

Machine Learning 2

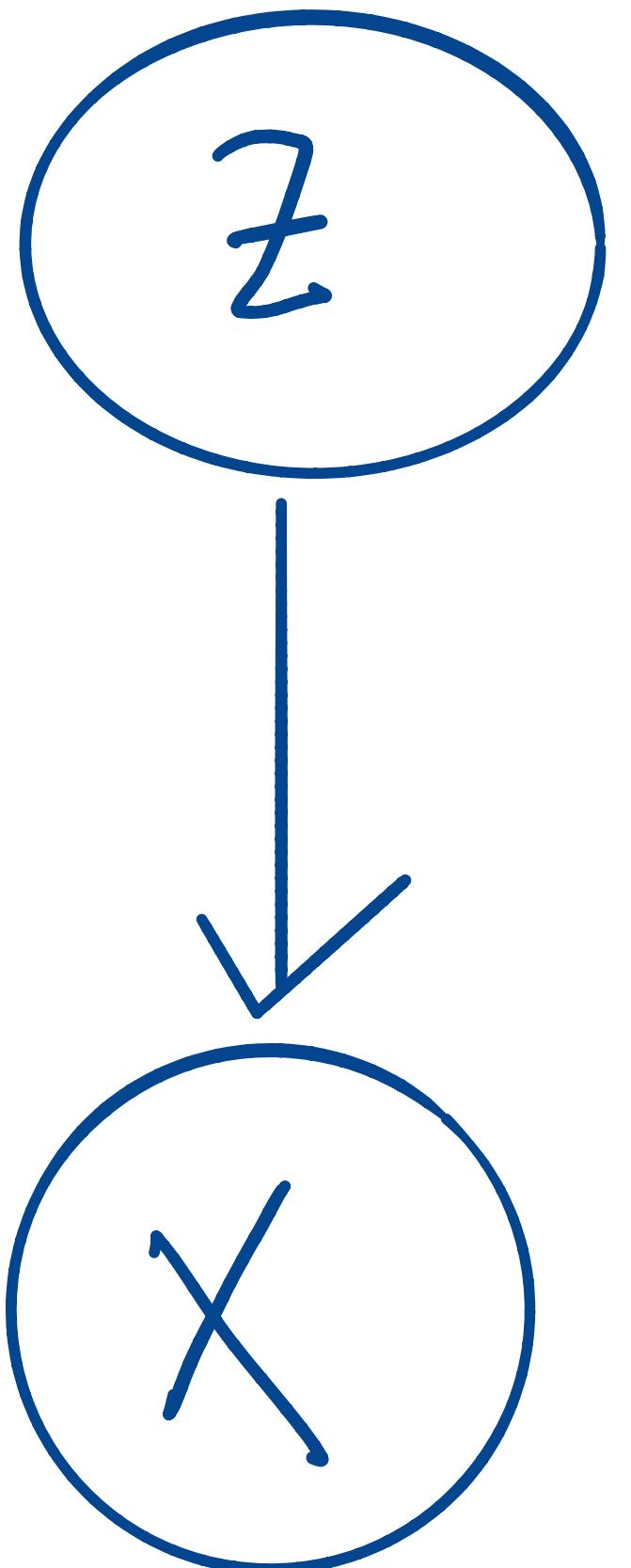
**Inference
- Motivation**

Patrick Forré

What is Inference?

- Let a joint probability distribution $p(x_1, \dots, x_M)$ be represented by the computer, e.g. probability tables, factorized forms, via parameters, etc.
- Inference in our context means to derive or compute a representation of certain properties of $p(x_1, \dots, x_M)$, e.g.
 - marginals distributions: $p(x_k, x_j)$,
 - conditionals distributions: $p(x_k, x_j | x_i, x_v)$
 - modes: $x^* \in \arg \max_x p(x)$
 - etc.

Example: Gaussian



$$p(z) = N(z | \mu, \Sigma)$$

$$p(x|z) = N(x | Az + b, S)$$

$$p(x, z) = p(x|z) \cdot p(z) \stackrel{?}{=} N\begin{pmatrix} z \\ x \end{pmatrix} | m, T$$

$$m = \begin{pmatrix} \mu \\ A\mu + b \end{pmatrix}, \quad T = \begin{pmatrix} \Sigma & \Sigma A^T \\ A\Sigma & S + A\Sigma A^T \end{pmatrix}$$

$$p(x) = N(x | A\mu + b, S + A\Sigma A^T)$$

joint ↓
 ↑ marginal

Example: Gaussian

$$N\left(\begin{pmatrix} z \\ x \end{pmatrix} \mid \begin{pmatrix} \mu \\ Ay+b \end{pmatrix}, \begin{pmatrix} \Sigma & \Sigma A^T \\ A\Sigma & A\Sigma A^T + S \end{pmatrix}\right)$$

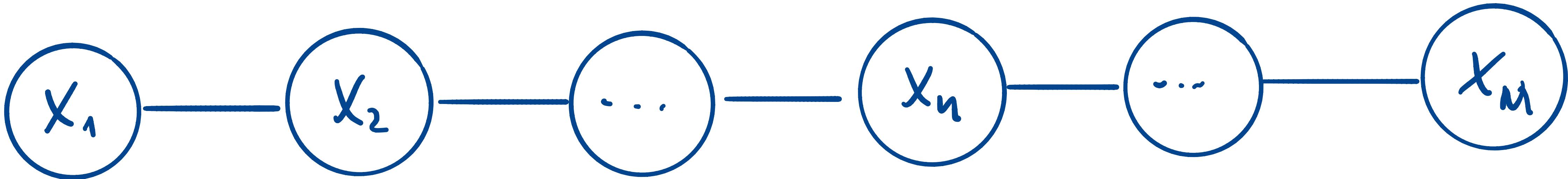
joint

$$p(z|x) = N(z \mid \mu + \Sigma A^T (A\Sigma A^T + S)^{-1} (x - Ay - b),$$
$$\Sigma - \Sigma A^T (A\Sigma A^T + S)^{-1} A \Sigma)$$

Inverse $\sim \mathcal{O}(n^3)$

(might be infeasible)

Example: Discrete Markov Chain



$$p(x) = \frac{1}{z} \cdot \prod_{i=1}^{M-1} \psi_{\{i, i+1\}}(x_i, x_{i+1})$$

each x_i with K classes

$$p(x_n) = \sum_{x_1, x_2, \dots, x_{n-1}, x_{n+1}, x_m} \sum_{\underbrace{x_1, x_2, \dots, x_{n-1}, x_{n+1}, x_m}_{K^{M-1} \text{ states}}} p(x_1, \dots, x_{n-1}, x_n, x_{n+1}, \dots, x_m) = \sum_{x_1} \sum_{x_n} \sum_{x_m} \frac{1}{z} \prod_i \psi_{\{i, i+1\}}(x_i, x_{i+1})$$

K^{M-1} states

to add up

Example: Discrete Markov Chain

$$P(x_n) = \sum_{x_1} \dots \sum_{x_M} \frac{1}{Z} \prod_{i=1}^M \psi_{\{i, i+1\}}(x_i, x_{i+1})$$
$$= \frac{1}{Z} \sum_{x_1} \sum_{x_2} \dots \psi_{\{1, 2\}}(x_1, x_2) \dots \psi_{\{n-1, n\}}(x_{n-1}, x_n) \sum_{x_{n+1}} \psi_{\{n, n+1\}}(x_n, x_{n+1}) \dots \left(\sum_{x_M} \psi_{\{M-1, M\}}(x_{M-1}, x_M) \right)$$

M times

K^2

$$\sim M \cdot K^2 \ll K^{M-1}$$

Different ways to do inference

- **Exact Inference** (desirable, but often computational infeasible)
 - Exact Inference in Graphical Models (exploiting sparse representations)
- **Approximate Inference**
 - **Variational Inference** (analytic or numerical approximations)
 - Expectation-Maximization (EM) Algorithm
 - Variational Bayes Mean Field Approximation
 - Variational Autoencoder (VAE)
 - **Monte Carlo Sampling Methods** (approximation creating and using (weighted) samples)
 - Importance, Slice, Rejection Sampling, etc.
 - Metropolis-Hastings MCMC

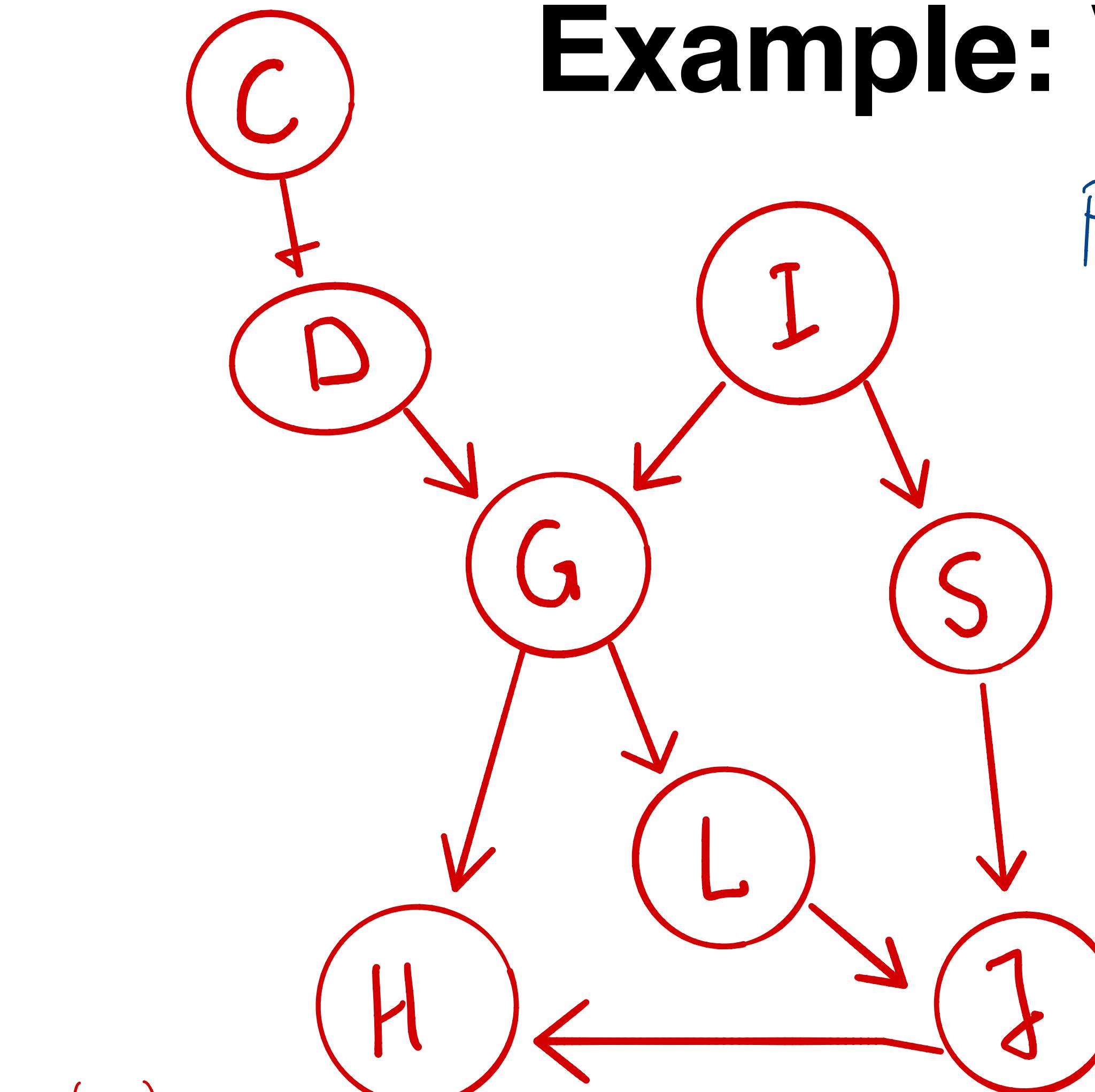
Machine Learning 2

Graphical Models

- Exact Inference**
- Variable Elimination Algorithm**

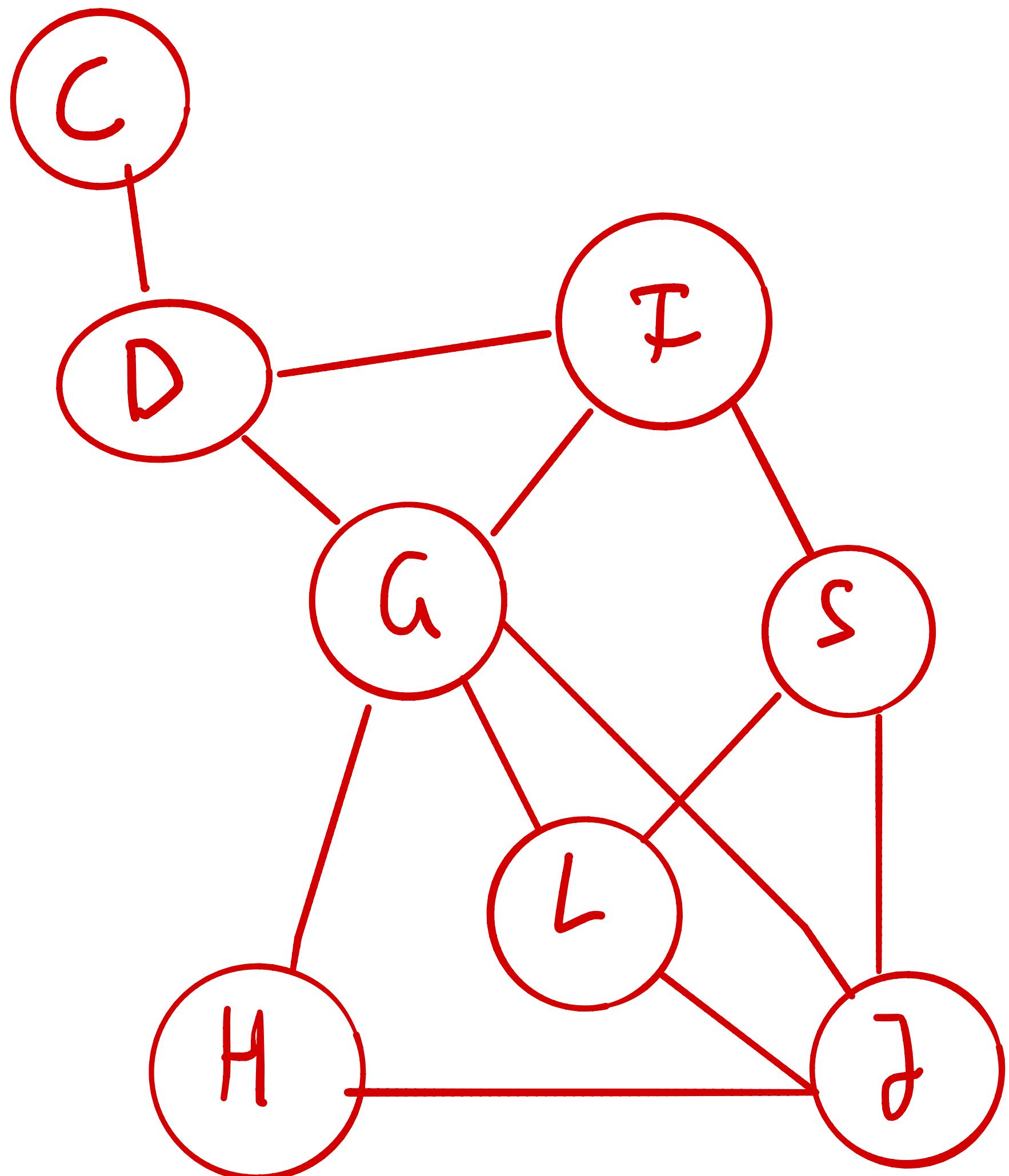
Patrick Forré

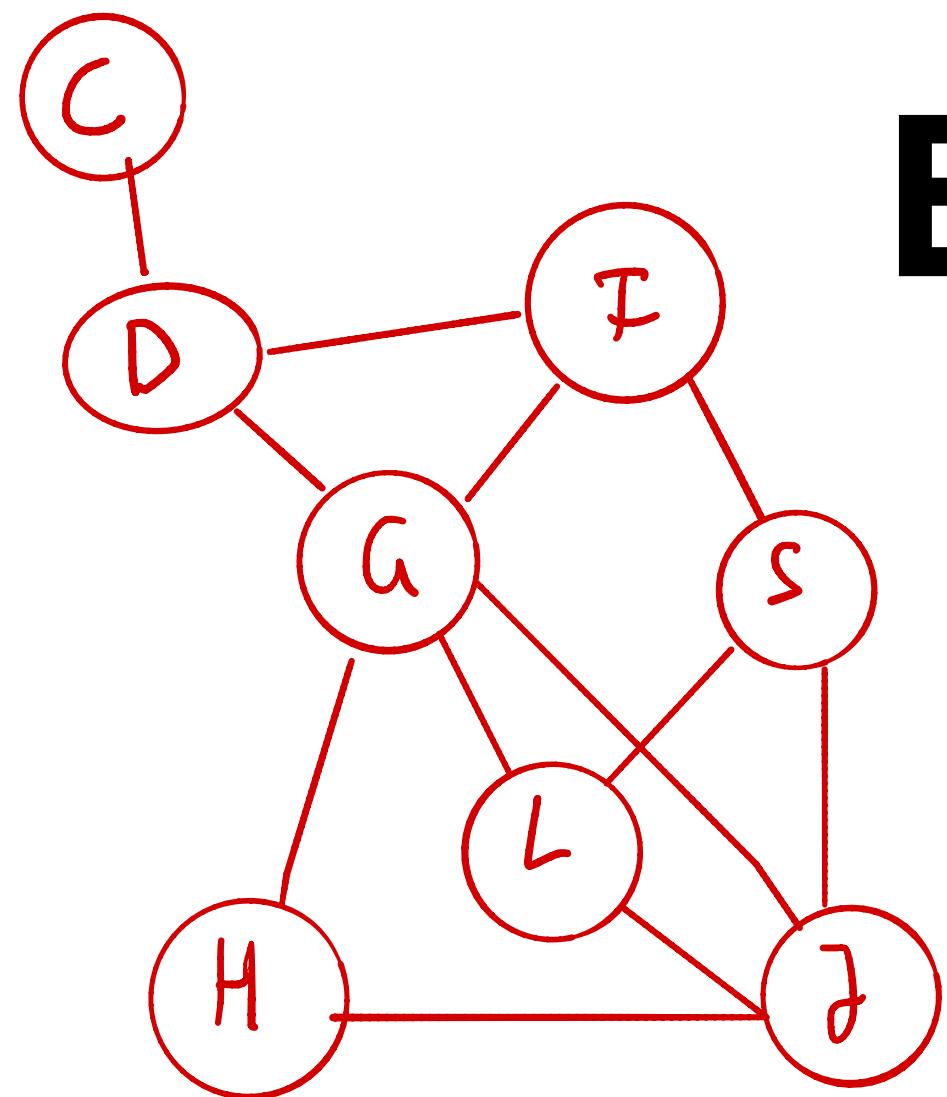
Example: Variable Elimination



$p(j)$

$$\begin{aligned}
 p(x) &= p(c) \cdot p(d|c) \cdot p(i) \cdot p(g|i,d) \cdot p(s|i) \cdot p(l|g) \cdot p(j|l,s) \cdot p(h|g,j). \\
 &= q(c) \cdot q(d|c) \cdot q(i) \cdot q(g,i,d) \cdot q(s|i) \cdot q(l,g) \cdot q(j,l,s) \cdot q(h|g,j)
 \end{aligned}$$





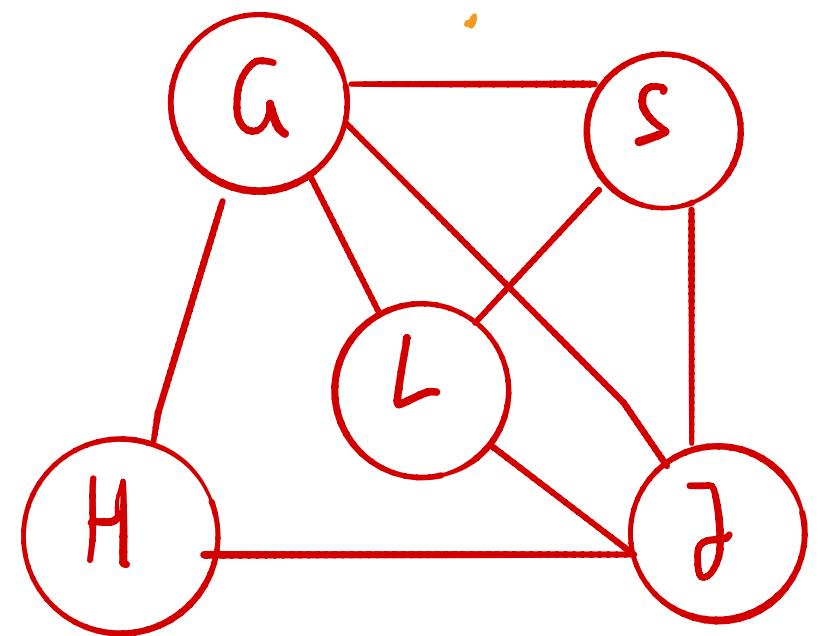
Example: Variable Elimination

$$P(X) = 4(C) \cdot 4(D, C) \cdot 4(I) \cdot 4(A, I, D) \cdot 4(S, I) \cdot 4(L, A) \cdot 4(J, L, S) \cdot 4(H, G, J)$$

$$P(j) = \sum_L \sum_S \sum_A \sum_H \sum_I \sum_D \sum_C P(X)$$

$$\begin{aligned}
 &= \sum_{L, S} 4(j, L, S) \cdot \sum_G 4(L, G) \sum_H 4(H, G, j) \cdot \sum_I 4(S, I) \cdot 4(I) \sum_D 4(A, I, D) \sum_C 4(C) 4(D, C). \\
 &= \sum_{L, S} 4(j, L, S) \cdot \sum_G 4(L, G) \sum_H 4(H, G, j) \cdot \sum_I 4(S, I) \cdot 4(I) \sum_D 4(A, I, D) \cdot \phi_1(D) \\
 &= \sum_{L, S} 4(j, L, S) \cdot \sum_G 4(L, G) \sum_H 4(H, G, j) \cdot \sum_I 4(S, I) \cdot 4(I) \cdot \phi_2(A, I) \\
 &= \sum_{L, S} 4(j, L, S) \cdot \sum_G 4(L, G) \sum_H 4(H, G, j) \cdot \phi_3(S, A)
 \end{aligned}$$

Example: Variable Elimination

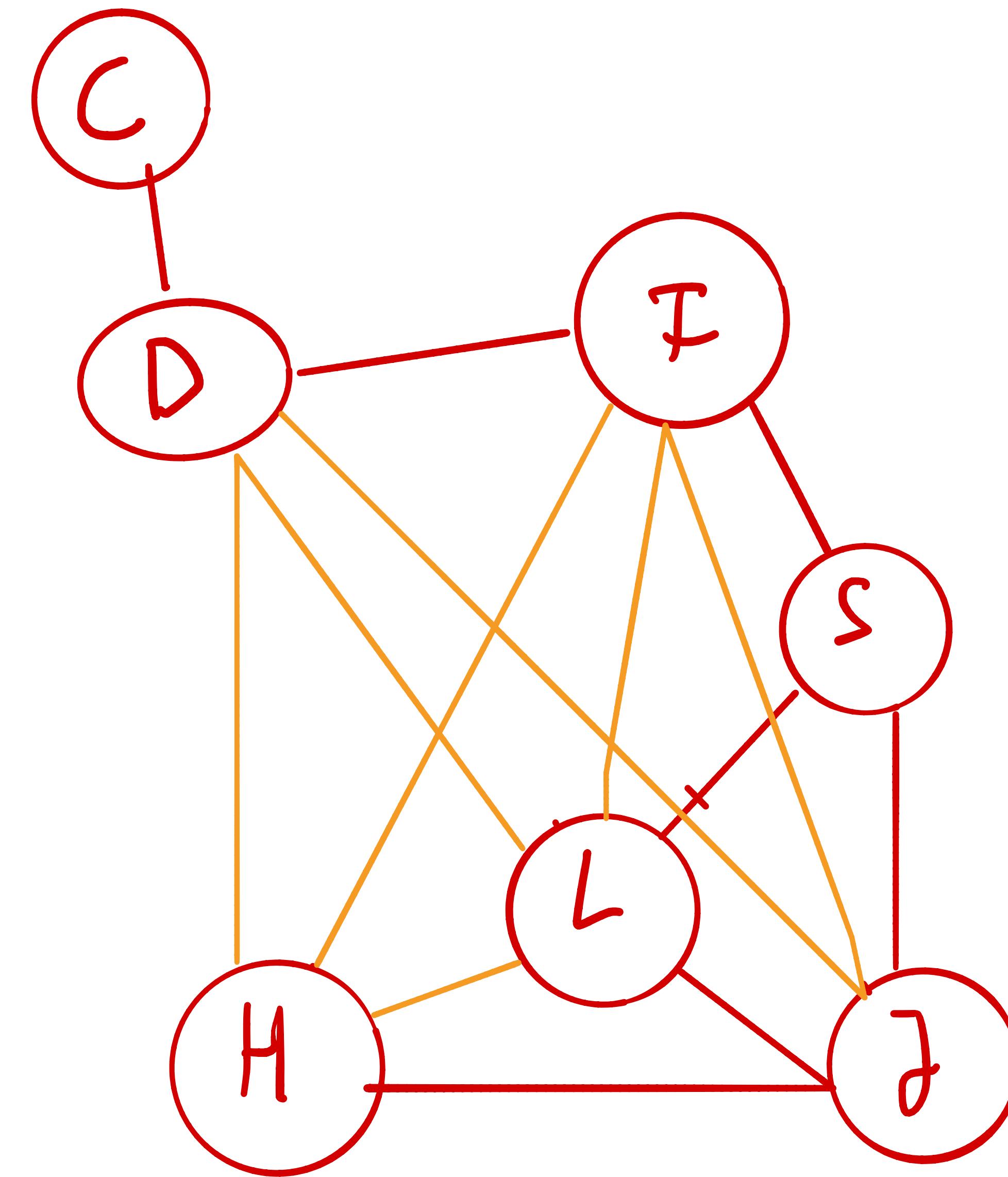
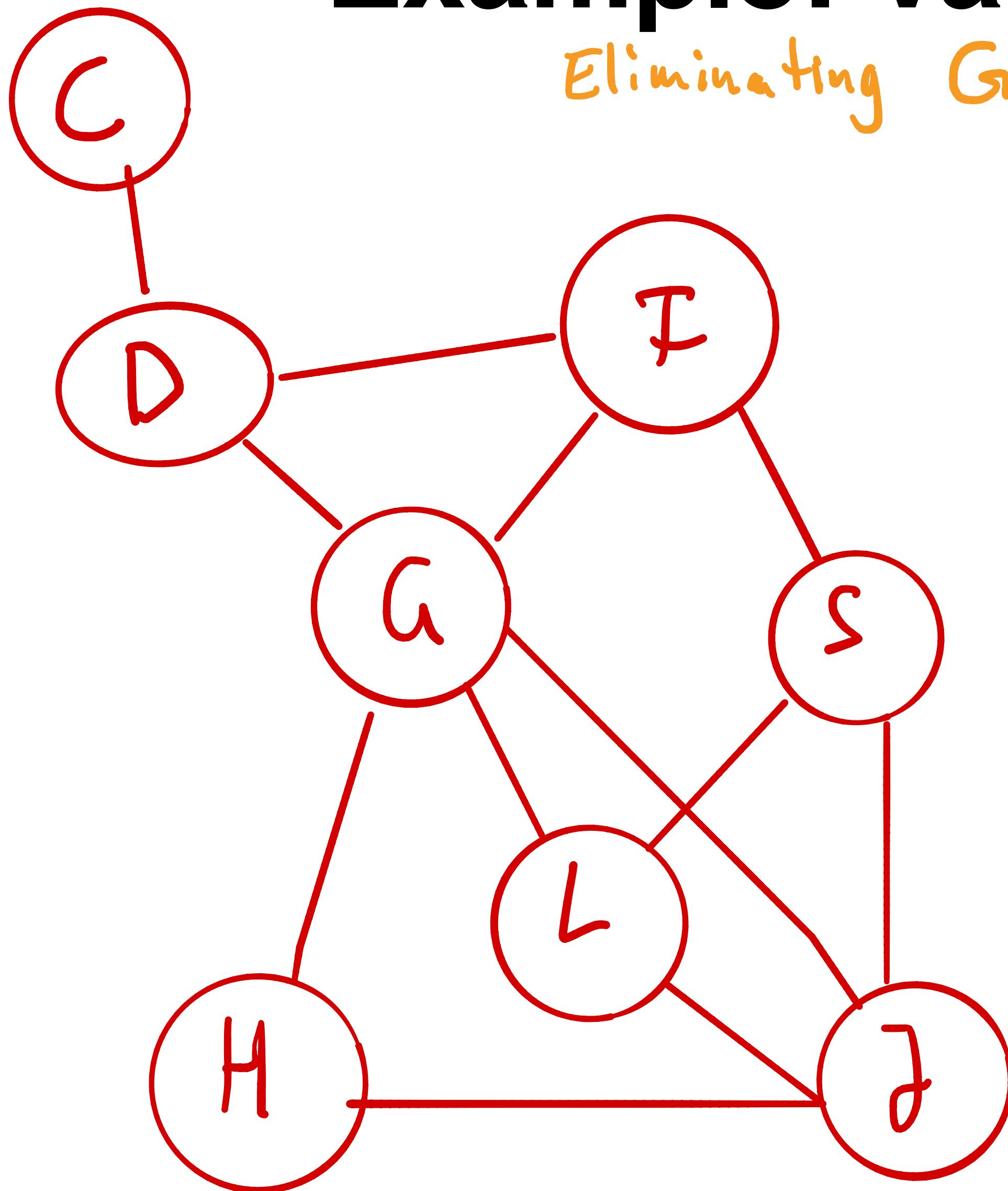


$$P(X) = 4(C) \cdot 4(D, C) \cdot 4(I) \cdot 4(A, I, D) \cdot 4(S, I) \cdot 4(L, A) \cdot 4(J, L, S) \cdot 4(H, G, J)$$

$$P(j) = \sum_L \sum_S \sum_A \sum_H \left(\sum_I \sum_D \sum_C P(X) \right)$$

$$\begin{aligned}
 P(j) &= \sum_{L, S} 4(j, L, S) \cdot \sum_G 4(L, G) \quad \phi_S(S, A) \underbrace{\sum_H 4(H, A, j)}_{\text{H}} \\
 &= \sum_{L, S} 4(j, L, S) \cdot \sum_G 4(L, G) \quad \phi_S(S, A) \cdot \underbrace{\phi_U(A, j)}_{\phi_S(L, S, j)} \\
 &= \sum_{L, S} 4(j, L, S) \\
 &= \phi_6(j)
 \end{aligned}$$

Example: Variable Elimination



Variable Elimination Algorithm

- Let (G, p) be BN or MRF. Goal: Compute $p(x_C)$
- Finding a good ‘elimination’ ordering of the nodes $V = \{v_1, \dots, v_M\}$:
 - Finding the optimal one is NP-hard (but Variable Elimination has exponential time complexity i.g.).
 - Heuristics for MRFs:
 - Least number of new edges created at each step.
 - Create smallest factor possible at each step.
 - Heuristics for BNs:
 - Iteratively eliminate child-less nodes not in C first. End up with $\text{Anc}^G(C)$, proceed from there.
 - Use heuristics from above.
- Eliminate / Marginalize variables out in that order.
 - by combining all involved factors by summing over all the variables values

Machine Learning 2

Graphical Models

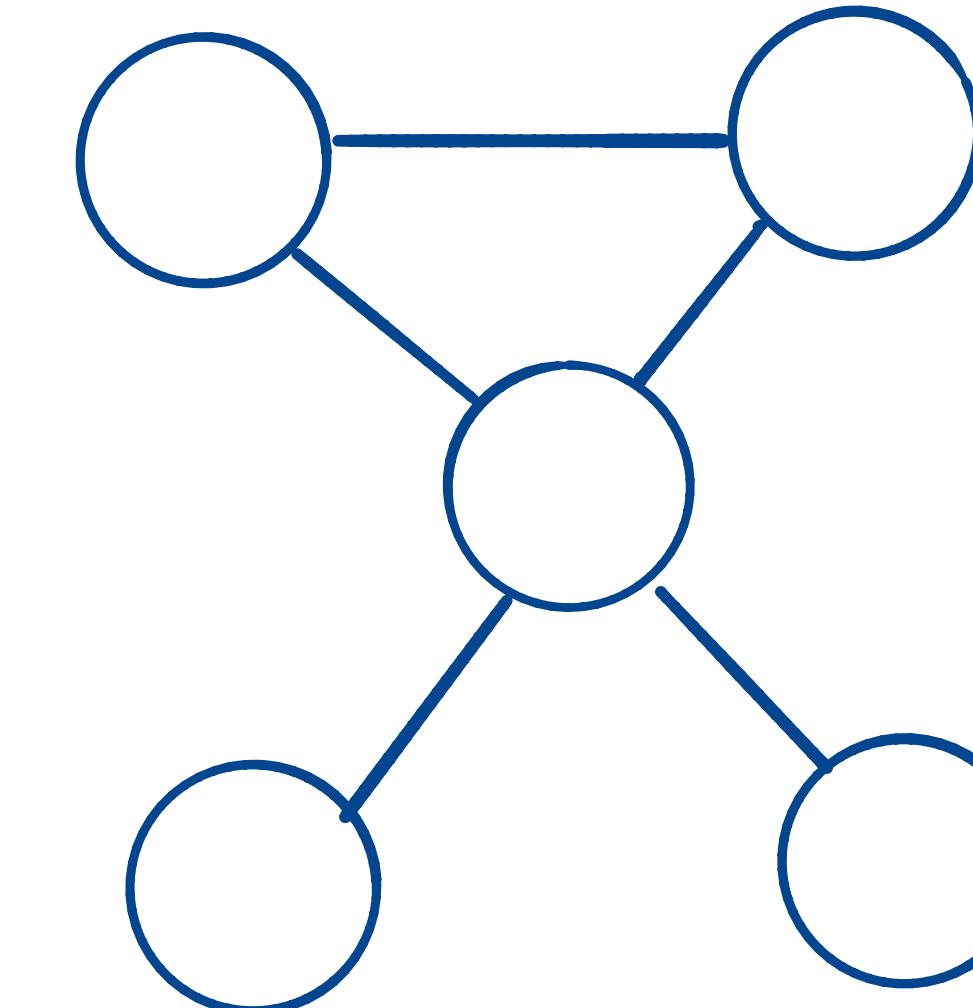
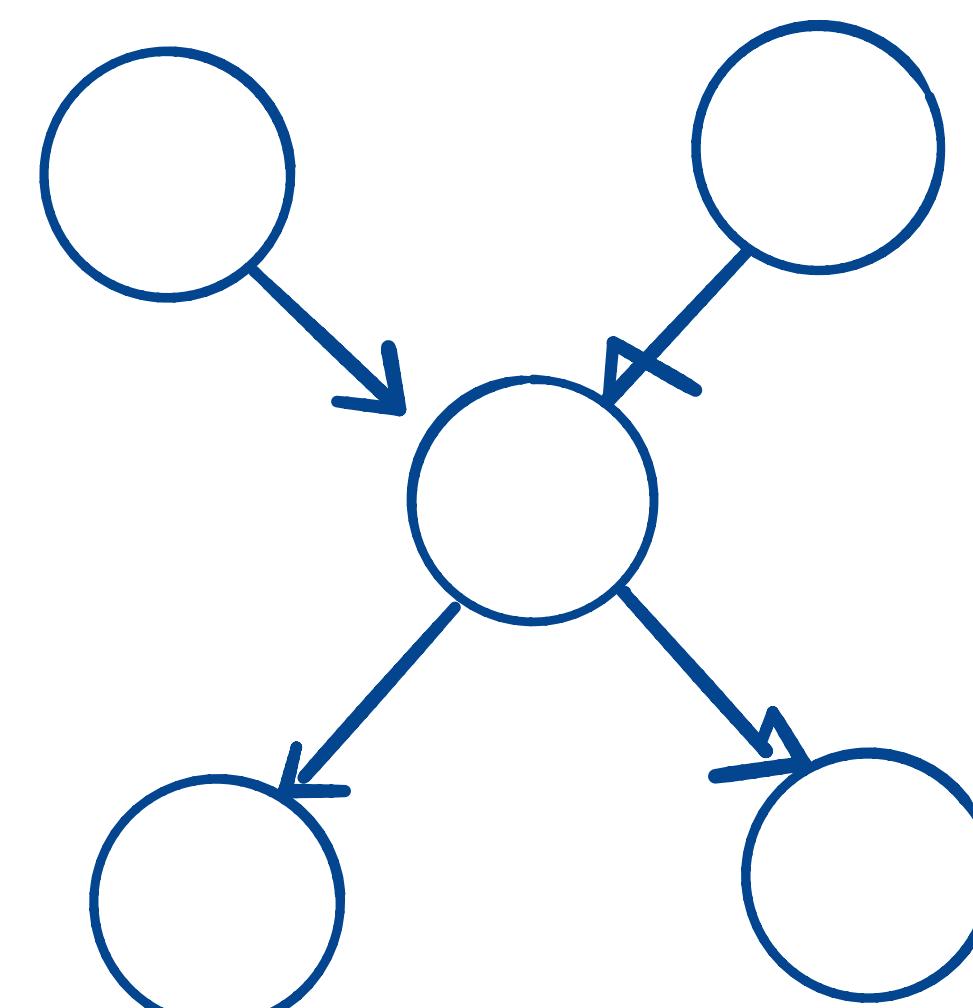
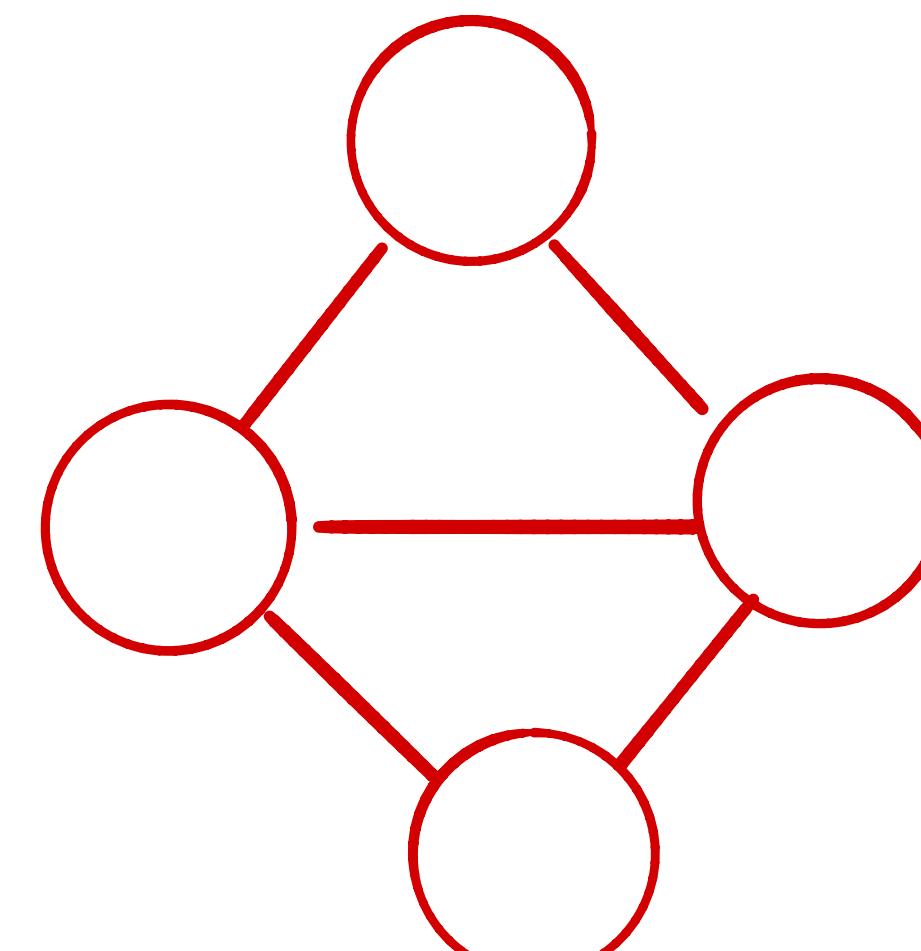
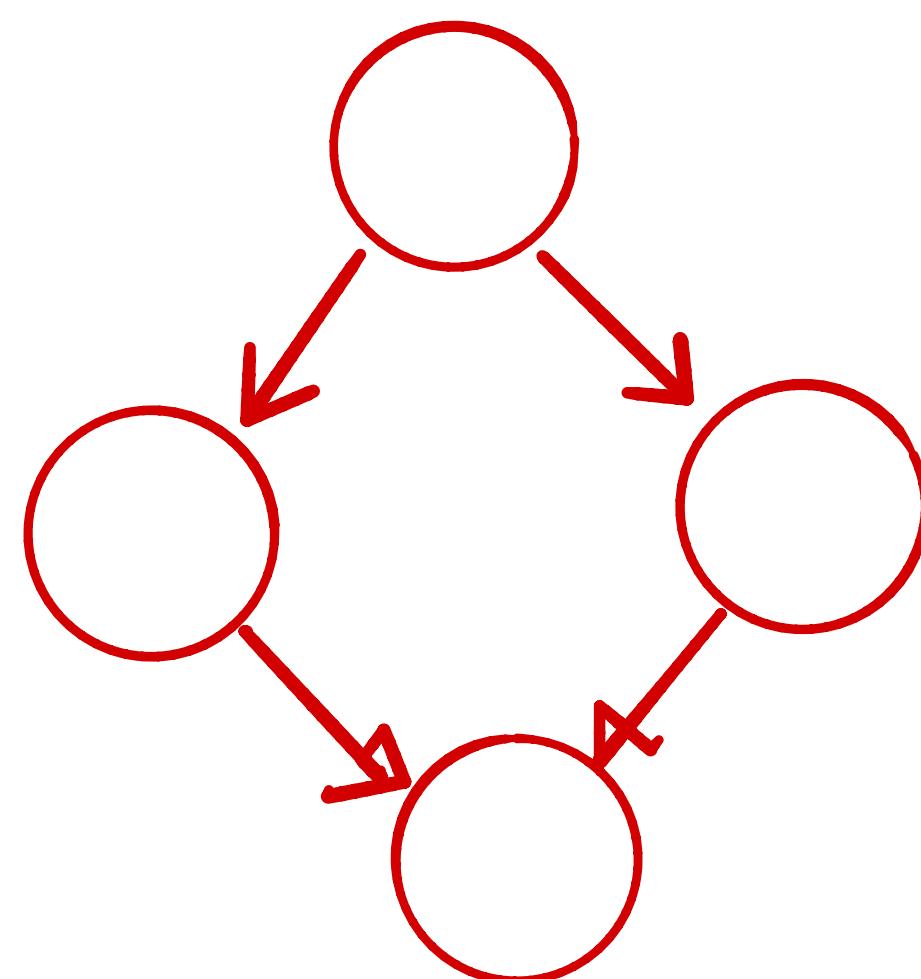
- Exact Inference**
- Factor Trees**

Patrick Forré

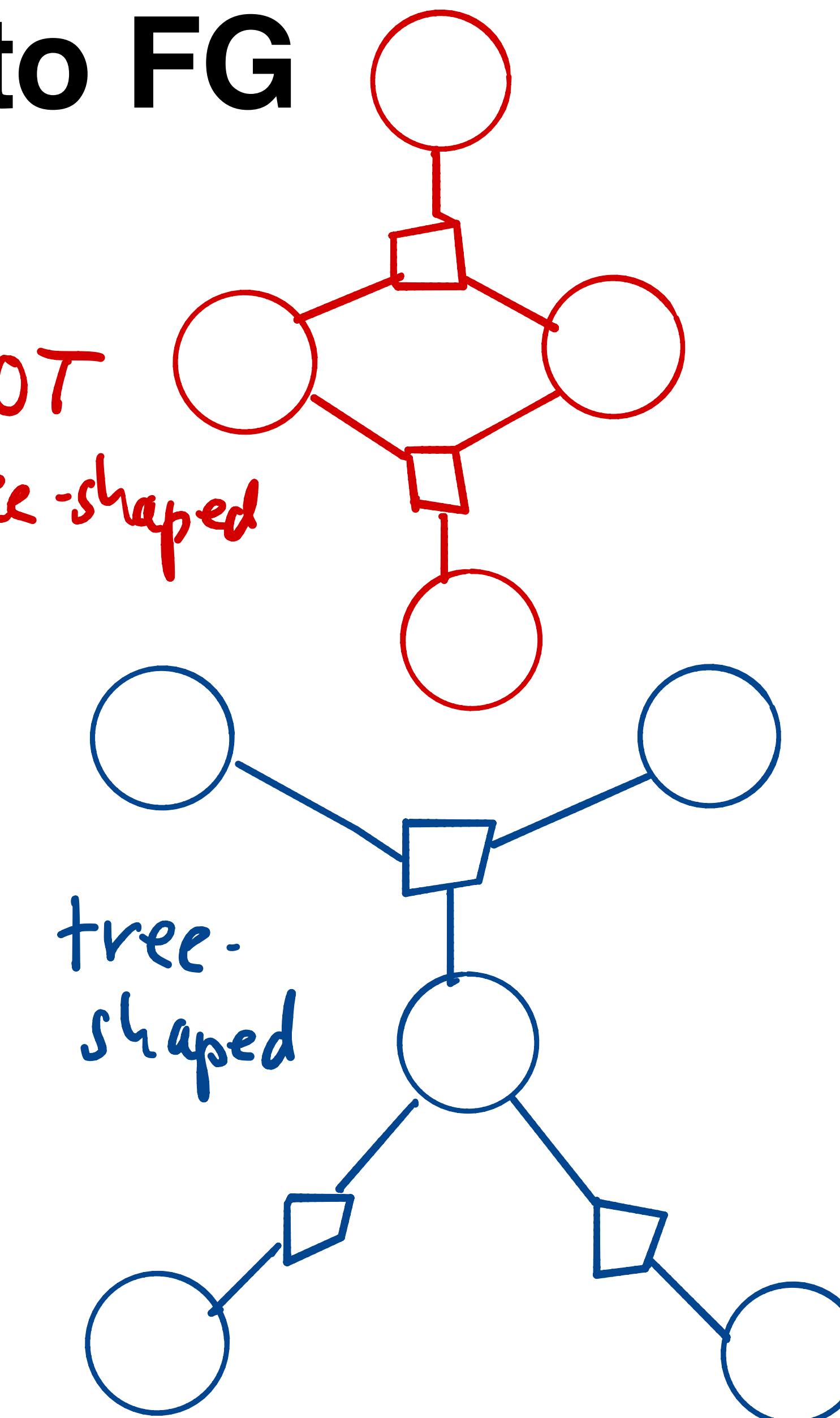
Tree-shaped Factor Graphs

- A Factor Graph $G = (V, F, E)$ is called tree-shaped or Factor Tree if it does not contain any (undirected) cycles.
- Since we can always convert BNs and MRFs to FGs everything said can be applied to BNs and MRFs as long as their FGs are tree-shaped.
needs to be checked !
- Note that a FG constructed from a BN ('acyclic') might NOT be tree-shaped in general.
- Note that MRFs with cycles might not have cycles anymore in their FGs.

Example: MRF to BN to FG



NOT
tree-shaped



Machine Learning 2

Graphical Models

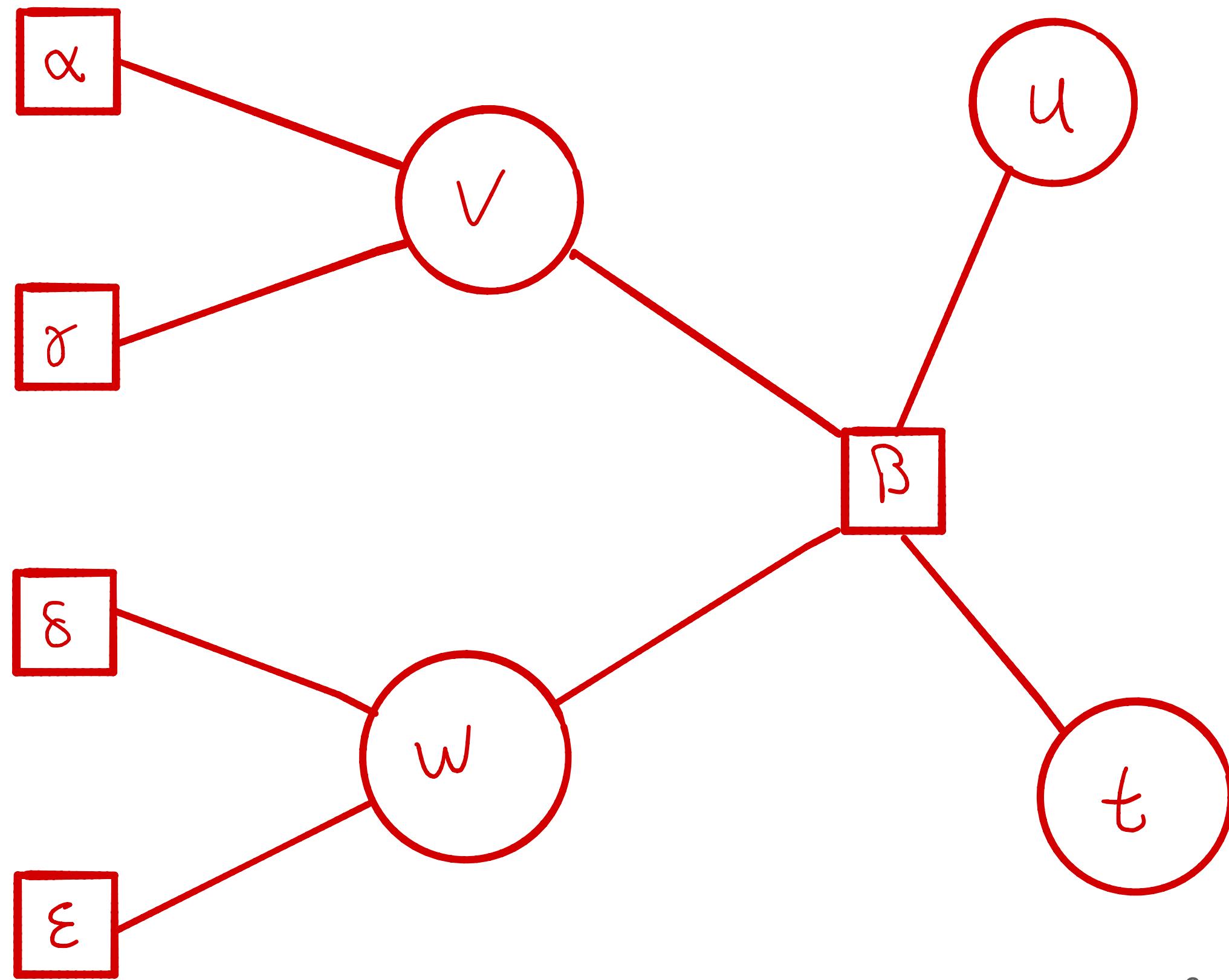
- Exact Inference**
- Sum-Product Algorithm**
- (in Factor Trees)**

Patrick Forré

Example: Factor Tree

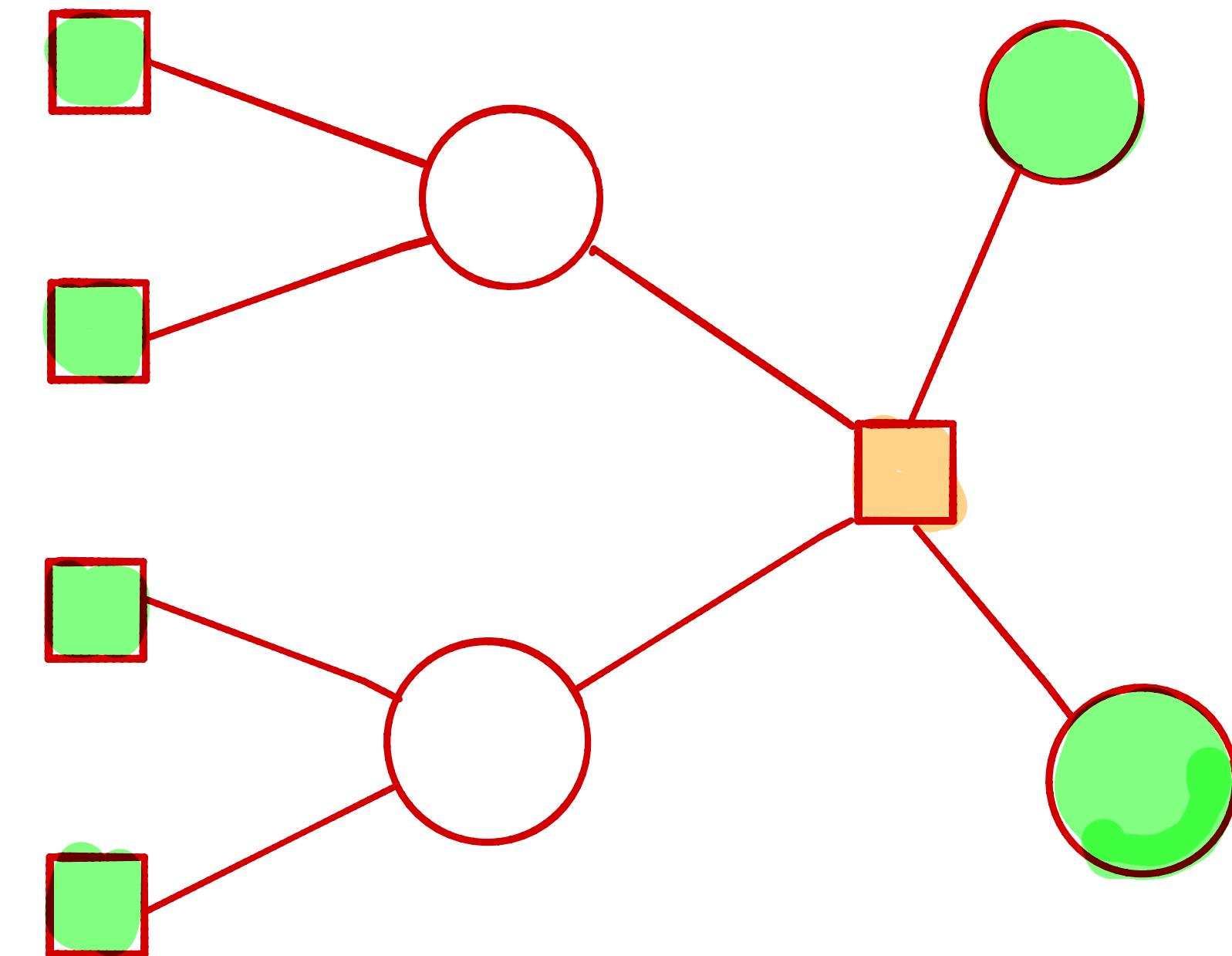
- $G = (V, F, E)$ Factor Tree with $p(x_V) = \frac{1}{Z} \prod_{\alpha \in F} \psi_\alpha(x_\alpha)$.

Goal: $p(x_r)$



Message Passing

- $G = (V, F, E)$ Factor Tree with $p(x_V) = \frac{1}{Z} \prod_{\alpha \in F} \psi_\alpha(x_\alpha)$.
- Goal: Find ~~Z~~ and all marginals $p(x_v)$ and $p(x_\alpha)$ for $v \in V$ and $\alpha \in F$.
- Pick a ‘Root’ node in the middle of G .
- Message Passing (only once forth and back):
 - Pass Messages from Leaves towards Root.
 - Pass Messages from Root back toward Leaves.
 - Locally compute ‘beliefs’ (marginals).



The Messages

- Leaf-variable-to-factor:

$$\mu_{u \rightarrow \beta}(x_u) = 1$$

- Non-leaf-variable-to-factor: $\mu_{v \rightarrow \beta}(x_v) =$

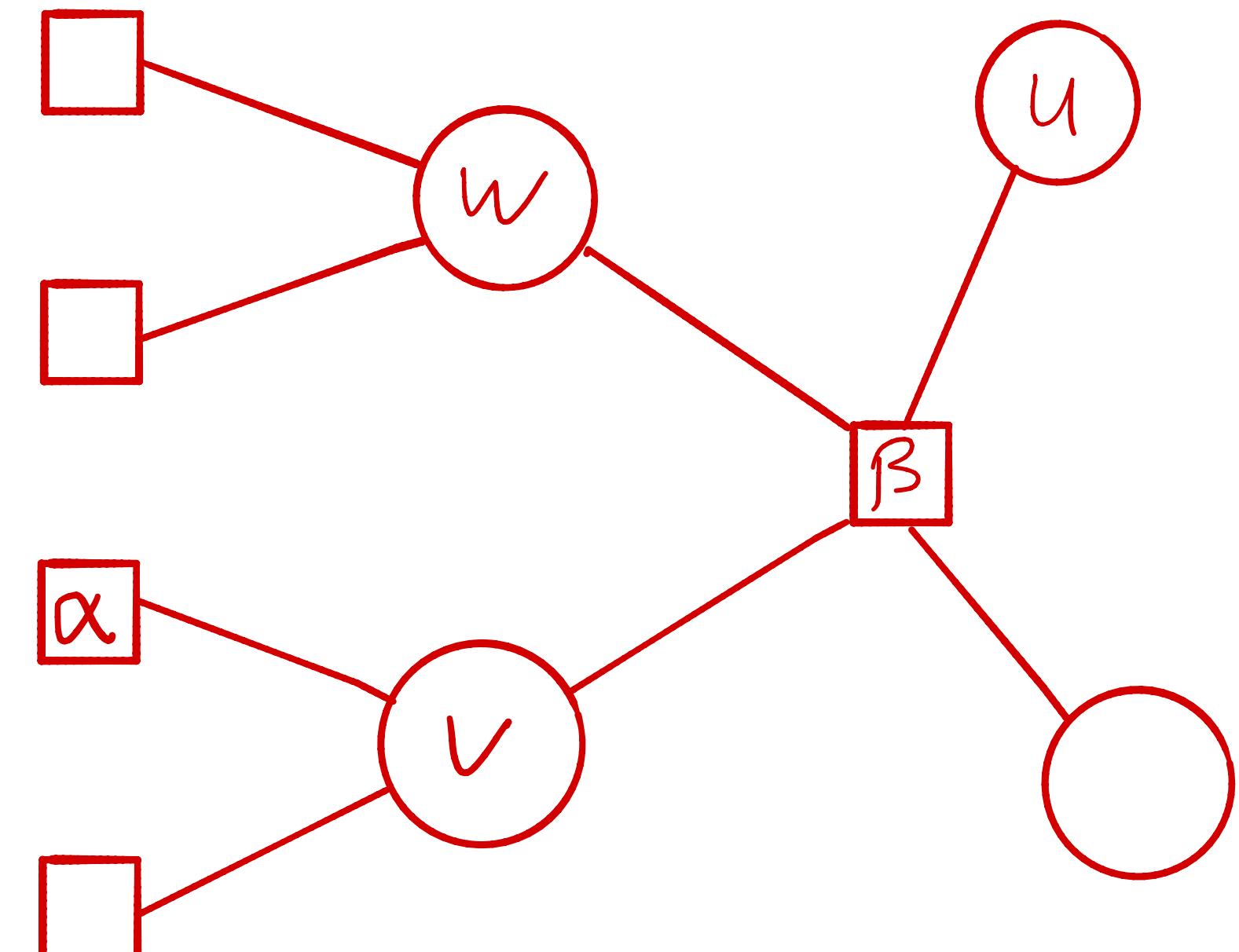
$$\prod_{\gamma \in \partial(v) \setminus \{\beta\}} \mu_{\gamma \rightarrow v}(x_v)$$

- Leaf-factor-to-variable:

$$\mu_{\alpha \rightarrow v}(x_v) = \psi_{\alpha}(x_v)$$

- Non-leaf-factor-to-variable: $\mu_{\beta \rightarrow w}(x_w) =$

$$\sum_{x_{\beta \setminus w}} \psi_{\beta}(x_{\beta}) \cdot \prod_{u \in \beta \setminus w} \mu_{u \rightarrow \beta}(x_u)$$



Computing the Beliefs

- **Variable beliefs** ($v \in V$):

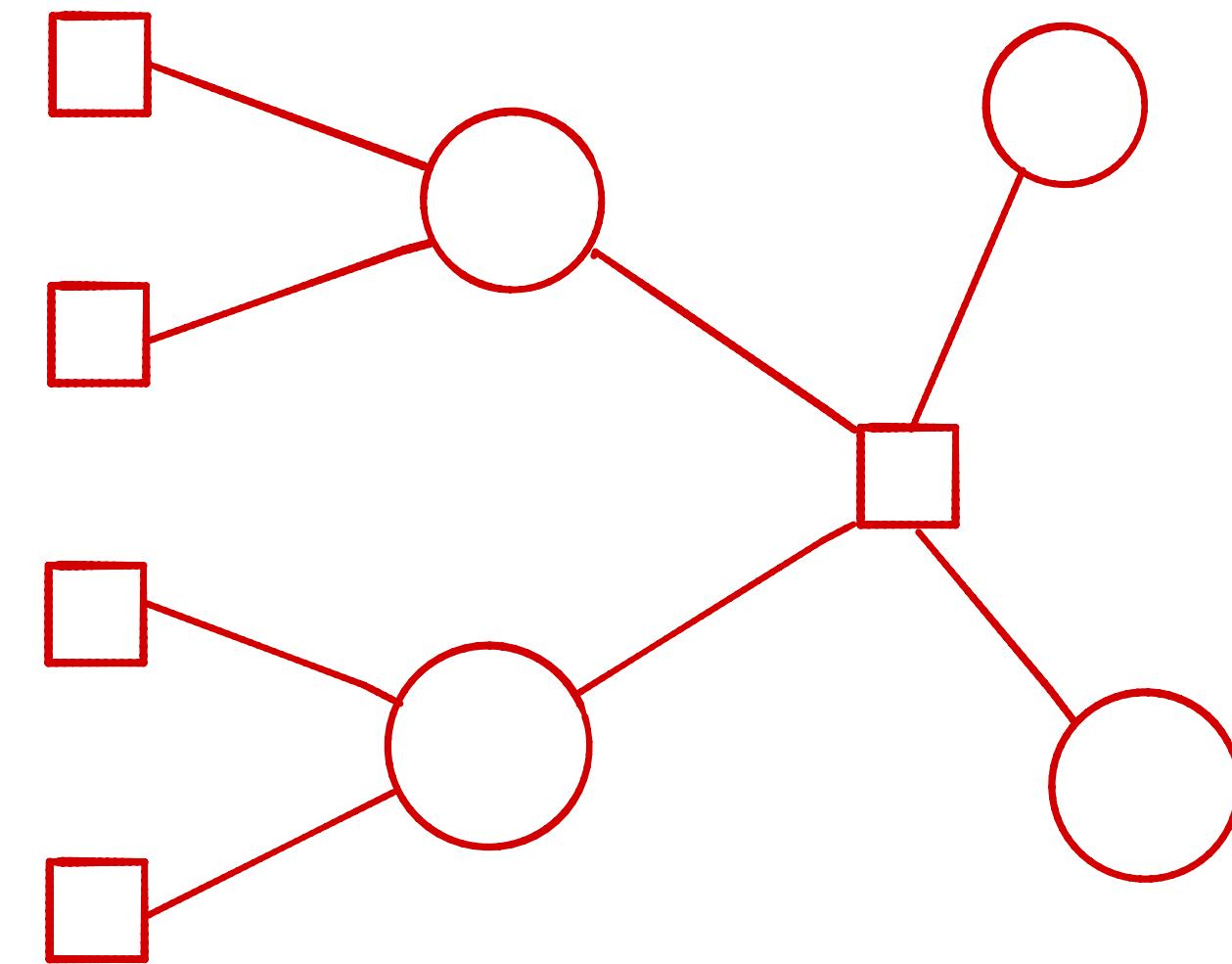
$$\bullet p(x_v) = \frac{1}{Z_v} \cdot \prod_{\alpha \in \partial(v)} M_{\alpha \rightarrow v}(x_v)$$

$$\bullet Z_v := \sum_{x_v} \prod_{\alpha \in \partial(v)} M_{\alpha \rightarrow v}(x_v)$$

- **Factor beliefs** ($\alpha \in F$):

$$\bullet p(x_\alpha) = \frac{1}{Z_\alpha} \psi_\alpha(x_\alpha) \cdot \prod_{v \in \partial(\alpha)} M_{v \rightarrow \alpha}(x_v)$$

$$\bullet Z_\alpha := \sum_{x_\alpha} \psi_\alpha(x_\alpha) \cdot \prod_{v \in \partial(\alpha)} M_{v \rightarrow \alpha}(x_v)$$



Exercise: Show that the ‘Beliefs’ are correct

- $G = (V, F, E)$ Factor Tree with $p(x_V) = \frac{1}{Z} \prod_{\alpha \in F} \psi_\alpha(x_\alpha)$.
- Show that the following equations after the Sum-Product-Algorithm has converged are actually correct:

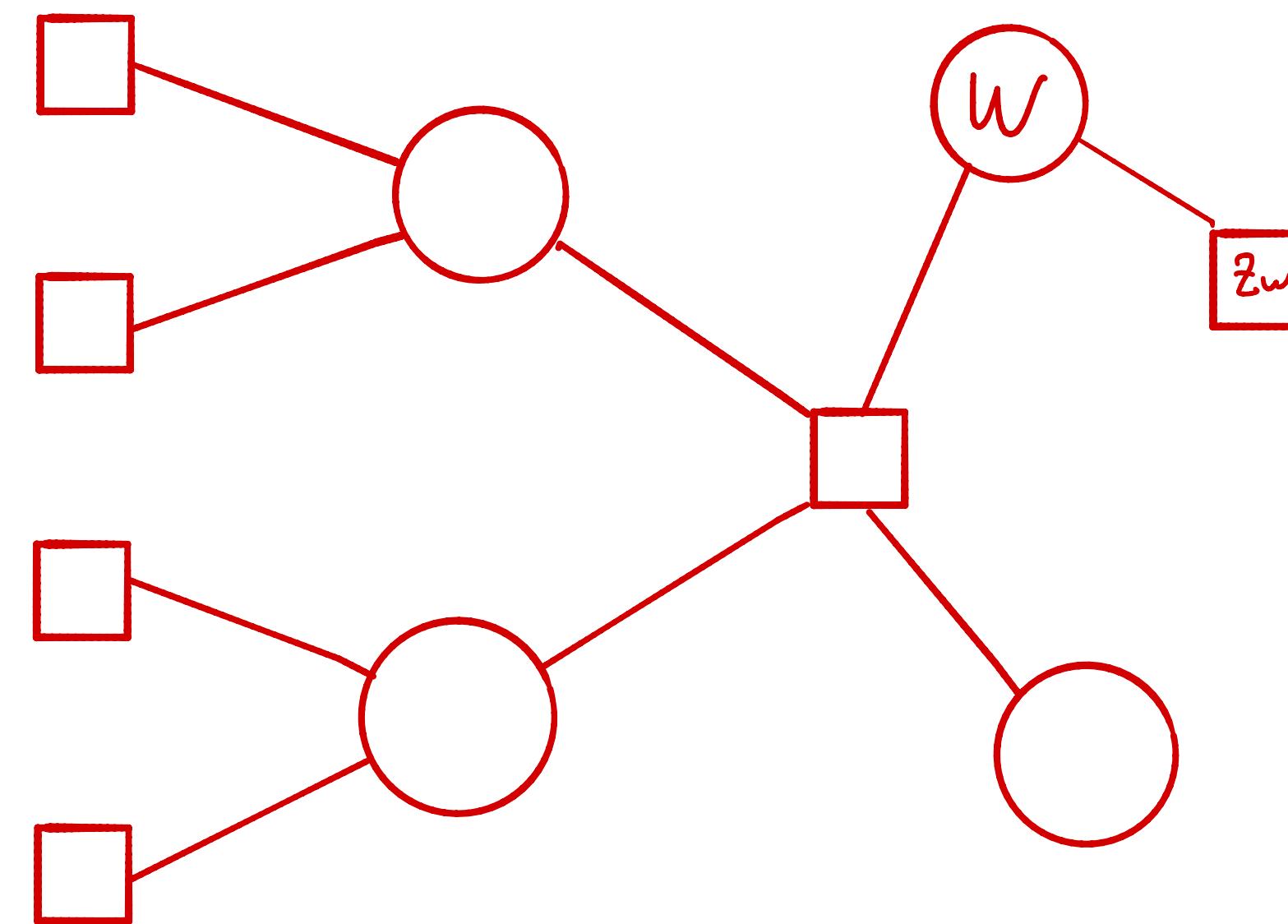
$$\bullet p(x_v) = \frac{1}{Z_v} \prod_{\alpha \in \partial(v)} \mu_{\alpha \rightarrow v}(x_v) \quad \text{with } Z_v := \sum_{x_v} \prod_{\alpha \in \partial(v)} \mu_{\alpha \rightarrow v}(x_v)$$
$$\bullet p(x_\alpha) = \frac{1}{Z_\alpha} \cdot \psi_\alpha(x_\alpha) \cdot \prod_{v \in \partial(\alpha)} \mu_{v \rightarrow \alpha}(x_v) \quad \text{with } Z_\alpha := \sum_{x_\alpha} \psi_\alpha(x_\alpha) \cdot \prod_{v \in \partial(\alpha)} \mu_{v \rightarrow \alpha}(x_v)$$

Conditioning via Evidence Factors

- $G = (V, F, E)$ Factor Tree with $p(x_V) = \frac{1}{Z} \prod_{\alpha \in F} \psi_\alpha(x_\alpha)$.
- Goal: Compute $p(x_v | x_C = z_C)$ for some fixed value z_C .
- For this introduce ‘evidence factors’: $\psi_w(x_w) := \mathbb{1}_{z_w}(x_w)$ for $w \in C$ into a new Factor graph:

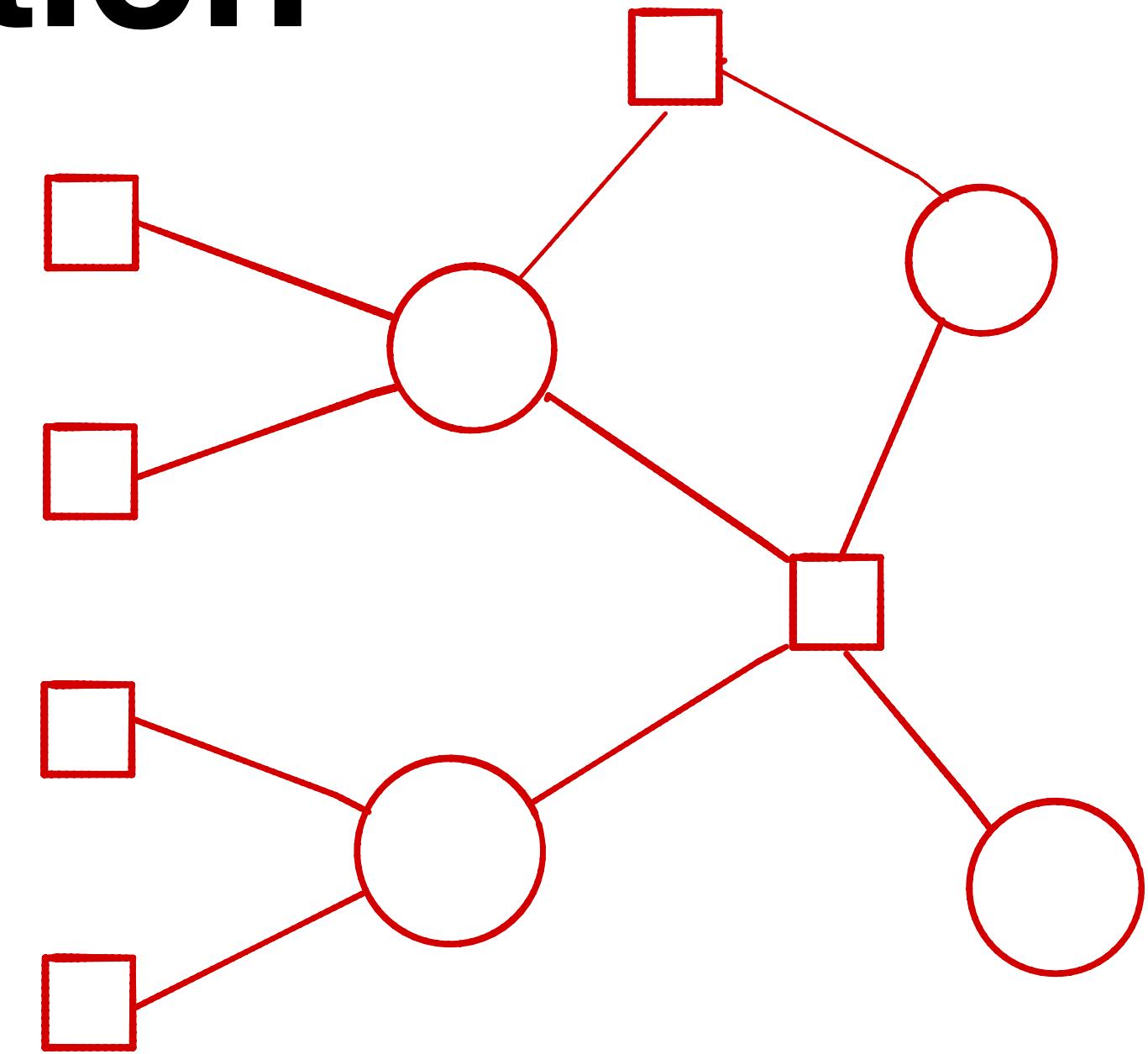
$$p(x_{V \setminus C} | z_C) \propto \prod_{\alpha \in F} \psi_\alpha(x_\alpha) \cdot \prod_{w \in C} \mathbb{1}_{z_w}(x_w)$$

- Run Sum-Product-Algorithm as before.
- Similarly for factor beliefs.



Loopy Belief Propagation

- $G = (V, F, E)$ Factor Graph (possibly non-tree),
$$p(x_V) = \frac{1}{Z} \prod_{\alpha \in F} \psi_\alpha(x_\alpha).$$
- Loopy Belief Propagation:
 - Initialize all messages to 1. Then until convergence do:
 - randomly draw node and use usual Sum-Product update rules.
 - Compute beliefs
 - Not many theoretical guarantees, not exact anymore (approximate inference)
 - In practice reported to converge quickly and approximation often good (but sometimes not).
 - Alternatively: Cluster nodes of a cycle together -> check Junction Tree Algorithm



Sum-Product Algorithm / Belief Propagation

- $G = (V, F, E)$ Factor Tree with $p(x_V) = \frac{1}{Z} \prod_{\alpha \in F} \psi_\alpha(x_\alpha)$. Pick a 'Root' inside the graph.

- Message Passing (once from Leaves to Root and back):

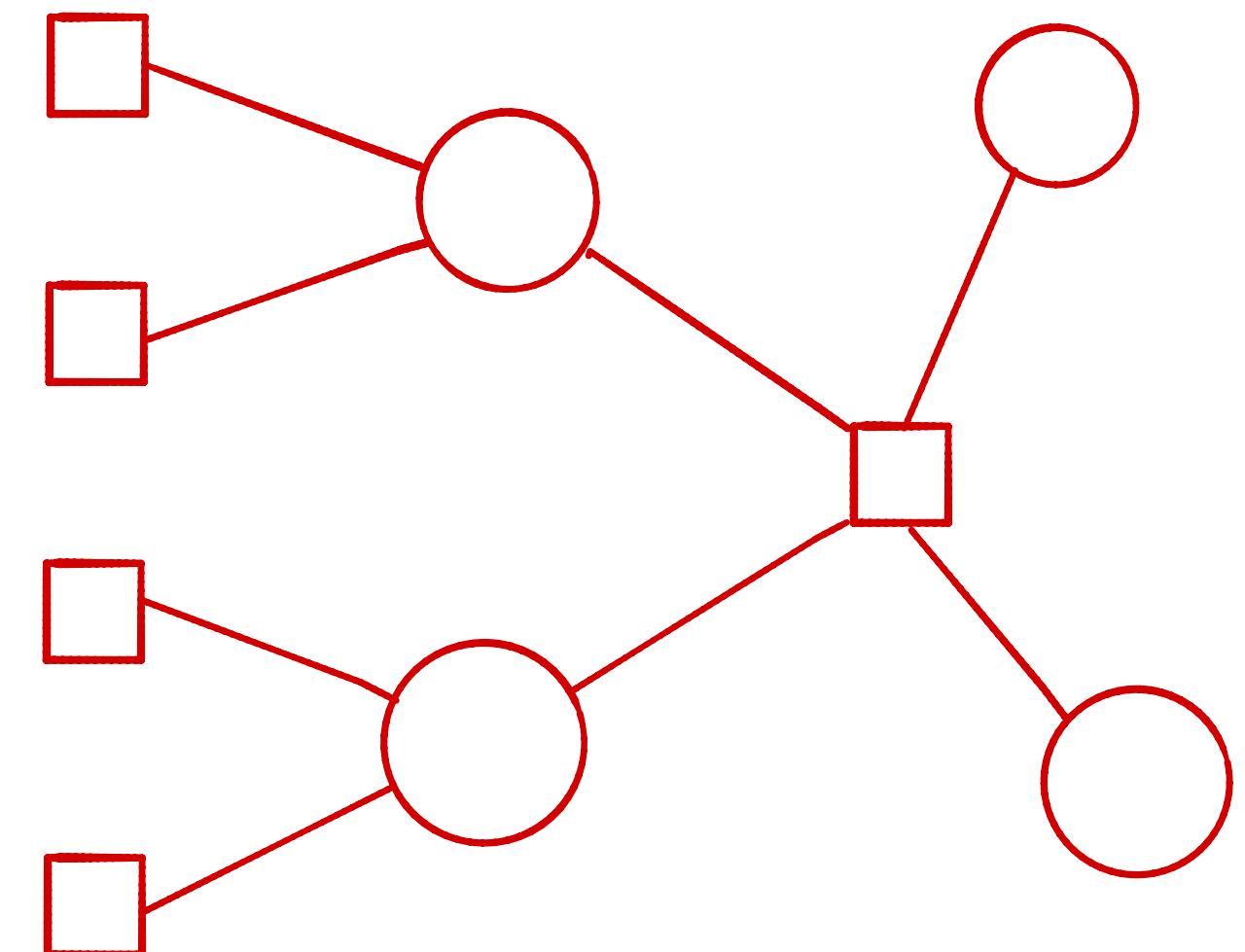
- Variable-to-factor: $\mu_{v \rightarrow \beta}(x_v) = \prod_{\gamma \in \partial(v) \setminus \{\beta\}} \mu_{\gamma \rightarrow v}(x_v)$ Leaf: $\mu_{u \rightarrow \beta}(x_u) = 1$

- Factor-to-variable: $\mu_{\beta \rightarrow w}(x_w) = \sum_{x_{\partial(\beta) \setminus \{w\}}} \psi_\beta(x_\beta) \prod_{u \in \partial(\beta) \setminus \{w\}} \mu_{u \rightarrow \beta}(x_u)$ Leaf: $\mu_{\alpha \rightarrow v}(x_v) = \psi_\alpha(x_v)$

- Compute 'Beliefs' / marginals via:

- $p(x_v) = \frac{1}{Z_v} \prod_{\alpha \in \partial(v)} \mu_{\alpha \rightarrow v}(x_v)$ with $Z_v := \sum_{x_v} \prod_{\alpha \in \partial(v)} \mu_{\alpha \rightarrow v}(x_v)$

- $p(x_\alpha) = \frac{1}{Z_\alpha} \cdot \psi_\alpha(x_\alpha) \cdot \prod_{v \in \partial(\alpha)} \mu_{v \rightarrow \alpha}(x_v)$ with $Z_\alpha := \sum_{x_\alpha} \psi_\alpha(x_\alpha) \cdot \prod_{v \in \partial(\alpha)} \mu_{v \rightarrow \alpha}(x_v)$



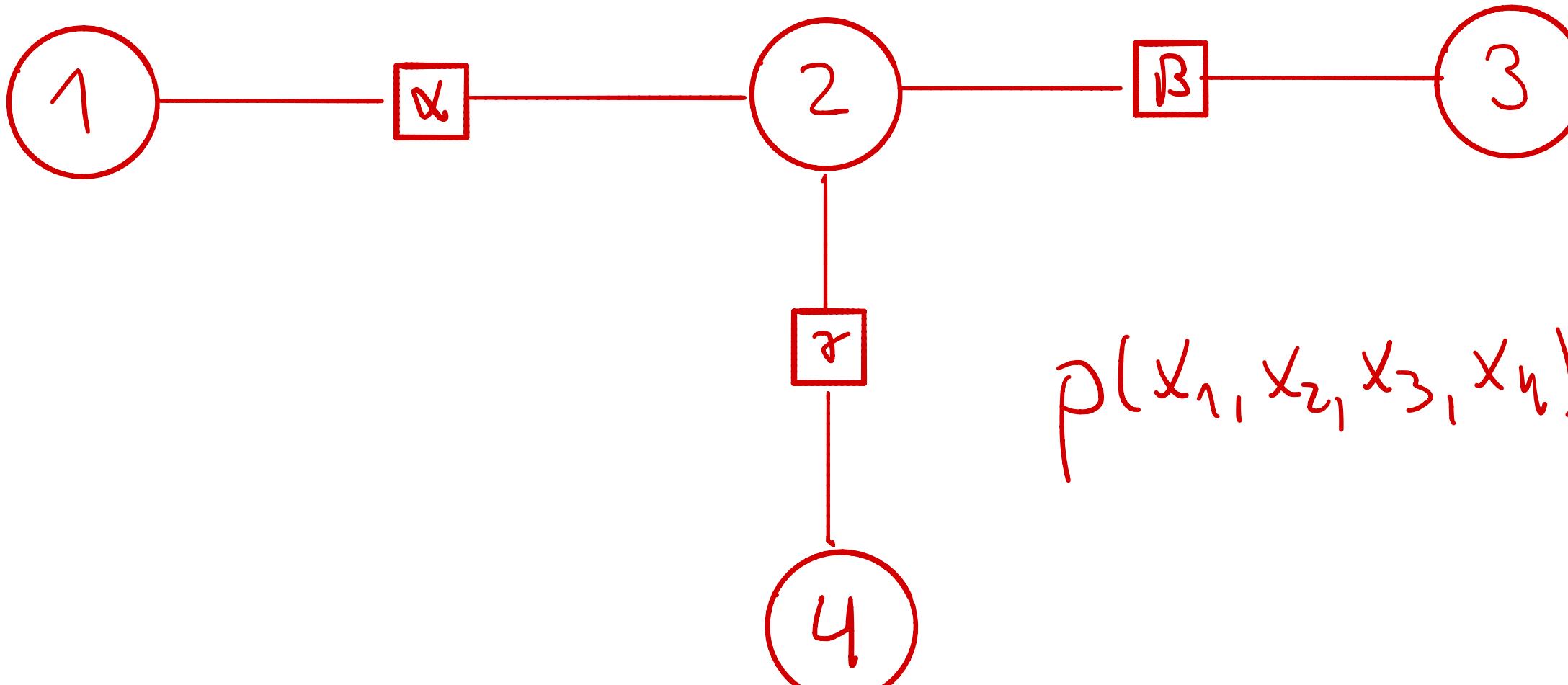
Machine Learning 2

Graphical Models

- Exact Inference**
- Sum-Product Algorithm**
- Example**

Patrick Forré

Example: Sum-Product Algorithm / Belief Propagation



$$p(x_1, x_2, x_3, x_4) = \frac{1}{Z} \cdot \psi_\alpha(x_1, x_2) \psi_\beta(x_2, x_3) \psi_\gamma(x_2, x_4)$$

$$\mu_{1 \rightarrow \alpha}(x_1) = 1$$

$$\mu_{\alpha \rightarrow 2}(x_2) = \sum_{x_1} \psi_\alpha(x_1, x_2) \cdot \overbrace{\mu_{1 \rightarrow \alpha}(x_1)}^{=} = \sum_{x_1} \psi_\alpha(x_1, x_2)$$

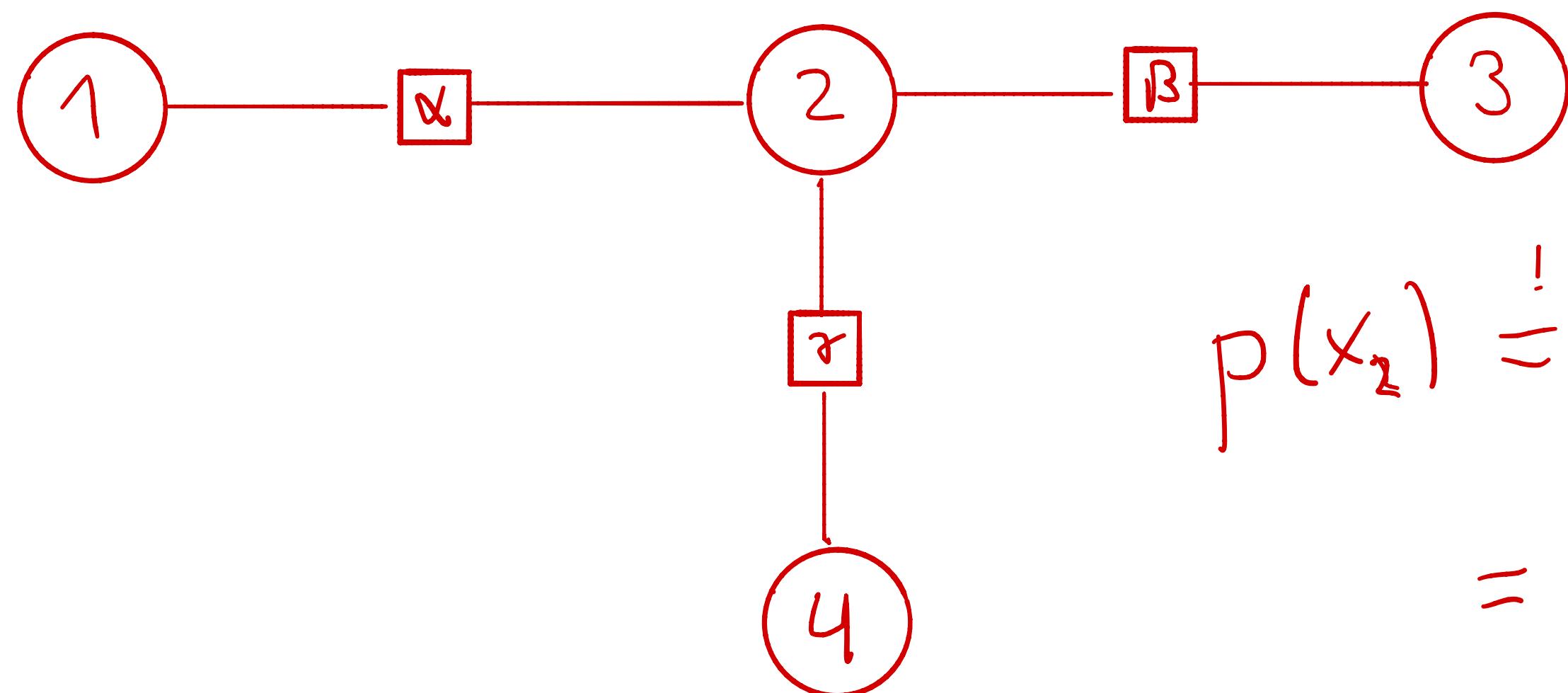
$$\mu_{4 \rightarrow 2}(x_2) = 1$$

$$\mu_{\gamma \rightarrow 2}(x_2) = \sum_{x_4} \psi_\gamma(x_2, x_4) \cdot \overbrace{\mu_{4 \rightarrow 2}(x_4)}^{=} = \sum_{x_4} \psi_\gamma(x_2, x_4)$$

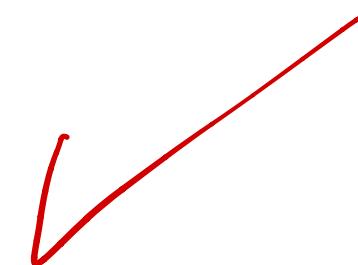
$$\mu_{3 \rightarrow \beta}(x_3) = 1$$

$$\mu_{\beta \rightarrow 2}(x_2) = \sum_{x_3} \psi_\beta(x_2, x_3) \cdot \overbrace{\mu_{3 \rightarrow \beta}(x_3)}^{=} = \sum_{x_3} \psi_\beta(x_2, x_3).$$

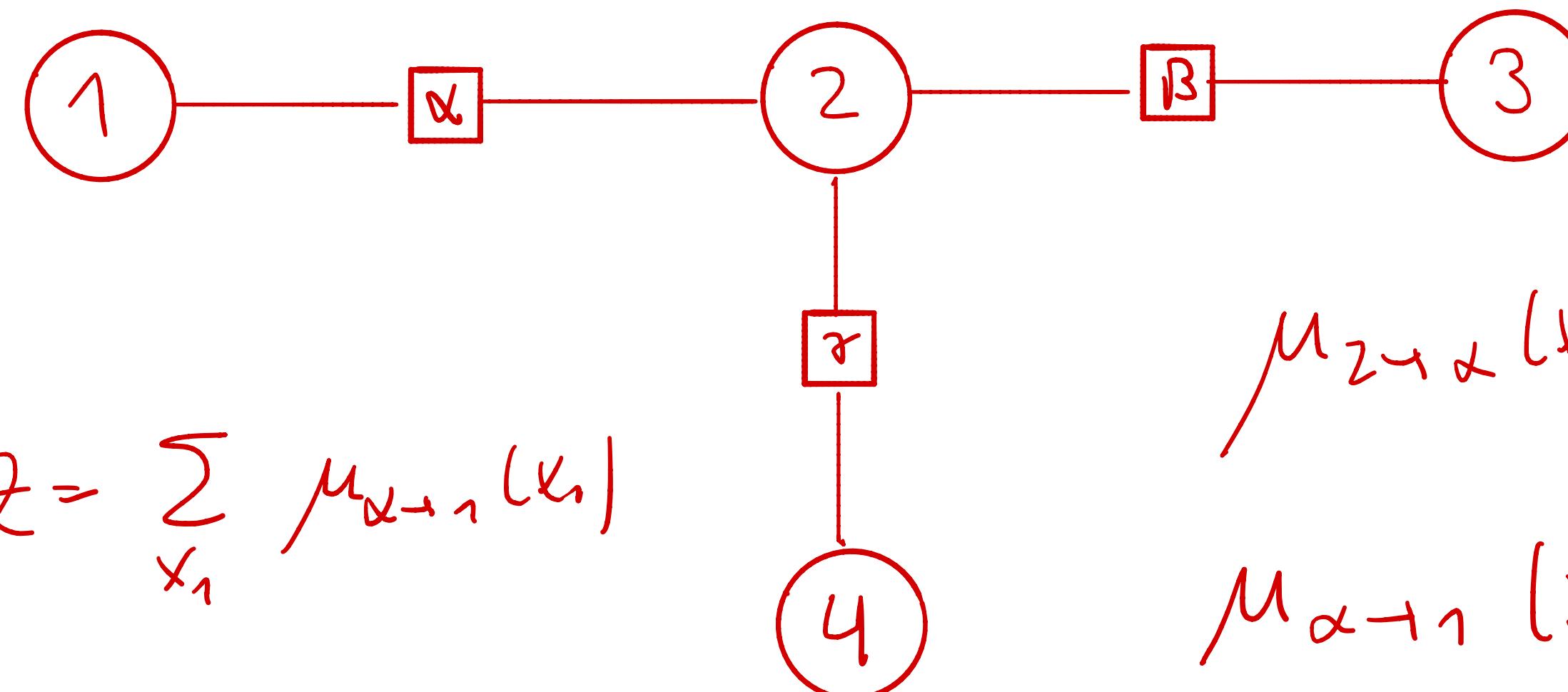
Example: Sum-Product Algorithm / Belief Propagation



$$\begin{aligned}
 p(x_2) &= \frac{1}{Z_2} \cdot \mu_{\alpha \rightarrow 2}(x_2) \cdot \mu_{\beta \rightarrow 2}(x_2) \cdot \mu_{\gamma \rightarrow 2}(x_2) \\
 &= \frac{1}{Z_2} \left(\sum_{x_1} \psi_\alpha(x_1, x_2) \right) \left(\sum_{x_3} \psi_\beta(x_2, x_3) \right) \left(\sum_{x_4} \psi_\gamma(x_2, x_4) \right) \\
 &= \frac{1}{Z_2} \sum_{x_1, x_3, x_4} \psi_\alpha(x_1, x_2) \cdot \psi_\beta(x_2, x_3) \cdot \psi_\gamma(x_2, x_4) \\
 &= \sum_{x_1, x_3, x_4} p(x_1, x_2, x_3, x_4)
 \end{aligned}$$



Example: Sum-Product Algorithm / Belief Propagation



$$Z = \sum_{x_1} \mu_{\alpha \rightarrow 1}(x_1)$$

$$\mu_{2 \rightarrow 2}(x_2) = \mu_{\gamma \rightarrow 2}(x_2) \cdot \mu_{\beta \rightarrow 2}(x_2)$$

$$\mu_{\alpha \rightarrow 1}(x_1) = \sum_{x_2} \psi_\alpha(x_1, x_2) \cdot \mu_{2 \rightarrow \alpha}(x_2)$$

$$P(x_1) = ? \quad \frac{1}{Z} \cdot \mu_{\alpha \rightarrow 1}(x_1)$$

$$= \frac{1}{Z} \sum_{x_2} \psi_\alpha(x_1, x_2) \cdot \mu_{2 \rightarrow \alpha}(x_2)$$

$$= \frac{1}{Z} \sum_{x_2} \psi_\alpha(x_1, x_2) \cdot \left(\sum_{x_n} \psi_\gamma(x_2, x_n) \right) \left(\sum_{x_3} \psi_\beta(x_2, x_3) \right)$$

$$= \frac{1}{Z} \sum_{x_2, x_3, x_4} \psi_\alpha(x_1, x_2) \psi_\beta(x_2, x_3) \cdot \psi_\gamma(x_2, x_4) = \sum_{x_2, x_3, x_4} P(x_1, x_2, x_3, x_4) \quad \checkmark$$

Machine Learning 2

Graphical Models

- Exact Inference**
- Sum-Product Algorithm**
- Summary**

Patrick Forré

Sum-Product Algorithm / Belief Propagation

- $G = (V, F, E)$ Factor Tree with $p(x_V) = \frac{1}{Z} \prod_{\alpha \in F} \psi_\alpha(x_\alpha)$. Pick a 'Root' inside the graph.

- Message Passing (once from Leaves to Root and back):

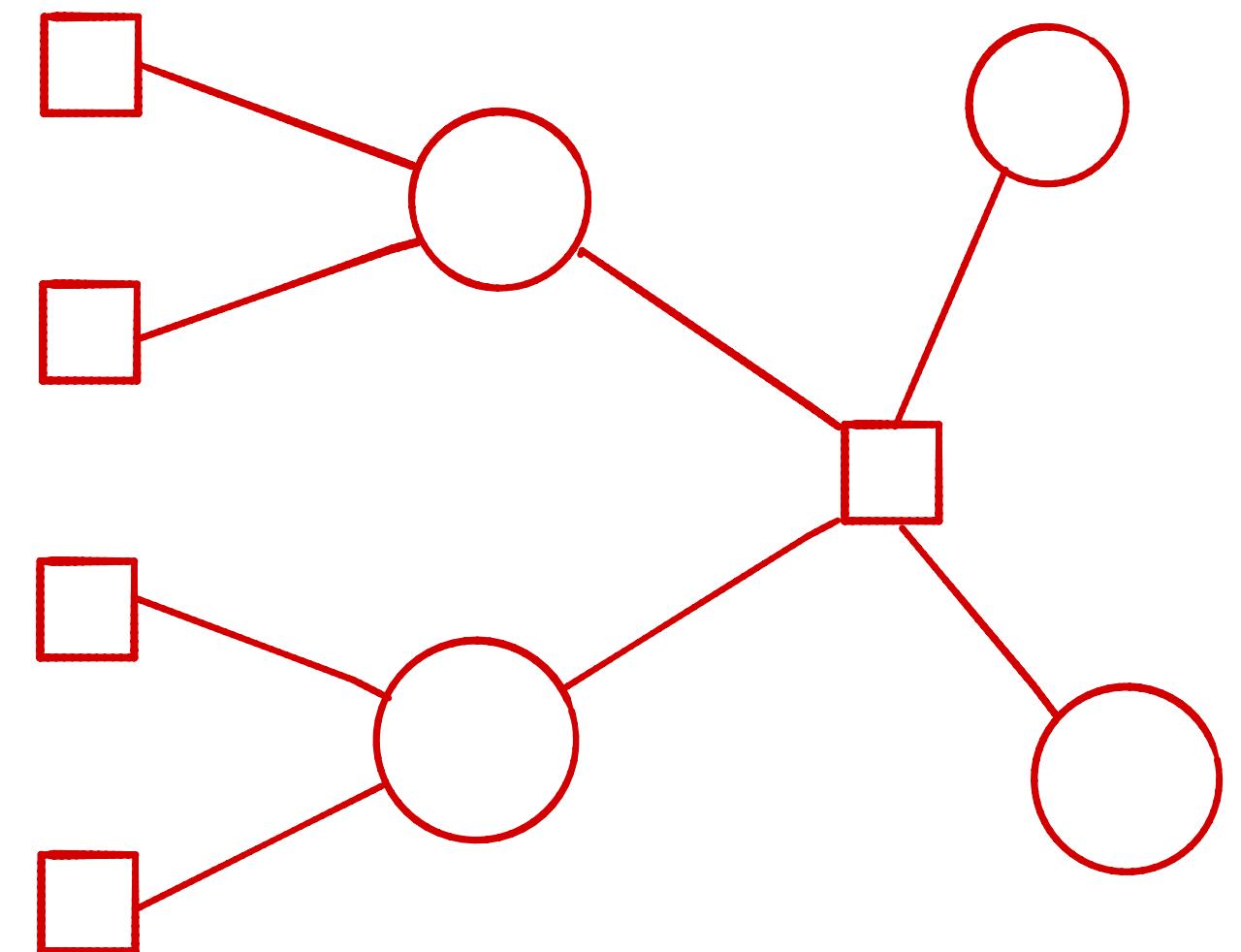
- Variable-to-factor: $\mu_{v \rightarrow \beta}(x_v) = \prod_{\gamma \in \partial(v) \setminus \{\beta\}} \mu_{\gamma \rightarrow v}(x_v)$ Leaf: $\mu_{u \rightarrow \beta}(x_u) = 1$

- Factor-to-variable: $\mu_{\beta \rightarrow w}(x_w) = \sum_{x_{\partial(\beta) \setminus \{w\}}} \psi_\beta(x_\beta) \prod_{u \in \partial(\beta) \setminus \{w\}} \mu_{u \rightarrow \beta}(x_u)$ Leaf: $\mu_{\alpha \rightarrow v}(x_v) = \psi_\alpha(x_v)$

- Compute 'Beliefs' / marginals via:

- $p(x_v) = \frac{1}{Z_v} \prod_{\alpha \in \partial(v)} \mu_{\alpha \rightarrow v}(x_v)$ with $Z_v := \sum_{x_v} \prod_{\alpha \in \partial(v)} \mu_{\alpha \rightarrow v}(x_v)$

- $p(x_\alpha) = \frac{1}{Z_\alpha} \cdot \psi_\alpha(x_\alpha) \cdot \prod_{v \in \partial(\alpha)} \mu_{v \rightarrow \alpha}(x_v)$ with $Z_\alpha := \sum_{x_\alpha} \psi_\alpha(x_\alpha) \cdot \prod_{v \in \partial(\alpha)} \mu_{v \rightarrow \alpha}(x_v)$



Remarks

- For computing **conditionals** $p(x_v | x_C = z_C)$ use **Evidence Factors**:

$$\psi_w(x_w) := \mathbf{1}_{z_w}(x_w).$$

- For **non-tree-shaped** factor graphs consider using **Loopy Belief Propagation**
 - Initialize messages to 1,
 - randomly draw node and update with usual Sum-Product update rules.
- Alternatively, use Junction Tree Algorithm to avoid cycles.