alarm

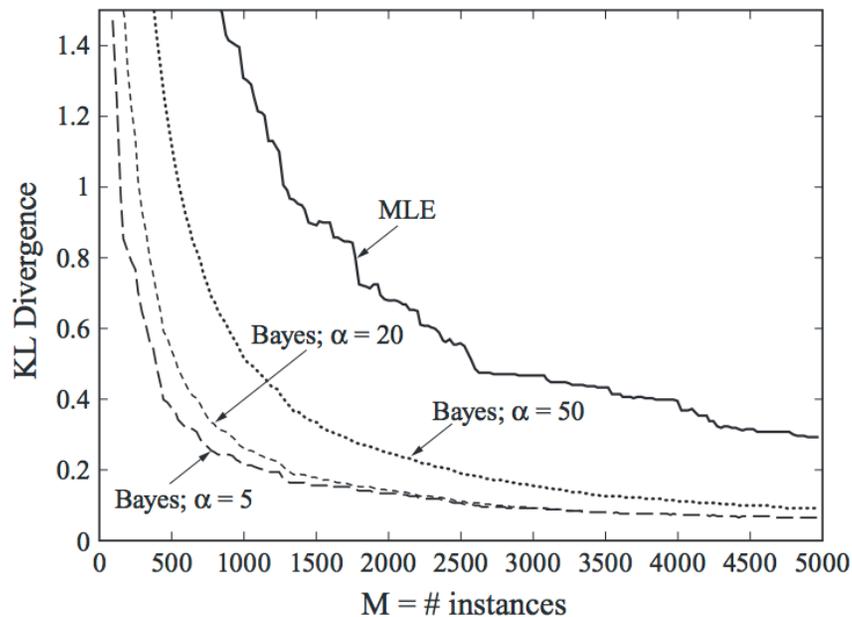
- Network Parameters
 - 37 nodes with table CPDs
 - 504 parameters (instead of 2^{54})
- Evaluate ability of parameter estimation method
 - To reconstruct network parameters from data
- Inputs for parameter estimation
 - Training samples
 - Obtained by sampling from specified network
 - Network structure

Learning curve for ICU alarm

Measuring the quality of the learned network

Relative entropy between generating distribution P^* and learned distribution P
Provides a global measure of extent to which learned distribution resembles true distribution

- Relative entropy to true model



Lowest error achieved

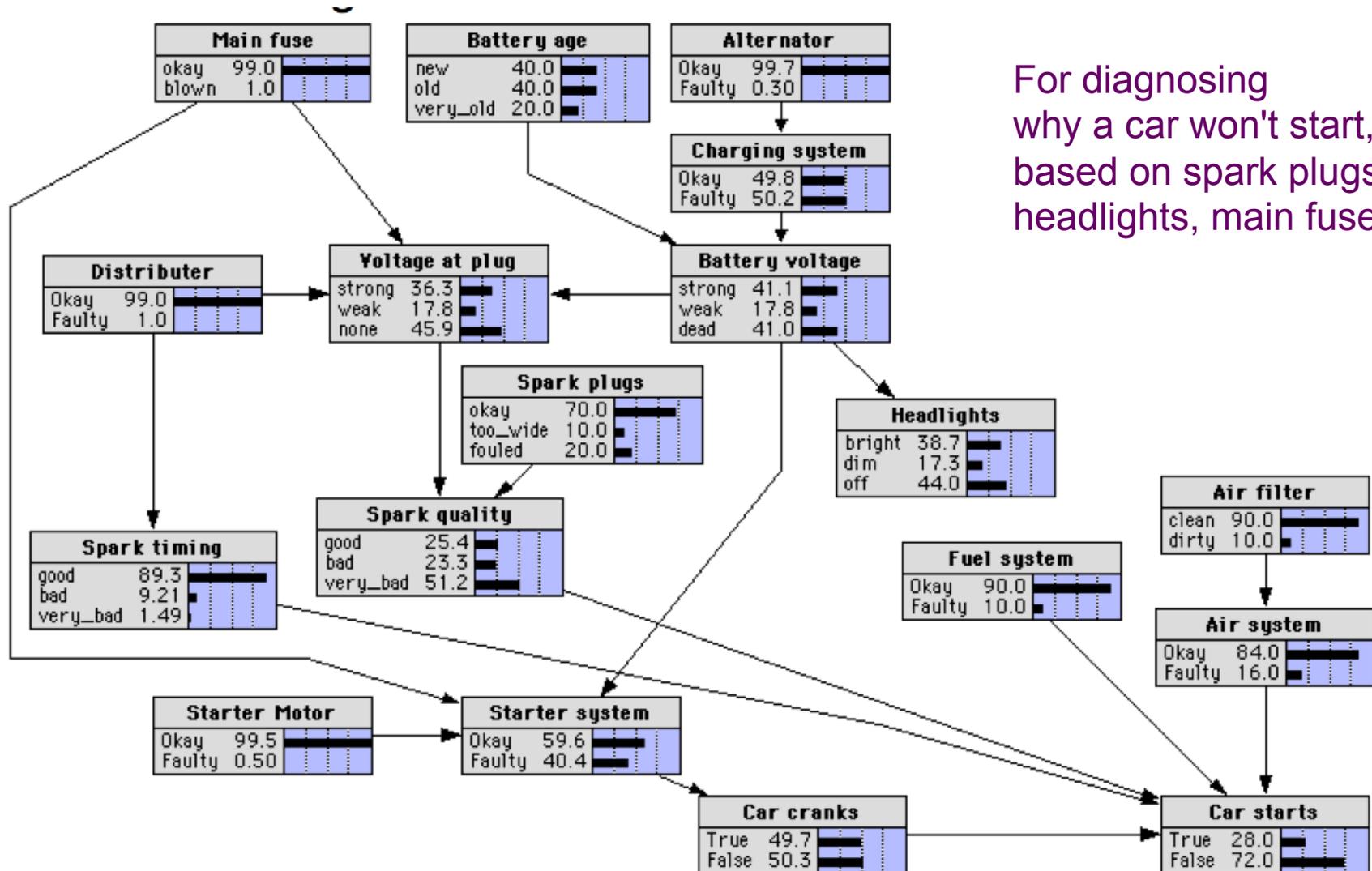
Very weak prior $\alpha=5$

Enough to provide smoothing

Higher values introduce a bias

Effect of high bias disappears
as samples increase

A Static Bayesian Network



For diagnosing why a car won't start, based on spark plugs, headlights, main fuse, etc.