

Readings: K&F 17.1, 17.2, 17.3, 17.4



Parameter Estimation

Lecture 9 – Apr 25, 2011
CSE 515, Statistical Methods, Spring 2011

Instructor: Su-In Lee
University of Washington, Seattle

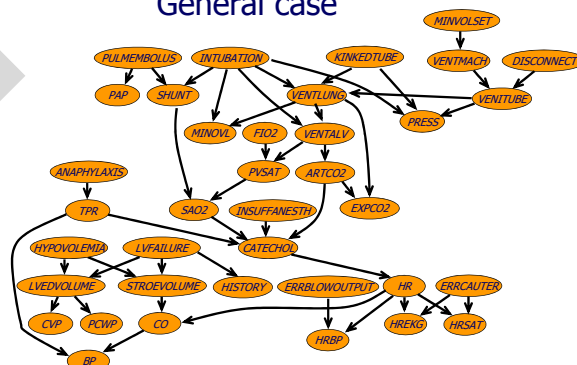
Parameter estimation

- Maximum likelihood estimation (MLE)
 - Parameter estimation based on observations
- Bayesian approach
 - Incorporate our prior knowledge

A single variable
Bayesian network



General case



Maximum Likelihood Estimator

- The *Coin* example – general case

- X: result of a coin toss (head or tail)
- Training data (instances) $D = \langle x[1], \dots, x[m] \rangle$ (M_H heads and M_T tails)
- Parameters: $P(X=h) = \theta$



- Goal:** find $\theta \in [0, 1]$ that predicts the data well

- "Predicts the data well" = likelihood of the data given θ

$$L(\theta : D) = P(D : \theta) = P(x[1], \dots, x[m] | \theta)$$

- MLE: Find θ maximizing likelihood

$$L(\theta : D) = \prod_{i=1}^m P(x[i] | \theta) = \prod_{i=1}^m P(x[i] | \theta) = \theta^{M_H} (1-\theta)^{M_T}$$

- Equivalent to maximizing log-likelihood

$$l(\theta : D) = \log P(D : \theta) = M_H \log \theta + M_T \log(1-\theta)$$

- Differentiating the log-likelihood and solving for θ , we get that the maximum likelihood parameter:

$$\frac{\partial l}{\partial \theta} \Big|_{\theta = \theta_{mle}} = 0$$

$$\theta_{mle} = \arg \max l(\theta : D) = \frac{M_H}{M_H + M_T}$$

3

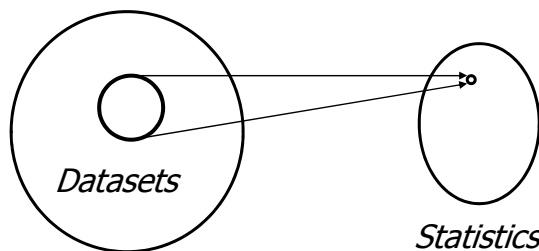
Sufficient Statistics

- For computing the parameter θ of the coin toss example, we only needed M_H and M_T since

$$L(\theta : D) = P(D : \theta) = \theta^{M_H} (1-\theta)^{M_T}$$

→ M_H and M_T are sufficient statistics

D: HHHT
D': THHH



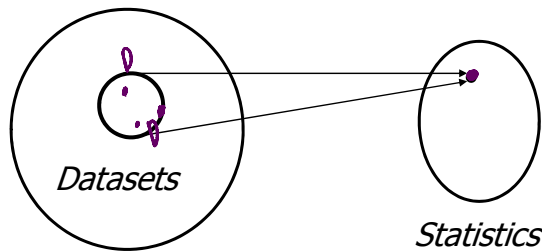
4

Sufficient Statistics

- A function $s(D)$ is a **sufficient statistic** from instances to a vector in \mathbb{R}^k if, for any two datasets D and D' and any $\theta \in \Theta$, we have

$$\sum_{x[i] \in D} s(x[i]) = \sum_{x[i] \in D'} s(x[i]) \Rightarrow L(D; \theta) = L(D'; \theta)$$

- We often refer to the tuple $\sum_{x[i] \in D} s(x[i])$ as the **sufficient statistics** of the data set D .
 - In coin toss experiment, M_0 and M_1 are sufficient statistics



5

Sufficient Statistics for Multinomial

- Y : multinomial, k values (e.g. result of a dice throw)
 $k=6$

- A **sufficient statistics** for a dataset D over Y is the tuple of counts $\langle M_1, \dots, M_k \rangle$ such that M_i is the number of times that the $Y=y^i$ in D

- Likelihood function:** $L(D; \theta) = \prod_{i=1}^k \theta_i^{M_i}$ where $\theta_i = P(Y = y^i)$

- MLE Principle:** Choose Θ that maximize $L(D; \Theta)$

- Multinomial MLE:** $\theta^i = \frac{M_i}{\sum_{i=1}^m M_i}$

6

Sufficient Statistic for Gaussian

- Gaussian distribution: $X \sim N(\mu, \sigma^2)$
 - Probability density function (pdf): $p(X) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$
 - Rewrite as $p(X) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}x^2 + \frac{x\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2}}$
- sufficient statistics for Gaussian: $\langle M(\sum_m x[m]), \sum_m x[m]^2 \rangle$

- MLE Principle: Choose Θ that maximize $L(D:\Theta)$

- Multinomial MLE: $\mu = \frac{1}{M} \sum_m x[m]$
 $\sigma = \sqrt{\frac{1}{M} \sum_m (x[m] - \mu)^2}$

7

MLE for Bayesian Networks

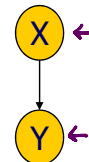
- Parameters
 - θ_{x0}, θ_{x1} θ_x
 - $\theta_{y0|x0}, \theta_{y1|x0}, \theta_{y0|x1}, \theta_{y1|x1}$ $\theta_{y|x}$

- Training data:
 - tuples $\langle x[m], y[m] \rangle$ $m=1, \dots, M$

- Likelihood function:

$$\begin{aligned}
 L(D:\theta) &= \prod_{m=1}^M P(x[m], y[m] : \theta) \\
 &= \prod_{m=1}^M P(x[m] : \theta_x) P(y[m] | x[m] : \theta_{y|x}) \\
 &= \left(\prod_{m=1}^M P(x[m] : \theta_x) \right) \left(\prod_{m=1}^M P(y[m] | x[m] : \theta_{y|x}) \right)
 \end{aligned}$$

X	
x^0	x^1
0.7	0.3



x	y	
	y^0	y^1
x^0	0.95	0.05
x^1	0.2	0.8

→ Likelihood decomposes into two separate terms, one for each variable ("decomposability of the likelihood function")

8

MLE for Bayesian Networks

- Terms further decompose by CPDs: $\langle x[m], y[m] \rangle$

$$\prod_{m=1}^M P(y[m] | x[m] : \theta) = \left[\prod_{m: x[m]=x^0} P(y[m] | x[m] : \theta_{Y|X}) \right] \left[\prod_{m: x[m]=x^1} P(y[m] | x[m] : \theta_{Y|X}) \right]$$

$$= \prod_{m: x[m]=x^0} P(y[m] | x[m] : \theta_{Y|X^0}) \prod_{m: x[m]=x^1} P(y[m] | x[m] : \theta_{Y|X^1})$$

- By sufficient statistics

$$\prod_{m: x[m]=x^1} P(y[m] | x[m] : \theta_{Y|X^1}) = \theta_{y^0|x^1}^{M[x^1, y^0]} \cdot \theta_{y^1|x^1}^{M[x^1, y^1]}$$

$\theta_{y^0|x^1}$ or $\theta_{y^1|x^1}$ depends on $y[m]$

where $M[x^1, y^1]$ is the number of data instances in which X takes the value x^1 and Y takes the value y^1

- MLE

$$\theta_{y^0|x^1} = \frac{M[x^1, y^0]}{M[x^1, y^0] + M[x^1, y^1]} = \frac{M[x^1, y^0]}{M[x^1]}$$

9

MLE for Bayesian Networks

- Likelihood for Bayesian network

$$L(\Theta : D) = \prod_m P(x[m] : \Theta)$$

$$= \prod_m \prod_i P(x_i[m] | Pa_i[m] : \Theta_i)$$

$$= \prod_i \left[\prod_m P(x_i[m] | Pa_i[m] : \Theta_i) \right]$$

$$= \prod_i L_i(\Theta_{x_i|Pa_i} : X_i, Pa_i)$$

Conditional likelihood or "Local likelihood"

→ if $\Theta_{x_i|Pa_i}$ are disjoint then MLE can be computed by maximizing each local likelihood separately

$$\theta_{x_i|Pa_i} = \theta_{x_i|Pa_i} \quad L_i(\theta_{x_i|Pa_i} : x_i, Pa_i)$$

10

MLE for Table CPD BayesNets

- Multinomial CPD

$$L_Y(D : \theta_{Y|X}) = \prod_m \theta_{y[m]|x[m]} \cdot \prod_{x \in \text{Val}(X)} \prod_{y \in \text{Val}(Y)} \theta_{y|x}^{M[x,y]}$$

Handwritten notes: $\{x \in \text{Val}(X) \rightarrow y \in \text{Val}(Y)\} \rightarrow D$, $\theta_{y|x}$ (circled), $\frac{PCT(x,y)}{\theta_{y|x}}$

- For each value $x \in X$ we get an independent multinomial problem where the MLE is

$$\theta_{y|x} = \frac{M[x, y^i]}{M[x]}$$

$\sum_Y \theta_{y|x} = 1$

$\theta_{\text{blue}} + \theta_{\text{green}} = 1$
 $\theta_{\text{blue}} \cdot (1 - \theta_{\text{blue}})$

MLE for Tree CPDs

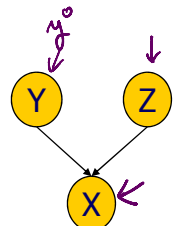
- Assume tree CPD with known tree structure

$$L(D : \theta_{X|Y,Z}) = \prod_{x,y,z} \theta_{x|y,z}^{M[x,y,z]}$$

$$= \prod_x \left(\theta_{x|y^0,z^0}^{M[x,y^0,z^0]} \cdot \theta_{x|y^0,z^1}^{M[x,y^0,z^1]} \cdot \theta_{x|y^1,z^0}^{M[x,y^1,z^0]} \cdot \theta_{x|y^1,z^1}^{M[x,y^1,z^1]} \right)$$

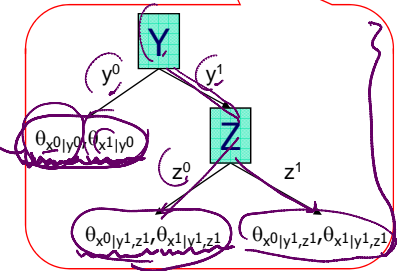
$$= \prod_x \left(\theta_{x|y^0}^{M[x,y^0,z^0]+M[x,y^0,z^1]} \cdot \theta_{x|y^1,z^0}^{M[x,y^1,z^0]} \cdot \theta_{x|y^1,z^1}^{M[x,y^1,z^1]} \right)$$

$$= \prod_x \left(\theta_{x|y^0}^{M[x,y^0,x]} \cdot \theta_{x|y^1,z^0}^{M[x,y^1,z^0,x]} \cdot \theta_{x|y^1,z^1}^{M[x,y^1,z^1,x]} \right)$$



Terms for $\langle y^0, z^0 \rangle$ and $\langle y^0, z^1 \rangle$ can be combined

Optimization can be done by leaves



MLE for Tree CPD BayesNets

- Tree CPD T , leaves l

$$\begin{aligned}
 L_Y(D; \theta_{Y|X}) &= \prod_{m=1}^m P(y[m] | \mathbf{x}[m]; \theta_{Y|X}) \\
 &= \prod_{m=1}^m \theta_{y[m] | l(\mathbf{x}[m])} \\
 &= \prod_{l \in \text{Leaves}(T)} \left[\prod_{y \in \text{Val}(Y)} \theta_{y|l}^{M[c_l, y]} \right]
 \end{aligned}$$

- For each value $l \in \text{Leaves}(T)$ we get an independent multinomial problem where the MLE is

$$\theta_{y^i|l} = \frac{M[c_l, y^i]}{M[c_l]} \quad M[c_l] = \sum_{x:l(x)=l} M[x, y^i]$$

13

Limitations of MLE

- A thumbtack is tossed 10 times, and comes out 'head' 3 of the 10 tosses \rightarrow Probability of head = 0.3 $\frac{3}{10}$
- A coin is tossed 10 times, and comes out 'head' 3 of the 10 tosses \rightarrow Probability of head = 0.3
- A coin is tossed 1,000,000 times, and comes out 'head' 300,000 of the 1,000,000 tosses \rightarrow Probability of head = 0.3
- Would you place the same bet on the next thumbtack toss as you would on the next coin toss?
- We need to incorporate **prior knowledge**
 - Prior knowledge should only be used as a guide

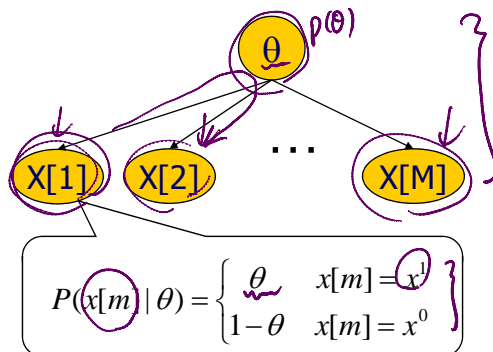
14

Bayesian Inference

Assumptions

- Given a fixed θ tosses are independent
- If θ is unknown tosses are not marginally independent – each toss tells us something about θ

- The following network captures our assumptions



15

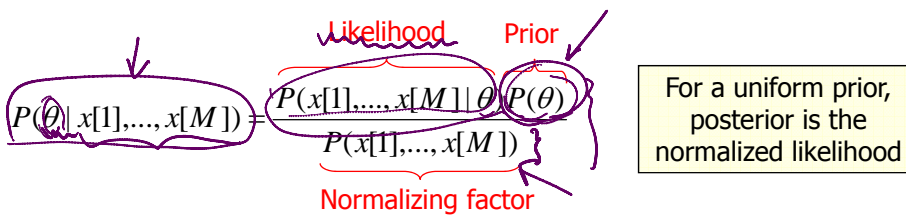
Bayesian Inference

Joint probabilistic model

$$\begin{aligned}
 P(x[1], \dots, x[M], \theta) &= P(x[1], \dots, x[M] | \theta) P(\theta) \\
 &= P(\theta) \prod_{i=1}^M P(x[i] | \theta) \\
 &= P(\theta) \theta^{M_1} (1 - \theta)^{M_0}
 \end{aligned}$$

Handwritten notes include $D: x[1] \dots x[M]$ and M_1 pointing to the product term, and M_0 pointing to the $(1-\theta)$ term. A small tree diagram shows θ as the parent of $x[1], x[2], \dots, x[M]$.

Posterior probability over θ



16

Bayesian Prediction

- Predict the data instance from the previous ones

$$\begin{aligned}
 &P(x[M+1] | x[1], \dots, x[M]) \\
 &\propto \int P(x[M+1], \theta | x[1], \dots, x[M]) d\theta \\
 &= \int P(x[M+1] | x[1], \dots, x[M], \theta) P(\theta | x[1], \dots, x[M]) d\theta \\
 &= \int P(x[M+1] | \theta) P(\theta | x[1], \dots, x[M]) d\theta
 \end{aligned}$$

$P(\theta) = 1$
 $\propto L(\theta; D) = \theta^{M_H} (1-\theta)^{M_T}$

- Solve for uniform prior $P(\theta) = 1$ (for $0 \leq \theta \leq 1$) and binomial variable

$$\frac{P(x[M+1] = x^1 | x[1], \dots, x[M])}{P(x[1], \dots, x[M])} = \frac{1}{\int_0^1 \theta \cdot \theta^{M_H} \cdot (1-\theta)^{M_T} d\theta}$$

"Bayesian estimate" $\rightarrow \frac{M_H + 1}{M_H + M_T + 2}$ "Imaginary counts"

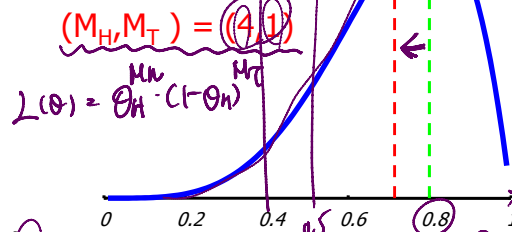
Example: Binomial Data

- Prior: uniform for θ in $[0, 1]$

$P(\theta) = 1$ $0 \leq \theta \leq 1$

$\rightarrow P(\theta | D)$ is proportional to the likelihood $L(D; \theta)$

$$P(\theta | x[1], \dots, x[M]) \propto P(x[1], \dots, x[M] | \theta)$$



- MLE for $P(X=H)$ is $4/5 = 0.8$

- Bayesian prediction is $5/7 = 0.71$

$$P(x[M+1] = H | D) = \int \theta \cdot P(\theta | D) d\theta = \frac{5}{7} = 0.7142 \dots$$

Dirichlet Priors

- A **Dirichlet prior** is specified by a set of (non-negative) hyper-parameters $\alpha_1, \dots, \alpha_k$ so that

$\theta = [\theta_1, \dots, \theta_k] \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k)$ if

$$p(\theta) = \frac{1}{Z} \prod_k \theta_k^{\alpha_k - 1} \quad \text{where} \quad \sum_k \theta_k = 1, \quad \Gamma(x) = \int_0^{\infty} t^{x-1} e^{-t} dt$$

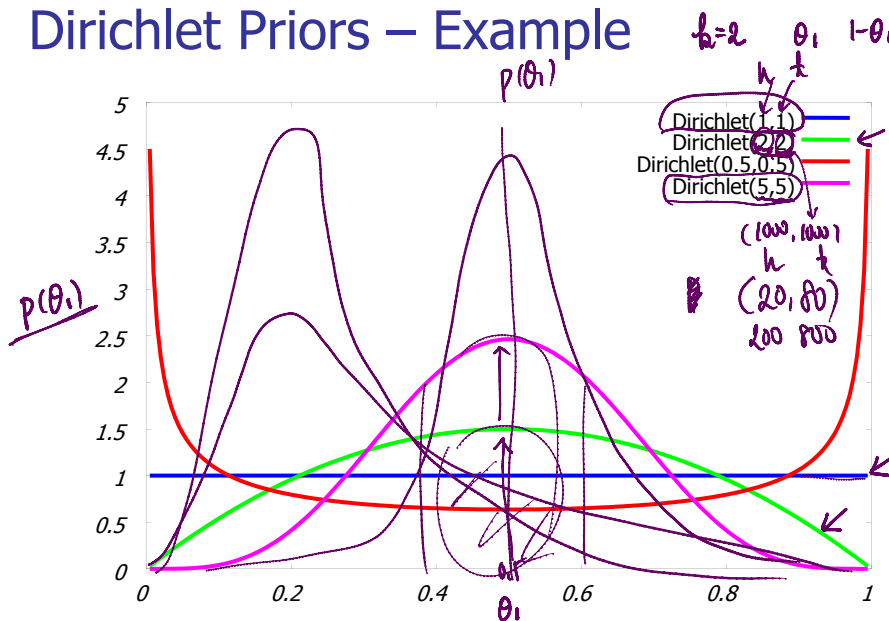
$$\int p(\theta) d\theta = 1 \quad \text{and} \quad Z = \frac{\prod_{i=1}^k \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^k \alpha_i)}$$

- Intuitively, hyper-parameters correspond to the number of imaginary counts before starting the coin toss experiment

$p(\theta)$ $h=2$ $\theta = [\theta_1, \theta_2]$
 α_1, α_2 $1-\theta_1$

19

Dirichlet Priors – Example



20

Dirichlet Priors

- Dirichlet priors have the property that the posterior is also Dirichlet

- Prior is $\text{Dir}(\alpha_1, \dots, \alpha_k)$ $p(\theta) = \frac{1}{Z} \prod_k \theta_k^{\alpha_k - 1}$
- Data counts are $\langle M_1, \dots, M_k \rangle$ D
- Posterior is $\text{Dir}(\alpha_1 + M_1, \dots, \alpha_k + M_k)$ $p(\theta | D) = \frac{1}{Z} \prod_k \theta_k^{\alpha_k + M_k - 1}$

- The hyperparameters $\alpha_1, \dots, \alpha_k$ can be thought of as "imaginary" counts from our prior experience

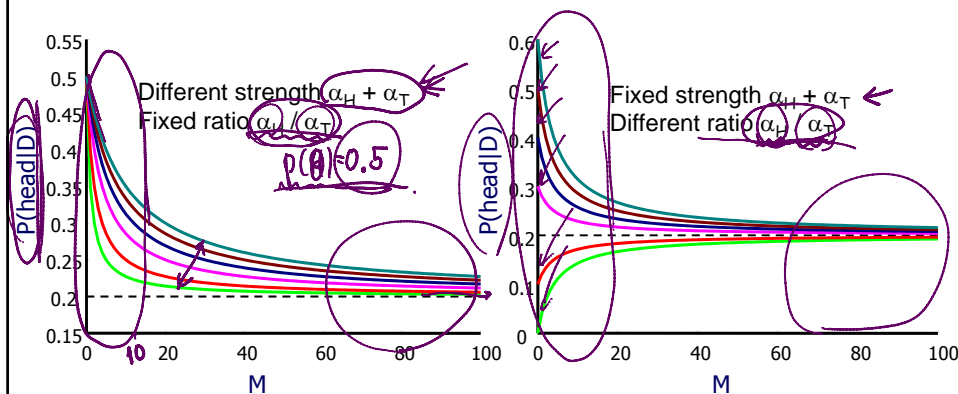
- Equivalent sample size = $\alpha_1 + \dots + \alpha_k$
 - The larger the equivalent sample size the more confident we are in our prior

21

Effect of Priors

- Prediction of $P(X=H)$ after seeing data with $M_H=0.2M$, $M_T=0.8M$ as a function of the sample size

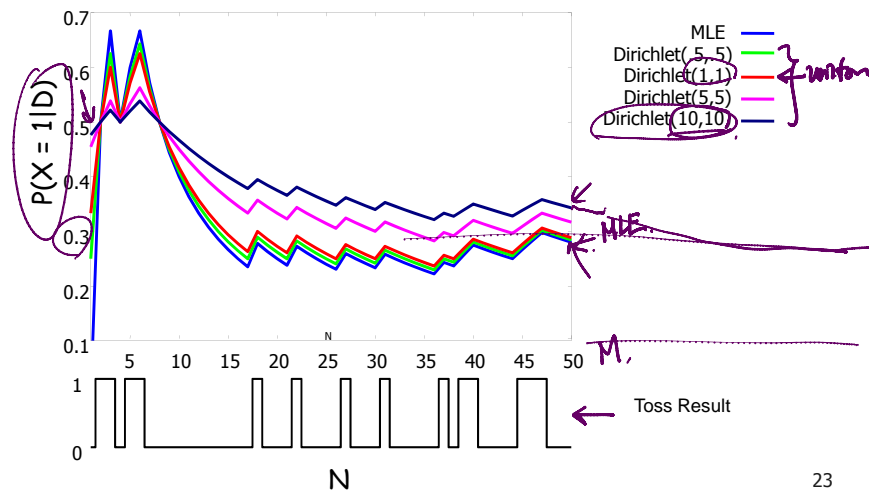
$$p(X=H) = 0.2$$



22

Effect of Priors (cont.)

- In real data, Bayesian estimates are less sensitive to noise in the data



General Formulation

- Joint distribution over D, θ

$$P(D, \theta) = P(D | \theta) P(\theta)$$

- Posterior distribution over parameters

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

- $P(D)$ is the marginal likelihood of the data

$$P(D) = \int_{\theta} P(D | \theta) P(\theta) d\theta$$

- As we saw, likelihood can be described compactly using sufficient statistics
- We want conditions in which posterior is also compact
 - E.g. Dirichlet priors

Conjugate Families

- A family of priors $P(\theta:\alpha)$ is **conjugate** to a model $P(\xi|\theta)$ if for any possible dataset D of i.i.d samples from $P(\xi|\theta)$ and choice of hyperparameters α for the prior over θ , there are hyperparameters α' that describe the posterior, i.e.,

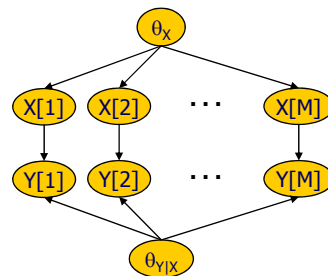
$$P(\theta:\alpha') \propto P(D|\theta)P(\theta:\alpha)$$

- Posterior has the same parametric form as the prior
 - Dirichlet prior is a **conjugate family** for the multinomial likelihood
- Conjugate families are useful since:
 - Many distributions can be represented with hyperparameters
 - They allow for sequential update within the same representation
 - In many cases we have closed-form solutions for prediction

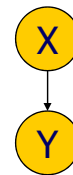
25

Bayesian Estimation in BayesNets

Bayesian network for parameter estimation



Bayesian network

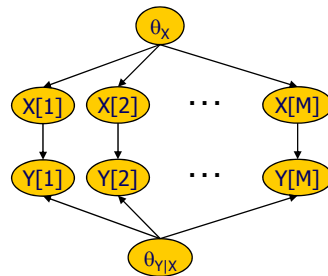


- **Instances are independent given the parameters**
 - $(x[m'], y[m'])$ are d-separated from $(x[m], y[m])$ given θ
- **Priors for individual variables are a priori independent**
 - Global independence of parameters $P(\theta) = \prod_i P(\theta_{X_i} | P_{\alpha(X_i)})$

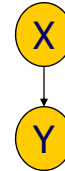
26

Bayesian Estimation in BayesNets

Bayesian network for parameter estimation



Bayesian network

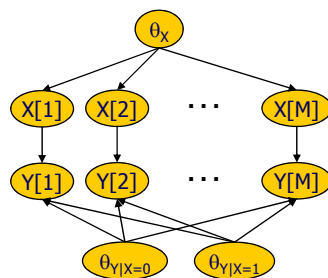


- **Posteriors of θ are independent given complete data**
 - Complete data d-separates parameters for different CPDs
 - $P(\theta_x, \theta_{Y|X} | D) = P(\theta_x | D)P(\theta_{Y|X} | D)$
 - As in MLE, we can solve each estimation problem separately

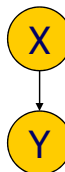
27

Bayesian Estimation in BayesNets

Bayesian network for parameter estimation



Bayesian network

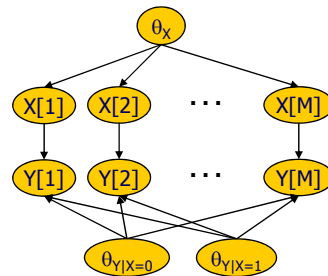


- **Posteriors of θ are independent given complete data**
 - Also holds for parameters within families
 - Note **context specific independence** between $\theta_{Y|X=0}$ and $\theta_{Y|X=1}$ when given both X and Y

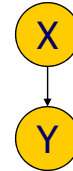
28

Bayesian Estimation in BayesNets

Bayesian network for parameter estimation



Bayesian network



- **Posteriors of θ can be computed independently**
 - For multinomial $\theta_{x_i|pa_i}$ posterior is Dirichlet with parameters $(\alpha_{x_i=1|pa_i} + M[X_i=1|pa_i]), \dots, (\alpha_{x_i=k|pa_i} + M[X_i=k|pa_i])$
 - $$P(X_i[M+1]=x_i | Pa_i[M+1]=pa_i, D) = \frac{\alpha_{x_i|pa_i} + M[x_i, pa_i]}{\sum \alpha_{x_i|pa_i} + M[x_i, pa_i]}$$

29

Assessing Priors for BayesNets

- We need the $\alpha(x_i, pa_i)$ for each node x_i
- We can use initial parameters Θ_0 as prior information
 - Need also an equivalent sample size parameter M'
 - Then, we let $\alpha(x_i, pa_i) = M' \cdot P(x_i, pa_i | \Theta_0)$
- This allows to update a network using new data
 - **Example network for priors**
 - $P(X=0)=P(X=1)=0.5$
 - $P(Y=0)=P(Y=1)=0.5$
 - $M'=1$
 - Note: $\alpha(x_0)=0.5 \alpha(x_0, y_0)=0.25$



30

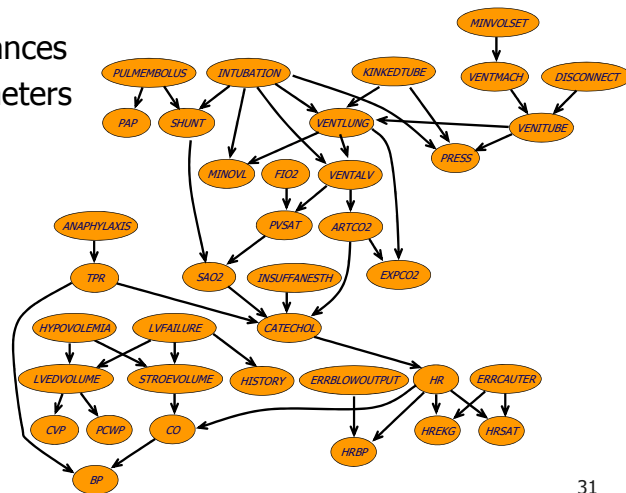
Case Study: ICU Alarm Network

- The "Alarm" network

- 37 variables

- Experiment

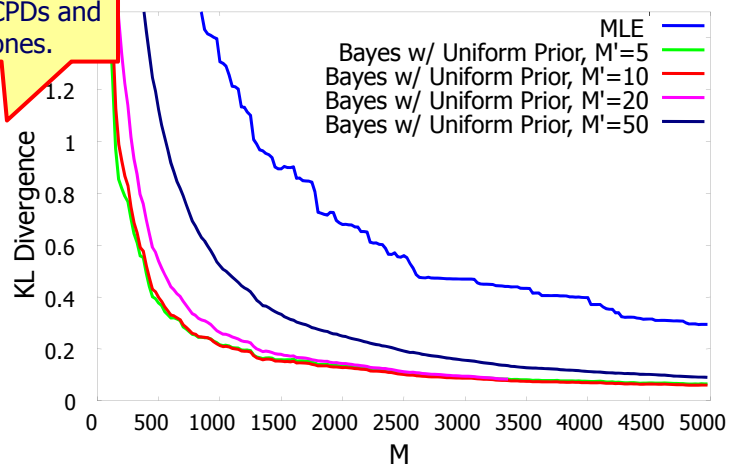
- Sample instances
 - Learn parameters
 - MLE
 - Bayesian



31

Case Study: ICU Alarm Network

The distance between the original CPDs and the learned ones.



- MLE performs worst
- Prior $M'=5$ provides best smoothing

32

Parameter Estimation Summary

- Estimation relies on **sufficient statistics**
 - For multinomials these are of the form $M[x_i, pa_i]$
 - Parameter estimation

$$\hat{\theta}_{x_i|pa_i} = \frac{M[x_i, pa_i]}{M[pa_i]} \quad P(x_i | pa_i, D) = \frac{\alpha_{x_i, pa_i} + M[x_i, pa_i]}{\alpha_{pa_i} + M[pa_i]}$$

MLE

Bayesian (Dirichlet)

- Bayesian methods also require choice of priors
- MLE and Bayesian are asymptotically equivalent
- Both can be implemented in an **online** manner by accumulating sufficient statistics

33

Acknowledgement

- These lecture notes were generated based on the slides from Prof Eran Segal.