

# Bayesian Parameter Estimation in Bayesian Networks

Sargur Srihari

[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics

## 1. Bayesian network

- where parameters are included as variable nodes

## 2. Global parameter independence

- Leads to global decomposition

## 3. Examples of Bayesian parameter estimation for Bayesian networks

### 1. ICU network

### 2. Text classification

- Naïve Bayes models
- Latent Dirichlet allocation

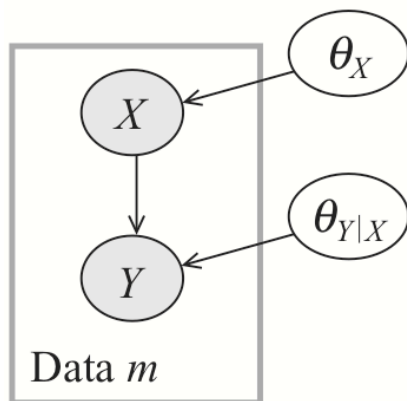
# Inclusion of Parameters as variables

- Bayesian Framework (Approach) requires
  - specify a joint distribution over unknown parameters  $\theta$  and data instances  $D$
- In turn, this joint distribution can be represented as a Bayesian network
  - instances and parameters of variables

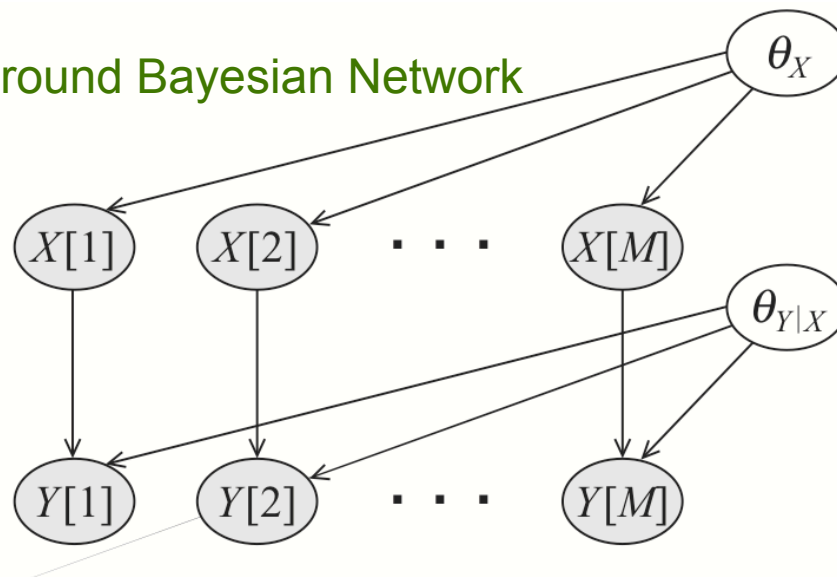
# Simple example with parameters

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

Plate Model



Ground Bayesian Network



# Global Parameter Independence

- If  $\mathcal{G}$  is a Bayesian network with parameters

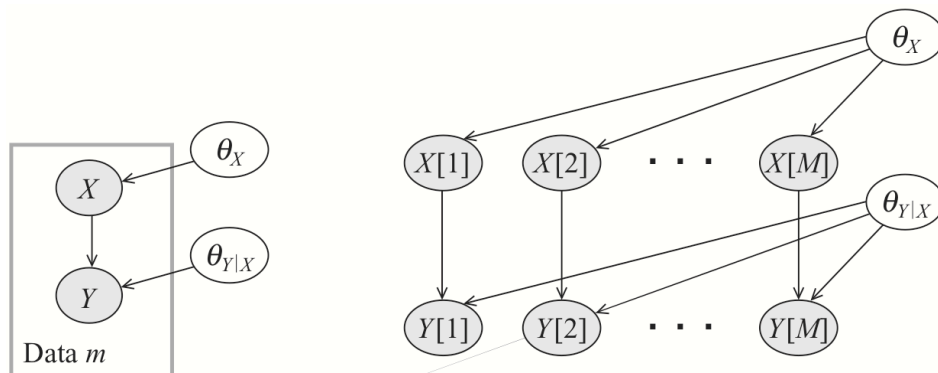
$$\theta = (\theta_{X_1|Pa_{X_1}}, \dots, \theta_{X_n|Pa_{X_n}})$$

- A prior  $P(\theta)$  satisfies global independence if

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

- In prior distribution over parameters, nodes are independent of each other

- Network with global parameter independence:



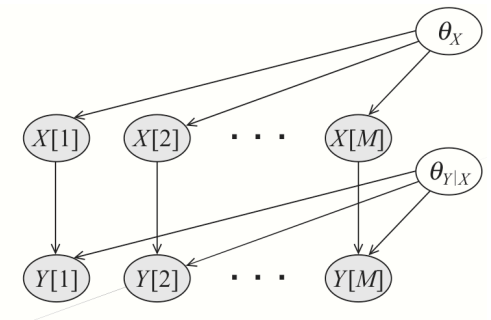
# Use of Global Parameter independence

- Commonly assumed
- But may not always be appropriate
  - Student Example
    - Takes two courses from the same instructor, grade distribution may be the same.
    - So the parameters are dependent

# Global Decomposition

- Assume global parameter independence
- Important conclusion arises
  - Parameters for different CPDs are D-separated
  - In  $X \rightarrow Y$  network
    - If  $x[m]$  and  $y[m]$  are observed for all  $m$  then  $\theta_X$  and  $\theta_{Y|X}$  are d-separated, i.e.,

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



## • Conclusion

- Given data set  $\mathcal{D}$  we can determine posterior over  $\theta_X$  independently of 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  $\theta$
- Need to specify prior  $P(\theta)$  over all possible parameterizations of network

# Posterior Distribution over Parameters

- Given data samples  $\mathcal{D}$

$$P(\theta | \mathcal{D}) = \frac{P(\mathcal{D} | \theta)P(\theta)}{P(\mathcal{D})}$$

- $P(\theta)$  is the Prior distribution
- $P(\mathcal{D} | \theta)$  is the probability of Data given a particular setting
  - Which is the likelihood function
- $P(\mathcal{D})$  Is the normalizing constant
  - Called Marginal likelihood

- Can prove that

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

Conclusion:

Posteriors can be determined independently

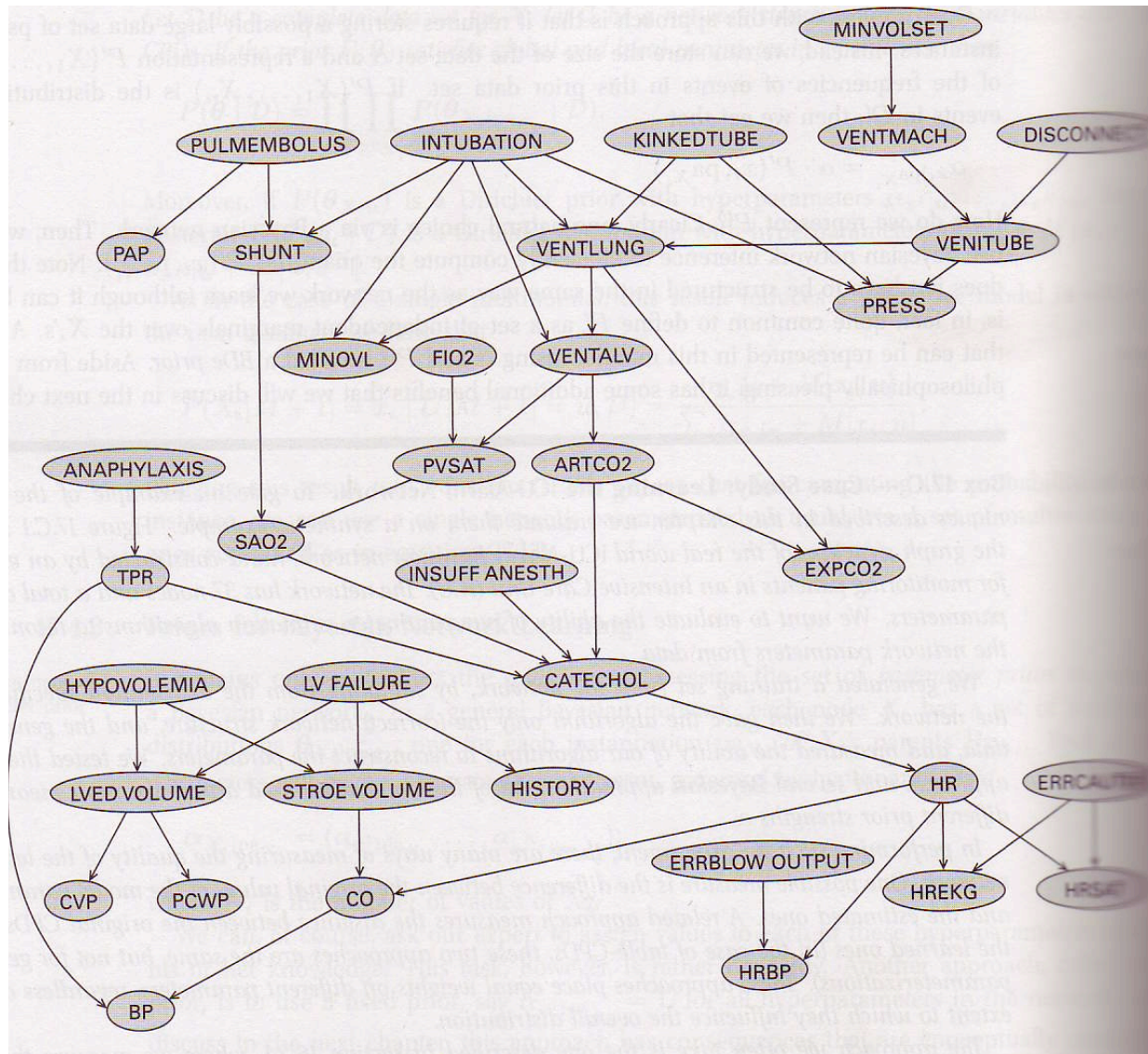
# Prediction problem

- 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

# ICU Alarm Bayesian Network



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

Benchmark for  
Bayesian Network  
Learning

# Parameters for ICU Alarm

- Network parameters
  - 37 nodes with table CPDs
  - 504 parameters
- Evaluate ability of parameter estimation
  - to reconstruct network parameters from data
- Inputs to Parameter Estimation
  1. Training samples
    - obtained by sampling from specified network
  2. Network Structure

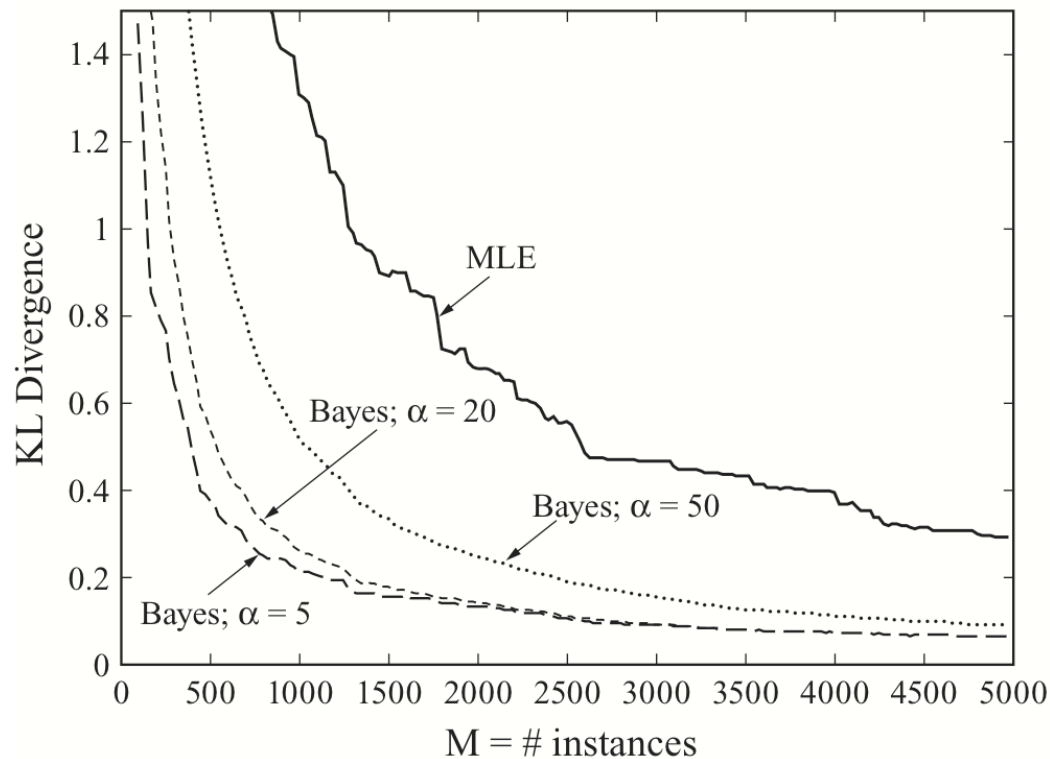
# Evaluation of Learning Methods

- Methods Evaluated
  - MLE approach
  - Several Bayesian approaches
    - Uniform prior mean
    - Different prior strengths  $\alpha=5,20,50$
- Measure of Quality
  1. Difference between original parameters and learned parameters
  2. Distance between original CPDs and estimated ones
  - Two are same in case of table CPDs but not general parameterizations

# 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



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

# Text Classification Networks

- Document to be put into a category
  - Assume document belongs to single category
  - E.g., “Sports”, “Economics”
- Bag of words model
- Distribution of bag in different categories
- Naïve Bayes model
  - There still are design choices affecting performance

# Random Variable Selection

- Remove extraneous characters
  - such as punctuation marks
- Remove stop words (content-free)
  - the, and,...
- Map words to canonical words
  - in pre-defined dictionary  $\mathcal{D}$ 
    - apples → apple
    - used → use
    - running → run

# Two approaches to define features

- Bernoulli Naïve Bayes model

No of features =  
No of words in D

- Binary attribute  $X_i$  indicates whether  $w_i \in D$  appears in document
  - Not how many times it appears in document
- Learn frequency over each document category
  - Probability that “ball” appears in “sports”

- Multinomial Naïve Bayes model

No of features =  
No of positions in doc

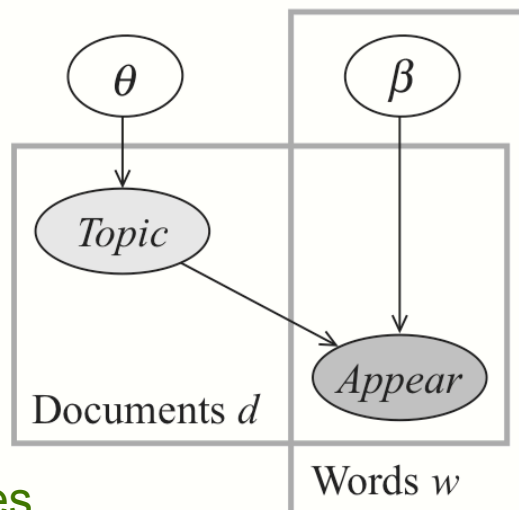
- Attributes describe specific sequence of words
  - attribute  $X_i$  indicates which D-word appears in  $i^{\text{th}}$  pos
    - $X_i$  takes on one of many values, one for each possible word
  - Simplifying assumption (parameter sharing)
    - Prob  $X_i=w$  is same as Prob  $X_j=w$

# Different Plate models for Text

Both associate document  $d$  with Topic  $T$

## Bernoulli Naïve Bayes

No of features =  
No of words in  $D$



Intersecting plates  
for documents  
and dictionary  
words

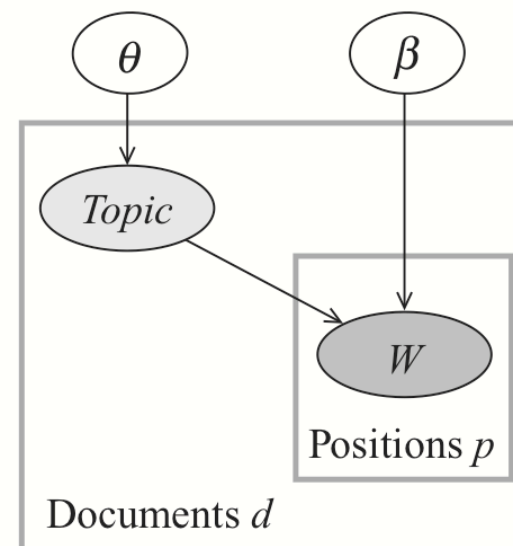
### Binary Attributes

$Appear(d,w)$  for document  $d$  and word  $w$   
takes value  $1$  if  $w$  appears in  $d$

Bernoulli parameter  $\beta_w[w]$  is different for  
different words

## Multinomial Naïve Bayes

No of features =  
No of positions in doc



Multinomial is  
generally more  
successful than  
Bernoulli

$W(d,p)$  for document  $d$  and position  $p$   
takes value of dictionary word

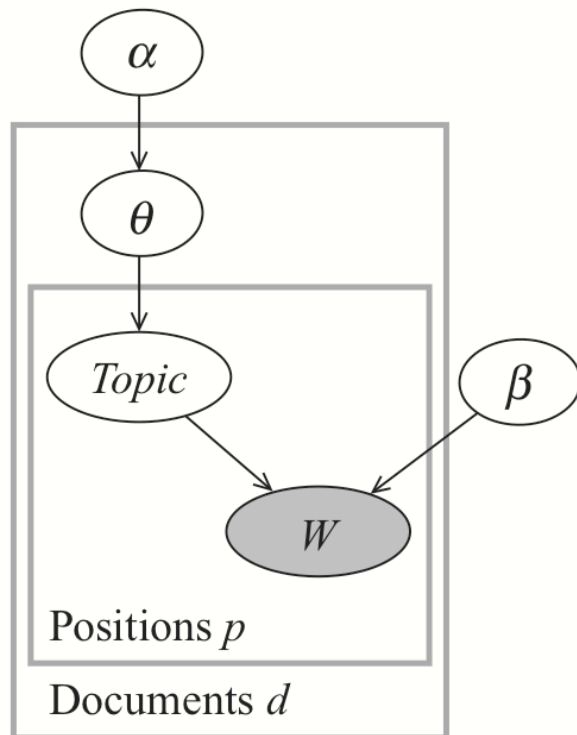
Parameter  $\beta_w$  is the same for all positions

# Parameter Estimation for Text

- In both models
  - Parameter estimated from data
  - Model used for classifying new documents
- Parameters measure prob. of word given topic
  - E.g., “bank” given “economics”
  - Bayesian parameter estimation avoids over-fitting
    - Especially ascribing zero to words not in training set
- With Bayesian estimation
  - can learn naïve Bayes using small corpus
    - Principal advantage of Naïve Bayes
  - More realistic language models are harder

# Latent Dirichlet Allocation

- Extends the multinomial naïve Bayes model
- Set of topics associated with set of multinomial distributions  $\theta_W$  over words



Document  $d$  is associated with a continuous mixture of topics defined using parameters  $\theta(d)$

Parameters selected independently from each document  $d$ , from a Dirichlet distribution parameterized for a set of hyper-parameters  $\alpha$

Word in position  $p$  of document  $d$  is selected by first selecting topic  $Topic(d,p)=t$  from mixture  $\theta(d)$  and then selecting specific dictionary word from Multinomial  $\beta$ , associated with topic  $t$

Avoids over-fitting problems with the two naïve Bayes models

# Summary of BN Parameter Estimation

- Examined parameter estimation for Bayesian networks
  - When data are complete
- Discussed two approaches
  - MLE and Bayesian