

# COMPSCI 688: Probabilistic Graphical Models

## Lecture 2: More Probability and Directed Graphical Models

Dan Sheldon

Manning College of Information and Computer Sciences  
University of Massachusetts Amherst

Partially based on materials by Benjamin M. Marlin (marlin@cs.umass.edu) and Justin Domke (domke@cs.umass.edu)

## Review

## Discrete Distributions

- ▶ Sample space  $\Omega$
- ▶ Atomic probability  $p(\omega)$  for all  $\omega \in \Omega$

$$p(\omega) \geq 0, \quad \sum_{\omega \in \Omega} p(\omega) = 1$$

- ▶ Events  $A \subseteq \Omega$  (only things that have probabilities!)

$$P(A) = \sum_{\omega \in A} p(\omega)$$

- ▶ Random variable  $X : \Omega \rightarrow \text{Val}(X)$  has probability mass function (PMF)

$$p_X(x) = P(X(\omega) = x) = P(X = x)$$

## Events vs Random Variables

- ▶ A random variable  $X$  is a mapping from  $\Omega$  to  $\text{Val}(X)$
- ▶ **But:** for any random variable  $X$ , we can also define the probability distribution with sample space  $\Omega = \text{Val}(X)$  and atomic probabilities  $p_X(x)$ . This is the **distribution** of  $X$ .
- ▶ If we only care about events involving  $X$ , it's easier to just define the distribution of  $X$  without using a different underlying probability space
- ▶ If we care about multiple random variables, we can similarly define their **joint distribution**

## Joint Distributions

5 / 37

## Random Variables and Data Sets

In ML and stats, probability distributions are defined over records described by multiple attributes modeled as random variables. This leads to joint distributions.

Gender	Blood Pressure	Cholesterol	Heart Disease
Male	Med	Low	No
Male	Hi	Hi	Yes
Male	Med	Med	Yes
Male	Med	Hi	No
Female	Med	Low	No
Male	Low	Med	No

6 / 37

## Joint Probability Distributions

- ▶ The *joint distribution* of random variables  $X_1, \dots, X_N$  is a probability distribution over their *canonical sample space*
- ▶ The *canonical sample space*  $\Omega$  of  $X_1, \dots, X_N$  is the Cartesian product of their domains  $\Omega = \text{Val}(X_1) \times \dots \times \text{Val}(X_N)$ .
- ▶ An element of  $\Omega$  is a joint assignment  $(x_1, \dots, x_N)$
- ▶ The joint probability mass function of  $X_1, \dots, X_N$  is

$$p(x_1, \dots, x_N) = P(X_1 = x_1, \dots, X_N = x_N)$$

7 / 37

## Joint Distributions: Heart Disease Example

**Example:** The joint distribution over random variables *Gender*, *BloodPressure*, *Cholesterol* and *HeartDisease* is given by a table like this:

Gender	BloodPressure	Cholesterol	HeartDisease	P
F	L	L	N	0.0127
F	L	L	Y	0.0007
F	L	M	N	0.0098
F	L	M	Y	0.0009
F	L	H	N	0.0087
F	L	H	Y	0.0010
...	...	...	...	...

8 / 37

## Random Vectors

- ▶ It's convenient to use vector-valued random variables  $\mathbf{X} = (X_1, \dots, X_N)$  (or "random vectors") and assignments  $\mathbf{x} = (x_1, \dots, x_N)$ :

$$P(\mathbf{X} = \mathbf{x}) = P(X_1 = x_1, \dots, X_N = x_N)$$

- ▶ The PMF is  $p_{\mathbf{X}}(\mathbf{x})$  or just  $p(\mathbf{x})$
- ▶ This is just notation: it means the same thing as a joint distribution over  $(X_1, \dots, X_N)$
- ▶ **Notation:** use  $\mathbf{X}_{-i}$  and  $\mathbf{x}_{-i}$  for vectors excluding  $X_i$  or  $x_i$

## Rules of Probability

## Marginal Distributions

- ▶ Suppose we have a joint distribution  $P(\mathbf{X} = \mathbf{x}, \mathbf{Y} = \mathbf{y})$ .
- ▶  $P(\mathbf{X} = \mathbf{x})$  is called a *marginal distribution*. How can we find  $P(\mathbf{X} = \mathbf{x})$ ?

$$\begin{aligned} P(\mathbf{X} = \mathbf{x}) &= \sum_{\mathbf{y} \in \text{Val}(\mathbf{Y})} P(\mathbf{X} = \mathbf{x}, \mathbf{Y} = \mathbf{y}) \\ &= \sum_{y_1 \in \text{Val}(Y_1)} \dots \sum_{y_M \in \text{Val}(Y_M)} P(X_1 = x_1, \dots, X_N = x_N, Y_1 = y_1, \dots, Y_M = y_M) \end{aligned}$$

## Marginal Distributions: Heart Disease Example

Given a joint distribution on  $G, BP, C, HD$ , we obtain the marginal probability  $P(G = M, BP = H, C = H)$  as follows:

$$\begin{aligned} P(G = M, BP = H, C = H) &= \sum_{h \in \{Y, N\}} P(G = M, BP = H, C = H, HD = h) \\ &= P(G = M, BP = H, C = H, HD = Y) \\ &\quad + P(G = M, BP = H, C = H, HD = N) \\ &= 0.050 + 0.005 \end{aligned}$$

Gender	BloodPressure	Cholesterol	HeartDisease	P
M	H	H	Y	0.050
M	H	H	N	0.005
M	H	M	Y	0.045
M	H	M	N	0.008
...	...	...	...	...

### Conditional Distributions

- ▶ Joint distributions are useful because we can use them to answer queries like “What is the probability that  $\mathbf{Y} = \mathbf{y}$  given that I observed  $\mathbf{X} = \mathbf{x}$ ?”:

$$P(\mathbf{Y} = \mathbf{y} | \mathbf{X} = \mathbf{x}) = \frac{P(\mathbf{X} = \mathbf{x}, \mathbf{Y} = \mathbf{y})}{P(\mathbf{X} = \mathbf{x})}$$

$$= \frac{P(\mathbf{X} = \mathbf{x}, \mathbf{Y} = \mathbf{y})}{\sum_{\mathbf{y} \in \text{Val}(\mathbf{Y})} P(\mathbf{X} = \mathbf{x}, \mathbf{Y} = \mathbf{y})}$$

- ▶ Write  $p(\mathbf{y} | \mathbf{x})$  to denote the PMF of  $\mathbf{Y}$  given  $\mathbf{X} = \mathbf{x}$

### Conditional Distributions: Heart Disease Example

$$P(HD = Y | G = M, BP = H, C = H) = \frac{P(G = M, BP = H, C = H, HD = Y)}{P(G = M, BP = H, C = H)}$$

$$= \frac{0.050}{0.050 + 0.005} = 0.91$$

Gender	BloodPressure	Cholesterol	HeartDisease	P
M	H	H	Y	0.050
M	H	H	N	0.005
M	H	M	Y	0.045
M	H	M	N	0.008
...	...	...	...	...

### Chain Rule

- ▶ By rearranging the definition of conditional probability, we get the chain rule:

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y} | \mathbf{x}) p(\mathbf{x})$$

- ▶ Applying the chain rule repeatedly to a random vector  $\mathbf{X}$  gives:

$$p(\mathbf{x}) = p(x_N | x_1, \dots, x_{N-1}) p(x_1, \dots, x_{N-1})$$

$$\vdots$$

$$= p(x_N | x_1, \dots, x_{N-1}) p(x_{N-1} | x_1, \dots, x_{N-2}) \dots p(x_3 | x_2, x_1) p(x_2 | x_1) p(x_1)$$

### Chain Rule: Heart Disease Example

We can apply the chain rule using any ordering of the variables:

$$p(g, bp, c, hd) = p(hd | c, bp, g) p(c | bp, g) p(bp | g) p(g)$$

$$p(g, bp, c, hd) = p(g | bp, c, hd) p(bp | c, hd) p(c | hd) p(hd)$$

$$p(g, bp, c, hd) = p(c | hd, g, bp) p(hd | g, bp) p(g | bp) p(bp)$$

## Bayes' Rule

- ▶ By using the definition of conditional probability twice, we obtain one of the most important equations in probability theory, Bayes' Rule:

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{y})} = \frac{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})}{p(\mathbf{y})}$$

- ▶ Bayes' rule lets us compute  $p(\mathbf{x}|\mathbf{y})$  from a joint distribution specified by  $p(\mathbf{x})$  and  $p(\mathbf{y}|\mathbf{x})$ .

## Conditional Independence

## Probabilistic Models

The solution to the problem of exponential-sized joint distributions is the use of **compact** probabilistic models.

- ▶ Bayesian networks achieve compactness by exploiting the chain rule and asserting (conditional) independence relations
- ▶ As a result, Bayesian networks can express high-dimensional distributions as products of simpler factors.

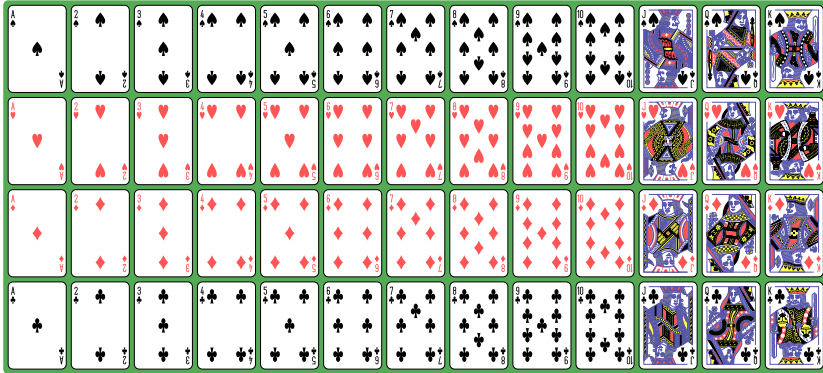
## Marginal Independence

$$\mathbf{X} \perp \mathbf{Y} \iff p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$$

$$\mathbf{X} \perp \mathbf{Y} \iff p(\mathbf{x}|\mathbf{y}) = p(\mathbf{x})$$

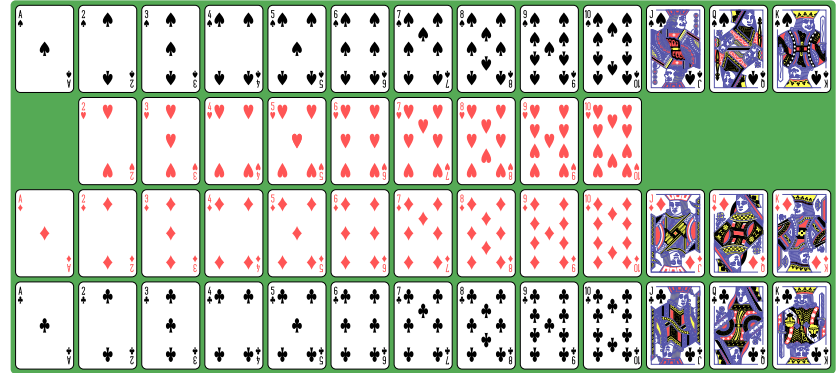
$$\mathbf{X} \perp \mathbf{Y} \iff p(\mathbf{y}|\mathbf{x}) = p(\mathbf{y})$$

### Card Example I



Draw a random card: is value  $\perp$  color? Yes

### Card Example II



What about with this deck? Is value  $\perp$  color? No

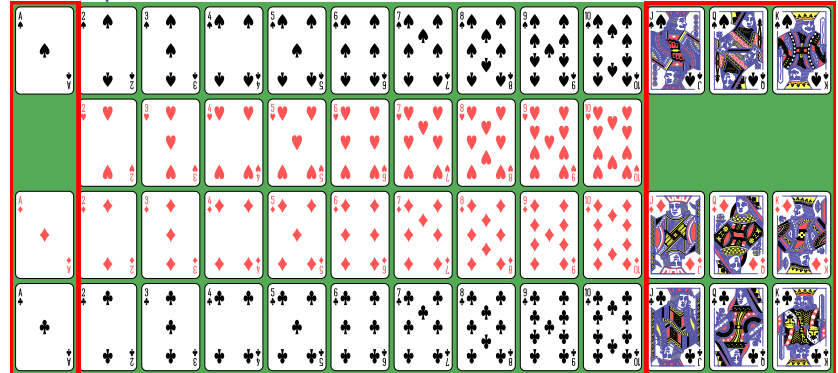
### Conditional Independence

$$X \perp Y | Z \iff p(y, x | z) = p(x | z)p(y | z)$$

$$X \perp Y | Z \iff p(x | y, z) = p(x | z)$$

$$X \perp Y | Z \iff p(y | x, z) = p(y | z)$$

### Card Example III



Is value  $\perp$  color | facecard? Yes

## Bayesian Networks

## Compactness from Independence

Suppose we have a joint distribution  $p(a, b, c)$  and we know that the independence relation  $C \perp A | B$  holds. How can we exploit this fact to simplify  $p(a, b, c)$ ?

$$p(a, b, c) = p(a)p(b|a)p(c|a, b) \quad \text{chain rule}$$

$$= p(a)p(b|a)p(c|b) \quad \text{conditional independence}$$

## Bayesian Networks: Main Idea

- ▶ The main idea of Bayesian networks is conceptually simple:
  1. Order the variables and apply the chain rule
  2. Drop some dependencies, which corresponds to conditional independence assumptions
- ▶ **Example:** variables  $G, C, HD, CP$ , assume: (1)  $G \perp C$ , (2)  $CP \perp G, C | HD$

## Bayesian Networks: Main Idea

- ▶ This idea has several consequences:
  - ▶ The variables can be arranged in a directed acyclic graph (DAG). (Sometimes interpreted causally, but beware.)
  - ▶ The distribution satisfies certain (local and global) conditional independence properties that can be derived from the graph
- ▶ We'll next introduce Bayesian networks formally and start discussing their properties

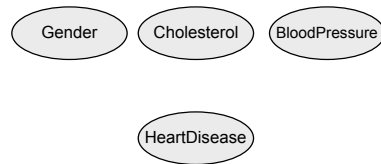
### Bayesian Networks: Nodes

Formally, a Bayesian network consists of a directed acyclic graph (DAG)  $\mathcal{G}$  and a joint distribution  $p(\mathbf{x}) = p(x_1, \dots, x_N)$  for random variables  $X_1, \dots, X_N$

The vertex set  $V$  has one node  $i$  for each random variable  $X_i$

**Warning:** it's also common to use the random variable itself, i.e.,  $X_i$  as the node

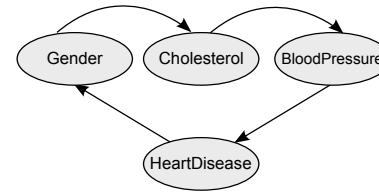
**Example:**



### Bayesian Networks: Edges

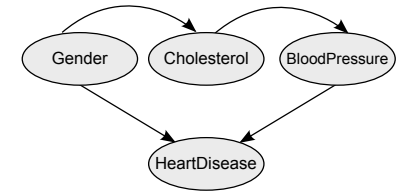
The DAG constraint means that  $\mathcal{G}$  can't contain any directed cycles  $i \rightarrow j \rightarrow \dots \rightarrow i$ .

**Example:**



**Not a valid DAG**  
Directed Cycle

**Example:**



**A valid DAG.**  
No directed cycle

### Bayesian Networks: Parents/Children

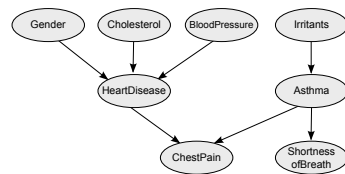
If there is a directed edge  $i \rightarrow j$ :

- ▶  $i$  is a *parent* of  $j$
- ▶  $j$  is a *child* of  $i$
- ▶ (sometimes:  $X_i$  is a parent of  $X_j$ , and so on)

Define

- ▶  $pa(i)$  = set of all parents of  $i$
- ▶  $ch(i)$  = set of all children of  $i$

**Example:**



$pa(CP) = \{HD, A\}$   
 $ch(A) = \{CP, SB\}$

### Bayesian Networks: Descendants/Non-Descendants

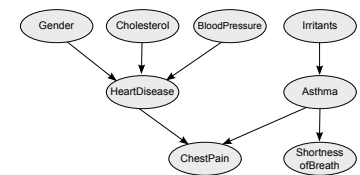
If there is a directed path from  $i$  to  $j$ :

- ▶  $j$  is a *descendant* of  $i$ .
- ▶ Else  $j$  is a *non-descendent* of  $i$ .

Define

- ▶  $de(i)$  = set of all descendants of  $i$
- ▶  $nd(i)$  = set of all non-descendants of  $i$

**Example:**



$de(I) = \{A, SB, CP\}$   
 $nd(BP) = \{G, C, I, A, SB\}$



### Bayesian Networks: Joint Distribution

The joint distribution implied by a Bayesian network is **factorized** into a product of local conditional probability distributions.

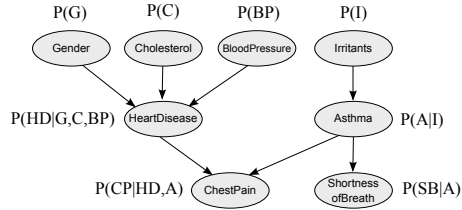


Figure 1: image

The joint distribution is the product of the conditional distributions:

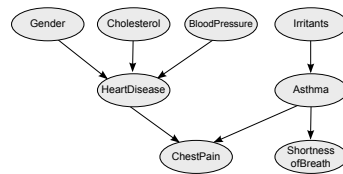
$$p(\mathbf{x}) = \prod_{i=1}^N p(x_i | \mathbf{x}_{pa(i)}).$$

### Bayesian Networks: CPDs and CPTs

- ▶ The individual factors  $p(x_i | \mathbf{x}_{pa(i)})$  in a Bayesian network are referred to as conditional probability distributions or CPDs.
- ▶ The CPD for node  $i$  must specify the probability that  $X_i$  takes any value  $x_i$  in its domain when conditioned on each joint assignment  $\mathbf{x}_{pa(i)}$  of its parents
- ▶ For discrete random variables, we can represent the CPD of each node using a look-up table called a conditional probability table or CPT.

### Bayesian Networks: CPT Example

hd	g	bp	ch	$p(hd g, bp, ch)$
No	M	Low	Low	0.95
Yes	M	Low	Low	0.05
No	F	Low	Low	0.99
Yes	F	Low	Low	0.01
⋮				



### Bayesian Networks: Storage Complexity

- ▶ What is the minimum amount of space needed to store the probability distribution for a single discrete random variable that takes  $V$  values?  $V - 1$
- ▶ How much space does it take to store the CPT for a binary-valued variable with  $D$  binary-valued parents?  $2^D$
- ▶ Suppose there are  $D$  binary variables connected in a chain  $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_D$ . What is the total storage cost?  $1 + 2(D - 1)$   
How large is the full joint?  $2^D - 1$

## Next Time

Next time, we'll discuss factorization and conditional independence in Bayesian networks.