

Bayesian Learning

Chapter 20.1-20.4

Some material adapted
from lecture notes by
Lise Getoor and Ron Parr

Naïve Bayes

Naïve Bayes

- Use Bayesian modeling
- Make the simplest possible independence assumption:
 - Each attribute is independent of the values of the other attributes, given the class variable
 - In our restaurant domain: Cuisine is independent of Patrons, *given* a decision to stay (or not)

Bayesian Formulation

- $p(C | F_1, \dots, F_n) = p(C) p(F_1, \dots, F_n | C) / P(F_1, \dots, F_n)$
 $= \alpha p(C) p(F_1, \dots, F_n | C)$
- Assume each feature F_i is *conditionally independent* of the other given the class C . Then:
 $p(C | F_1, \dots, F_n) = \alpha p(C) \prod_i p(F_i | C)$
- Estimate each of these conditional probabilities from the observed **counts** in the training data:
 $p(F_i | C) = N(F_i \wedge C) / N(C)$
 - One subtlety of using the algorithm in practice: When your estimated probabilities are zero, ugly things happen
 - The fix: Add one to every count (aka “Laplacian smoothing”—they have a different name for *everything!*)

Naive Bayes: Example

$$\begin{aligned} p(\text{Wait} \mid \text{Cuisine, Patrons, Rainy?}) &= \\ &= \alpha \cdot p(\text{Wait}) \cdot p(\text{Cuisine} \mid \text{Wait}) \cdot p(\text{Patrons} \mid \text{Wait}) \cdot p(\text{Rainy?} \mid \text{Wait}) \\ &= \frac{p(\text{Wait}) \cdot p(\text{Cuisine} \mid \text{Wait}) \cdot p(\text{Patrons} \mid \text{Wait}) \cdot p(\text{Rainy?} \mid \text{Wait})}{p(\text{Cuisine}) \cdot p(\text{Patrons}) \cdot p(\text{Rainy?})} \end{aligned}$$

We can estimate all of the parameters ($p(F)$ and $p(C)$) just by counting from the training examples

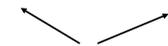
Naive Bayes: Analysis

- Naive Bayes is amazingly easy to implement (once you understand the bit of math behind it)
- Remarkably, naive Bayes can outperform many much more complex algorithms—it's a baseline that should pretty much always be used for comparison
- Naive Bayes can't capture interdependencies between variables (obviously)—for that, we need Bayes nets!

Learning Bayesian Networks

Bayesian learning: Bayes' rule

- Given some model space (set of hypotheses h_i) and evidence (data D):
 - $P(h_i \mid D) = \alpha P(D \mid h_i) P(h_i)$
- We assume that observations are independent of each other, given a model (hypothesis), so:
 - $P(h_i \mid D) = \alpha \prod_j P(d_j \mid h_i) P(h_i)$
- To predict the value of some unknown quantity, X (e.g., the class label for a future observation):
 - $P(X \mid D) = \sum_i P(X \mid D, h_i) P(h_i \mid D) = \sum_i P(X \mid h_i) P(h_i \mid D)$



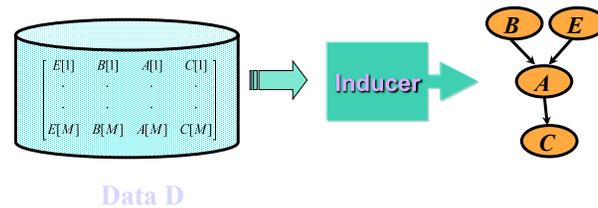
These are equal by our independence assumption

Bayesian learning

- We can apply Bayesian learning in three basic ways:
 - BMA (Bayesian Model Averaging):** Don't just choose one hypothesis; instead, make predictions based on the weighted average of all hypotheses (or some set of best hypotheses)
 - MAP (Maximum A Posteriori) hypothesis:** Choose the hypothesis with the highest *a posteriori* probability, given the data
 - MLE (Maximum Likelihood Estimate):** Assume that all hypotheses are equally likely *a priori*; then the best hypothesis is just the one that maximizes the likelihood (i.e., the probability of the data given the hypothesis)
- MDL (Minimum Description Length) principle:** Use some encoding to model the complexity of the hypothesis, and the fit of the data to the hypothesis, then minimize the overall description of $h_i + D$

Learning Bayesian networks

- Given training set
- Find B that best matches $D \quad D = \{x[1], \dots, x[M]\}$
 - model selection
 - parameter estimation



Learning Bayesian Networks

- Describe a BN by specifying its (1) structure and (2) CPD tables
- Both can be learned from data, but
 - learning structure is much harder than learning parameters
 - learning when some nodes are hidden, or with missing data harder still
- We have four cases:

Structure	Observability	Method
Known	Full	Maximum Likelihood Estimation
Known	Partial	EM (or gradient ascent)
Unknown	Full	Search through model space
Unknown	Partial	EM + search through model space

Parameter estimation

- Assume known structure
- Goal: estimate BN parameters θ
 - entries in local probability models, $P(X | \text{Parents}(X))$
- A parameterization θ is good if it is likely to generate the observed data:

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

i.i.d. samples

- Maximum Likelihood Estimation (MLE) Principle: Choose θ^* so as to maximize L

Parameter estimation II

- The likelihood **decomposes** according to the structure of the network
 - we get a separate estimation task for each parameter
- The MLE (maximum likelihood estimate) solution:
 - for each value x of a node X
 - and each instantiation \mathbf{u} of $Parents(X)$

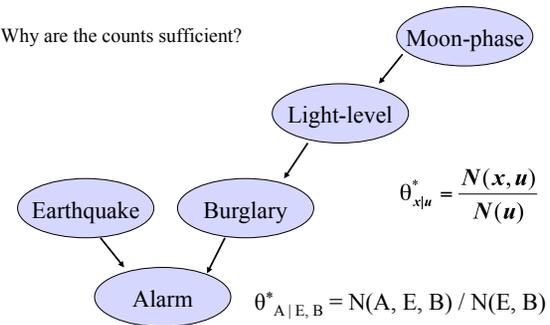
$$\theta_{x|\mathbf{u}}^* = \frac{N(\mathbf{x}, \mathbf{u})}{N(\mathbf{u})}$$

← sufficient statistics

- Just need to collect the counts for every combination of parents and children observed in the data
- MLE is equivalent to an assumption of a uniform prior over parameter values

Sufficient statistics: Example

- Why are the counts sufficient?



Model selection

Goal: Select the best network structure, given the data

Input:

- Training data
- Scoring function

Output:

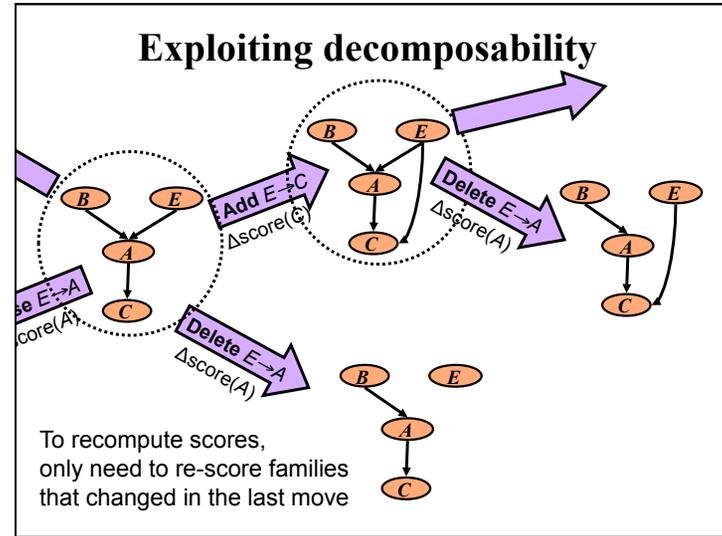
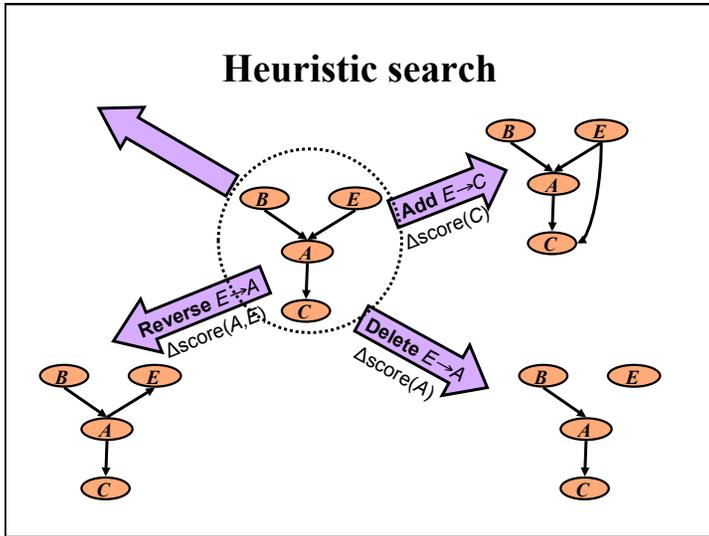
- A network that maximizes the score

Structure selection: Scoring

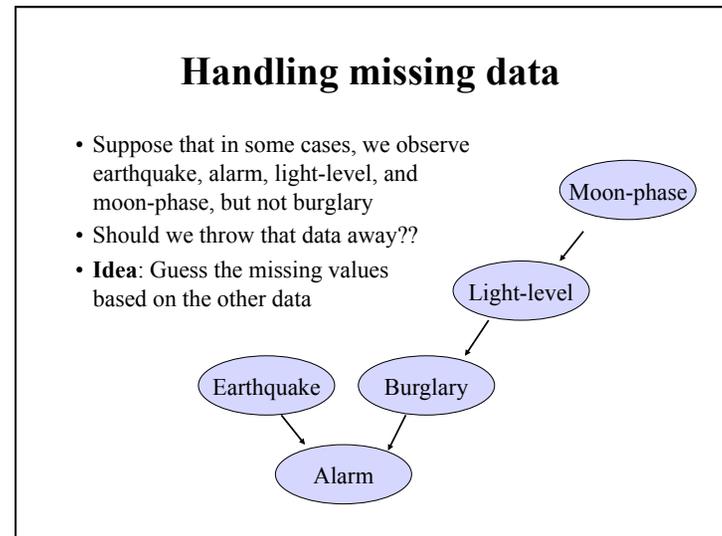
- Bayesian: prior over parameters and structure
 - get balance between model complexity and fit to data as a byproduct
- Score (G:D) = log P(G|D) \propto log [P(D|G) P(G)]
 - Marginal likelihood* (points to P(D|G))
 - Prior* (points to P(G))
- Marginal likelihood just comes from our parameter estimates
- Prior on structure can be any measure we want; typically a function of the network complexity

Same key property: Decomposability

$$\text{Score}(\text{structure}) = \sum_i \text{Score}(\text{family of } X_i)$$



- ### Variations on a theme
- **Known structure, fully observable:** only need to do parameter estimation
 - **Unknown structure, fully observable:** do heuristic search through structure space, then parameter estimation
 - **Known structure, missing values:** use expectation maximization (EM) to estimate parameters
 - **Known structure, hidden variables:** apply adaptive probabilistic network (APN) techniques
 - **Unknown structure, hidden variables:** too hard to solve!



EM (expectation maximization)

- **Guess** probabilities for nodes with **missing values** (e.g., based on other observations)
- **Compute the probability distribution** over the missing values, given our guess
- **Update the probabilities** based on the guessed values
- **Repeat** until convergence

EM example

- Suppose we have observed Earthquake and Alarm but not Burglary for an observation on November 27
- We estimate the CPTs based on the *rest* of the data
- We then estimate $P(\text{Burglary})$ for November 27 from those CPTs
- Now we recompute the CPTs as if that estimated value had been observed
- Repeat until convergence!

