# CSC 4020 Fundamentals of Machine Learning:
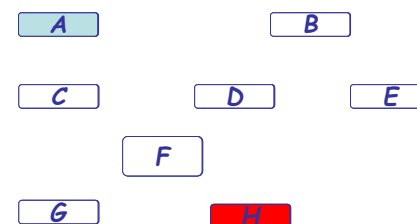## Bayesian Networks

Baoyuan Wu

April 14

# Representing Multivariate Distribution

- **Representation: what is the joint probability dist. on multiple variables?**

$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8,)$$

  - **How many state configurations in total? --- $2^8$**
  - **Are they all needed to be represented?**
  - **Do we get any scientific/medical insight?**

- **Factored representation: the chain-rule**

$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$$
$$= P(X_1)P(X_2 \mid X_1)P(X_3 \mid X_1, X_2)P(X_4 \mid X_1, X_2, X_3)P(X_5 \mid X_1, X_2, X_3, X_4)P(X_6 \mid X_1, X_2, X_3, X_4, X_5)$$
$$P(X_7 \mid X_1, X_2, X_3, X_4, X_5, X_6)P(X_8 \mid X_1, X_2, X_3, X_4, X_5, X_6, X_7)$$

  - **This factorization is true for any distribution and any variable ordering**
  - **Do we save any parameterization cost?**

- **If $X_i$'s are independent: ($P(X_i|\cdot)= P(X_i)$)**

$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$$
$$= P(X_1)P(X_2)P(X_3)P(X_4)P(X_5)P(X_6)P(X_7)P(X_8) = \prod P(X_i)$$

- **What do we gain?**
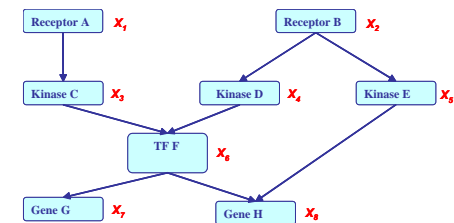- **What do we lose?**

A    B

C    D    E

F

G    H

# Two types of GMs

- **Directed edges give causality relationships (Bayesian Network or Directed Graphical Model):**

$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$

$= P(X_1) P(X_2) P(X_3 / X_1) P(X_4 / X_2) P(X_5 / X_2)$
$P(X_6 / X_3, X_4) P(X_7 / X_6) P(X_8 / X_5, X_6)$

| Receptor A | $X_1$ | | Receptor B | $X_2$ |

Kinase C $X_3$  Kinase D $X_4$  Kinase E $X_5$

TF F $X_6$

Gene G $X_7$  Gene H $X_8$
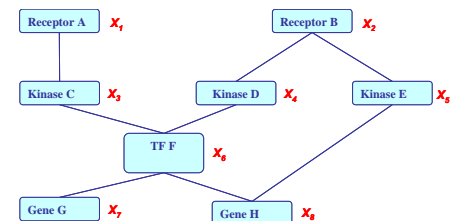
- **Undirected edges simply give correlations between variables (Markov Random Field or Undirected Graphical model):**
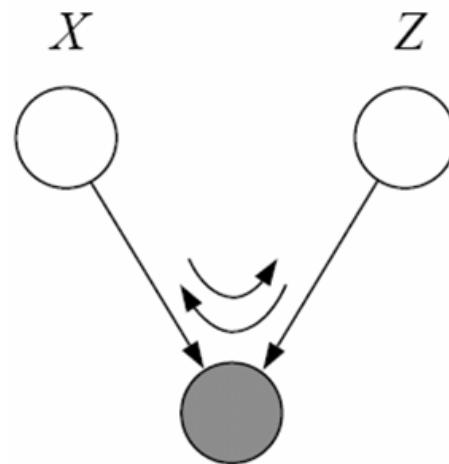
$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$

$= 1/Z \exp\{E(X_1)+E(X_2)+E(X_3, X_1)+E(X_4, X_2)+E(X_5, X_2)$
$+ E(X_6, X_3, X_4)+E(X_7, X_6)+E(X_8, X_5, X_6)\}$

Receptor A $X_1$   Receptor B $X_2$

Kinase C $X_3$  Kinase D $X_4$  Kinase E $X_5$

TF F $X_6$

Gene G $X_7$  Gene H $X_8$

# Representation of directed GM

# Example: The Dishonest Casino

A casino has two dice:
- Fair dice
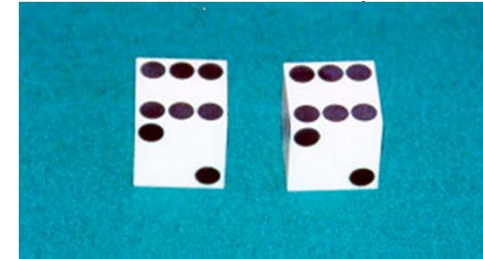  P(1) = P(2) = P(3) = P(5) = P(6) = 1/6
- Loaded dice
  P(1) = P(2) = P(3) = P(5) = 1/10
  P(6) = 1/2

Casino player switches back-&-forth between fair and loaded dice once every 20 turns

**<u>Game:</u>**
1. You bet $1
2. You roll (always with a fair die)
3. Casino player rolls (maybe with fair die, maybe with loaded die)
4. Highest number wins $2
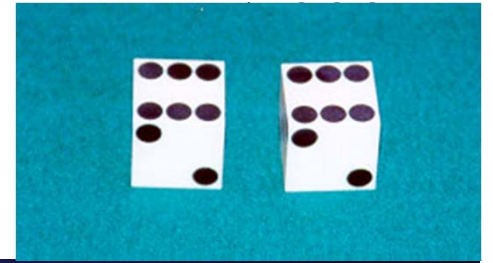
# Puzzles regarding the dishonest casino

**GIVEN:** A sequence of rolls by the casino player

124552646214614613613666166466163661636616361651561511514612356 2344

**QUESTION**

- How likely is this sequence, given our model of how the casino works?
  - This is the **EVALUATION** problem

- What portion of the sequence was generated with the fair die, and what portion with the loaded die?
  - This is the **DECODING** question

- How "loaded" is the loaded die? How "fair" is the fair die? How often does the casino player change from fair to loaded, and back?
  - This is the **LEARNING** question

# Knowledge Engineering

- **Picking variables**
  - Observed
  - Hidden

- **Picking structure**
  - CAUSAL
  - Generative
  - Coupling

- **Picking Probabilities**
  - Zero probabilities
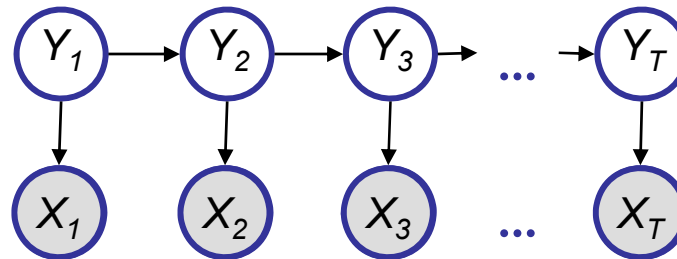  - Orders of magnitudes
  - Relative values

# Hidden Markov Model

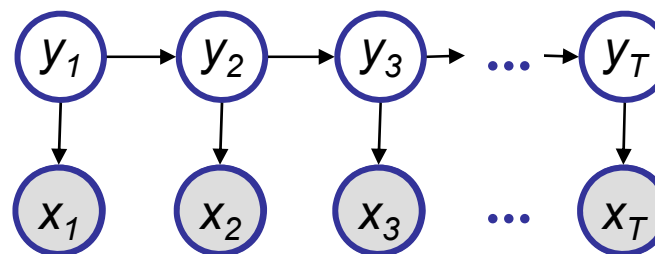**The underlying source:**

Speech signal
genome function

dice

**The sequence:**

Phonemes
DNA sequence
sequence of rolls

# Probability of a parse

- Given a sequence $\mathbf{x} = x_1\ldots\ldots x_T$
  and a parse $\mathbf{y} = y_1, \ldots\ldots, y_T$,
- To find how likely is the parse:
  (given our HMM and the sequence)



$$p(\mathbf{x}, \mathbf{y}) \quad = p(x_1\ldots\ldots x_T, y_1, \ldots\ldots, y_T) \qquad \text{(Joint probability)}$$

$$= p(y_1)\, p(x_1 \mid y_1)\, p(y_2 \mid y_1)\, p(x_2 \mid y_2) \ldots p(y_T \mid y_{T-1})\, p(x_T \mid y_T)$$

$$= p(y_1)\, \mathrm{P}(y_2 \mid y_1) \ldots p(y_T \mid y_{T-1}) \times p(x_1 \mid y_1)\, p(x_2 \mid y_2) \ldots p(x_T \mid y_T)$$

$$= p(y_1, \ldots\ldots, y_T)\, p(x_1\ldots\ldots x_T \mid y_1, \ldots\ldots, y_T)$$

- Marginal probability: $\quad p(\mathbf{x}) = \sum_{\mathbf{y}} p(\mathbf{x},\mathbf{y}) = \sum_{y_1} \sum_{y_2} \cdots \sum_{y_N} \pi_{y_1} \prod_{t=2}^{T} a_{y_{t-1},y_t} \prod_{t=1}^{T} p(x_t \mid y_t)$
- Posterior probability: $\quad p(\mathbf{y} \mid \mathbf{x}) = p(\mathbf{x},\mathbf{y}) / p(\mathbf{x})$

- We will learn how to do this explicitly (polynomial time)

# Bayesian Network:

- A BN is a directed graph whose nodes represent the random variables and whose edges represent direct influence of one variable on another.

- It is a data structure that provides the skeleton for representing **a joint distribution** compactly in a **factorized** way;

- It offers a compact representation for **a set of conditional independence assumptions** about a distribution;

- We can view the graph as encoding a generative sampling process executed by nature, where the value for each variable is selected by nature using a distribution that depends only on its parents. In other words, each variable is a stochastic function of its parents.

# Bayesian Network: Factorization Theorem

- **Theorem:**

  Given a DAG, The most general form of the probability distribution that is consistent with the graph factors according to "node given its parents":

  $$P(\mathbf{X}) = \prod_{i=1:d} P(X_i \mid \mathbf{X}_{\pi_i})$$
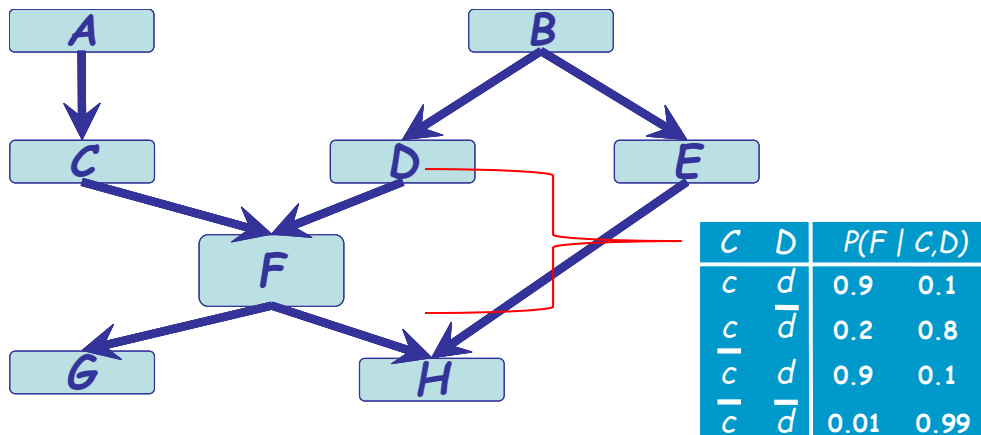
  where $\mathbf{X}_{\pi_i}$ is the set of parents of $X_i$, $d$ is the number of nodes (variables) in the graph.



$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$$

$$= P(X_1)\, P(X_2)\, P(X_3 \mid X_1)\, P(X_4 \mid X_2)\, P(X_5 \mid X_2)$$
$$P(X_6 \mid X_3, X_4)\, P(X_7 \mid X_6)\, P(X_8 \mid X_5, X_6)$$

# Specification of a directed GM

- There are two components to any GM:
  - the *qualitative* specification
  - the *quantitative* specification



| C | D | P(F \| C,D) | |
|---|---|---|---|
| c | d | 0.9 | 0.1 |
| c | d̄ | 0.2 | 0.8 |
| c̄ | d | 0.9 | 0.1 |
| c̄ | d̄ | 0.01 | 0.99 |

# Qualitative Specification

- Where does the qualitative specification come from?

  - Prior knowledge of causal relationships

  - Prior knowledge of modular relationships

  - Assessment from experts

  - Learning from data

  - We simply link a certain architecture (e.g. a layered graph)

  - …

# Local Structures & Independencies

- Common parent
  - Fixing B decouples A and C

    "given the level of gene B, the levels of A and C are independent"
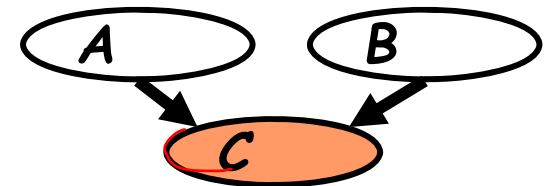
- Cascade
  - Knowing B decouples A and C

    "given the level of gene B, the level gene A provides no extra prediction value for the level of gene C"

- V-structure
  - Knowing C couples A and B

    because A can "explain away" B w.r.t. C

    "If A correlates to C, then chance for B to also correlate to B will decrease"

- The language is compact, the concepts are rich!

# Common parent



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

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

$$a \not\!\perp\!\!\!\perp b \mid \emptyset$$
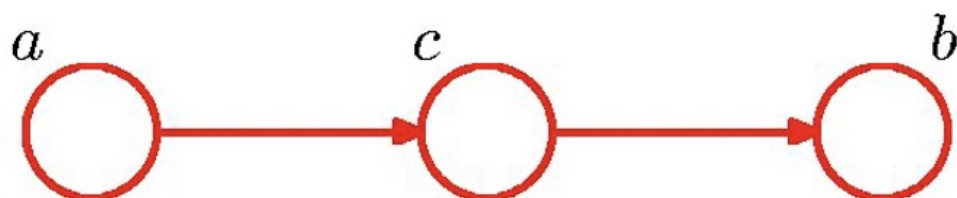
# Common parent



$$p(a, b|c) = \frac{p(a, b, c)}{p(c)}$$
$$= p(a|c)p(b|c)$$
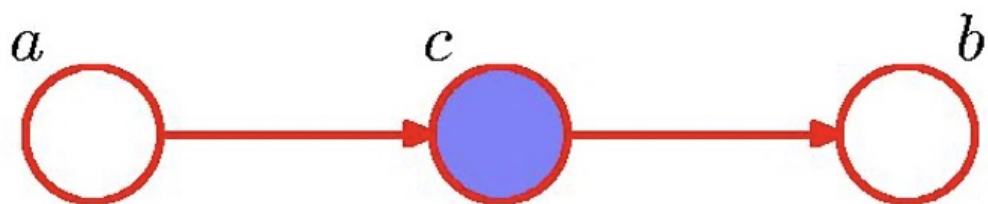
$$a \perp\!\!\!\perp b \mid c$$

# Chain



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

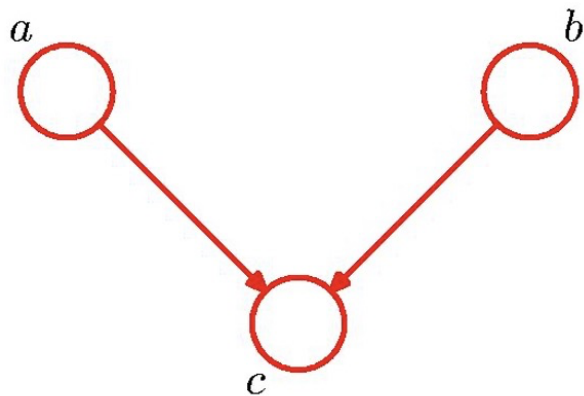$$p(a, b) = p(a) \sum_c p(c|a)p(b|c) = p(a)p(b|a)$$

$$a \not\!\perp\!\!\!\perp b \mid \emptyset$$

# Chain



$$p(a, b|c) = \frac{p(a, b, c)}{p(c)}$$

$$= \frac{p(a)p(c|a)p(b|c)}{p(c)}$$

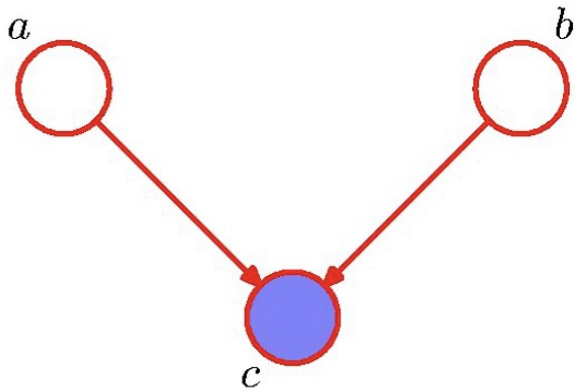$$= p(a|c)p(b|c)$$

$$a \perp\!\!\!\perp b \mid c$$

# V-structure



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

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

$$a \perp\!\!\!\perp b \mid \emptyset$$

# V-structure



$$p(a, b | c) = \frac{p(a, b, c)}{p(c)}$$

$$= \frac{p(a)p(b)p(c|a, b)}{p(c)}$$
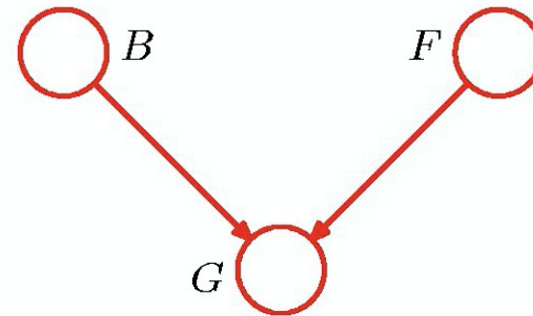
$$a \not\perp b \mid c$$

# One example: "Am I out of fuel?"

$$p(G = 1 | B = 1, F = 1) = 0.8$$
$$p(G = 1 | B = 1, F = 0) = 0.2$$
$$p(G = 1 | B = 0, F = 1) = 0.2$$
$$p(G = 1 | B = 0, F = 0) = 0.1$$
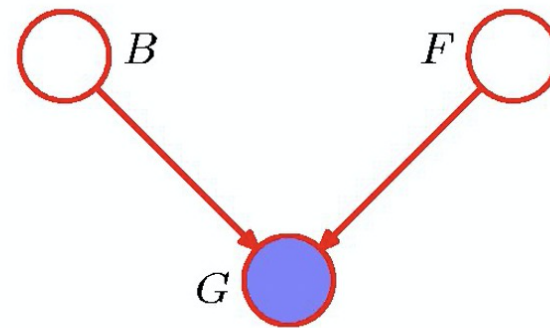


$$p(B = 1) = 0.9$$
$$p(F = 1) = 0.9$$

and hence

$$p(F = 0) = 0.1$$

$B$ = Battery (0=flat, 1=fully charged)
$F$ = Fuel Tank (0=empty, 1=full)
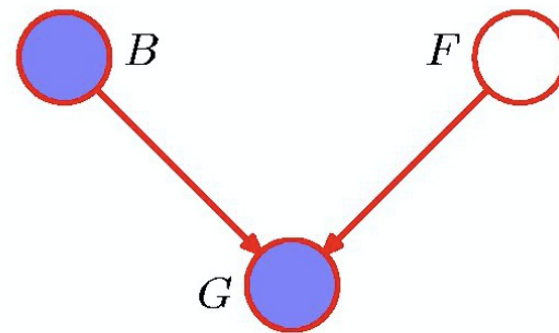$G$ = Fuel Gauge Reading
   (0=empty, 1=full)

# One example: "Am I out of fuel?"



$$p(F = 0|G = 0) = \frac{p(G = 0|F = 0)p(F = 0)}{p(G = 0)}$$

$$\simeq 0.257$$

Probability of an empty tank increased by observing $G = 0$.

# One example: "Am I out of fuel?"



$$p(F = 0 | G = 0, B = 0) = \frac{p(G = 0 | B = 0, F = 0)p(F = 0)}{\sum_{F \in \{0,1\}} p(G = 0 | B = 0, F)p(F)}$$

$$\simeq 0.111$$

Probability of an empty tank reduced by observing $B = 0$.
This referred to as "explaining away".