

CS262a Bayesian Networks

© Tingfeng Xia @ UCLA

Winter Quarter, 2022

This work is licensed under a [Creative Commons “Attribution-NonCommercial-ShareAlike 4.0 International”](#) license.



Preface

This document is intentionally kept short, and is not a complete summary of what is covered in this course. Excerpt from online:

The objective of this class is to provide an in-depth exposition of knowledge representation, reasoning, and machine learning under uncertainty using the framework of Bayesian networks. Both theoretical underpinnings and practical considerations will be covered, with a special emphasis on constructing and learning graphical models, and on various exact and approximate inference algorithms. Additional topics include logical approaches to probabilistic inference, compilation techniques, sensitivity analysis, undirected graphical models, and statistical relational learning.

Instructor: Professor Adnan Darwiche

Book: Adnan Darwiche. *Modeling and Reasoning with Bayesian Networks*. Cambridge University Press 2009.

1 Propositional Logic

1.1 Principle Logical Forms

Inconsistent. Something that never holds; $\text{Mods}(\cdot) = \emptyset$; $\Pr(\alpha) = 0$

Valid. Something that always holds; $\text{Mods}(\cdot) = \Omega$; $\Pr(\alpha) = 1$

Equivalent. $\text{Mods}(\alpha) = \text{Mods}(\beta)$

Mutual Exclusive. $\text{Mods}(\alpha) \cap \text{Mods}(\beta) = \emptyset$

Exhaustive. $\text{Mods}(\alpha) \cup \text{Mods}(\beta) = \Omega$

Entailment / Implication. $\alpha \models \beta \triangleq \text{Mods}(\alpha) \subseteq \text{Mods}(\beta)$

1.2 Equivalent Forms

- $\text{Mods}(\alpha \wedge \beta) = \text{Mods}(\alpha) \cap \text{Mods}(\beta)$
- $\text{Mods}(\alpha \vee \beta) = \text{Mods}(\alpha) \cup \text{Mods}(\beta)$
- $\text{Mods}(\neg\alpha) = \overline{\text{Mods}(\alpha)}$

1.3 Instantiation Agreement

Two instantiation, each of which can cover a subset of different variables, are said to be compatible with each other if they agree on all common variables. Denoted as $\mathbf{x} \sim \mathbf{y}$.

1.4 Information Theory

Entropy.

$$\text{ENT}(X) = - \sum_x \text{Pr}(x) \log \text{Pr}(x) \quad (1.1)$$

where $0 \log 0 = 0$ by convention. With a higher entropy, we say that it is more chaotic.

Conditional Entropy.

$$\text{ENT}(X|Y) = \sum_y \text{Pr}(y) \text{ENT}(X|y) \quad \text{where} \quad \text{ENT}(X|y) = - \sum_x \text{Pr}(x|y) \log \text{Pr}(x|y) \quad (1.2)$$

Conditioning never increases the entropy, i.e.

$$\text{ENT}(X|Y) \leq \text{ENT}(X) \quad (1.3)$$

Mutual Information

$$\text{MI}(X; Y) = \sum_{x,y} \text{Pr}(x,y) \log \frac{\text{Pr}(x,y)}{\text{Pr}(x)\text{Pr}(y)} \quad (1.4)$$

$$= \text{ENT}(X) - \text{ENT}(X|Y) \quad (1.5)$$

$$= \text{ENT}(Y) - \text{ENT}(Y|X) \quad (1.6)$$

Conditional Mutual Information

$$MI(X; Y|Z) = \sum_{x,y,z} \Pr(x, y, z) \log \frac{\Pr(x, y|z)}{\Pr(x|z)\Pr(y|z)} \quad (1.7)$$

$$= ENT(X|Z) - ENT(X|Y, Z) \quad (1.8)$$

$$= ENT(Y|Z) - ENT(Y|X, Z) \quad (1.9)$$

2 Probability Calculus

2.1 Bayesian Conditioning

Bayesian Condition is specified by the formula

$$\Pr(\alpha|\beta) = \frac{\Pr(\alpha \wedge \beta)}{\Pr(\beta)} \quad (2.1)$$

In particular, not to be confused with Bayesian inference (to be added later).

2.2 Independence and Notations

Independence.

$$\alpha \perp\!\!\!\perp \beta \iff \Pr(\alpha|\beta) = \Pr(\alpha) \vee \Pr(\beta) = 0 \quad (2.2)$$

$$\iff \Pr(\alpha \wedge \beta) = \Pr(\alpha)\Pr(\beta) \quad (2.3)$$

Conditional Independence.

$$(\alpha \perp\!\!\!\perp \beta)|\gamma \iff \Pr(\alpha|\beta \wedge \gamma) = \Pr(\alpha|\gamma) \vee \Pr(\beta \wedge \gamma) = 0 \quad (2.4)$$

$$\iff \Pr(\alpha \wedge \beta|\gamma) = \Pr(\alpha|\gamma)\Pr(\beta|\gamma) \vee \Pr(\gamma) = 1 \quad (2.5)$$

Set Independence

$$I_{Pr}(X, Z, Y) \iff (x \perp\!\!\!\perp y)|z, \quad \forall x, y, z \in X, Y, Z \quad (2.6)$$

3 Bayesian Networks

3.1 Soft Evidence

3.1.1 All Things Considered Method

We normalize / rescale w according to new evidence.

$$\Pr'(w) = \begin{cases} \frac{\Pr'(\beta)}{\Pr(\beta)} \Pr(w) & \text{if } w \models \beta \\ \frac{\Pr'(\neg\beta)}{\Pr(\neg\beta)} \Pr(w) & \text{if } w \models \neg\beta \end{cases} \quad (3.1)$$

The closed form is called the Jefferey's Rule.

Jeffery's Rule

$$\Pr'(\alpha) = q\Pr(\alpha|\beta) + (1 - q)\Pr(\alpha|\neg\beta) \quad (3.2)$$

Jeffery's Rule - General Case

$$\Pr'(\alpha) = \sum_{i=1}^n \Pr'(\beta_i) \Pr(\alpha|\beta_i) \quad (3.3)$$

3.1.2 Nothing-else Considered Method

Odds.

$$O(\beta) = \frac{\Pr(\beta)}{\Pr(\neg\beta)} \quad (3.4)$$

Bayes Factor

$$k = \frac{O'(\beta)}{O(\beta)} = \frac{\Pr'(\beta)/\Pr'(\neg\beta)}{\dots} \quad (3.5)$$

from where we can expand and organize

$$\Pr'(\beta) = \frac{k\Pr(\beta)}{k\Pr(\beta) + \Pr(\neg\beta)} \quad (3.6)$$

Closed Form Solution.

$$\Pr'(\alpha) = \frac{k\Pr(\alpha \wedge \beta) + \Pr(\alpha \wedge \neg\beta)}{k\Pr(\beta) + \Pr(\neg\beta)} \quad (3.7)$$

3.2 Noisy Sensors

$$O'(\beta) = \underbrace{\frac{1-f_n}{f_p}}_{k^+} O(\beta) \quad O'(\beta) = \underbrace{\frac{f_n}{1-f_p}}_{k^-} O(\beta) \quad (3.8)$$

3.3 Markov Assumptions

$$\text{Markov}(G) = \{I_{\text{Pr}}(V, \text{Parents}(V), \text{ND}(V))\}_V \quad (3.9)$$

where ND means non-descendants, and includes all nodes except for $V, \text{Parents}(V)$ and $\text{Descendants}(V)$ (all the way till leaf)

3.4 Graphoid Axioms

Symmetry.

$$I_{\text{Pr}}(X, Z, Y) \iff I_{\text{Pr}}(Y, Z, X) \quad (3.10)$$

Decomposition.

$$I_{\text{Pr}}(X, Z, Y \cup W) \implies I_{\text{Pr}}(X, Z, Y) \wedge I_{\text{Pr}}(X, Z, W) \quad (3.11)$$

Weak Union.

$$I_{\text{Pr}}(X, Z, Y \cup W) \implies I_{\text{Pr}}(X, Z \cup Y, W) \quad (3.12)$$

Contraction.

$$I_{\text{Pr}}(X, Z, Y) \wedge I_{\text{Pr}}(X, Z \cup Y, W) \implies I_{\text{Pr}}(X, Z, Y \cup W) \quad (3.13)$$

Triviality.

$$I_{\text{Pr}}(X, Z, \emptyset) \quad (3.14)$$

3.5 Positive Graphoid Axioms

... includes everything from Graphoid Axioms (Section 3.4) and in addition has

Intersection.

$$I_{\text{Pr}}(X, Z \cup W, Y) \wedge I_{\text{Pr}}(X, Z \cup Y, W) \implies I_{\text{Pr}}(X, Z, Y \cup W) \quad (3.15)$$

3.6 D-seperation Linear Prune Theorem

3.7 D-seperation Properties

Soundness.

$$\text{dsep}_G(X, Z, Y) \implies I_{\text{Pr}}(X, Z, Y) \quad (3.16)$$

(Weak) Completeness. There exists a parametrization Θ that for every DAG G such that

$$I_{\text{Pr}}(X, Z, Y) \iff \text{dsep}_G(X, Z, Y) \quad (3.17)$$

4 Inference by Factor Elimination

4.1 Elimination Trees

Variables. $\text{vars}(i)$ denotes the variables mentioned at node i . $\text{vars}(i, j)$ denotes all variables mentioned in nodes to the i -side of the graph (inclusive). Hence, it holds that $\text{vars}(i) \subseteq \text{vars}(i, j)$.

Separators.

$$S_{ij} \triangleq \text{vars}(i, j) \cap \text{vars}(j, i) \quad (4.1)$$

Clusters.

$$C_i \triangleq \text{vars}(i) \cup \bigcup_j S_{ij} \quad (4.2)$$

5 Inference by Conditioning

5.1 Run time Comparison

5.1.1 Variable Elimination (VE)

Let w be the width of the tree, n denote the number of variables, and $|Q|$ as query variable size.

Time Complexity. $O(n \exp(w))$

Space Complexity.

$$O(n \exp(w) + n \exp(|Q|)) \equiv O(n \exp(w)), \quad \text{if } |Q| < \infty \quad (5.1)$$

5.1.2 Message Passing

One Message. The cost to pass one message is

$$O(\exp(|C_i|)) \quad (5.2)$$

where $|C_i|$ is the cluster size.

Every Message.

$$O(\exp(w)) \quad (5.3)$$

what is this again?

where w is the width of the elimination tree.

Amount of Messages. The total amount of messages is

$$O(2(m-1)) \quad \text{where } m = |V| \quad (5.4)$$

since we have a tree structure, meaning we have exactly $(m-1)$ edges and each edge can have forward backward each once.

All Messages - All Cluster Marginals. The total time to pass all messages and compute all cluster marginals is

$$O(m \exp(w)) \quad \text{or} \quad O(n \exp(w)), \text{ for } O(n) \text{ edges.} \quad (5.5)$$

5.1.3 Polytree / Belief Propagation

Runtime. Define k as the max number of parents in the poly tree, then k is the same as the width of elimination tree. Let, also, n denote the number of nodes in the polytree. The algorithm has runtime

$$O(n \exp(k)) \quad (5.6)$$

5.1.4 Cut Set Conditioning

Time and Space (Total) Complexity.

$$O(n \exp(k)) \quad \text{where } n = |N| \text{ and } k \text{ is width} \quad (5.7)$$

5.1.5 Any Space Recursive Cut Set

	no cache	all cache	$\Delta no \rightarrow all$
space	$O(wn)$	$O(n \exp(w))$	\uparrow
time	$O(n \exp(w \log n))$	$O(n \exp(w))$	\downarrow

6 Compiling Bayesian Networks

6.1 Network Polynomials

The network polynomial is a summation over all instantiations of a network,

$$f \triangleq \sum_z \prod_{\theta_{x|u} \sim z} \theta_{x|u} \prod_{\lambda_x \sim z} \lambda_x \quad (6.1)$$

6.2 AC Properties

AC Size. of an AC is defined as the number of edges in the circuit.

AC Complexity. is the size of smallest AC that represents the network polynomial.

Decomposable. At each \star node, we need

$$\text{vars}(AC_A) \cap \text{vars}(AC_B) = \emptyset \quad (6.2)$$

Deterministic. At each $+$ node, we require at most one positive input is non-zero for all *complete instantiation*.

Smooth. At each $+$ node, we require

$$\text{vars}(AC_A) = \text{vars}(AC_B) \quad (6.3)$$

AC for Marginals. requires decomposable and smooth. This guarantees that sub-circuits are of complete variable instantiations.

AC for Marginals and MPE. requires all three above: decomposable, deterministic, and smooth. The additional determinism guarantees a 1-to-1 mapping between sub-circuits and complete variable instantiations.

6.3 AC Derivative Probabilistic Implications

$$\frac{\partial f}{\partial \lambda_{\mathbf{x}}}(\mathbf{e}) = \Pr(\mathbf{x}, \mathbf{e} - X) \quad (6.4)$$

and

$$\theta_{\mathbf{x}|\mathbf{u}} \frac{\partial f}{\partial \theta_{\mathbf{x}|\mathbf{u}}}(\mathbf{e}) = \Pr(\mathbf{x}, \mathbf{u}, \mathbf{e}) \quad (6.5)$$