

Readings:

K&F: 3.3, 3.4, 16.1, 16.2, 16.3, 16.4

Learning P-maps Param. Learning

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

September 24th, 2008

10-708 – ©Carlos Guestrin 2006-2008

1

Perfect maps (P-maps)

- I-maps are not unique and often not simple enough
- Define “simplest” G that is I-map for P
 - A BN structure G is a perfect map for a distribution P if $I(P) = I(G)$
- Our goal:
 - Find a perfect map!
 - Must address equivalent BNs

10-708 – ©Carlos Guestrin 2006-2008

2

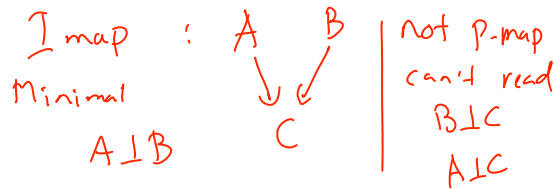
Inexistence of P-maps 1

- XOR (this is a hint for the homework)

$$A = B \text{ XOR } C$$

$$\begin{array}{l|l} A \perp B & \neg A \perp B \perp C \\ B \perp C & \neg A \perp C \perp B \\ C \perp A & \neg B \perp C \perp A \end{array}$$

P-MAP?
extra credit



10-708 - ©Carlos Guestrin 2006-2008

3

Obtaining a P-map

- Given the independence assertions that are true for P
- Assume that there exists a perfect map G^*
 - Want to find G^*
- Many structures may encode same independencies as G^* , when are we done?
 - Find all equivalent structures simultaneously!

10-708 - ©Carlos Guestrin 2006-2008

4

I-Equivalence

- Two graphs G_1 and G_2 are **I-equivalent** if $I(G_1) = I(G_2)$
- **Equivalence class** of BN structures
 - Mutually-exclusive and exhaustive partition of graphs

- How do we characterize these equivalence classes?

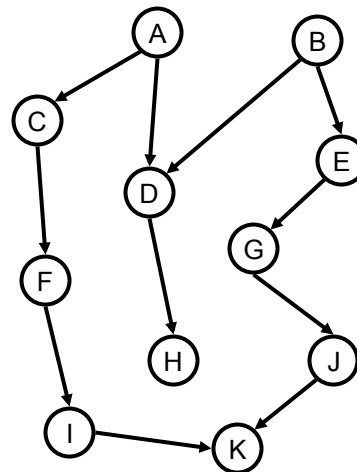
10-708 – ©Carlos Guestrin 2006-2008

5

Skeleton of a BN

- **Skeleton** of a BN structure G is an **undirected graph** over the same variables that has an edge $X-Y$ for every $X \rightarrow Y$ or $Y \rightarrow X$ in G

- (Little) **Lemma**: Two I-equivalent BN structures must have the same skeleton

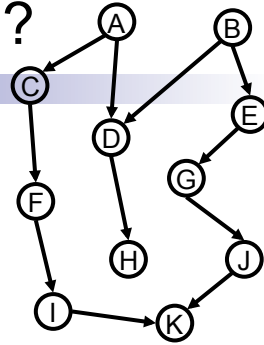


10-708 – ©Carlos Guestrin 2006-2008

6

What about V-structures?

- V-structures are key property of BN structure



- **Theorem:** If G_1 and G_2 have the same skeleton and V-structures, then G_1 and G_2 are I-equivalent

10-708 – ©Carlos Guestrin 2006-2008

7

Same V-structures not necessary

- **Theorem:** If G_1 and G_2 have the same skeleton and V-structures, then G_1 and G_2 are I-equivalent
- Though sufficient, same V-structures not necessary

10-708 – ©Carlos Guestrin 2006-2008

8

Immoralities & I-Equivalence

- Key concept not V-structures, but “immoralities” (unmarried parents ☺)
 - $X \rightarrow Z \leftarrow Y$, with no arrow between X and Y
 - Important pattern: X and Y independent given their parents, but not given Z
 - (If edge exists between X and Y, we have *covered* the V-structure)
- **Theorem:** G_1 and G_2 have the same skeleton and immoralities if and only if G_1 and G_2 are I-equivalent

10-708 – ©Carlos Guestrin 2006-2008

9

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

10-708 – ©Carlos Guestrin 2006-2008

10

Identifying the skeleton 1

- When is there an edge between X and Y?

- When is there no edge between X and Y?

Identifying the skeleton 2

- Assume d is max number of parents (d could be n)

- For each X_i and X_j
 - $E_{ij} \leftarrow \text{true}$
 - For each $\mathbf{U} \subseteq \mathbf{X} - \{X_i, X_j\}$, $|\mathbf{U}| \leq d$
 - Is $(X_i \perp X_j \mid \mathbf{U})$?
 - $E_{ij} \leftarrow \text{false}$
 - If E_{ij} is true
 - Add edge X – Y to skeleton

Identifying immoralities

- Consider $X - Z - Y$ in skeleton, when should it be an immorality?
- Must be $X \rightarrow Z \leftarrow Y$ (immorality):
 - When X and Y are **never independent** given \mathbf{U} , if $Z \in \mathbf{U}$
- Must **not** be $X \rightarrow Z \leftarrow Y$ (not immorality):
 - When there exists \mathbf{U} with $Z \in \mathbf{U}$, such that X and Y are **independent** given \mathbf{U}

10-708 – ©Carlos Guestrin 2006-2008

13

From immoralities and skeleton to BN structures

- Representing BN equivalence class as a **partially-directed acyclic graph (PDAG)**
- **Immoralities force direction on some other BN edges**
- Full (polynomial-time) procedure described in reading

10-708 – ©Carlos Guestrin 2006-2008

14

What you need to know

- Minimal I-map
 - every P has one, but usually many
- Perfect map
 - better choice for BN structure
 - not every P has one
 - can find one (if it exists) by considering I-equivalence
 - Two structures are I-equivalent if they have same skeleton and immoralities

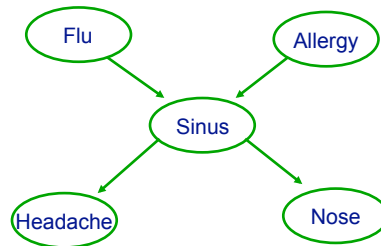
Announcements

- Recitation tomorrow
 - Don't miss it!
- No class on Monday ☹

Review

■ Bayesian Networks

- Compact representation for probability distributions
- Exponential reduction in number of parameters
- Exploits independencies



■ Next – Learn BNs

- parameters
- structure

10-708 – ©Carlos Guestrin 2006-2008

17

Thumbtack – Binomial Distribution

- $P(\text{Heads}) = \theta$, $P(\text{Tails}) = 1 - \theta$

■ Flips are i.i.d.:

- Independent events
- Identically distributed according to Binomial distribution

- Sequence D of α_H Heads and α_T Tails

$$P(D | \theta) = \theta^{\alpha_H} (1 - \theta)^{\alpha_T}$$

10-708 – ©Carlos Guestrin 2006-2008

18

Maximum Likelihood Estimation

- **Data:** Observed set D of α_H Heads and α_T Tails
- **Hypothesis:** Binomial distribution
- Learning θ is an optimization problem
 - What's the objective function?
- MLE: Choose θ that maximizes the probability of observed data:

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

10-708 – ©Carlos Guestrin 2006-2008

19

Your first learning algorithm

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

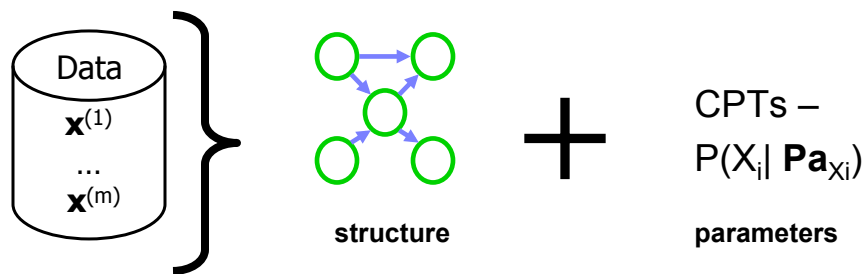
- Set derivative to zero: $\frac{d}{d\theta} \ln P(\mathcal{D} | \theta) = 0$

10-708 – ©Carlos Guestrin 2006-2008

20

Learning Bayes nets

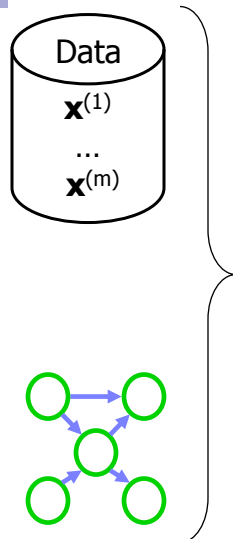
	Known structure	Unknown structure
Fully observable data		
Missing data		



10-708 – ©Carlos Guestrin 2006-2008

21

Learning the CPTs



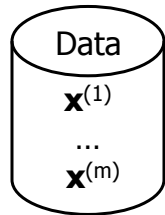
For each discrete variable X_i

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

10-708 – ©Carlos Guestrin 2006-2008

22

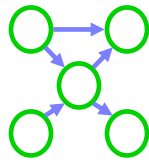
Learning the CPTs



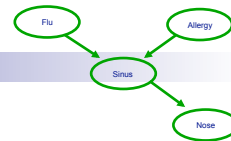
For each discrete variable X_i

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

WHY????????????



Maximum likelihood estimation (MLE) of BN parameters – example



- Given structure, log likelihood of data:
 $\log P(\mathcal{D} | \theta_{\mathcal{G}}, \mathcal{G})$

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

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

10-708 – ©Carlos Guestrin 2006-2008

25

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 2006-2008

26

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 2006-2008

29

Bayesian Learning for Thumbtack

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

- Likelihood function is simply Binomial:

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

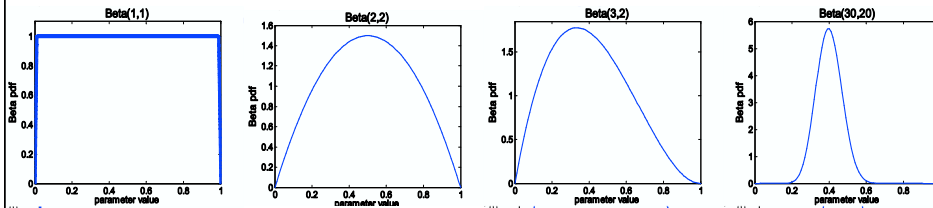
- What about prior?
 - Represent expert knowledge
 - Simple posterior form
- Conjugate priors:
 - Closed-form representation of posterior (more details soon)
 - **For Binomial, conjugate prior is Beta distribution**

10-708 – ©Carlos Guestrin 2006-2008

30

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

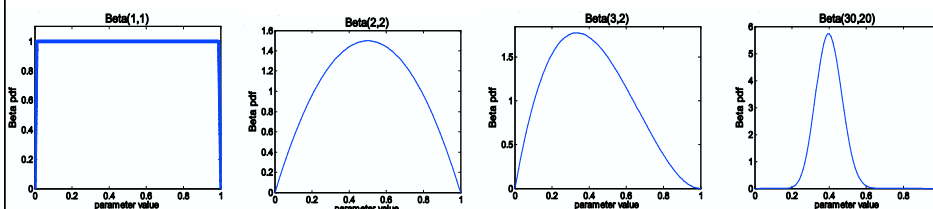
10-708 – ©Carlos Guestrin 2006-2008

31

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 2006-2008

32

Conjugate prior

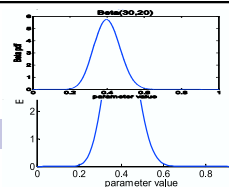
- Prior: $Beta(\alpha_H, \alpha_T)$
- Data: m_H heads and m_T tails (binomial likelihood)
- Posterior distribution:
$$P(\theta | \mathcal{D}) \sim Beta(m_H + \alpha_H, m_T + \alpha_T)$$
- Given likelihood function $P(D|\theta)$
- (Parametric) prior of the form $P(\theta|\alpha)$ is **conjugate** to likelihood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

10-708 – ©Carlos Guestrin 2006-2008

33

Using Bayesian posterior

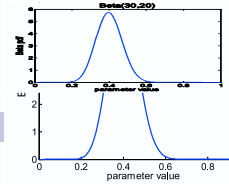
- Posterior distribution:
$$P(\theta | \mathcal{D}) \sim Beta(m_H + \alpha_H, m_T + \alpha_T)$$
- Bayesian inference:
 - No longer single parameter:
$$E[f(\theta)] = \int_0^1 f(\theta)P(\theta | \mathcal{D})d\theta$$
 - Integral is often hard to compute



10-708 – ©Carlos Guestrin 2006-2008

34

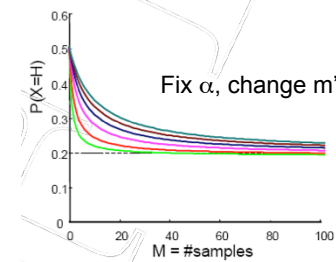
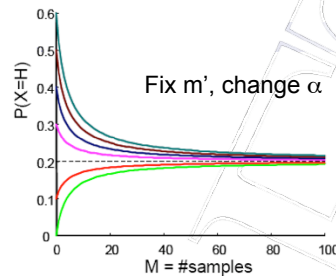
Bayesian prediction of a new coin flip



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

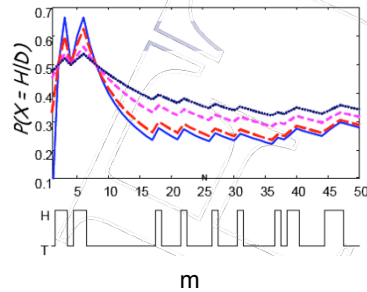
Asymptotic behavior and equivalent sample size

- 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 α_H, α_T , or
 - m' (equivalent sample size) and α
 - $$E[\theta] = \frac{m_H + \alpha m'}{m_H + m_T + m'}$$



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 2006-2008

37

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 2006-2008

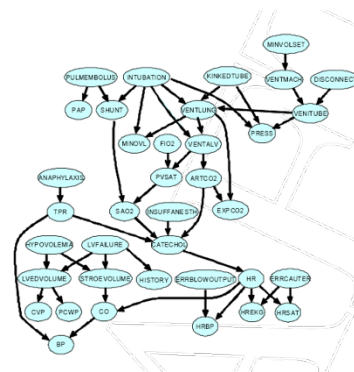
38

Bayesian learning for two-node BN

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

Very important assumption on prior: Global parameter independence

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

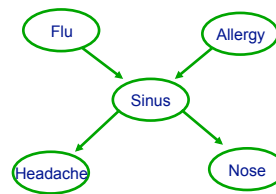


Global parameter independence, d-separation and local prediction

- 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 2006-2008

41

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 2006-2008

42

Priors for BN CPTs

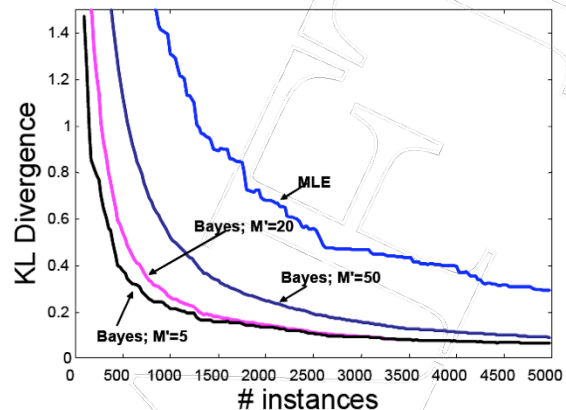
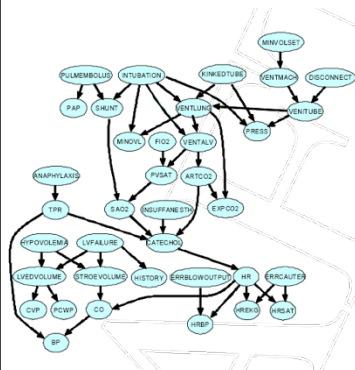
(more when we talk about structure learning)

- Consider each CPT: $P(X|\mathbf{U}=\mathbf{u})$
- Conjugate prior:
 - Dirichlet($\alpha_{X=1|\mathbf{U}=\mathbf{u}}, \dots, \alpha_{X=k|\mathbf{U}=\mathbf{u}}$)
- More intuitive:
 - “prior data set” D' with m' equivalent sample size
 - “prior counts”:
 - prediction:

10-708 – ©Carlos Guestrin 2006-2008

43

An example



10-708 – ©Carlos Guestrin 2006-2008

44

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