

# Bayesian Networks

David Rosenberg

New York University

April 21, 2016

# Introduction

# Probabilistic Reasoning

- **Represent** system of interest by a set of random variables

$$(X_1, \dots, X_d).$$

- Suppose by research or ML, we get a joint probability distribution

$$p(x_1, \dots, x_d).$$

- We'd like to be able to do “**inference**” on this model – essentially, answer queries:
  - ① What is the most likely of value  $X_1$ ?
  - ② What is the most likely of value  $X_1$ , given we've observed  $X_2 = 1$ ?
  - ③ Distribution of  $(X_1, X_2)$  given observation of  $(X_3 = x_3, \dots, X_d = x_d)$ ?

## Example: Medical Diagnosis

- Variables for each **symptom**
  - fever, cough, fast breathing, shaking, nausea, vomiting
- Variables for each **disease**
  - pneumonia, flu, common cold, bronchitis, tuberculosis
- Diagnosis is performed by **inference** in the model:

$$p(\text{pneumonia} = 1 \mid \text{cough} = 1, \text{fever} = 1, \text{vomiting} = 0)$$

- The QMR-DT (Quick Medical Reference - Decision Theoretic) has
  - 600 diseases
  - 4000 symptoms

# Discrete Probability Distribution Review

# Some Notation

- This lecture we'll only be considering **discrete** random variables.
- Capital letters  $X_1, \dots, X_d, Y$ , etc. denote **random variables**.
- Lower case letters  $x_1, \dots, x_n, y$  denote the values taken.
- Probability that  $X_1 = x_1$  and  $X_2 = x_2$  will be denoted

$$\mathbb{P}(X_1 = x_1, X_2 = x_2).$$

- We'll generally write things in terms of the probability mass function:

$$p(x_1, x_2, \dots, x_d) := \mathbb{P}(X_1 = x_1, X_2 = x_2, \dots, X_d = x_d)$$

# Representing Probability Distributions

- Let's consider the case of discrete random variables.
- Conceptually, everything can be represented with probability tables.
- Variables
  - Temperature  $T \in \{\text{hot}, \text{cold}\}$
  - Weather  $W \in \{\text{sun}, \text{rain}\}$

$t$	$p(t)$
hot	0.5
cold	0.5

$w$	$p(w)$
sun	0.6
rain	0.4

- These are the **marginal** probability distributions.
- To do reasoning, we need the **joint probability distribution**.

# Joint Probability Distributions

- A joint probability distribution for  $T$  and  $W$  is given by

$t$	$w$	$p(t, w)$
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

- With the joint we can answer question such as
  - $p(\text{sun} \mid \text{hot})=?$
  - $p(\text{rain} \mid \text{cold})=?$



# Representing Joint Distributions

- Consider random variables  $X_1, \dots, X_d \in \{0, 1\}$ .
- How many parameters do we need to represent the joint distribution?
- Joint probability table has  $2^d$  rows.
- For QMR-DT, that's  $2^{4600} > 10^{1000}$  rows.
- That's not going to happen.
- Only  $\sim 10^{80}$  atoms in the universe.
- Having exponentially many parameters is a problem for
  - storage
  - computation (inference is summing over exponentially many rows)
  - statistical estimation / learning

# How to Restrict the Complexity?

- Restrict the space of probability distributions
- We will make various **independence** assumptions.
- Extreme assumption:  $X_1, \dots, X_d$  are **mutually independent**.

## Definition

Discrete random variables  $X_1, \dots, X_d$  are **mutually independent** if their joint probability mass function (PMF) factorizes as

$$p(x_1, x_2, \dots, x_d) = p(x_1)p(x_2) \cdots p(x_d).$$

- Note: We usually just write **independent** for “mutually independent”.
- How many parameters to represent the joint distribution, assuming independence?

# Assume Full Independence

- How many parameters to represent the joint distribution?
- Say  $p(X_i = 1) = \theta_i$ , for  $i = 1, \dots, d$ .
- **Clever representation:** Since  $x_i \in \{0, 1\}$ , we can write

$$\mathbb{P}(X_i = x_i) = \theta_i^{x_i} (1 - \theta_i)^{1-x_i}.$$

- Then by independence,

$$p(x_1, \dots, x_d) = \prod_{i=1}^d \theta_i^{x_i} (1 - \theta_i)^{1-x_i}$$

- How many parameters?
- $d$  parameters needed to represent the joint.

# Conditional Interpretation of Independence

- Suppose  $X$  and  $Y$  are independent, then

$$p(x | y) = p(x).$$

- Proof:

$$\begin{aligned} p(x | y) &= \frac{p(x, y)}{p(y)} \\ &= \frac{p(x)p(y)}{p(y)} = p(x). \end{aligned}$$

- With full independence, we have no relationships among variables.
- Information about one variable says nothing about any other variable.
  - Would mean diseases don't have symptoms.

# Conditional Independence

- Consider 3 events:
  - 1  $W = \{\text{The grass is wet}\}$
  - 2  $S = \{\text{The road is slippery}\}$
  - 3  $R = \{\text{It's raining}\}$
- These events are certainly **not** independent.
  - Raining ( $R$ )  $\implies$  Grass is wet AND The road is slippery ( $W \cap S$ )
  - Grass is wet ( $W$ )  $\implies$  More likely that the road is slippery ( $S$ )
- Suppose we know that **it's raining**.
  - Then, we learn that **the grass is wet**.
  - Does this tell us anything new about whether **the road is slippery**?
- Once we know  $R$ , then  $W$  and  $S$  become independent.
- This is called **conditional independence**, and we'll denote it as

$$W \perp S \mid R.$$

# Conditional Independence

## Definition

We say  $W$  and  $S$  are **conditionally independent** given  $R$ , denoted

$$W \perp S \mid R,$$

if the conditional joint factorizes as

$$p(w, s \mid r) = p(w \mid r)p(s \mid r).$$

Also holds when  $W$ ,  $S$ , and  $R$  represent **sets of random variables**.

- Can have conditional independence without independence.
- Can have independence without conditional independence.

## Example: Rainy, Slippery, Wet

- Consider 3 events:
  - 1  $W = \{\text{The grass is wet}\}$
  - 2  $S = \{\text{The road is slippery}\}$
  - 3  $R = \{\text{It's raining}\}$
- Represent joint distribution as

$$\begin{aligned}
 p(w, s, r) &= p(w, s | r)p(r) && \text{(no assumptions so far)} \\
 &= p(w | r)p(s | r)p(r) && \text{(assuming } W \perp S | R)
 \end{aligned}$$

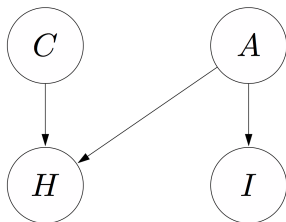
- How many parameters to specify the joint?
  - $p(w | r)$  requires two parameters: one for  $r = 1$  and one for  $r = 0$ .
  - $p(s | r)$  requires two.
  - $p(r)$  requires one parameter,
- Full joint: 7 parameters. Conditional independence: 5 parameters.  
Full independence: 3 parameters.

# Bayesian Networks



# Bayesian Networks: Introduction

- Bayesian Networks are
  - used to specify joint probability distributions that
  - have a particular factorization.



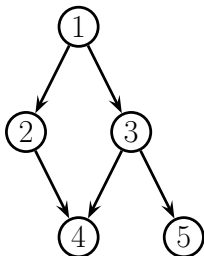
$$p(c, h, a, i) = p(c)p(a) \times p(h | c, a)p(i | a)$$

- With practice, one can read conditional independence relationships directly from the graph.

# Directed Graphs

A **directed graph** is a pair  $G = (\mathcal{V}, \mathcal{E})$ , where

- $\mathcal{V} = \{1, \dots, d\}$  is a set of **nodes** and
- $\mathcal{E} = \{(s, t) \mid s, t \in \mathcal{V}\}$  is a set of **directed edges**.



$$\text{Parents}(5) = \{3\}$$

$$\text{Parents}(4) = \{2, 3\}$$

$$\text{Children}(3) = \{4, 5\}$$

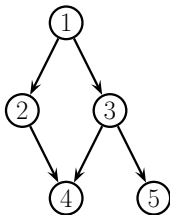
$$\text{Descendants}(1) = \{2, 3, 4, 5\}$$

$$\text{NonDescendants}(3) = \{1, 2\}$$

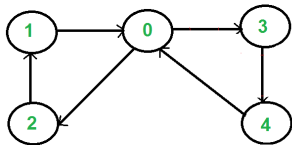
# Directed Acyclic Graphs (DAGs)

A **DAG** is a directed graph with **no directed cycles**.

DAG



Not a DAG



Every DAG has a **topological ordering**, in which parents have lower numbers than their children.

<http://www.geeksforgeeks.org/wp-content/uploads/SCC1.png> and KPM Figure 10.2(a).

# Bayesian Networks

## Definition

A **Bayesian network** is a

- DAG  $G = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V} = \{1, \dots, d\}$ , and
- a corresponding set of random variables  $X = \{X_1, \dots, X_d\}$

where

- the joint probability distribution over  $X$  factorizes as

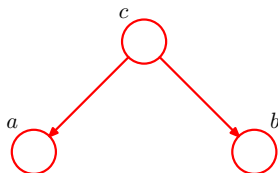
$$p(x_1, \dots, x_d) = \prod_{i=1}^d p(x_i \mid x_{\text{Parents}(i)}).$$

Bayesian networks are also known as

- **directed graphical models**, and
- **belief networks**.

# Conditional Independencies

# Bayesian Networks: “A Common Cause”



$$p(a, b, c) = p(c)p(a | c)p(b | c)$$

Are  $a$  and  $b$  independent? ( $c$ =Rain,  $a$ =Slippery,  $b$ =Wet?)

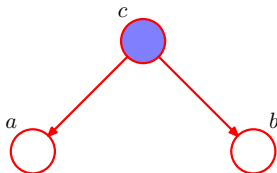
$$p(a, b) = \sum_c p(c)p(a | c)p(b | c),$$

which in general will not be equal to  $p(a)p(b)$ .

---

From Bishop's *Pattern recognition and machine learning*, Figure 8.15.

# Bayesian Networks: “A Common Cause”



$$p(a, b, c) = p(c)p(a | c)p(b | c)$$

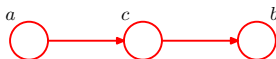
Are  $a$  and  $b$  independent, conditioned on observing  $c$ ? ( $c$ =Rain,  $a$ =Slippery,  $b$ =Wet?)

$$\begin{aligned} p(a, b | c) &= p(a, b, c) / p(c) \\ &= p(a | c)p(b | c) \end{aligned}$$

So  $a \perp b | c$ .

From Bishop's *Pattern recognition and machine learning*, Figure 8.16.

# Bayesian Networks: “An Indirect Effect”



$$p(a, b, c) = p(a)p(c | a)p(b | c)$$

Are  $a$  and  $b$  independent? (Note: This is a **Markov chain**)  
(e.g.  $a$ =raining,  $c$ =wet ground,  $b$ =mud on shoes)

$$\begin{aligned} p(a, b) &= \sum_c p(a, b, c) \\ &= p(a) \sum_c p(c | a)p(b | c) \end{aligned}$$

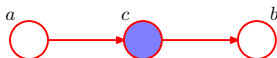
So doesn't factorize, thus not independent, in general.

---

From Bishop's *Pattern recognition and machine learning*, Figure 8.17.



# Bayesian Networks: “An Indirect Effect”



$$p(a, b, c) = p(a)p(c | a)p(b | c)$$

Are  $a$  and  $b$  independent after observing  $c$ ?  
 (e.g.  $a$ =raining,  $c$ =wet ground,  $b$ =mud on shoes)

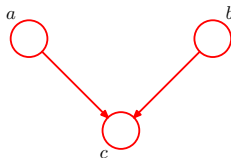
$$\begin{aligned}
 p(a, b | c) &= p(a, b, c) / p(c) \\
 &= p(a)p(c | a)p(b | c) / p(c) \\
 &= p(a | c)p(b | c)
 \end{aligned}$$

So  $a \perp b | c$ .

---

From Bishop's *Pattern recognition and machine learning*, Figure 8.18.

# Bayesian Networks: “A Common Effect”



$$p(a, b, c) = p(a)p(b)p(c | a, b)$$

Are  $a$  and  $b$  independent? ( $a$ =course difficulty,  $b$ =knowledge,  $c$ = grade)

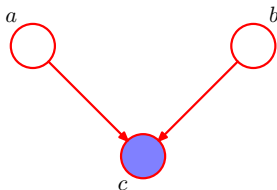
$$\begin{aligned}
 p(a, b) &= \sum_c p(a)p(b)p(c | a, b) \\
 &= p(a)p(b) \sum_c p(c | a, b) \\
 &= p(a)p(b)
 \end{aligned}$$

So  $a \perp b$ .

---

From Bishop's *Pattern recognition and machine learning*, Figure 8.19.

# Bayesian Networks: “A Common Effect” or “V-Structure”



$$p(a, b, c) = p(a)p(b)p(c | a, b)$$

Are  $a$  and  $b$  independent, given observation of  $c$ ? ( $a$ =course difficulty,  $b$ =knowledge,  $c$ = grade)

$$p(a, b | c) = p(a)p(b)p(c | a, b)/p(c)$$

which does not factorize into  $p(a | c)p(b | c)$ , in general.

---

From Bishop's *Pattern recognition and machine learning*, Figure 8.20.

# Conditional Independence from Graph Structure

- In general, given 3 sets of nodes  $A$ ,  $B$ , and  $C$
- How can we determine whether

$$A \perp B \mid C?$$

- There is a purely graph-theoretic notion of “**d-separation**” that is equivalent to conditional independence.
- Suppose we have observed  $C$  and we want to do inference on  $A$ .
- We could ignore any evidence collected about  $B$ , where  $A \perp B \mid C$ .
- See KPM Section 10.5.1 for details.

# Markov Blanket

- Suppose we have a very large Bayesian network.
- We're interested in a single variable  $A$ , which we cannot observe.
- To get maximal information about  $A$ , do we have to observe all other variables?
- No! We only need to observe the **Markov blanket** of  $A$ :

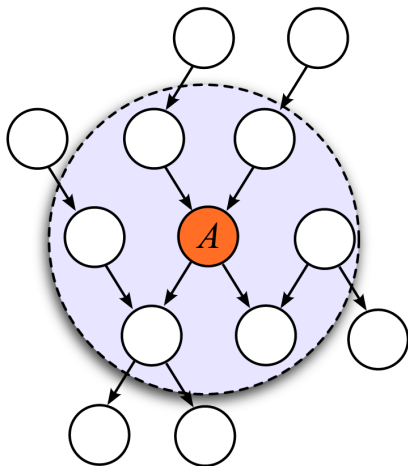
$$p(A \mid \text{all other nodes}) = p(A \mid \text{MarkovBlanket}(A)).$$

- In a Bayesian network, the Markov blanket of  $A$  consists of
  - the parents of  $A$
  - the children of  $A$
  - the “co-parents” of  $A$ , i.e. the parents of the children of  $A$

(See KPM Sec. 10.5.3 for details.)

# Markov Blanket

Markov Blanket of  $A$  in a Bayesian Network:



From [http://en.wikipedia.org/wiki/Markov\\_blanket](http://en.wikipedia.org/wiki/Markov_blanket): "Diagram of a Markov blanket" by Laughsinthestocks - Licensed under CC0 via Wikimedia Commons

## When to use Bayesian Networks?

# Bayesian Networks

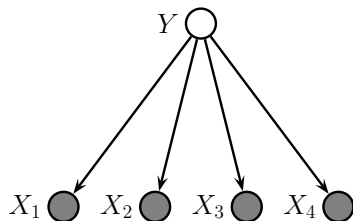
- Bayesian Networks are great when
  - you know something about the relationships between your variables, or
  - you will routinely need to make inferences with incomplete data.
- Challenges:
  - The naive approach to inference doesn't work beyond small scale.
  - Need more sophisticated algorithm:
    - exact inference
    - approximate inference



# Naive Bayes

# Binary Naive Bayes: A Generative Model for Classification

- $\mathcal{X} = \left\{ (X_1, X_2, X_3, X_4) \in \{0, 1\}^4 \right\}$        $\mathcal{Y} = \{0, 1\}$  be a class label.
- Consider the Bayesian network depicted below:



- BN structure implies joint distribution factors as:

$$p(x_1, x_2, x_3, x_4, y) = p(y)p(x_1 | y)p(x_2 | y)p(x_3 | y)p(x_4 | y)$$

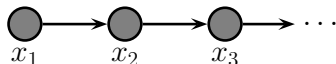
- Features  $X_1, \dots, X_4$  are independent given the class label  $Y$ .

KPM Figure 10.2(a).

# Markov Models

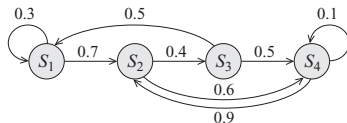
# Markov Chain Model

- A Markov chain model has structure:



$$p(x_1, x_2, x_3, \dots) = p(x_1)p(x_2 | x_1)p(x_3 | x_2)\dots$$

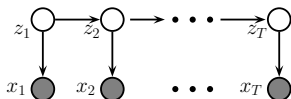
- Conditional distributions  $p(x_i | x_{i-1})$  is called the **transition model**.
- When conditional distribution independent of  $i$ , called **time-homogeneous**.
- 4-state transition model for  $X_i \in \{S_1, S_2, S_3, S_4\}$ :



KPM Figure 10.3(a) and Koller and Friedman's *Probabilistic Graphical Models* Figure 6.04.

# Hidden Markov Model

- A hidden Markov model (HMM) has structure:



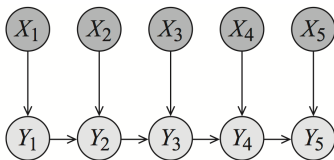
$$p(x_1, z_1, x_2, z_2, x_3, z_3, \dots) = p(z_1) \underbrace{\prod_{t=2}^T p(z_t | z_{t-1})}_{\text{Transition Model}} \underbrace{\prod_{t=1}^T p(x_t | z_t)}_{\text{Observation Model}}$$

- At deployment time, we typically only observe  $X_1, \dots, X_T$ .
- Want to infer  $Z_1, \dots, Z_T$ .
- e.g. Want to most likely sequence  $(Z_1, \dots, Z_T)$ . (Use **Viterbi algorithm**.)

KPM Figure 10.4

# Maximum Entropy Markov Model

- A maximum entropy Markov model (MEMM) has structure:



$$p(y_1 \dots, y_5 | x) = \underbrace{p(y_0) \prod_{t=1}^5 p(y_t | y_{t-1}, x)}_{\text{Conditional Transition Model}}$$

- At deployment time, we only observe  $X_1, \dots, X_T$ .
- This is a **conditional model**. (And not a generative model).