# Bayesian Param. Learning

# Bayesian Structure Learning

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

October 6th, 2008

---

# Decomposable score

- Log data likelihood

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

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

for MLE    Fam Score $(X_i | Pa_{X_i} : D) = m \hat{I}(X_i; Pa_{X_i}) - m \hat{H}(X_i)$

# Chow-Liu tree learning algorithm 1

- For each pair of variables $X_i, X_j$
  - Compute empirical distribution:
  $$\bar{P}(x_i, x_j) \overset{\text{MLE}}{=} \frac{\text{Count}(x_i, x_j)}{m}$$
  - Compute mutual information:
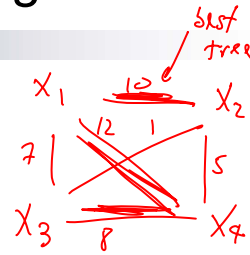  $$\hat{I}(X_i, X_j) = \sum_{x_i, x_j} \hat{P}(x_i, x_j) \log \frac{\hat{P}(x_i, x_j)}{\hat{P}(x_i)\hat{P}(x_j)}$$
- Define a graph
  - Nodes $X_1, \ldots, X_n$    $w_{ij}$
  - Edge (i,j) gets weight $\hat{I}(X_i, X_j)$

*[handwritten annotations: best tree; diagram with $X_1$ —10— $X_2$, $X_3$ —8— $X_4$ with edge weights 12, 1, 7, 5; find Maximum Spanning tree; $\max_{\text{trees}} \text{Score(tree)} = \sum_{ij} I(X_i, X_j) = \sum_{ij} w_{ij} \Rightarrow$ best tree BN]*
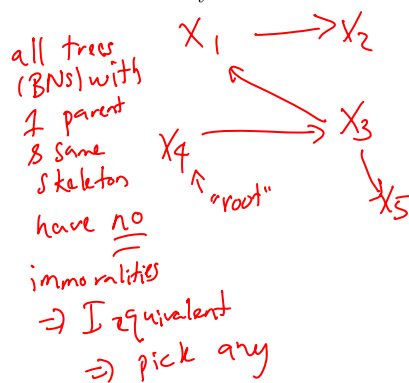
---

# Chow-Liu tree learning algorithm 2

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

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

*[handwritten annotations: using Chow-Liu OPTIMAL tree BN; all trees (BNs) with 1 parent & same skeleton have no immoralities ⇒ I equivalent ⇒ pick any; diagram $X_1 \rightarrow X_2$, $X_4$ "root", $X_4 \rightarrow X_3$, $X_3 \rightarrow X_5$]*
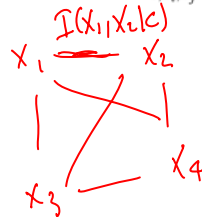
# Can we extend Chow-Liu 1
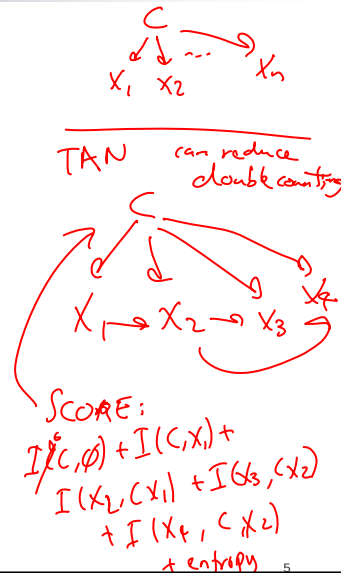
- Tree augmented naïve Bayes (TAN)
  [Friedman et al. '97]
  - ☐ Naïve Bayes model overcounts, because correlation between features not considered
  - ☐ Same as Chow-Liu, but score edges with:

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

# Can we extend Chow-Liu 2

- (Approximately learning) models with tree-width up to *k*
  - ☐ [Chechetka & Guestrin '07]
  - ☐ But, $O(n^{2k+6})$

# 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 $O(N^{2k+6})$)

---

# Maximum likelihood score overfits!

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

- Information never hurts:

$$\uparrow I(X_i, Pa_{X_i}) = H(X_i) - H(X_i \mid Pa_{X_i})$$

$$H(A \mid B) \leq H(A \mid C) \qquad C \subseteq B$$

the more parents the higher $I(X_i, Pa_{X_i})$

- Adding a parent always increases score!!!

$$MLE \Rightarrow \qquad \text{Complete Graph}$$

# Bayesian score

- **Prior distributions:**
  - Over structures
  - Over parameters of a structure
- Posterior over structures given data:

*(handwritten annotations)*

nots LD

$P(D|G, \theta_G)$

prior over graphs, e.g.
$P(G) \propto e^{-c |\text{number of edges}|}$

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

prior over CPT parameters

$= \dfrac{\int_{\theta_G} P(D|G,\theta_G) P(\theta_G|G) P(G) d\theta_G}{P(D)}$

prior over Graphs

$$\log P(\mathcal{G} \mid D) = \log P(\mathcal{G}) + \log \int_{\theta_{\mathcal{G}}} P(D \mid \mathcal{G}, \theta_{\mathcal{G}}) P(\theta_{\mathcal{G}} | \mathcal{G}) d\theta_{\mathcal{G}}$$

posterior

$+ \text{constant} \leftarrow \log P(D)$

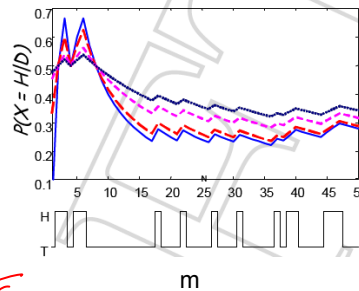---

# Can we really trust MLE?

- **What is better?**
  - 3 heads, 2 tails     $\theta_{MLE} = \dfrac{3}{5}$
  - 30 heads, 20 tails    $\theta_{MLE} = \dfrac{3}{5}$
  - $3 \times 10^{23}$ heads, $2 \times 10^{23}$ tails    $\theta_{MLE} = \dfrac{3}{5}$



- Many possible answers, we need distributions over possible parameters

# Bayesian Learning

- Use Bayes rule:

*(handwritten: likelihood, prior, posterior)*

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

- Or equivalently:

$$P(\theta \mid \mathcal{D}) \;\propto\; P(\mathcal{D} \mid \theta)P(\theta)$$

---

# Bayesian Learning for Thumbtack

*(handwritten: posterior, likelihood, prior)*

$$P(\theta \mid \mathcal{D}) \;\propto\; P(\mathcal{D} \mid \theta)P(\theta)$$

- Likelihood function is simply Binomial:

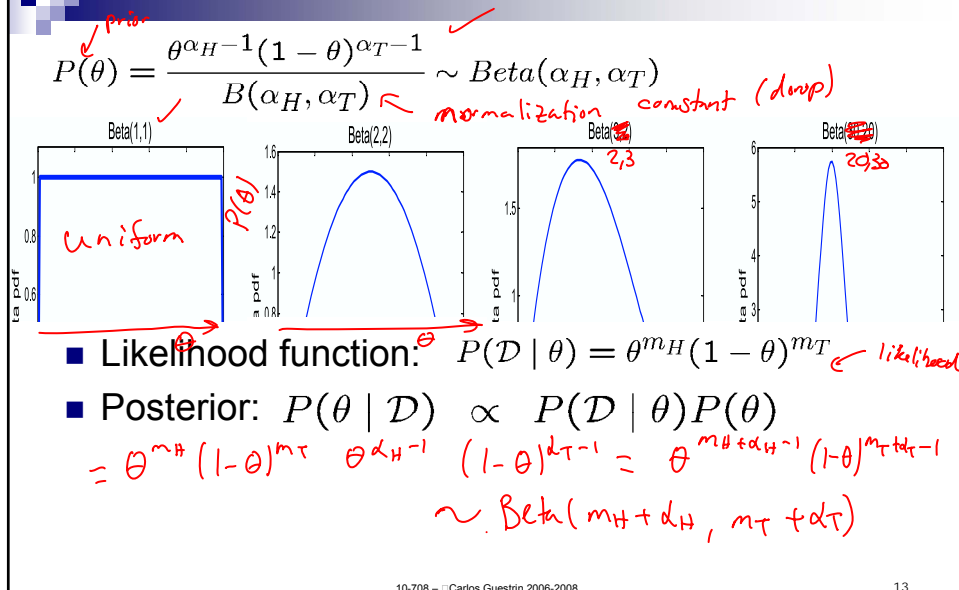$$P(\mathcal{D} \mid \theta) = \theta^{m_H}(1 - \theta)^{m_T}$$

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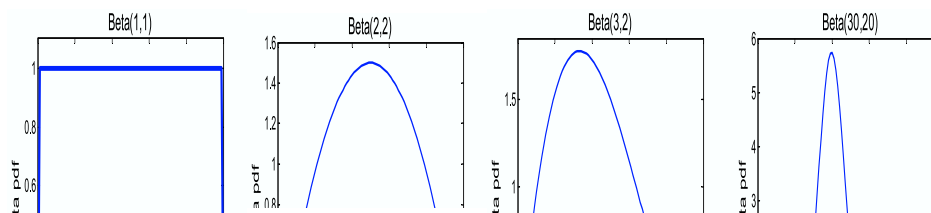  - **For Binomial, conjugate prior is Beta distribution**

# Beta prior distribution – P(θ)

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

*[handwritten annotations: "prior", "normalization constant (drop)", "uniform", "P(θ)"]*



Beta(1,1)   Beta(2,2)   Beta(3,2)   Beta(30,20)

- Likelihood function: $P(\mathcal{D} \mid \theta) = \theta^{m_H}(1-\theta)^{m_T}$ *[handwritten: likelihood]*
- Posterior: $P(\theta \mid \mathcal{D}) \propto P(\mathcal{D} \mid \theta)P(\theta)$

*[handwritten:]* $= \theta^{m_H}(1-\theta)^{m_T} \theta^{\alpha_H - 1}(1-\theta)^{\alpha_T - 1} = \theta^{m_H + \alpha_H - 1}(1-\theta)^{m_T + \alpha_T - 1}$

$\sim Beta(m_H + \alpha_H, m_T + \alpha_T)$

# Posterior distribution

- Prior: $Beta(\alpha_H, \alpha_T)$
- Data: $m_H$ heads and $m_T$ tails

- Posterior distribution:

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



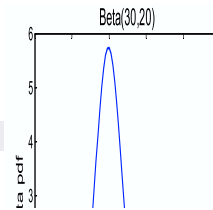Beta(1,1)   Beta(2,2)   Beta(3,2)   Beta(30,20)

# Conjugate prior

- Prior: $Beta(\alpha_H, \alpha_T)$
- Data: $m_H$ heads and $m_T$ tails (binomial likelihood)
- Posterior distribution:

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

- Given likelihood function P($D$|θ)

- (Parametric) prior of the form P(θ|α) is **conjugate** to likelihood function if posterior is of the same parametric family, and can be written as:
  - P(θ|α'), for some new set of parameters α'

# Using Bayesian posterior

Beta(30,20)

- Posterior distribution:

$$P(\theta \mid \mathcal{D}) \sim Beta(m_H + \alpha_H, m_T + \alpha_T)$$
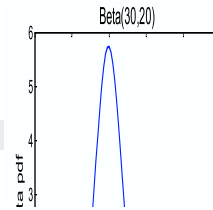
- Bayesian inference:
  - No longer single parameter:

  $$E[f(\theta)] = \int_0^1 f(\theta) P(\theta \mid \mathcal{D}) d\theta$$

  *posterior*

  *utility*

  - Integral is often hard to compute

  *often → mean parameter*
  *mode parameter*

# Bayesian prediction of a new coin flip

Beta(30,20)

- Prior: $\text{Beta}(\alpha_H, \alpha_T)$
- Observed $m_H$ heads, $m_T$ tails, what is probability of m+1 flip is heads?

Posterior $\text{Beta}(\alpha_H + m_H, \alpha_T + m_T)$

$$P(m+1 \text{ flip} = \text{heads} \mid m_H, m_T)$$

$$= \int_\theta P(m+1 \text{ flip} = H \mid \theta) P(\theta \mid m_H, m_T) \, d\theta$$

$$= \int_\theta \theta \, P(\theta \mid m_H, m_T) \, d\theta \equiv \text{mean} = \frac{\alpha_H + m_H}{\alpha_H + m_H + \alpha_T + m_T}$$

$$\text{Beta}(\alpha_H + m_H, \alpha_T + m_T)$$

# Asymptotic behavior and equivalent sample size

Prior $\text{Beta}(\alpha m', (1-\alpha)m')$
$m' \leftarrow$ equivalent Sample Size

- Beta prior equivalent to extra thumbtack flips:

$\alpha_H = \alpha m'$
$\alpha_T = (1-\alpha)m'$

$$E[\theta] = \frac{m_H + \alpha_H}{m_H + \alpha_H + m_T + \alpha_T}$$

Fix m', change $\alpha$

$m \to \infty$

- As $m \to \infty$, prior is "forgotten"
- **But, for small sample size, prior is important!**
- **Equivalent sample size**:
  - Prior parameterized by $\alpha_H, \alpha_T$, or
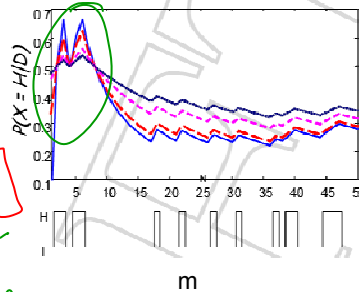  - m' (equivalent sample size) and $\alpha$

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

bigger m'

Fix $\alpha$, change m'

$m \to \infty$

# Bayesian learning corresponds to smoothing

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

$$= \frac{m}{m+m'}\left[\frac{m_H}{m}\right] + \frac{m'}{m+m'}\left[\frac{\alpha m'}{m'}\right]$$

MLE estimate

$\alpha$ prior "mean"

m

- m=0 ⇒ prior parameter
- m→∞ ⇒ MLE

$$Beta(\alpha_H, \alpha_T) \leftarrow mode \quad \frac{\alpha_H - 1}{\alpha_H + \alpha_T - 2}$$

---

# Bayesian learning for multinomial

$m_i \in$ # observations of class, value or side i

- What if you have a k sided coin???
- Likelihood function if **multinomial**:
  - $P(D \mid \theta_1, \ldots, \theta_k) = \theta_1^{m_1} \theta_2^{m_2} \cdots \theta_k^{m_k}$
  - $\sum \theta_i = 1, \quad \theta_i \geq 0$
- **Conjugate** prior for multinomial is **Dirichlet**:
  - $\theta \sim Dirichlet(\alpha_1, \ldots, \alpha_k) \sim \prod_i \theta_i^{\alpha_i - 1}$
  
  $\alpha_i \geq 0$
- **Observe** $m$ data points, $m_i$ from assignment i, **posterior**:

$$P(\theta_1 \ldots \theta_k \mid m_1 \ldots m_k) \propto P(m_1 \ldots m_k \mid \theta_1 \ldots \theta_k) \, P(\theta)$$

$$\equiv Dirichlet(\alpha_1 + m_1, \alpha_2 + m_2, \ldots, \alpha_k + m_k)$$

- **Prediction**:

$$E[\theta_i] = \frac{m_i + \alpha_i}{\sum_j (m_j + \alpha_j)}$$

# Bayesian learning for two-node BN

- Parameters $\theta_X$, $\theta_{Y|X}$

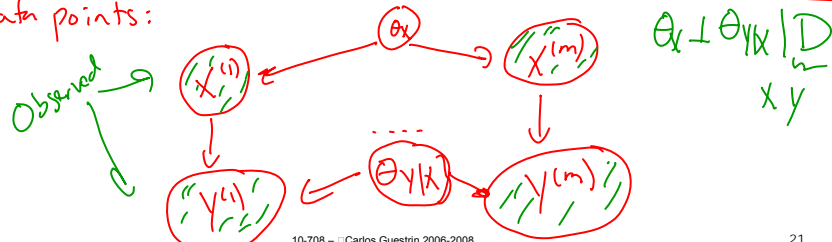  $X \longrightarrow Y$
  $P(X) \qquad P(Y|X)$

- Priors:
  - $P(\theta_X)$: $Dirichlet(\alpha_{X=1}, \alpha_{X=2}, \cdots, \alpha_{X=x})$
  - $P(\theta_{Y|X})$: for each value $X=x$
    a set of parameters $\theta_{Y|X=x}$

  $P(\theta_{Y|X=x}) \equiv Dirichlet(\alpha_{Y=1|X=x}, \cdots, \alpha_{Y=k|X=x})$

m data points:

$\theta_X \perp \theta_{Y|X} \mid D$

Observed

$X^{(1)}$ ← $\theta_X$ → $X^{(m)}$

$X\ Y$

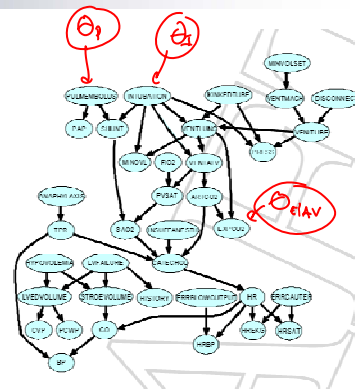$Y^{(1)}$ ← $\theta_{Y|X}$ → $Y^{(m)}$

---

# Very important assumption on prior: Global parameter independence

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

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

$\theta_1 \quad \theta_2$

$\theta_{clav}$

# Global parameter independence, d-separation and local prediction
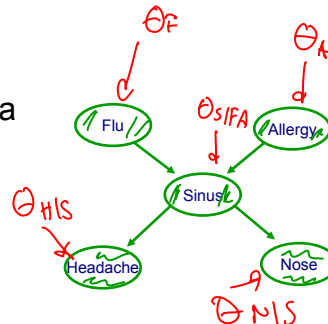
- Independencies in **meta BN**:

  *add prior vars to the BN*

  $$P(\theta) = P(\theta_F) \, P(\theta_A) \, P(\theta_{S|FA}) \, P(\theta_{N|S}) \, P(\theta_{H|S})$$

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

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

  *params indep. given data*

  $\theta_F$  $\theta_A$

  $\theta_{S|FA}$

  $\theta_{H|S}$

  Flu  Allergy

  Sinus

  Headache  Nose

  $\theta_{N|S}$

---

# 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:

# Priors for BN CPTs
(more when we talk about structure learning)

- Consider each CPT: P(X|**U=u**)
- Conjugate prior:
  - Dirichlet($\alpha_{X=1|U=u}, \ldots, \alpha_{X=k|U=u}$) $\equiv Dirichlet\left(Count'(X=1, U=u) \ldots, Count'(X=k, U=u)\right)$
- More intuitive:
  - "prior data set" $D'$ with m' equivalent sample size
  - "prior counts": $Count'(X=x, U=u)$ or $m' \cdot P'(X=x, U=u)$
  - prediction:

$$E[\theta_{X=x|U=u}] = \frac{Count(X=x, U=u) + Count'(X=x, U=u)}{Count(U=u) + Count'(U=u)}$$