Probabilistic Graphical Models

Lecture 4 – Learning Bayesian Networks

> CS/CNS/EE 155 Andreas Krause

Announcements

- Another TA: Hongchao Zhou
- Please fill out the questionnaire about recitations
- Homework 1 out. Due in class Wed Oct 21
- Project proposals due Monday Oct 19

Representing the world using BNs







True distribution P' with cond. ind. I(P')

Bayes net (G,P) with I(P)

- Want to make sure that $I(P) \subseteq I(P')$
- Need to understand CI properties of BN (G,P)

Factorization Theorem



G is an **I-map** of P (independence map)



True distribution P can be represented exactly as Bayesian network (G,P) $P(X_1, ..., X_n) = \prod P(X_i | \mathbf{Pa}_{X_i})$

Additional conditional independencies

- BN specifies joint distribution through conditional parameterization that satisfies Local Markov Property I_{loc}(G) = {(X_i ⊥ Nondescendants_{Xi} | Pa_{Xi})}
- But we also talked about additional properties of CI
 - Weak Union, Intersection, Contraction, ...
- Which additional CI does a particular BN specify?
 - All CI that can be derived through algebraic operations

proving Cl is very cumbersome!!

Is there an easy way to find all independences of a BN just by looking at its graph??

Examples



 $A \perp F$ $A \perp F \mid C$ $A \perp P \mid C \mid X$

Active trails

- An undirected path in BN structure G is called active trail for observed variables O

 {X₁,...,X_n}, if for every consecutive triple of vars X,Y,Z on the path
 - $X \rightarrow Y \rightarrow Z$ and Y is unobserved ($Y \notin O$)
 - $X \leftarrow Y \leftarrow Z$ and Y is unobserved ($Y \notin O$)
 - $X \leftarrow Y \rightarrow Z$ and Y is unobserved ($Y \notin O$)
 - X → Y ← Z and Y or any of Y's descendants is observed
- Any variables X_i and X_j for which [‡] active trail for observations O are called d-separated by O
 We write d-sep(X_i;X_i | O)
- Sets A and B are d-separated given O if d-sep(X,Y | O) for all X∈A, Y∈B. Write d-sep(A; B | O)

Soundness of d-separation

Have seen: P factorizes according to G \leftarrow I_{loc}(G) <= I(P)</p>

• Define I(G) = {(X \perp Y | Z): d-sep_G(X;Y |Z)}

Theorem: Soundness of d-separation
 P factorizes over G → I(G) ⊆ I(P)

- Hence, d-separation captures only true independences
- How about I(G) = I(P)?

Completeness of d-separation

- Theorem: For "almost all" distributions P that factorize over G it holds that I(G) = I(P)
 - "almost all": except for a set of distributions with measure
 0, assuming only that no finite set of distributions has
 measure > 0

Algorithm for d-separation

- How can we check if $X \perp Y \mid Z$?
 - Idea: Check every possible path connecting X and Y and verify conditions
 - Exponentially many paths!!! 🙁
- Linear time algorithm:
 Find all nodes reachable from X
 - 1. Mark **Z** and its ancestors
 - 2. Do breadth-first search starting from X; stop if path is blocked
 - Have to be careful with implementation details (see reading)



Representing the world using BNs





True distribution P' with cond. ind. I(P')

Bayes net (G,P) with I(P)

- Want to make sure that $I(P) \subseteq I(P')$
- Ideally: I(P) = I(P')
- Want BN that exactly captures independencies in P'!

Minimal I-map

 Graph G is called minimal I-map if it's an I-map, and if any edge is deleted
 no longer I-map.

Uniqueness of Minimal I-maps

Is the minimal I-Map unique?



Perfect maps

- Minimal I-maps are easy to find, but can contain many unnecessary dependencies.
- A BN structure G is called P-map (perfect map) for distribution P if I(G) = I(P)
- Does every distribution P have a P-map?

I-Equivalence

- Two graphs G, G' are called I-equivalent if I(G) = I(G')
- I-equivalence partitions graphs into equivalence classes



Skeletons of BNs



I-equivalent BNs must have same skeleton

Immoralities and I-equivalence

 A V-structure X → Y ← Z is called immoral if there is no edge between X and Z ("unmarried parents")

immorality

Theorem: I(G) = I(G') ⇔ G and G' have the same skeleton and the same immoralities.

Today: Learning BN from data

- Want P-map if one exists
- Need to find
 - Skeleton
 - Immoralities

Identifying the skeleton

- When is there an edge between X and Y?
- If r(X + Y)? X
 If r(X + Y) evolution else
 If r(X + Y) evolution else
 If r(X + Y) evolution else
 When is there no edge between X and Y?
 If r(X + Y) evolution
 When is there no edge between X and Y?
 If r(X + Y) evolution
 If r(X + Y)

Algorithm for identifying the skeleton

Not necessarily practical ...

Identifying immoralities

Q

When is X – Z – Y an immorality?

X (X 12/12) X

• Immoral \Leftrightarrow for all **U**, $Z \in \mathbf{U}$: \neg (X \perp Y | **U**)

From skeleton & immoralities to BN Structures

Represent I-equivalence class as partially-directed acyclic graph (PDAG)

How do I convert PDAG into BN?

Have to be careful when orienting edges to avoid cyclos Polytime also in reading

Testing independence

- So far, assumed that we know I(P'), i.e., all independencies associated with true dist. P'
- Often, access to P' only through sample data (e.g., sensor measurements, etc.)
- Given vars X, Y, Z, want to test whether $X \perp Y \mid Z$ $X \perp Y \mid Z \iff \widehat{L}(x; Y \mid Z) = 0$ $\sum_{z} \sum_{xy} e^{(x,y)} \log \frac{P(x,y)}{P(x|Z)} e^{(y)}$ Fit at $e^{p(x,y,z)}$ from data $\Rightarrow \widehat{P}(x,y,z)$ $Compute \widehat{\Gamma}(x; Y|Z) = \sum_{xyZ} \widehat{P}(x,y,z) \log \frac{\widehat{P}(x,y|Z)}{\widehat{P}(x|Z)\widehat{P}(y|Z)}$ Test whether $\widehat{\Gamma}(x; Y|Z) < E$

Next topic: Learning BN from Data

- Two main parts:
 - Learning structure (conditional independencies)
 - Learning parameters (CPDs)

Parameter learning

- Suppose X is Bernoulli distribution (coin flip) with unknown parameter P(X=H) =θ.
- Given training data $D = \{x^{(1)}, \dots, x^{(m)}\}$ (e.g., H H T H H H T T H T H H H..) how do we estimate θ ? $P(D(6) = O^{M_H} \cdot ((-6)^{M_T})$ ny: #of H in J nr: - - T - $\hat{O} = argmax P(D|\Theta) = argmax \log P(D|\Theta)$ = argmax M_H log θ + M_T log(1- θ) Want dlag P(0(6) ! 0.

Maximum Likelihood Estimation

- Given: data set D
- **Hypothesis**: data generated i.i.d. from binomial distribution with $P(X = H) = \theta$
- Optimize for θ which makes D most likely:

Solving the optimization problem $\log P(D(\theta) = m_{H} \log \theta + m_{T} \log (1-\theta)$ $\frac{d}{d\theta} - r = \frac{m_H}{A} - \frac{m_T}{1-\theta} = 0$ $m_{\pm}(1-\theta) = m_{T}\theta$ $m_{H} = \theta(m_{T} + m_{H})$ $=) \Theta = \frac{n_H}{n_T + n_H} = \frac{n_H}{m}$

= Count (X=H)

Learning general BNs

| | Known structure | Unknown structure |
|------------------|-----------------|-------------------|
| Fully observable | Easy i | hard 2. |
| Missing data | hard 3. (EM) | very hand (last) |

Estimating CPDs

Given data D = {(x₁,y₁),...,(x_n,y_n)} of samples from X,Y, want to estimate P(X | Y)

P(X=x |Y=y) = Oxig



MLE for Bayes nets

$$log P(D|\theta) = log TT TT P(X_{i}^{(l)} Pa_{i}^{(l)}; \theta) \qquad Parameton$$

$$= \sum_{l i} \sum_{i} log P(X_{i}^{(l)} | Pa_{i}^{(l)}; \theta_{X_{i}^{(l)}} | Pa_{i}^{(l)}; \theta_{X_{$$

$$\frac{\partial}{\partial \theta_{x_i}(Pa_i)} \log P(D|\theta) = \sum_{j \in \mathbb{N}} \sum_{\substack{\substack{i \in \mathbb{N} \\ i \in \mathbb{N}}}} \sum_{\substack{\substack{\substack{i \in \mathbb{N} \\ i \in \mathbb{N}}}} \sum_{\substack{\substack{\substack{i \in \mathbb{N} \\ i \in \mathbb{N} \\ i \in \mathbb{N}}}} \sum_{\substack{\substack{\substack{i \in \mathbb{N} \\ i \in \mathbb{N} \\ i \in \mathbb{N}}}} \sum_{\substack{\substack{\substack{i \in \mathbb{N} \\ i \in \mathbb$$

Algorithm for BN MLE

Given BN structure G

Learning general BNs

| | Known structure | Unknown structure |
|------------------|---|-------------------|
| Fully observable | Easy! () Get CPDs by counting (for M(E) | ??? |
| Missing data | Hard (EM) | Very hard (later) |

Structure learning

- Two main classes of approaches:
- Constraint based
 - Search for P-map (if one exists):
 - Identify PDAG
 - Turn PDAG into BN (using algorithm in reading)
 - Key problem: Perform independence tests
- - Define scoring function (e.g., likelihood of data)
 - Think about structure as parameters
 - More common; can solve simple cases exactly

MLE for structure learning

For fixed structure, can compute likelihood of data

$$\log P(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G}) = \sum_{\ell} \sum_{i} \log P(X_{i} = x_{i}^{(\ell)} \mid \mathbf{Pa}_{i} = \mathbf{pa}_{i}^{(\ell)} \stackrel{\circ}{\rho} \stackrel{\circ}{\sigma}_{i}^{(\ell)})$$

$$\stackrel{\mathsf{M}(\mathsf{E}}{=} \sum_{i} \sum_{X_{i}} \sum_{\mathsf{P}^{\mathsf{a}_{i}}} \mathcal{C}_{\mathsf{ourt}} \stackrel{\circ}{\leftarrow} (X_{i}, \mathsf{Pa}_{i}) \stackrel{\circ}{\rho} \stackrel{\circ}{\rho} \stackrel{\circ}{\left(\mathsf{pa}_{i}\right)}$$

$$= m \sum_{i} \sum_{X_{i}} \sum_{\mathsf{pa}_{i}} \stackrel{\circ}{\rho} \stackrel{\circ}{\left(X_{i}, \mathsf{pa}_{i}\right)} \stackrel{\circ}{\log} \frac{\stackrel{\circ}{\rho} \stackrel{\circ}{\left(X_{i}, \mathsf{pa}_{i}\right)}{\stackrel{\circ}{\rho} \left(\mathsf{pa}_{i}\right)} \stackrel{\circ}{\rho} \stackrel{\circ}{\left(\mathsf{pa}_{i}\right)} \stackrel{\circ}{\rho} \stackrel{\circ}{\left(\mathsf{pa}_{i}\right)}$$

$$= m \sum_{i} \sum_{X_{i}} \sum_{\mathsf{pa}_{i}} \stackrel{\circ}{\rho} \stackrel{\circ}{\left(X_{i}, \mathsf{pa}_{i}\right)} \stackrel{\circ}{\rho} \stackrel{\circ}{\left(\mathsf{pa}_{i}\right)} \stackrel{\circ}{\rho} \stackrel{\circ}$$

Decomposable score

Log-data likelihood

$$\log \widehat{P}(\mathcal{D} \mid \theta_{\mathcal{G}}, \mathcal{G}) = m \sum_{i} \widehat{I}(X_{i}, \mathbf{Pa}_{i}) - m \sum_{i} \widehat{H}(X_{i})$$
independent of graph

- MLE score decomposes over families of the BN (nodes + parents)
- Score(G ; D) = \sum_{i} FamScore(X_i | Pa_i; D)
- Can exploit for computational efficiency!

Finding the optimal MLE structure

Log-likelihood score:

Score(
$$\mathcal{G}; \mathcal{D}$$
) = $\sum_{i} \widehat{I}(X_i, \mathbf{Pa}_i)$

- Want G^{*} = argmax_G Score(G ; D)
- Lemma: $G \subseteq G' \rightarrow Score(G; D) \leq Score(G'; D)$

"Information never huts" RV X, A CB H(×IA) Z H(×IB) I(X;A) = H(×I-H(×IA) => I(X;B) = I(×;A)

Finding the optimal MLE structure

Optimal solution for MLE is always the fully connected graph!!! ⁽³⁾

Non-compact representation; Overfitting!!

Solutions:

- Priors over parameters / structures (later)
- Constraint optimization (e.g., bound #parents)

Constraint optimization of BN structures

- Theorem: for any fixed d ≥ 2, finding the optimal BN (w.r.t. MLE score) is NP-hard
- What about d=1??
- Want to find optimal tree!

Finding the optimal tree BN

Scoring function

$$Score(\mathcal{G}; \mathcal{D}) = \sum_{i} \widehat{I}(X_i, \mathbf{Pa}_i)$$

Scoring a tree

$$\begin{array}{c} (\mathcal{B} \longrightarrow \mathcal{B} \longrightarrow \mathcal{B}) \\ \tilde{I}(\mathcal{Z};\mathcal{X}) + \tilde{I}(\mathcal{Y};\mathcal{X}) \\ \mathcal{B}/c \cdot \tilde{I}(\mathcal{X};\mathcal{Y}) = \tilde{I}(\mathcal{Y},\mathcal{X}) \end{array} \end{array}$$

Some shelotor => same scorp

Finding the optimal tree skeleton

- Can reduce to following problem:
- Given graph G = (V,E), and nonnegative weights w_e for each edge e=(X_i,X_j)
 In our case: w_e = I(X_i,X_i)
- Want to find tree T \subseteq E that maximizes $\sum_{e \in T} w_e$
- Maximum spanning tree problem!
- Can solve in time O(|E| log |E|)!

Chow-Liu algorithm

For each pair X_i, X_i of variables compute

$$\widehat{P}(x_i, x_j) = \frac{\operatorname{Count}(x_i, x_j)}{m}$$

Compute mutual information

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

- Define complete graph with weight of edge (X_i,X_i) given by the mutual information
- Find maximum spanning tree
 skeleton
- Orient the skeleton using breadth-first search

Generalizing Chow-Liu

- Tree-augmented Naïve Bayes Model [Friedman '97]
- If evidence variables are correlated, Naïve Bayes models can be overconfident
- Key idea: Learn optimal tree for conditional distribution P(X₁,...,X_n | Y)
- Can do optimally using Chow-Liu (homework! ③)

Tasks

- Subscribe to Mailing list <u>https://utils.its.caltech.edu/mailman/listinfo/cs155</u>
- Select recitation times
- Read Koller & Friedman Chapter 17.1-17.3, 18.1-2, 18.4.1
- Form groups and think about class projects. If you have difficulty finding a group, email Pete Trautman