# 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

---

# Review

---

## Discrete Distributions

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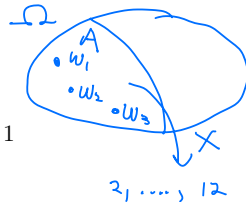- Atomic probability $p(\omega)$ for all $\omega \in \Omega$

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

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

$$P(A) = \sum_{\omega \in A} p(\omega) \quad \subset \quad p(\omega_1) + p(\omega_2) + p(\omega_3)$$

- Random variable $X : \Omega \to \mathrm{Val}(X)$ has probability mass function (PMF)

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

$p(x)$

---

## Events vs Random Variables

- A random variable $X$ is a a mapping from $\Omega$ to $\mathrm{Val}(X)$

- **But**: for any random variable $X$, we can also define the probability distribution with sample space $\Omega = \mathrm{Val}(X)$ 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

---

## 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.

*Attributes/variables*

$X_1$ *Gender* | $X_2$ *Blood Pressure* | $X_3$ *Cholesterol* | $X_4$ *Heart Disease*

| 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 |

*record*

---

## Joint Probability Distributions

▶ The *joint distribution* of random variables $X_1, \ldots, X_N$ is a probability distribution over their *canonical sample space*

▶ The *canonical sample space* $\Omega$ of $X_1, \ldots, X_N$ is the Cartesian product of their domains $\Omega = \mathrm{Val}(X_1) \times \ldots \times \mathrm{Val}(X_N)$.

▶ An element of $\Omega$ is a joint assignment $(x_1, \ldots, x_N)$

▶ The joint probability mass function of $X_1, \ldots, X_N$ is

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

$$P(X = a, Y = b) = P([X_1 = a] \cap [X_2 = b]) \quad \text{``and''}$$

---

## 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 |
| ... | ... | ... | ... | ... |

*all possible records*

$\Omega$

*exponential size in # variables*

## Random Vectors

$p(X_1, ..., X_N)$   $p(x)$   $x = (x_1, ..., x_N)$

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

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

▶ 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 $(X_1, ... , X_N)$

▶ **Notation**: use $\mathbf{X}_{-i}$ and $\mathbf{x}_{-i}$ for vectors excluding $X_i$ or $x_i$

$$\left(X_1, ..., X_{i-1}, X_{i+1}, ..., X_N\right)$$

---
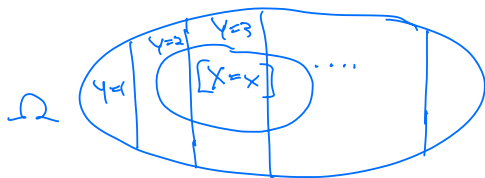
# Rules of Probability

---

## Marginal Distributions

$$(X_1, ..., X_N, Y_1, ..., Y_M)$$

▶ 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 $P(\mathbf{X} = \mathbf{x})$?

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

|  | $Y$ 1 | 2 | $P(x = x)$ |
|---|---|---|---|
| $X$ 1 | .2 | .3 | .5 |
| 2 | .2 | .3 | .5 |
|  | .4 | .6 |  |

---

## 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:

$$P(G = M, BP = H, C = H) = 0.050 + 0.005 = 0.055$$

| 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 $\mathbf{X} = \mathbf{x}$?":

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

*joint*
*marginal*

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

*free* *fixed at observed value*

$p_{X,Y}(x,y) = P(X=x, Y=y)$

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

$p_{Y|X}(y|x)$   $P(Y=y|X=1) \; \forall y$   $p_{X|Y}(x|y)$ $p_{Y|X}(y|x)$

---

## 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

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

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

$P(X=x, Y=y) = P(X=x)P(Y=y|X=x)$

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

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

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

$\vdots = p(x_N|x_1,...,x_{N-1})p(x_{N-1}|x_1,...,x_{N-2})p(x_1,...,x_{N-2})$

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

$$= \prod_{i=1}^{N} p(x_i | x_1, \cdots, x_{i-1})$$

---

## Chain Rule: Heart Disease Example

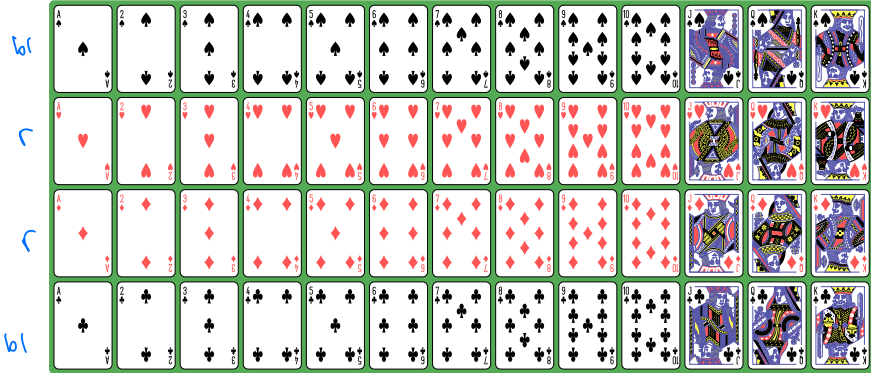We can apply the chain rule using any ordering of the variables:

1 2 3 4
$$p(g, bp, c, hd) = p(hd|c, bp, g)p(c|bp, g)p(bp|g)p(g)$$
4 3 2 1
$$p(g, bp, c, hd) = p(g|bp, c, hd)p(bp|c, hd)p(c|hd)p(hd)$$
2 1 4 3
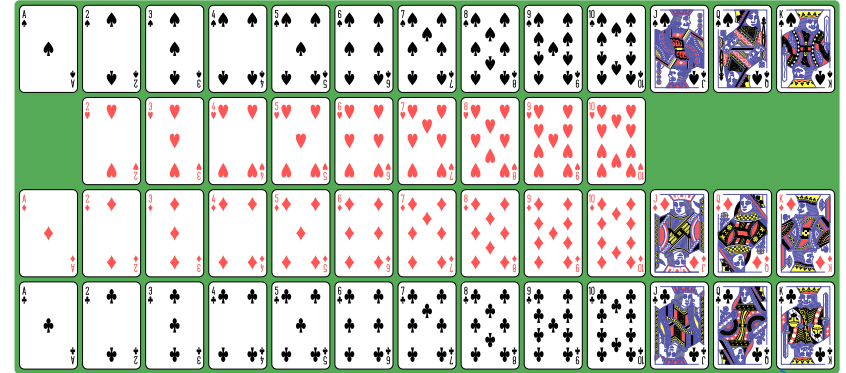$$p(g, bp, c, hd) = p(c|hd, g, bp)p(hd|g, bp)p(g|bp)p(bp)$$

## Bayes' Rule

hidden | observed/evidence | "world model" | measurement model

$p(y|x)$

$x \rightarrow y$

$p(x|y)$

$p(x)$ $\qquad$ $p(y|x)$

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

$p(\text{hidden}|\text{observed})$

posterior $\longrightarrow$ $p(\mathbf{x}|\mathbf{y}) = \dfrac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{y})} = \dfrac{p(\mathbf{y}|\mathbf{x})p(\mathbf{x})}{p(\mathbf{y})} = \dfrac{p(y|x)p(x)}{\sum_x p(y|x)p(x)}$

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

$p(y) = \sum_x p(x, y) = \sum_x p(y|x)p(x)$

## 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

$P(X=x, Y=y) = P(X=x)P(Y=y)$

"$X$ is indept. of $Y$"

$X=1$ $\quad$ $X=2$

$Y=1$

$Y=2$

$\Omega$

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

$\searrow p(x)p(y|x)$

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

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

## Slide 21/37

### Card Example I

$$P(value = 2 \mid col = r) = \frac{1}{13} \qquad P(value = 2) = \frac{1}{13}$$

bl

r

r

bl

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

## Slide 22/37

### Card Example II

$$P(val = J \mid col = r) \neq P(val = J)$$

What about with this deck? Is value $\perp$ color?

$$P(col = r \mid val \in \{2, 4, 6, 8\})$$

## Slide 23/37

### Conditional Independence

$$P(X = x, Y = y, Z = z)$$

X is conditionally independent of Y given Z if

$$\mathbf{X} \perp \mathbf{Y} \mid \mathbf{Z} \iff p(\mathbf{y}, \mathbf{x} \mid \mathbf{z}) = p(\mathbf{x} \mid \mathbf{z}) p(\mathbf{y} \mid \mathbf{z}) \qquad \forall x, y, z$$

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

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

"Given z, knowing y tells you nothing about the dist. of x"

## Slide 24/37

### Card Example III

facecard = yes: value $\perp$ color
facecard = no: value $\perp$ color

Is value $\perp$ color | facecard?

$$P(val = J \mid col = bl, face = y) = \frac{1}{4}$$
$$P(val = J \mid face = y) = \frac{1}{4}$$

# 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 $p(a, b, c)$?
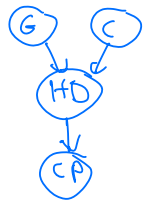
chain rule $\quad p(a,b,c) = p(a)\, p(b|a)\, p(c|a,b)$ $\quad$ $a \rightarrow b \rightarrow c$

$\parallel$

conditional indep. $\quad p(a,b,c) = p(a)\, p(b|a)\, p(c|b)$ $\quad$ $a \rightarrow b \rightarrow c$

---

## Bayesian Networks: Main Idea

$G \rightarrow C \rightarrow HD \rightarrow CP$

▶ 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

$CP \perp G | HD$

▶ **Example**: variables $G, C, HD, CP$, assume: (1) $G \perp C$, (2) $CP \perp (G, C) HD$

1. $p(g, c, hd, cp) = p(g)\, p(c|g)\, p(hd|g,c)\, p(cp|g, c, hd)$

$\parallel$ $\qquad\qquad\qquad\qquad$ $\parallel$

2. $p(g, c, hd, cp) = p(g)\, p(c)\, p(hd|g,c)\, p(cp|hd)$
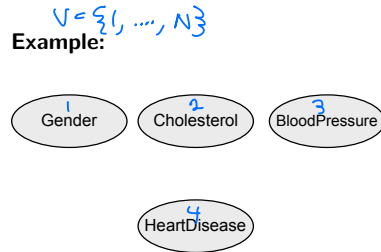
---

## 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

$\mathcal{G}, \; p(x)$

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 $X_1, \dots, X_N$

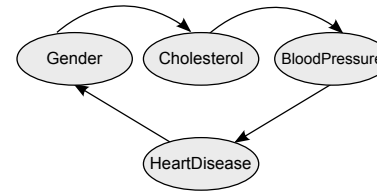The vertex set $V$ has one node $i$ for each random variable $X_i$

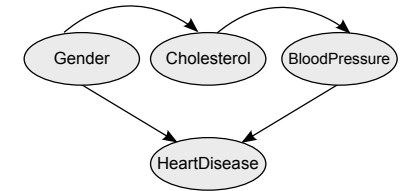**Warning**: it's also common to use the random variable itself, i.e., $X_i$ as the node

**Example:**

$V = \{1, \dots, N\}$

## Bayesian Networks: Edges

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

**Example:**



**Example:**



**Not a valid DAG**
Directed Cycle

**A valid DAG.**
No directed cycle

## Bayesian Networks: Parents/Children
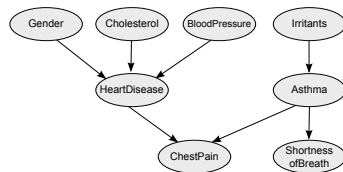
If there is a directed edge $i \to j$:
- $i$ is a *parent* of $j$
- $j$ is a *child* of $i$
- (sometimes: $X_i$ is a parent of $X_j$, and so on)

Define
- $\mathrm{pa}(i)$ = set of all parents of $i$
- $\mathrm{ch}(i)$ = set of all children of $i$

**Example:**



$$\mathrm{pa}(CP) = \{HD, A\}$$
$$\mathrm{ch}(A) = \{CP, SB\}$$
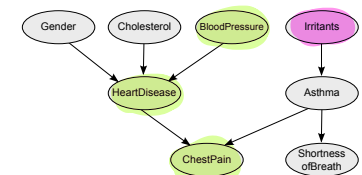
## Bayesian Networks: Descendants/Non-Descendants

If there is a directed path from $i$ to $j$:
- $j$ is a *descendant* of $i$.
- Else $j$ is a *non-descendant* of $i$.

Define
- $\mathrm{de}(i)$ = set of all descendants of $i$
- $\mathrm{nd}(i)$ = set of all non-descendants of $i$

**Example:**



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

# Bayesian Networks: Joint Distribution

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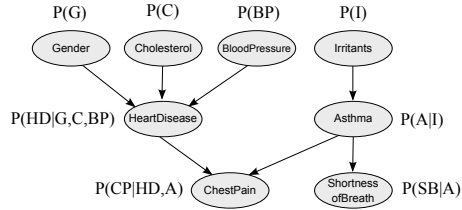


Figure 1: image

$P(g,c,bp,...,sb) = P(g)P(c)P(bp)P(i)P(hd|g,c,bp)P(a|i)P(cp|hd,a)P(sb|a)$

The joint distribution is the product of the conditional distributions:
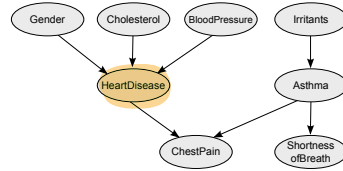
$p(\mathbf{x}) = \prod_{i=1}^{N} p(x_i \mid \mathbf{x}_{\mathbf{pa}(i)})$.

# Bayesian Networks: CPDs and CPTs

▶ The individual factors $p(x_i \mid \mathbf{x}_{\mathbf{pa}(i)})$ 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 $\mathbf{x}_{\mathbf{pa}(i)}$ 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

$p(hd,g \mid bp)$

$2^D$

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

exponential in (#parents + 1)

# 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? $p(a \mid b_1,...,b_D)$ $2^D \cdot (V-1) = 2^D$ $V=2$

▶ Suppose there are $D$ binary variables connected in a chain $X_1 \to X_2 \to ... \to X_D$. What is the total storage cost? $2\cdot(D-1)+1 = 2D-1$ How large is the full joint? $2^D$

# Next Time

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