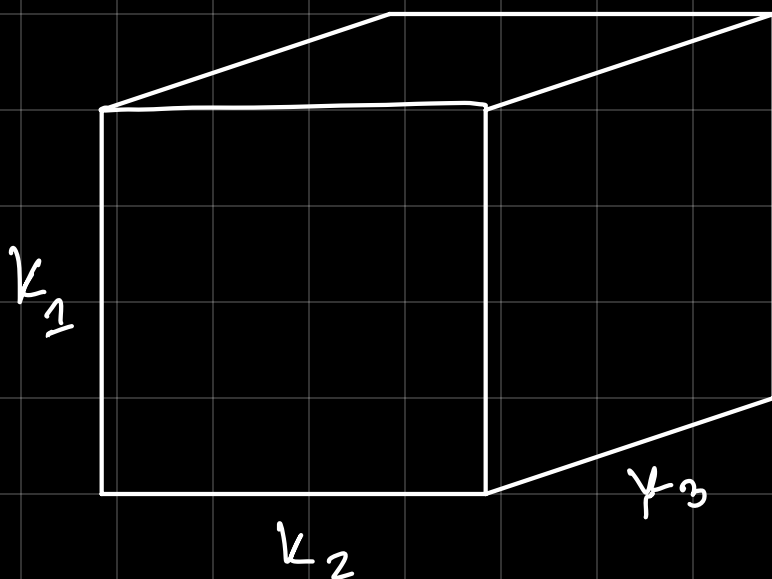# probabilistic graphical models (PGMs)

- model $p(X_1 \ldots X_n)$

- e.g., discrete $X_i \in \{1 \ldots k_i\}$, $i = 1 \ldots n$

- notation:
$$p(x_1 \ldots x_n) \equiv P(X_1 = x_1, \ldots, X_n = x_n)$$

- $p(x_1 \ldots x_n)$ stored in a ~~probability~~ table

- e.g., $n = 3$



- $$\sum_{x_1=1}^{k_1} \sum_{x_2=1}^{k_2} \sum_{x_3=1}^{k_3} P(X_1 = x_1, X_2 = x_2, X_3 = x_3) = 1$$

- $k_1 \times k_2 \times k_3$ entries

- <u>exponential</u> in $n$

- e.g. $k^n$ if $k_1 = k_2 \ldots = k_n = k$

- So what? Consider a conditional distribution (e.g., a posterior)...

$$P(x_1 \cdots x_{n-1} \mid X_n = x_n) = \frac{P(x_1 \cdots x_n)}{\underbrace{\sum_{x_1=1}^{k} \cdots \sum_{x_{n-1}=1}^{k} P(x_1 \cdots x_{n-1}, x_n)}_{k^{n-1} \text{ summands}}}$$
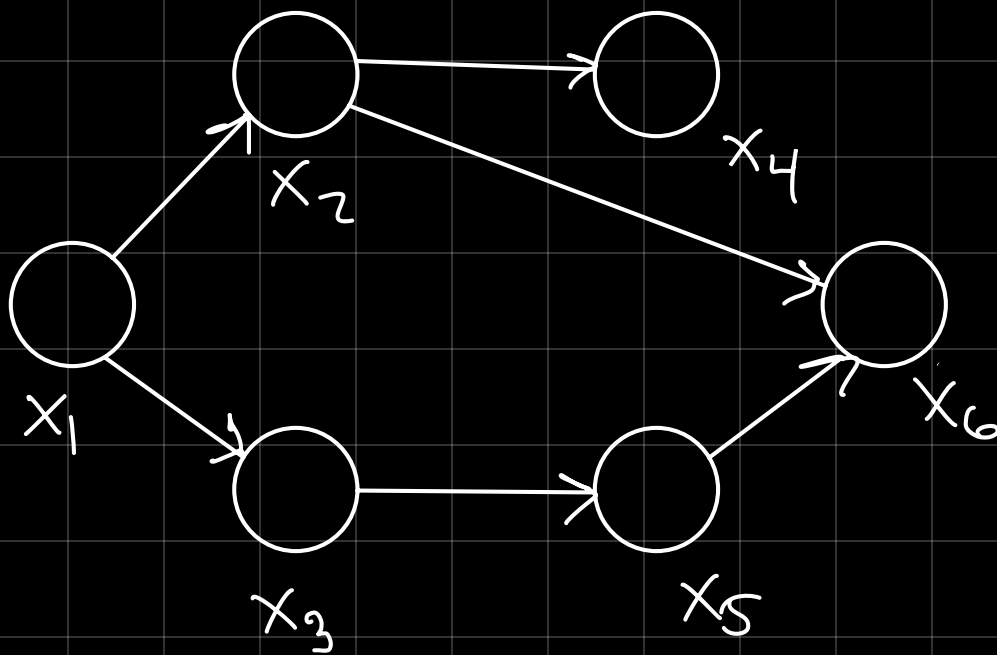
- This becomes <u>intractable</u> quickly

- e.g., a single binary trial $x_i \in \{0,1\}$ in all $n = 140$ search cells

  - $2^{140}$ cells of $P(x_1 \cdots x_n)$

  - modern processor can compute $\sim 10^9$ flops/sec

  - so summing $2^{140}$ cells would take $\sim 10^{33}$ seconds $\approx 10^{26}$ years !! (age of universe: $10^{10}$ years)

# Directed graphical models (DGMs)

- The reason there were so many cells in $P(x_1 \ldots x_n)$ is because we did not account for any <u>conditional independence</u> structure.

  - e.g. say $X_1 \perp\!\!\!\perp X_2 \perp\!\!\!\perp \ldots \perp\!\!\!\perp X_n$

    $\hookrightarrow P(x_1 \ldots x_n) = \prod_{i=1}^{n} P(x_i)$

    we only need to store $\sum_{i=1}^{n} K_i$ cells
    (as opposed to $\prod_{i=1}^{n} K_i$)

- DGMs provide a formal language to describe the set of conditional independencies in a joint distribution.

- We are already used to defining models (i.e. joint distributions) in terms of <u>forward sampling algorithms</u>

  - e.g.) $X_1 \sim P(x_1)$

    $X_2 \sim P(X_2 | X_1)$

    $X_3 \sim P(X_3 | X_1)$

    $X_4 \sim P(X_4 | X_2)$

    $X_5 \sim P(X_5 | X_3)$

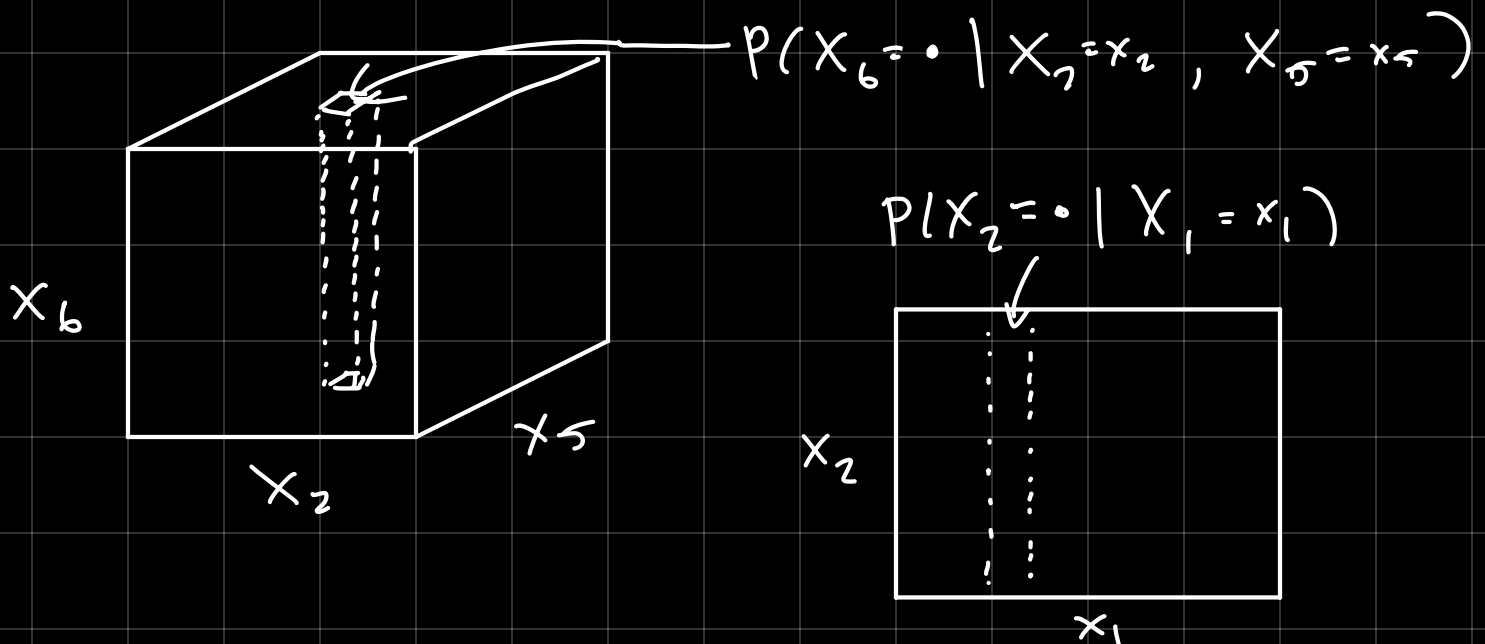    $X_6 \sim P(X_6 | X_5, X_2)$

- We can represent this algorithm graphically:



- Both the algorithm and the graph represent a joint dist. that <u>factorizes</u> as

$$P(x_1, \ldots x_6) = P(x_1) P(x_2 | x_1) \ldots P(x_6 | x_2, x_5)$$

- Each <u>factor</u> e.g. $P(x_6 | x_2, x_5)$ comes from a <u>local probability table</u> (LPT)

$$P(X_6 = \bullet | X_2 = x_2, X_5 = x_5)$$



$$P(X_2 = \bullet | X_1 = x_1)$$

- DGM is a directed acyclic graph (DAG)

    - nodes ≡ random variables
    - edges ≡ "parenthood"

        - $\pi_i \triangleq$ parents$(x_i)$
        - e.g. $\pi_6 = \{2, 5\}$

- It is defined WRT a specific topological ordering of the variables. (in this case $x_1, x_2, \ldots, x_6$)

- The joint distribution is:
$$P(x_1, \ldots x_n) = \prod_{i=1}^{n} P(x_i \mid x_{\pi_i})$$

- Each $P(x_i \mid x_{\pi_i})$ is an LPT

- Total # of cells : $\sum_{i=1}^{n} k_i \prod_{j \in \pi_i} k_j$

- e.g. $k_1 = \ldots = k_n = 2$

    $\rightarrow 2^n$ vs. $\sum_{i=1}^{n} 2^{|\pi_i|}$

    (exponential in $n$ vs. exponential in $|\pi_i|$)

# Graph separation

- A subset of the conditional indep. relations are encoded directly by graph seperation

- Define non-parent ancestors:

$$V_i \triangleq \text{ancestors}(x_i) \setminus \text{parents}(x_i)$$

- Graph seperation encodes:

$$X_i \perp\!\!\!\perp X_{V_i} \mid X_{\pi_i}$$

- e.g., show $X_5 \perp\!\!\!\perp X_1 \mid X_3$

$$P(x_5 \mid x_3, x_1) = \cfrac{P(x_5, x_3, x_1)}{\sum_{x_5} P(x_5, x_3, x_1)}$$

$$P(x_5, x_3, x_1) = P(x_5 \mid x_3) P(x_3 \mid x_1) P(x_1)$$

$$\searrow = \frac{P(x_5 \mid x_3) P(x_3 \mid x_1) P(x_1)}{\sum_{x_5} P(x_5 \mid x_3) P(x_3 \mid x_1) P(x_1)}$$

$$= \frac{P(x_5 \mid x_3) \cancel{P(x_3 \mid x_1)} \cancel{P(x_1)}}{\cancel{P(x_3 \mid x_1)} \cancel{P(x_1)} \underbrace{\sum_{x_5} P(x_5 \mid x_3)}_{"1}}$$

$$= P(x_5 \mid x_3) \checkmark$$

- Using graph seperation, we know:

  - $X_4 \perp\!\!\!\perp X_2 \mid X_1$
  - $X_5 \perp\!\!\!\perp X_1 \mid X_3$
  - $X_6 \perp\!\!\!\perp X_5, X_3 \mid X_1$

- Are these the only conditional indepencies among $X_1 \dots X_6$ implied by the graph? No.

- Why? It only reflects graph seperation for a <u>single topological ordering</u> $X_1 \dots X_6$

- e.g. $P(X_1 \dots X_6) = P(X_6) P(X_4 \mid X_6) \dots$

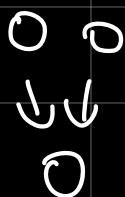- Nevertheless, a single DGM implies all conditional independencies via "<u>d-seperation</u>".

# D-separation

- "directional" separation  - Pearl (1988)

- 3 simple DAGs

  ① chain $\circ \to \circ \to \circ$

  ② tree $\circ \overset{\to \circ}{\underset{\to \circ}{}}$

  ③ V-structure
  $$\circ \quad \circ$$
  $$\downarrow \downarrow$$
  $$\circ$$

- shading $\equiv$ conditioning

  e.g. $\underset{X}{\circ} \to \underset{Y}{\bullet} \to \underset{Z}{\circ} \equiv P(X, Z \mid Y)$

## ① Chain

$$\underset{X}{\circ} \to \underset{Y}{\circ} \to \underset{Z}{\circ} \qquad Z \perp\!\!\!\perp X \mid Y$$

(example of a Markov assumption)
future $\perp\!\!\!\perp$ past | present

② Tree

$P(x, y, z)$

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



$X \perp\!\!\!\perp Z$ ?

$P(x|z) = \dfrac{\sum_y p(y) \, p(x|y) \, p(z|y)}{\sum_x \sum_y p(y) \, p(x|y) \, p(z|y)}$

no.

$X \perp\!\!\!\perp Z \mid Y$ ?

$P(x|z,y) = \dfrac{p(y) \, p(x|y) \, p(z|y)}{\sum_x p(y) \, p(x|y) \, p(z|y)}$

$= \dfrac{\cancel{p(y)} \, p(x|y) \, \cancel{p(z|y)}}{\cancel{p(y)} \, \cancel{p(z|y)}}$
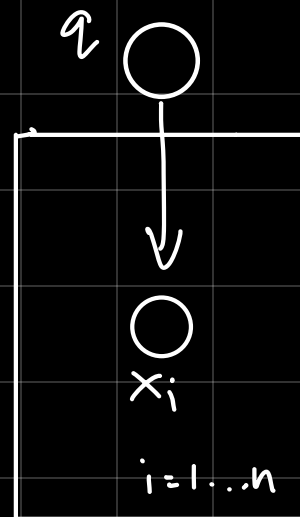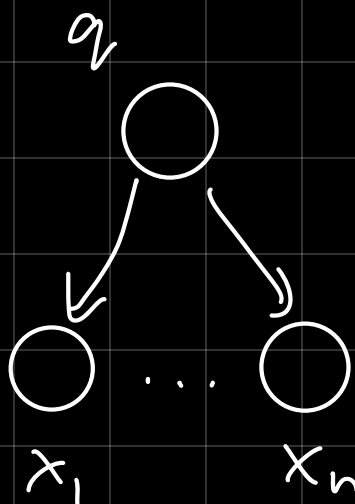
$= p(x|y) \checkmark$

# Example: repeated iid sampling

$q \sim \text{Beta}(\alpha, \beta)$

$X_i \overset{iid}{\sim} \text{Bern}(q)$

$P(X_1, \ldots X_n \mid q)$
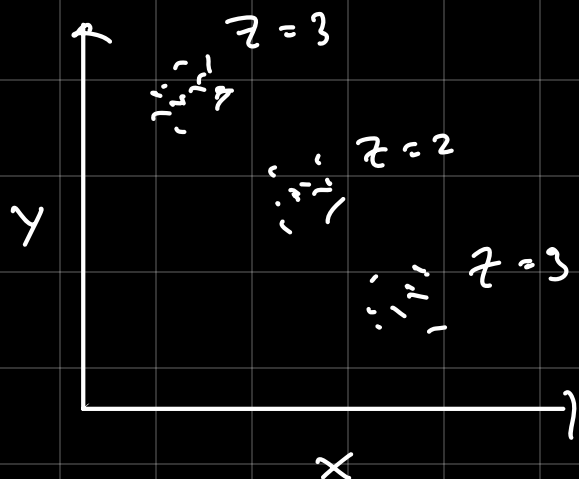
$= \prod_{i=1}^{n} P(X_i \mid q)$

Plates denote repetition

$P(X_1, \ldots X_n) = \underbrace{P(X_1)}_{\text{Bern}(\frac{\alpha}{\alpha+\beta})} \underbrace{P(X_2 \mid X_1)}_{\text{Bern}(\frac{\alpha + X_1}{\alpha+\beta+1})} \cdots$
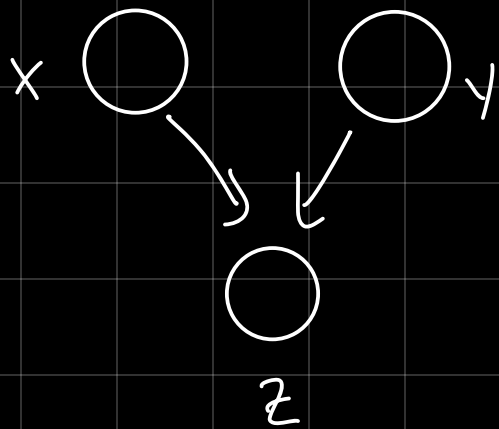
cond. indep $\implies$ marginal $\underline{\text{dependence}}$

# Example: confounding

$X \not\perp\!\!\!\perp Y \mid z = 3$

$z = 3$

$z = 2$

$z = 3$

Y

X

③ V-structure



$X$     $Y$

$Z$

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

$$X \perp\!\!\!\perp Y \ ?$$

$$P(x \mid y) = \frac{\sum_{z} P(x) \cancel{P(y)} \cancel{P(z \mid x, y)}}{\sum_{x} \sum_{z} \cancel{P(x)} \cancel{P(y)} \cancel{P(z \mid x, y)}}$$

$$= P(x) \ \checkmark$$



$$X \perp\!\!\!\perp Y \mid Z \ ?$$

$$P(x \mid y) = \frac{\sum_{z} P(x) P(y) P(z \mid x, y)}{\sum_{x} \sum_{z} P(x) P(y) P(z \mid x, y)}$$

No.

e.g. "explaining away"

$X$ = covid

$y$ = common cold

$Z$ = coughing

• $X \perp\!\!\!\perp Y \mid Z$ does not come from graph separation...

# Bayes ball

- a general algorithm to determine if
$$X_A \perp\!\!\!\perp X_B \mid X_C$$

  for subsets $A, B, C$ of nodes in a DGM.

- based on <u>reachability</u> : if a ball <u>cannot</u> bounce from $X_A$ to $X_B$ when $X_C$ is observed, then $X_A \perp\!\!\!\perp X_B \mid X_C$

- The Rules

  ✓ → ○ →          ✗ → ⊛ →

  ✓ ← ○ →          ✗ ← ⊛ →
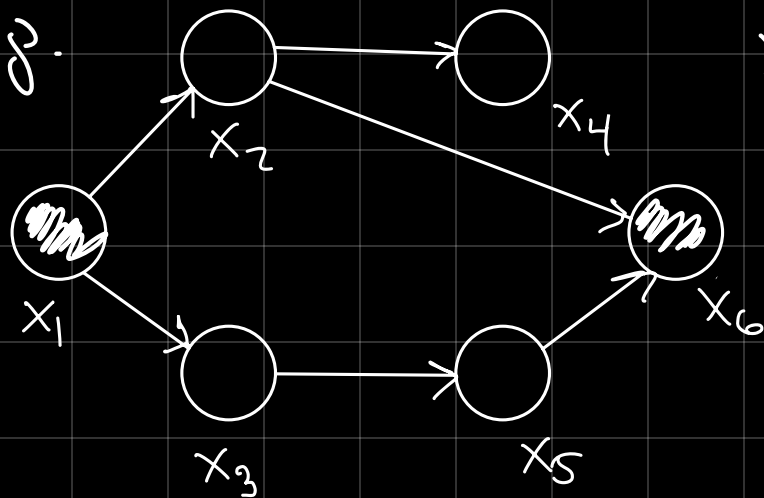
  ✗ → ○ ←          ✓ → ⊛ ←

e.g.



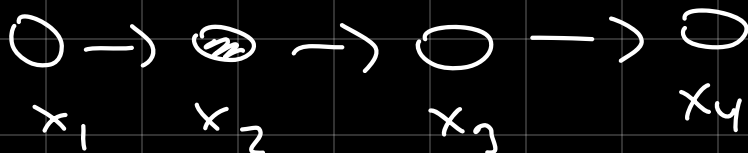$X_1 \perp\!\!\!\perp X_6 \mid X_2, X_3$ ?

yes

e.g.



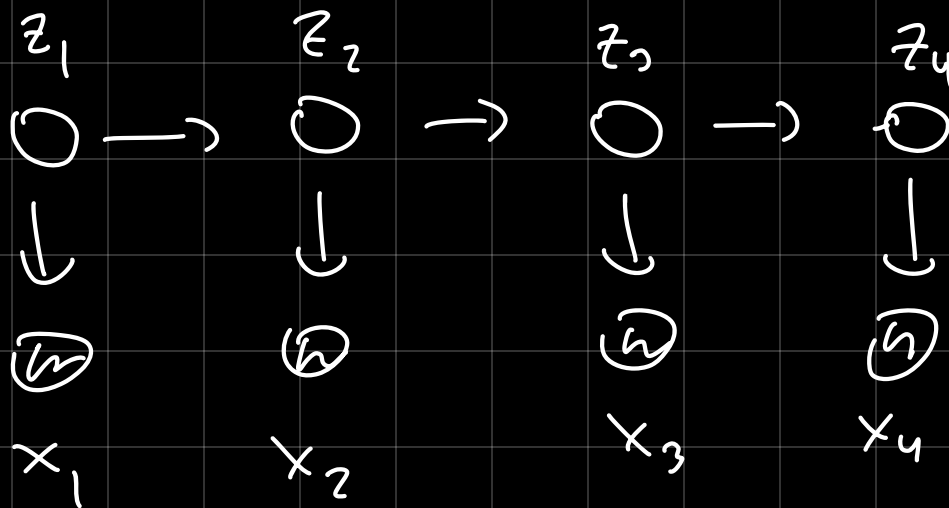$X_2 \perp\!\!\!\perp X_3 \mid X_1, X_6$ ?

no

e.g.



Markov chain

$X_3 \perp\!\!\!\perp X_1 \mid X_2$   (by design)

$X_4 \perp\!\!\!\perp X_1 \mid X_2$   (by Bayes ball)

e.g.



Hidden Markov model (HMM)

$X_1 \perp\!\!\!\perp X_4 \mid X_2, X_3$ ?

no.

# Theorem (Hammersley - Clifford)

- $G = (V, E)$ is a DAG over nodes $V = \{x_1 \cdots x_n\}$

- $S_1 = \{p : p \text{ respects } G\}$

  all joint dists $p = p(x_1 \cdots x_n)$ that respect all cond. independencies implied by $G$

- $S_2 = \{P_{G, \Phi} \text{ for all } \overline{\Phi}\}$
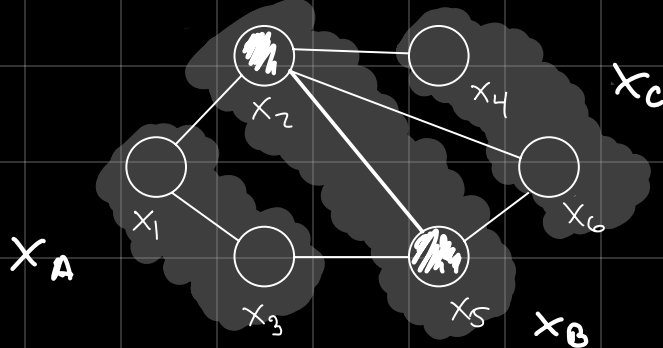
  where $\overline{\Phi}$ is a value for the LPTs in $G$

- Thm : $S_1 = S_2$

# Undirected graphical models

- "Markov random fields" (MRFs)

- $X_A \perp\!\!\!\perp X_C \mid X_B$    IFF    $X_B$ graph-separates $X_A, X_C$



- DGMs are acyclic and have an _ordering_; they therefore define the joint via LPTs and the chain rule

- How do UGMs parametrize the joint?

- Recall: $P(x_1 \ldots x_n) = \dfrac{f(x_1 \ldots x_n)}{Z}$     ← "kernel"
  
  "normalizer"

- $f(x_1 \ldots x_n) = \prod_{c \in C} \psi_c(x_c)$

- $C$ are all _maximal cliques_ in $G$

- $\psi_c(\cdot)$ is the _potential function_ for $X_c$

- $\psi_c(x_c) > 0$ is the _non-negative_ potential for configuration $X_c = x_c$

- These are like the LPTs, but they are not normalized

- **Maximal Cliques:**
  - connected components
  - $C = \{ \underset{c=1}{\{1,3\}}, \underset{2}{\{1,2\}}, \underset{3}{\{2,5,6\}}, \underset{4}{\{2,4\}} \}$
  - notice $x_1$ appears in $c=1$ and $c=2$
  - $\ell_1(x_1, x_3)$, $\ell_2(x_1, x_2)$, $\cdots$
  - $\ell_c$ measures the agreement of a clique

- $Z = \sum_{x_1} \cdots \sum_{x_n} \prod_c \ell(x_c)$   "partition function"

- Hard to compute

- "Energy"   $\ell(x_c) \triangleq \exp\left(- \underbrace{H_c(x_c)}\right)$

$$P(x_1 \cdots x_n) \triangleq \frac{1}{Z} \prod_{c \in C} \exp(- H_c(x_c))$$

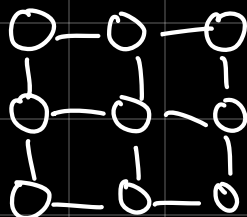$$= \exp\left(- \underbrace{\sum_c H(x_c)} - \log Z\right)$$

$$= H(x_1 \cdots x_n)$$

"Boltzmann distribution"

- low-energy configurations are more probable.

e.g., Ising model          $H(x_1 \cdots x_n) = \underset{i,j}{\sum} \underbrace{x_i x_j c_{ij}}_{\ell(x_i, x_j)}$

$x_{ij} \in \{-1, 1\}$  spin

## DGMs          vs.          UGMs

- graph → independence          dependence → graph
- good for generative          • good for complex
                                     dependence

- not all __families__ can be expressed by both
  DGMs and UGMs

X ○        ○ Y

Y ↓

Ⓜ

Z

$X \perp\!\!\!\perp Y$ and $X \not\perp\!\!\!\perp Y \mid Z$

(no way to express this
with an UGM)

vs.



W    X

Y    Z

$Y \perp\!\!\!\perp X \mid W, Z$

$W \perp\!\!\!\perp Z \mid X, Y$

(no way to express
with DGM)

Note for a given LPT:

$$P(x = 1 \mid y = 1, z) = p$$

$$P(x = 1 \mid y = 0, z) = p \quad \forall z$$

$$\Downarrow$$

$$X \perp\!\!\!\perp Y \mid Z$$