

# Bayesian Networks (Structure) Learning

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

November 12<sup>th</sup>, 2007

©2005-2007 Carlos Guestrin

1

## Information-theoretic interpretation of maximum likelihood 1

$\log \prod_{i=1}^n P(x^{(i)} | \theta_G, G) = \sum_i \log P(x^{(i)} | \theta_G, G)$   
 $\log ab = \log a + \log b$

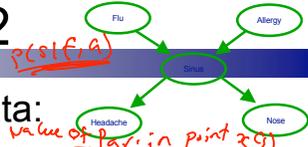
$\log P(\mathcal{D} | \theta_G, G) =$

$\log \prod_{i=1}^n P(x^{(i)} | \theta_G, G) = \sum_i \log P(x^{(i)} | \theta_G, G)$   
 $= \sum_i \log P(f^{(i)}) P(a^{(i)}) P(s^{(i)} | f^{(i)}, a^{(i)}) P(h^{(i)} | s^{(i)}) P(n^{(i)} | s^{(i)})$   
 $= \underbrace{\left[ \sum_i \log P(f^{(i)}) \right]}_{P(F)} + \underbrace{\left[ \sum_i \log P(a^{(i)}) \right]}_{P(A)} + \underbrace{\left[ \sum_i \log P(s^{(i)} | f^{(i)}, a^{(i)}) \right]}_{P(S|F,A)} + \underbrace{\left[ \sum_i \log P(h^{(i)} | s^{(i)}) \right]}_{P(H|S)} + \underbrace{\left[ \sum_i \log P(n^{(i)} | s^{(i)}) \right]}_{P(N|S)}$

one per CPT  $P(x_i | \text{Pa}x_i)$   
 only this term changes about  $P(A|F)$

$G:$

# Information-theoretic interpretation of maximum likelihood 2



Given structure, log likelihood of data:

$$\log P(\mathcal{D} | \theta_G, \mathcal{G}) = \sum_{j=1}^m \sum_{i=1}^n \log P(X_i = x_i^{(j)} | \text{Pa}_{X_i} = x^{(j)} [\text{Pa}_{X_i}])$$

*data points nodes* (pointing to  $\sum_{i=1}^n$ )  
*point j* (pointing to  $\sum_{j=1}^m$ )  
*value of Pa<sub>X<sub>i</sub></sub> in point z<sup>(j)</sup>* (pointing to  $x^{(j)}$ )

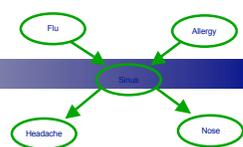
$$= \sum_{i=1}^n \sum_{j=1}^m \log P(X_i = x_i^{(j)} | \text{Pa}_{X_i} = z^{(j)} [\text{Pa}_{X_i}])$$

*for point j where F=t, A=f, S=t (Count(F=t, A=f, S=t))*

$$= \sum_{i=1}^n \sum_{x_i} \sum_u \frac{\text{count}(X_i = x_i, \text{Pa}_{X_i} = u)}{m} \log P(X_i = x_i | \text{Pa}_{X_i} = u)$$

*MLE estimate* (pointing to  $\hat{P}(X_i = x_i, \text{Pa}_{X_i} = u)$ )

# Information-theoretic interpretation of maximum likelihood 3



Given structure, log likelihood of data:

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

*hat because MLE* (pointing to  $\hat{P}$ )  
*MLE* (pointing to  $\hat{P}(x_i, \text{Pa}_{x_i, \mathcal{G}})$ )  
*tried to pick parents s.t. H(X<sub>i</sub> | Pa<sub>X<sub>i</sub></sub>) small; little uncertainty about X<sub>i</sub> | Pa<sub>X<sub>i</sub></sub>* (pointing to the log term)

$$= m \sum_i H(X_i | \text{Pa}_{X_i})$$

$$= m \sum_i (I(X_i, \text{Pa}_{X_i}) - H(X_i))$$

$$= H(A) - H(A|B)$$

*H(A|B) = - \sum\_{a,b} P(a,b) \cdot \log P(a|b)*  
*I(A, B) = - \sum\_{a,b} P(a,b) \log \frac{P(a,b)}{P(a)P(b)}*  
*= H(A) - H(A|B)*

## Decomposable score

- Log data likelihood

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

- Decomposable score:

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

## How many trees are there?

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

## Scoring a tree 1: equivalent trees


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

## Scoring a tree 2: similar trees


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

## 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:

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

- Nodes  $X_1, \dots, X_n$
- Edge (i,j) gets weight

$$\hat{I}(X_i, X_j)$$

## Chow-Liu tree learning algorithm 2

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

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

## 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:

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

## Can we extend Chow-Liu 2



- (Approximately learning) models with tree-width up to  $k$ 
  - [Checheta & 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  $O(N^{2k+6})$ )

## Scoring general graphical models – Model selection problem

What's the best structure?



$\langle x_1^{(1)}, \dots, x_n^{(1)} \rangle$   
...  
 $\langle x_1^{(m)}, \dots, x_n^{(m)} \rangle$



**The more edges, the fewer independence assumptions, the higher the likelihood of the data, but will overfit...**

## Maximum likelihood overfits!

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

- Information never hurts:
  
  
  
  
  
  
  
  
  
  
- Adding a parent always increases score!!!

## Bayesian score avoids overfitting

- Given a structure, distribution over parameters

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

- Difficult integral: use Bayes information criterion (BIC) approximation (equivalent as M! 1)

$$\log P(D | \mathcal{G}) \approx \log P(D | \mathcal{G}, \theta_{\mathcal{G}}) - \frac{\text{NumberParams}(\mathcal{G})}{2} \log M + \mathcal{O}(1)$$

- Note: regularize with MDL score
- Best BN under BIC still NP-hard

## Structure learning for general graphs

- In a tree, a node only has one parent
- **Theorem:**
  - The problem of learning a BN structure with at most  $d$  parents is **NP-hard for any (fixed)  $d \geq 2$**
- Most structure learning approaches use heuristics
  - Exploit score decomposition
  - (Quickly) Describe two heuristics that exploit decomposition in different ways

## Learn BN structure using local search

Starting from  
Chow-Liu tree

**Local search,**  
possible moves:

- Add edge
- Delete edge
- Invert edge

**Score using BIC**

## What you need to know about learning BNs

### ■ Learning BNs

- Maximum likelihood or MAP learns parameters
- Decomposable score
- Best tree (Chow-Liu)
- Best TAN
- Other BNs, usually local search with BIC score

## Unsupervised Learning Clustering K-means

Machine Learning – 10701/15781

Carlos Guestrin

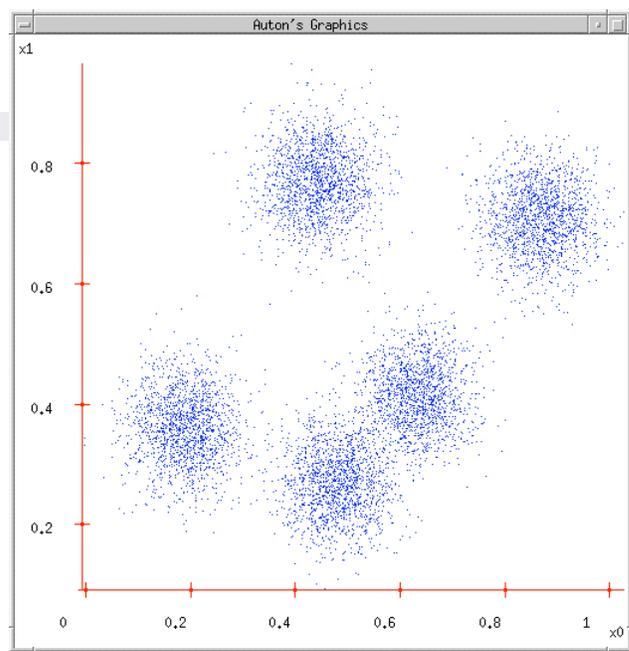
Carnegie Mellon University

November 12<sup>th</sup>, 2007

©2005-2007 Carlos Guestrin

**20**

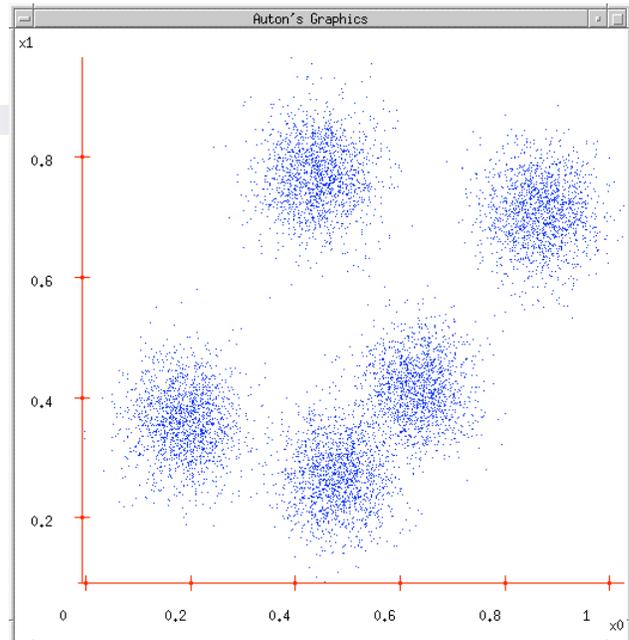
## Some Data



## K-means

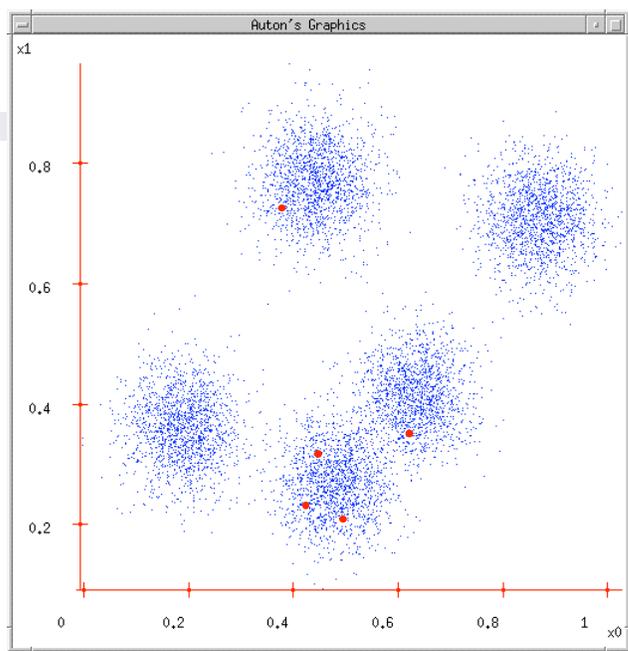


1. Ask user how many clusters they'd like.  
(e.g.  $k=5$ )



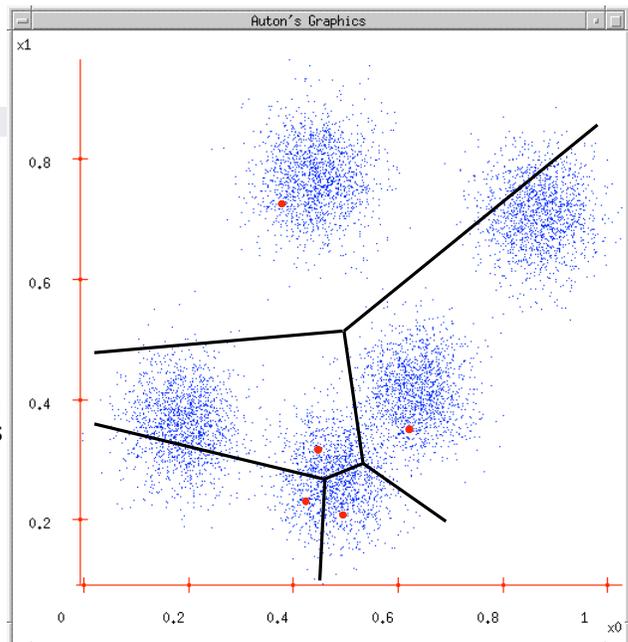
# K-means

1. Ask user how many clusters they'd like. (e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations



# K-means

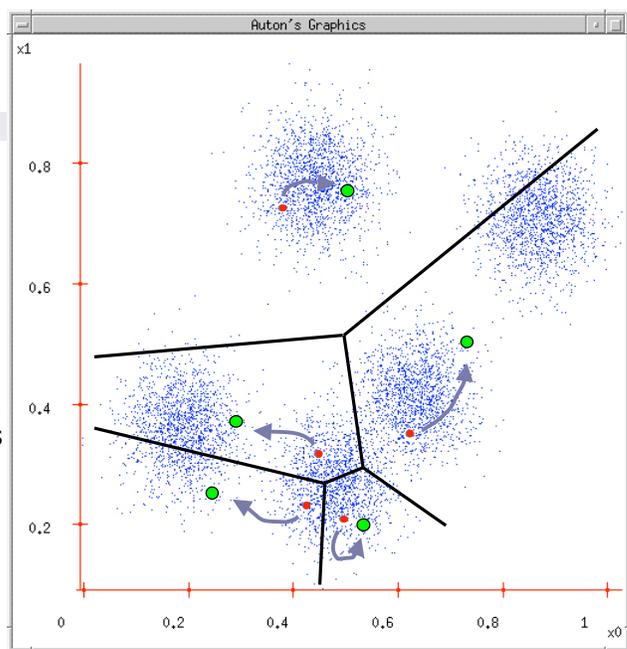
1. Ask user how many clusters they'd like. (e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to. (Thus each Center "owns" a set of datapoints)



# K-means



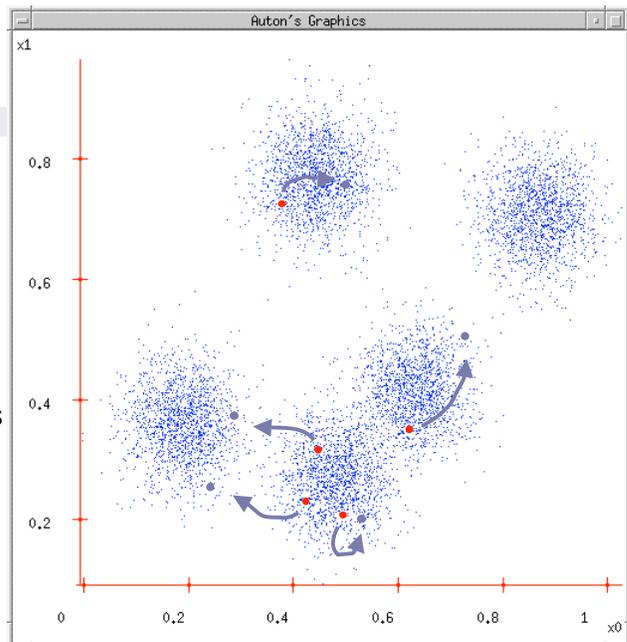
1. Ask user how many clusters they'd like. (e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns



# K-means



1. Ask user how many clusters they'd like. (e.g.  $k=5$ )
2. Randomly guess  $k$  cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns...
5. ...and jumps there
6. ...Repeat until terminated!



# K-means

- Randomly initialize  $k$  centers

- $\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$

- **Classify:** Assign each point  $j \in \{1, \dots, m\}$  to nearest center:

- $C^{(t)}(j) \leftarrow \arg \min_i \|\mu_i - x_j\|^2$

- **Recenter:**  $\mu_i$  becomes centroid of its point:

- $\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: C(j)=i} \|\mu - x_j\|^2$

- Equivalent to  $\mu_i \leftarrow$  average of its points!

# What is K-means optimizing?

- Potential function  $F(\mu, C)$  of centers  $\mu$  and point allocations  $C$ :

- $F(\mu, C) = \sum_{j=1}^m \|\mu_{C(j)} - x_j\|^2$

- Optimal K-means:

- $\min_{\mu} \min_C F(\mu, C)$

## Does K-means converge??? Part 1

- Optimize potential function:

$$\min_{\mu} \min_C F(\mu, C) = \min_{\mu} \min_C \sum_{i=1}^k \sum_{j:C(j)=i} \| \mu_i - x_j \|^2$$

- Fix  $\mu$ , optimize C

## Does K-means converge??? Part 2

- Optimize potential function:

$$\min_{\mu} \min_C F(\mu, C) = \min_{\mu} \min_C \sum_{i=1}^k \sum_{j:C(j)=i} \| \mu_i - x_j \|^2$$

- Fix C, optimize  $\mu$

# Coordinate descent algorithms


$$\min_{\mu} \min_C F(\mu, C) = \min_{\mu} \min_C \sum_{i=1}^k \sum_{j:C(j)=i} \|\mu_i - x_j\|^2$$

- Want:  $\min_a \min_b F(a,b)$
- Coordinate descent:
  - fix a, minimize b
  - fix b, minimize a
  - repeat
- Converges!!!
  - if F is bounded
  - to a (often good) local optimum
    - as we saw in applet (play with it!)
  
- K-means is a coordinate descent algorithm!