

Readings:

K&F: 14.1, 14.2, 14.3, 14.4, 15.1, 15.2, 15.3.1, 15.4.1

## Parameter Learning 2

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

September 27<sup>th</sup>, 2006

1

### Your first learning algorithm

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} \ln P(\mathcal{D} | \theta) \\ &= \arg \max_{\theta} \ln \theta^{\alpha_H} (1 - \theta)^{\alpha_T}\end{aligned}$$

$$\ln a^b = b \ln a$$

$$\ln a \cdot b = \ln a + \ln b$$

$$\frac{\partial}{\partial \theta} \ln \theta = \frac{1}{\theta}$$

$$\frac{\partial}{\partial \theta} \ln(1-\theta) = \frac{-1}{1-\theta}$$

■ Set derivative to zero:

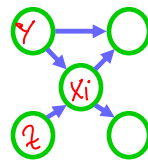
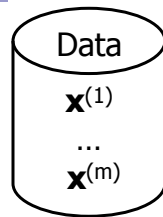
$$\frac{d}{d\theta} \ln P(\mathcal{D} | \theta) = 0$$

$$\begin{aligned}& \frac{\partial}{\partial \theta} \ln[\theta^{\alpha_H} (1-\theta)^{\alpha_T}] \\ &= \frac{\partial}{\partial \theta} [\alpha_H \ln \theta + \alpha_T \ln(1-\theta)] = \frac{\partial}{\partial \theta} \alpha_H \ln \theta + \frac{\partial}{\partial \theta} \alpha_T \ln(1-\theta) \\ &= \frac{\alpha_H}{\theta} - \frac{\alpha_T}{1-\theta} = 0 \Rightarrow \theta = \frac{\alpha_H}{\alpha_H + \alpha_T}\end{aligned}$$

10-708 – ©Carlos Guestrin 2006

2

# Learning the CPTs



For each discrete variable  $X_i$

$$P(X_i | \text{Pa}_{X_i}) = P(x_i | y, z)$$

$$\underset{\text{MLE}}{\approx} P(x_i = x_i | y = y, z = z) \underset{\text{MLE}}{\approx} \frac{\text{Count}(X_i = x_i, Y = y, Z = z)}{\text{Count}(Y = y, Z = z)}$$

$$\text{MLE: } P(X_i = x_i | X_j = x_j) = \frac{\text{Count}(X_i = x_i, X_j = x_j)}{\text{Count}(X_j = x_j)}$$

10-708 - ©Carlos Guestrin 2006

3

## Maximum likelihood estimation (MLE) of BN parameters – General case

- Data:  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$
- Restriction:  $\mathbf{x}^{(i)}[\text{Pa}_{X_i}] \rightarrow$  assignment to  $\text{Pa}_{X_i}$  in  $\mathbf{x}^{(i)}$
- Given structure, log likelihood of data:

$$\log P(\mathcal{D} | \theta_{\mathcal{G}}, \mathcal{G}) = \log \prod_i P(x_i = x_i^{(i)} | \text{Pa}_{X_i} = x^{(i)}[\text{Pa}_{X_i}])$$

10-708 - ©Carlos Guestrin 2006

4

## Taking derivatives of MLE of BN parameters – General case

$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G}) = \sum_{j=1}^m \sum_{i=1}^n \log P\left(X_i = x_i^{(j)} \mid \mathbf{Pa}_{X_i} = \mathbf{x}^{(j)} [\mathbf{Pa}_{X_i}]\right)$$

10-708 – ©Carlos Guestrin 2006

5

## General MLE for a CPT

- Take a CPT:  $P(X|\mathbf{U})$
- Log likelihood term for this CPT
  
- Parameter  $\theta_{X=x|\mathbf{U}=\mathbf{u}}$  :

$$\text{MLE: } P(X = x \mid \mathbf{U} = \mathbf{u}) = \theta_{X=x|\mathbf{U}=\mathbf{u}} = \frac{\text{Count}(X = x, \mathbf{U} = \mathbf{u})}{\text{Count}(\mathbf{U} = \mathbf{u})}$$

10-708 – ©Carlos Guestrin 2006

6

# Announcements

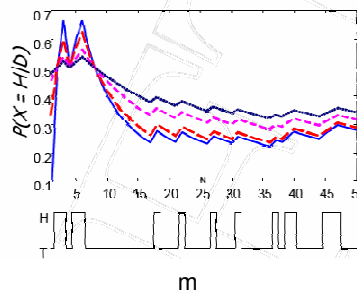
- Late homeworks:
  - 3 late days for the semester
    - one late day corresponds to 24 hours! (i.e., 3 late days due Saturday by noon)
    - Give late homeworks to Monica Hopes, Wean Hall 4619
      - If she is not in her office, time stamp (date and time) your homework, sign it, and put it under her door
  - After late days are used up:
    - Half credit within 48 hours
    - Zero credit after 48 hours
  - All homeworks **must be handed in**, even for zero credit
- Homework 2 out later today
- Recitation tomorrow:
  - review perfect maps, parameter learning

10-708 - ©Carlos Guestrin 2006

7

## Can we really trust MLE?

- What is better?
  - 3 heads, 2 tails
  - 30 heads, 20 tails
  - $3 \times 10^{23}$  heads,  $2 \times 10^{23}$  tails
- Many possible answers, we need distributions over possible parameters



10-708 - ©Carlos Guestrin 2006

8

# Bayesian Learning

- Use Bayes rule:

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

- Or equivalently:

$$P(\theta | \mathcal{D}) \propto P(\mathcal{D} | \theta)P(\theta)$$

10-708 – ©Carlos Guestrin 2006

9

# Bayesian Learning for Thumbtack

$$P(\theta | \mathcal{D}) \propto P(\mathcal{D} | \theta)P(\theta)$$

- Likelihood function is simply Binomial:

$$P(\mathcal{D} | \theta) = \theta^{m_H} (1 - \theta)^{m_T}$$

- What about prior?

- ☐ Represent expert knowledge
- ☐ Simple posterior form

- Conjugate priors:

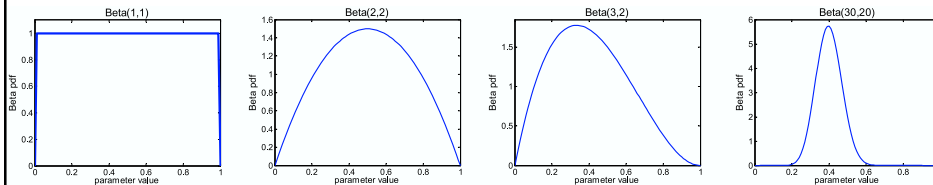
- ☐ Closed-form representation of posterior (more details soon)
- ☐ **For Binomial, conjugate prior is Beta distribution**

10-708 – ©Carlos Guestrin 2006

10

## Beta prior distribution – $P(\theta)$

$$P(\theta) = \frac{\theta^{\alpha_H-1}(1-\theta)^{\alpha_T-1}}{B(\alpha_H, \alpha_T)} \sim \text{Beta}(\alpha_H, \alpha_T)$$



- Likelihood function:  $P(\mathcal{D} | \theta) = \theta^{m_H}(1-\theta)^{m_T}$
- Posterior:  $P(\theta | \mathcal{D}) \propto P(\mathcal{D} | \theta)P(\theta)$

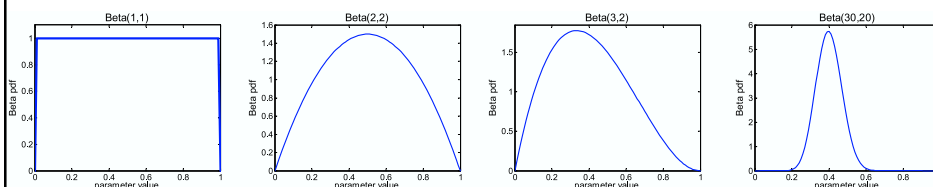
10.708 – ©Carlos Guestrin 2006

11

## Posterior distribution

- Prior:  $\text{Beta}(\alpha_H, \alpha_T)$
- Data:  $m_H$  heads and  $m_T$  tails
- Posterior distribution:

$$P(\theta | \mathcal{D}) \sim \text{Beta}(m_H + \alpha_H, m_T + \alpha_T)$$



10.708 – ©Carlos Guestrin 2006

12

# Conjugate prior

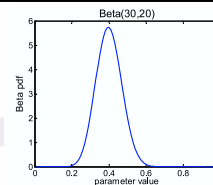
- Prior:  $Beta(\alpha_H, \alpha_T)$
- Data:  $m_H$  heads and  $m_T$  tails (binomial likelihood)
- Posterior distribution:  
$$P(\theta | \mathcal{D}) \sim Beta(m_H + \alpha_H, m_T + \alpha_T)$$
- Given likelihood function  $P(D|\theta)$
- (Parametric) prior of the form  $P(\theta|\alpha)$  is **conjugate** to likelihood function if posterior is of the same parametric family, and can be written as:
  - $P(\theta|\alpha')$ , for some new set of parameters  $\alpha'$

10-708 - ©Carlos Guestrin 2006

13

# Using Bayesian posterior

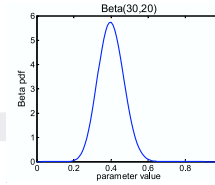
- Posterior distribution:  
$$P(\theta | \mathcal{D}) \sim Beta(m_H + \alpha_H, m_T + \alpha_T)$$
- Bayesian inference:
  - No longer single parameter:  
$$E[f(\theta)] = \int_0^1 f(\theta) P(\theta | \mathcal{D}) d\theta$$
  - Integral is often hard to compute



10-708 - ©Carlos Guestrin 2006

14

# Bayesian prediction of a new coin flip



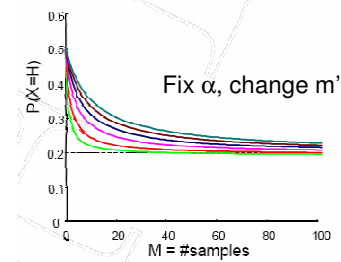
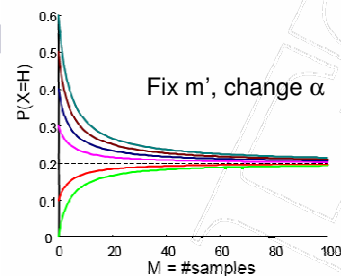
- Prior:
- Observed  $m_H$  heads,  $m_T$  tails, what is probability of  $m+1$  flip is heads?

10-708 – ©Carlos Guestrin 2006

15

## Asymptotic behavior and equivalent sample size

- Beta prior equivalent to extra thumbtack flips:
  - $$E[\theta] = \frac{m_H + \alpha_H}{m_H + \alpha_H + m_T + \alpha_T}$$
- As  $m \rightarrow \infty$ , prior is “forgotten”
- **But, for small sample size, prior is important!**
- **Equivalent sample size:**
  - Prior parameterized by  $\alpha_H, \alpha_T$ , or
  - $m'$  (equivalent sample size) and  $\alpha$
  - $$E[\theta] = \frac{m_H + \alpha m'}{m_H + m_T + m'}$$



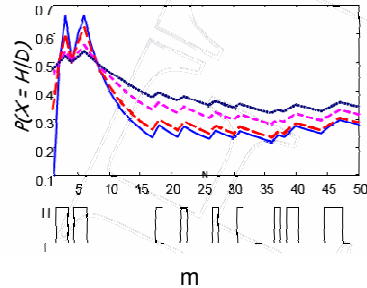
10-708 – ©Carlos Guestrin 2006

16



## Bayesian learning corresponds to smoothing

$$E[\theta] = \frac{m_H + \alpha m'}{m_H + m_T + m'}$$



- $m=0 \Rightarrow$  prior parameter
- $m \rightarrow \infty \Rightarrow$  MLE

10-708 - ©Carlos Guestrin 2006

17

## Bayesian learning for multinomial

- What if you have a  $k$  sided coin???
- Likelihood function if **multinomial**:
  - 
  -
- **Conjugate** prior for multinomial is **Dirichlet**:
  - $\theta \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k) \sim \prod_i \theta_i^{\alpha_i - 1}$
- **Observe**  $m$  data points,  $m_i$  from assignment  $i$ , **posterior**:
- **Prediction**:

10-708 - ©Carlos Guestrin 2006

18

## Bayesian learning for two-node BN

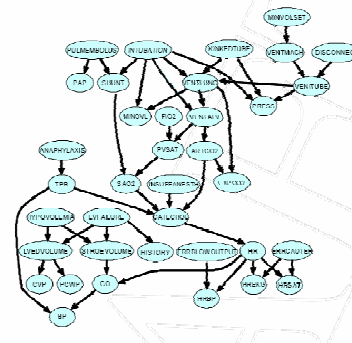
- Parameters  $\theta_X, \theta_{Y|X}$
- Priors:
  - $P(\theta_X)$ :
  - $P(\theta_{Y|X})$ :

10-708 – ©Carlos Guestrin 2006

19

## Very important assumption on prior: Global parameter independence

- **Global parameter independence:**
  - Prior over parameters is product of prior over CPTs



10-708 – ©Carlos Guestrin 2006

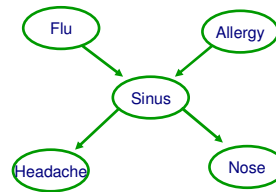
20

## Global parameter independence, d-separation and local prediction

- Independencies in **meta BN**:

- **Proposition:** For fully observable data  $D$ , if prior satisfies global parameter independence, then

$$P(\theta \mid \mathcal{D}) = \prod_i P(\theta_{X_i} \mid \text{Pa}_{X_i} \mid \mathcal{D})$$



10-708 – ©Carlos Guestrin 2006

21

## Within a CPT

- Meta BN including CPT parameters:

- Are  $\theta_{Y|X=t}$  and  $\theta_{Y|X=f}$  d-separated given  $D$ ?
- Are  $\theta_{Y|X=t}$  and  $\theta_{Y|X=f}$  independent given  $D$ ?
  - Context-specific independence!!!
- Posterior decomposes:

10-708 – ©Carlos Guestrin 2006

22



## What you need to know about parameter learning

- MLE:
  - score decomposes according to CPTs
  - optimize each CPT separately
- Bayesian parameter learning:
  - motivation for Bayesian approach
  - Bayesian prediction
  - conjugate priors, equivalent sample size
  - Bayesian learning  $\Rightarrow$  smoothing
- Bayesian learning for BN parameters
  - Global parameter independence
  - Decomposition of prediction according to CPTs
  - Decomposition within a CPT

10-708 – ©Carlos Guestrin 2006

25

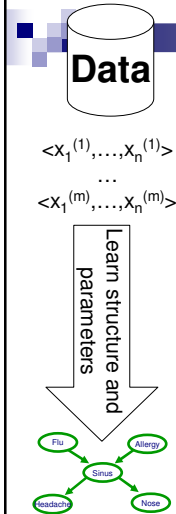
## Where are we with learning BNs?

- Given structure, estimate parameters
  - Maximum likelihood estimation
  - Bayesian learning
- What about learning structure?

10-708 – ©Carlos Guestrin 2006

26

# Learning the structure of a BN



## ■ Constraint-based approach

- BN encodes conditional independencies
- Test conditional independencies in data
- Find an I-map

## ■ Score-based approach

- Finding a structure and parameters is a density estimation task
- Evaluate model as we evaluated parameters
  - Maximum likelihood
  - Bayesian
  - etc.

10-708 – ©Carlos Guestrin 2006

27

# Remember: Obtaining a P-map?

- Given the independence assertions that are true for  $P$ 
  - Obtain skeleton
  - Obtain immoralities
- From skeleton and immoralities, obtain every (and any) BN structure from the equivalence class

## ■ Constraint-based approach:

- Use Learn PDAG algorithm
- Key question: **Independence test**

10-708 – ©Carlos Guestrin 2006

28

# Independence tests

- Statistically difficult task!
- Intuitive approach: **Mutual information**

$$I(X_i, X_j) = \sum_{x_i, x_j} P(x_i, x_j) \log \frac{P(x_i, x_j)}{P(x_i)P(x_j)}$$

- Mutual information and independence:
  - $X_i$  and  $X_j$  independent if and only if  $I(X_i, X_j)=0$
- Conditional mutual information:

10.708 – ©Carlos Guestrin 2006

29

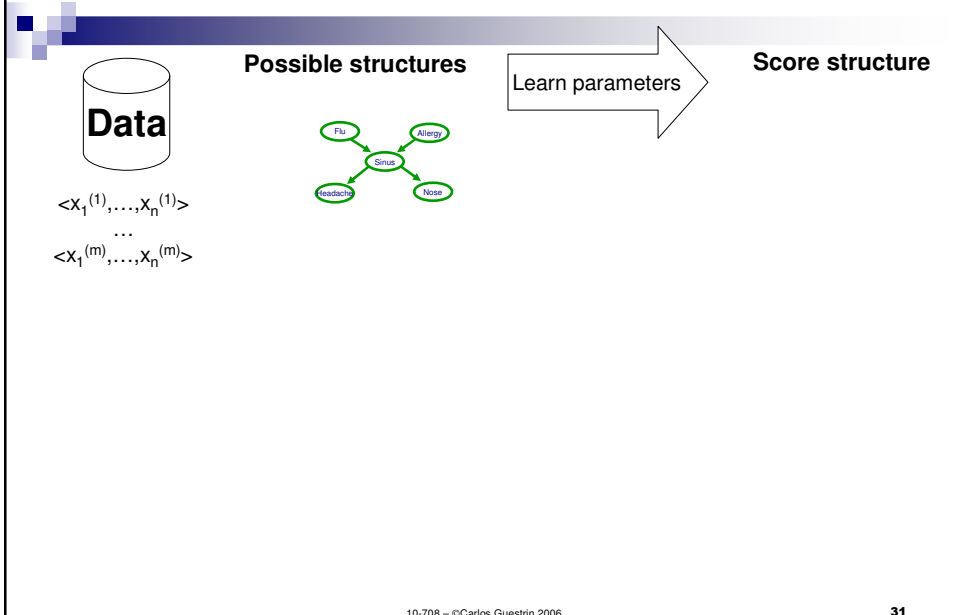
# Independence tests and the constraint based approach

- Using the data  $D$ 
  - Empirical distribution:  $\hat{P}(x_i, x_j) = \frac{\text{Count}(x_i, x_j)}{m}$
  - Mutual information:  $\hat{I}(X_i, X_j) = \sum_{x_i, x_j} \hat{P}(x_i, x_j) \log \frac{\hat{P}(x_i, x_j)}{\hat{P}(x_i)\hat{P}(x_j)}$
  - Similarly for conditional MI
- Use learning PDAG algorithm:
  - When algorithm asks:  $(X \perp Y | \mathbf{U})$ ?
- Many other types of independence tests
  - See reading...

10.708 – ©Carlos Guestrin 2006

30

# Score-based approach



## Information-theoretic interpretation of maximum likelihood

- Given structure, log likelihood of data:

$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G}) = \sum_{j=1}^m \sum_{i=1}^n \log P\left(X_i = x_i^{(j)} \mid \mathbf{Pa}_{X_i} = \mathbf{x}^{(j)}[\mathbf{Pa}_{X_i}]\right)$$

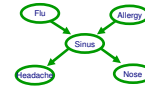




## Information-theoretic interpretation of maximum likelihood 2

- Given structure, log likelihood of data:

$$\log \hat{P}(\mathcal{D} \mid \theta, \mathcal{G}) = m \sum_i \sum_{x_i, \mathbf{Pa}_{x_i, \mathcal{G}}} \hat{P}(x_i, \mathbf{Pa}_{x_i, \mathcal{G}}) \log \hat{P}(x_i \mid \mathbf{Pa}_{x_i, \mathcal{G}})$$



10-708 – ©Carlos Guestrin 2006

33

## Decomposable score

- Log data likelihood

$$\log \hat{P}(\mathcal{D} \mid \theta, \mathcal{G}) = m \sum_i \hat{I}(x_i, \mathbf{Pa}_{x_i, \mathcal{G}}) - M \sum_i \hat{H}(X_i)$$

- Decomposable score:

- Decomposes over families in BN (node and its parents)
- Will lead to significant computational efficiency!!!
- $\text{Score}(G : D) = \sum_i \text{FamScore}(X_i \mid \mathbf{Pa}_{X_i} : D)$

10-708 – ©Carlos Guestrin 2006

34

## How many trees are there?

**Nonetheless – Efficient optimal algorithm finds best tree**

10-708 – ©Carlos Guestrin 2006

35

## Scoring a tree 1: I-equivalent trees

$$\log \hat{P}(\mathcal{D} \mid \theta, \mathcal{G}) = M \sum_i \hat{I}(x_i, \text{Pa}_{x_i, \mathcal{G}}) - M \sum_i \hat{H}(X_i)$$

10-708 – ©Carlos Guestrin 2006

36

## Scoring a tree 2: similar trees

$$\log \hat{P}(\mathcal{D} \mid \theta, \mathcal{G}) = M \sum_i \hat{I}(x_i, \text{Pa}_{x_i, \mathcal{G}}) - M \sum_i \hat{H}(X_i)$$

10-708 – ©Carlos Guestrin 2006

37

## Chow-Liu tree learning algorithm 1

- For each pair of variables  $X_i, X_j$

- Compute empirical distribution:

$$\hat{P}(x_i, x_j) = \frac{\text{Count}(x_i, x_j)}{m}$$

- Compute mutual information:

$$\hat{I}(X_i, X_j) = \sum_{x_i, x_j} \hat{P}(x_i, x_j) \log \frac{\hat{P}(x_i, x_j)}{\hat{P}(x_i) \hat{P}(x_j)}$$

- Define a graph

- Nodes  $X_1, \dots, X_n$
- Edge  $(i, j)$  gets weight  $\hat{I}(X_i, X_j)$

10-708 – ©Carlos Guestrin 2006

38

## Chow-Liu tree learning algorithm 2

$$\log \hat{P}(\mathcal{D} \mid \theta, \mathcal{G}) = M \sum_i \hat{I}(x_i, \text{Pa}_{x_i, \mathcal{G}}) - M \sum_i \hat{H}(X_i)$$

### ■ Optimal tree BN

- Compute maximum weight spanning tree
- Directions in BN: pick any node as root, breadth-first-search defines directions

10-708 – ©Carlos Guestrin 2006

39

## Can we extend Chow-Liu 1

### ■ Tree augmented naïve Bayes (TAN)

[Friedman et al. '97]

- Naïve Bayes model overcounts, because correlation between features not considered
- Same as Chow-Liu, but score edges with:

$$\hat{I}(X_i, X_j \mid C) = \sum_{c, x_i, x_j} \hat{P}(c, x_i, x_j) \log \frac{\hat{P}(x_i, x_j \mid c)}{\hat{P}(x_i \mid c) \hat{P}(x_j \mid c)}$$

10-708 – ©Carlos Guestrin 2006

40

## Can we extend Chow-Liu 2

- (Approximately learning) models with tree-width up to  $k$ 
  - [Narasimhan & Bilmes '04]
  - But,  $O(n^{k+1})\dots$ 
    - and more subtleties

10-708 – ©Carlos Guestrin 2006

41

## What you need to know about learning BN structures so far

- Decomposable scores
  - Maximum likelihood
  - Information theoretic interpretation
- Best tree (Chow-Liu)
- Best TAN
- Nearly best  $k$ -treewidth (in  $O(N^{k+1})$ )

10-708 – ©Carlos Guestrin 2006

42