

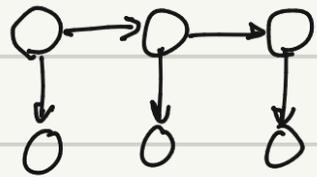
PRML Chapter 8 Graphical Model

Merit

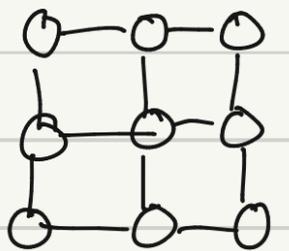
- 1) provide a simple way to visualize a probabilistic model
- 2) Insights into the properties of the model
- 3) complex computations can be expressed in terms of graphical manipulations.

Graphical Model

→ Bayesian Network.
useful for expressing causal relation between random variables



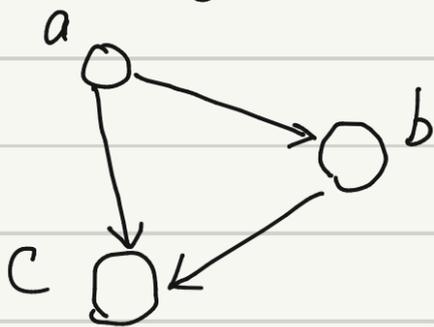
→ Markov Random Field
useful for expressing soft constraints between random variables.



8.1 Bayesian Networks

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

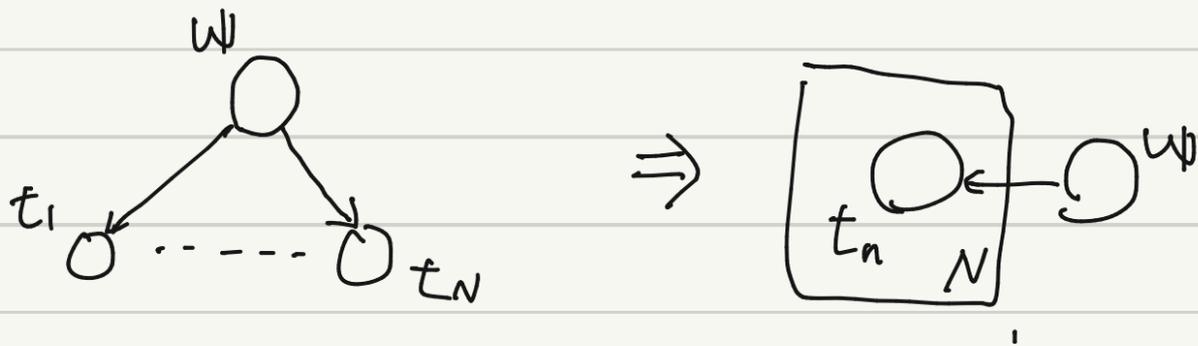
generalize



Absolutely, decomposition is not unique

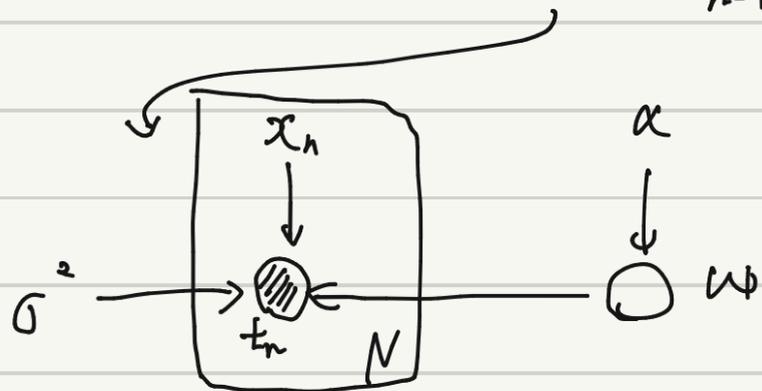
$$p(x_1, \dots, x_k) = p(x_k | x_1, \dots, x_{k-1}) \dots p(x_2 | x_1) p(x_1)$$

8.1.1 Polynomial regression



if we want stochastic variables to be explicit.

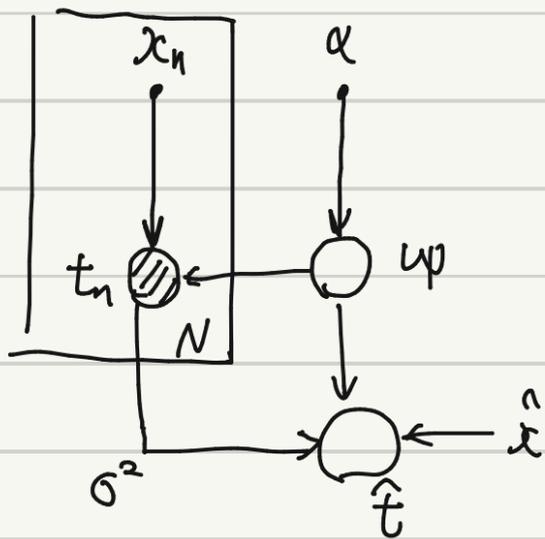
$$p(\mathbf{t}, \mathbf{w} | \mathbf{x}, \alpha, \sigma^2) = p(\mathbf{w} | \alpha) \prod_{n=1}^N p(t_n | w, x_n, \sigma^2)$$



⊗: represents observed values.

purpose: predict \hat{t}

$$p(\hat{t}, \mathbf{t}, \mathbf{w} | \hat{\mathbf{x}}, \mathbf{x}, \alpha, \sigma^2) = \left[\prod_{n=1}^N p(t_n | x_n, w, \sigma^2) \right] p(\mathbf{w} | \alpha) p(\hat{t} | \hat{\mathbf{x}}, w, \sigma^2)$$



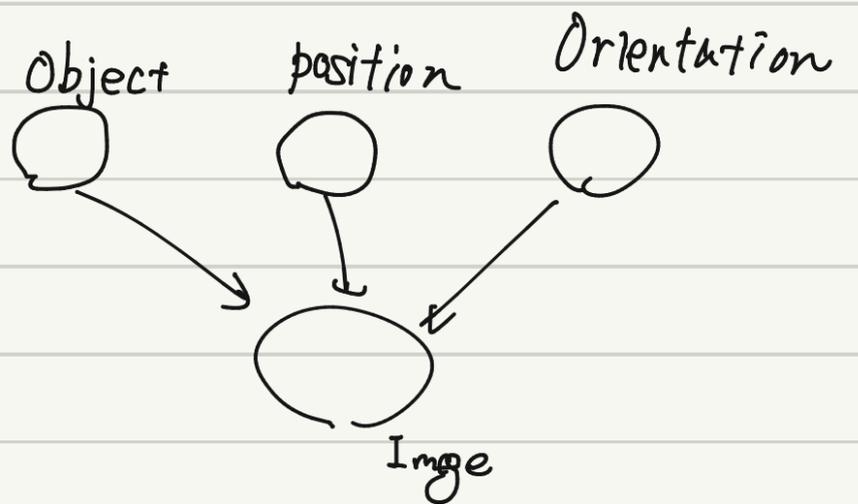
8.1.2 Generative models

Sometimes, we want to draw samples from given probability distribution

↳ (discuss later in Chapter 11)

• Ancestral sampling

ex) images object creation



• The graphical model : represents the causal process

↓
called "generative models"

8.1.3 Discrete variables

• We want to choose models so that parents and childs are conjugate.

• probability distribution for a single discrete variable x having K possible states

$$p(x | \mu) = \prod_{k=1}^K \mu_k^{x_k}$$

(governed by μ) ($\sum_k \mu_k = 1$)

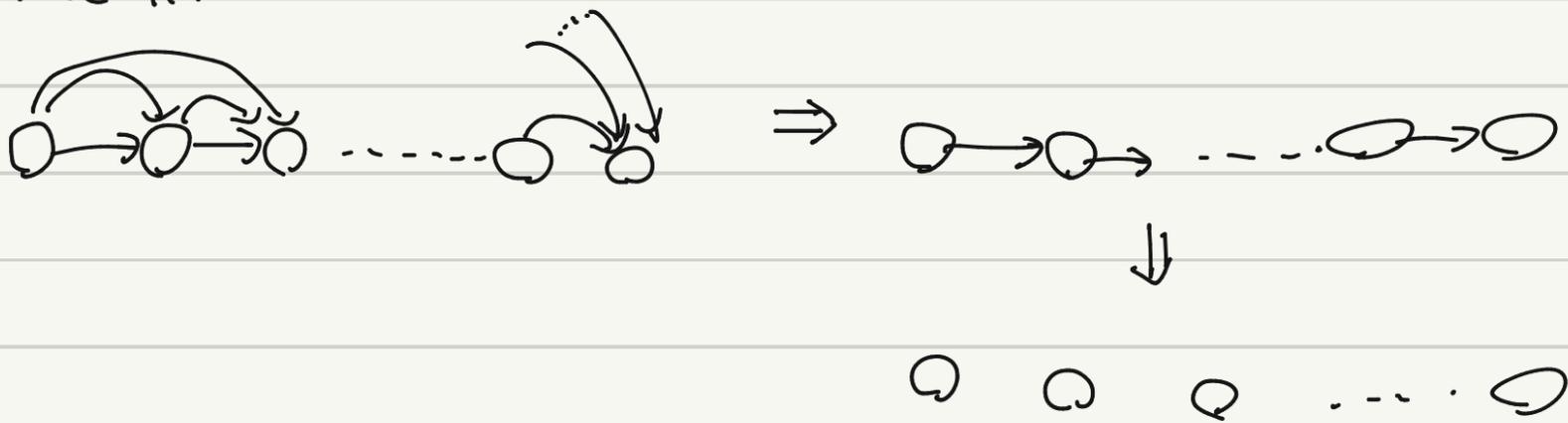
• Suppose that we have M discrete variable, each of which has k states

parameters = $k^M - 1$

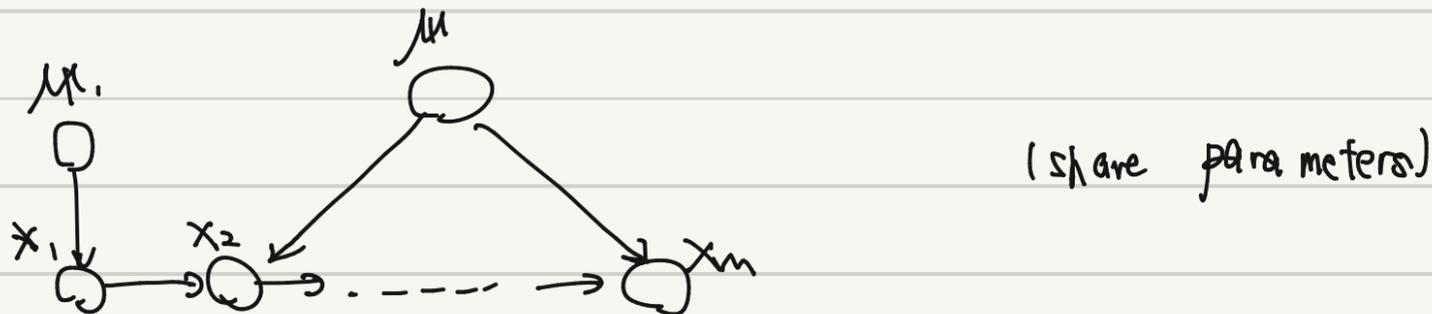
if they were independent, $M(k-1)$ want to reduce.

• How to reduce the number of independent parameters

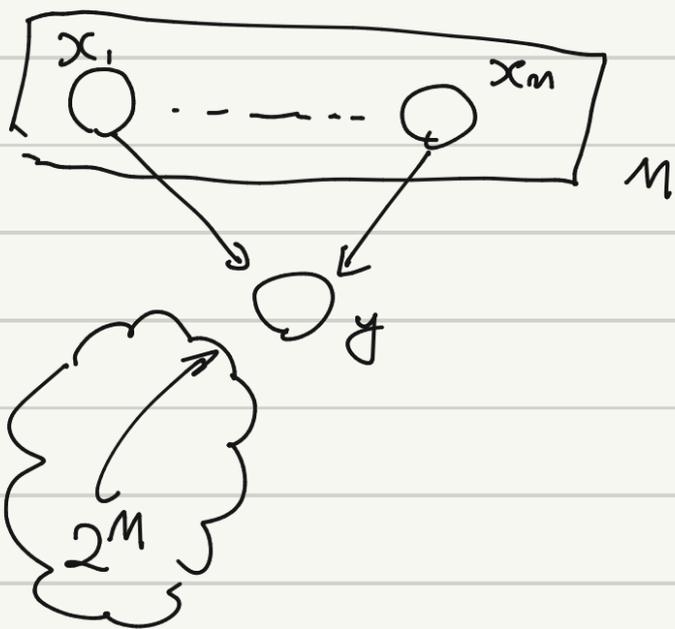
1. delete link



2. Sharing



3. use parameterized models



$$\begin{aligned}
 & 2^M \\
 & \downarrow \\
 & p(y|x_1, \dots, x_m) = \sigma \left(w_0 + \sum_{i=1}^m w_i x_i \right) \\
 & = \sigma(w^T x) \\
 & M+1 \\
 & (w_{\frac{y}{x}}; s(M+1) \text{ dimensional vector})
 \end{aligned}$$

8.1.4 Linear-Gaussian models

- how a multivariate Gaussian can be expressed as a graphical model which corresponds to a linear-Gaussian model

$$p(x_i | pa_i) = \mathcal{N} \left(x_i \mid \sum_{j \in pa_i} w_{ij} x_j + b_i, \mathcal{U}_i \right)$$

$pa_i =$ parent node

$$\ln p(\mathcal{X}) = \sum_{i=1}^D \ln p(x_i | pa_i)$$

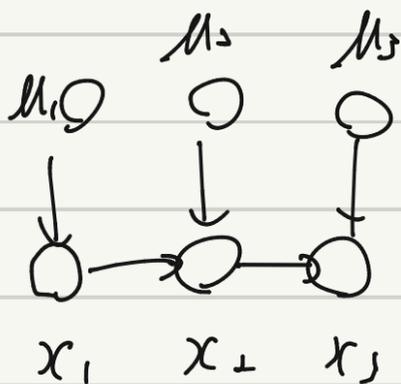
$$= - \sum_i \frac{1}{2\mathcal{U}_i} \left(x_i - \sum_{j \in pa_i} w_{ij} x_j - b_i \right)^2 + C$$

↑
quadratic function

↓ $p(x)$ is a multivariate Gaussian.

★

$$\begin{cases} E[x_i] = \sum_{j \in pa_i} w_{ij} E[x_j] + b_i \\ \text{cov}[x_i, x_j] = \sum_{k \in pa_i} w_{jk} \text{cov}[x_i, x_k] + I_{ij} \mathcal{U}_j \end{cases}$$



8.2 Conditional Independence

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

a is conditionally independent of b given c

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

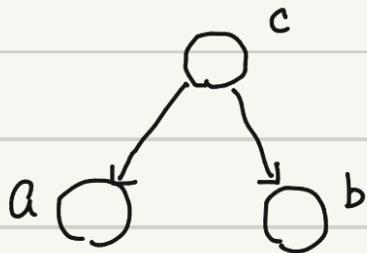
$$a \perp\!\!\!\perp b | c \quad p(a | b, c) = p(a | c)$$

• How to examine independency?

↳ judge from the shape of graph.

8.2.1 Three example graphs

1. tail-to-tail



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

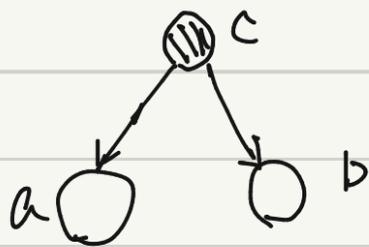
↓ marginalize with respect to c

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

↳ can't factorize into $p(a) \cdot p(b)$

↳ $a \not\perp\!\!\!\perp b | \emptyset$

dependence



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

↓ condition on the variable c

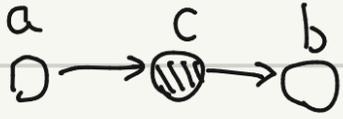
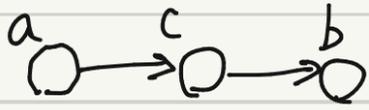
$$p(a, b | c) = \frac{p(a, b, c)}{p(c)} = p(a | c) p(b | c)$$

↳ obtain $a \perp\!\!\!\perp b | c$

conditional independence

(if c is observed...)

2. head-to-tail



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

↓ marginalize

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

↳ $a \perp\!\!\!\perp b | \emptyset$

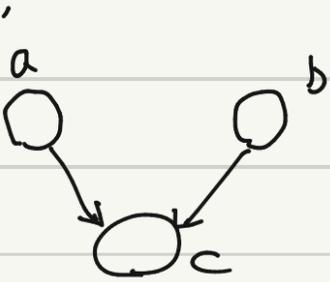
↓ condition on the variable c

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

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

↳ $a \perp\!\!\!\perp b | c$

3. head-to-head



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

↓ marginalize

$$p(a, b) = p(a) p(b) \rightarrow a \perp\!\!\!\perp b | \emptyset$$

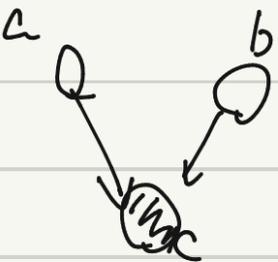
independent

↓ condition on the variable c

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

↳ $a \perp\!\!\!\perp b | c$

* explaining way

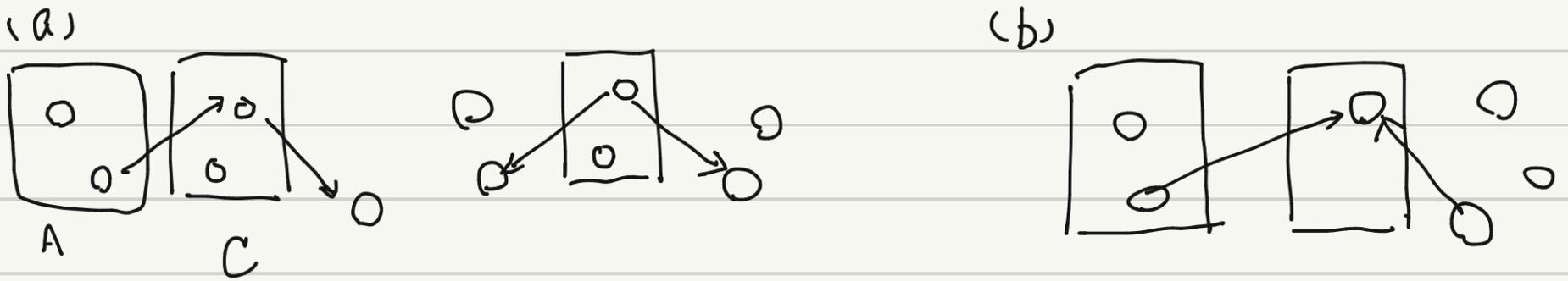


if b is observed, probability of a is changed

8.2.2 d -separation

Any path is said to be blocked if it includes a node such that

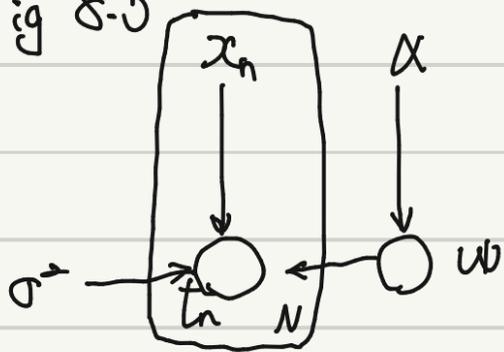
- (a) the arrows on the path meet either head-to-tail or tail-to-tail at the node, and node is in C
- (b) the arrows meet head-to-head at the node, and neither the node, nor any of its descendants, is in C



If All paths are blocked, A is said to be d -separated from B by C.

$$A \perp\!\!\!\perp B \mid C$$

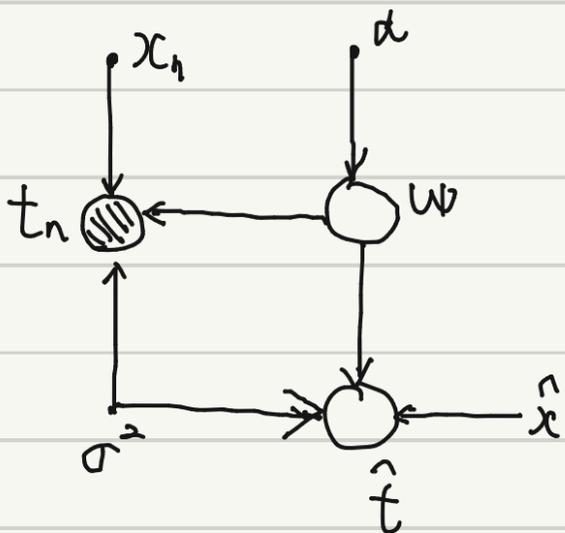
fig 8.5



α and σ^2 behave as observed nodes.
are tail to tail.

↳ play no role in d -separation.

Fig 8.7

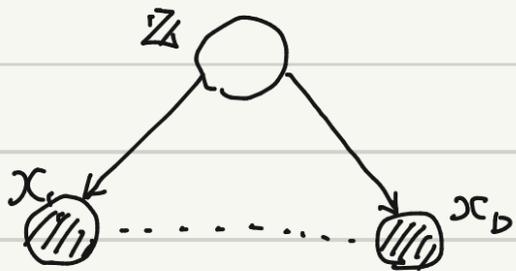


w : tail-to-tail

$$\hat{t} \perp\!\!\!\perp t_n \mid w$$

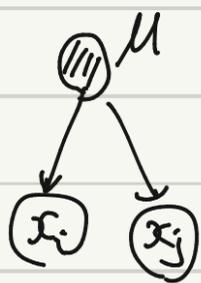
↳ conditioned on the w , \hat{t} is independent of t_n

Naive Bayes



conditioned on Z , x are assumed to be independent.

The graph - represents a specific decomposition of a joint probability distribution into a product of conditional probabilities
 - expresses a set of conditional independence statements obtained through d-separation criterion.

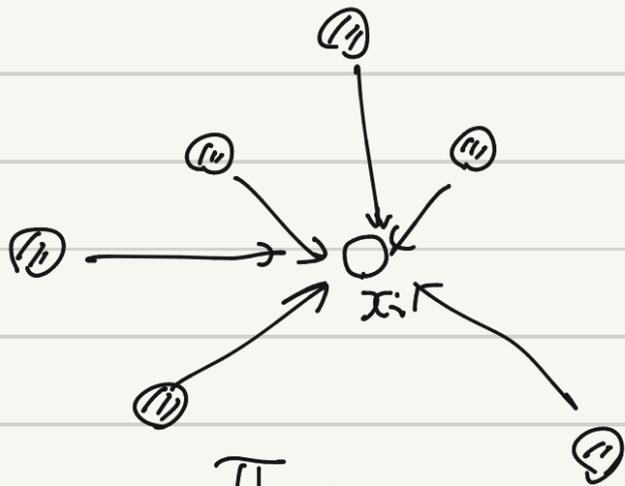


- μ is observed (we already knew μ (probabilistic distribution))

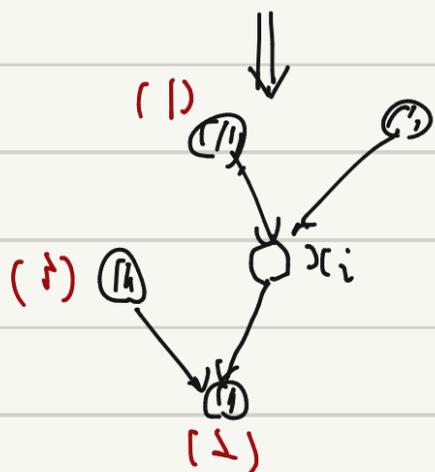
x_i and x_j are i.i.d (same probabilistic distribution but independent)

Markov blanket

consider the situation where every node except x_i was observed



$$p(x_i | x_{j \neq i}) = \frac{\prod_k p(x_k | p_{a_k})}{\int \prod_k p(x_k | p_{a_k}) dx_i}$$

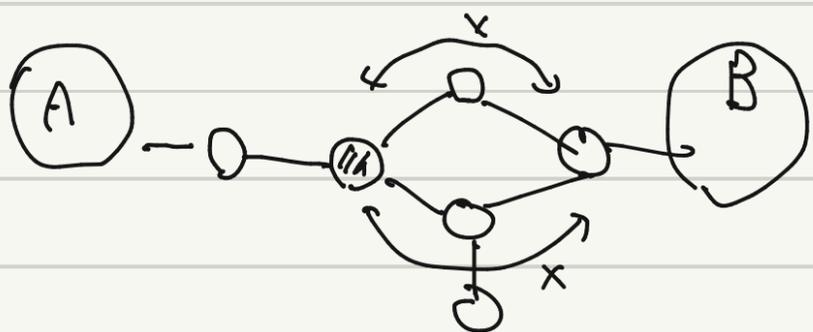


$(1) \sim (3)$ are depend on x_i

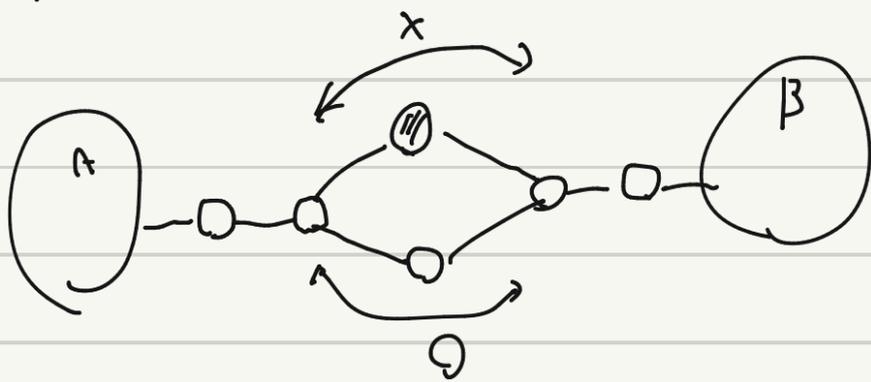
$\{(1)\} + \{(2)\} + \{(3)\} + \{x_i\} = \text{Markov blanket}$

8.3 Markov Random Fields

8.3.1 Conditional independence properties

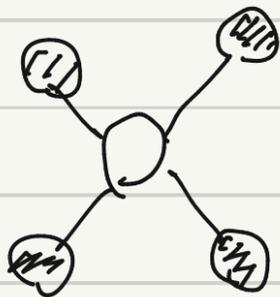


$A \perp\!\!\!\perp B \mid C$



$A \not\perp\!\!\!\perp B \mid C$

Markov blanket



8.3.2 Factorization properties

How to express the joint distribution $p(x)$ as a product of functions defined over sets of variables

Conditional independence properties can be expressed as
$$p(x_i, x_j \mid X_{\setminus \{i, j\}}) = p(x_i \mid X_{\setminus \{i, j\}}) \cdot p(x_j \mid X_{\setminus \{i, j\}})$$

let x_i and x_j not appear in the same factor

clique

a subset of the nodes in a graph such that there exists a link between all pairs of nodes in the subset.

Let C be a clique.

$$p(x) = \frac{1}{Z} \prod_C \psi_C(x_C)$$

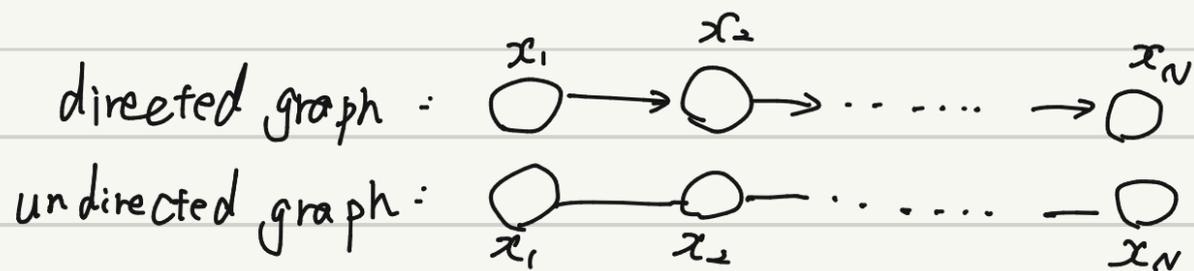
↖ hard to calculate

How can we define ψ_C ?

↳ viewing potential function as expressing which configurations of the local variables are preferred to others.

$$\psi_C(x_C) = \exp\{-E(x_C)\} \quad (\text{cause } \forall x_C, \psi_C(x_C) > 0)$$

8.3.4 Relation to directed graph



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

$$p(x) = \sum \psi_{1,2}(x_1, x_2) \psi_{2,3}(x_2, x_3) \dots \psi_{N-1,N}(x_{N-1}, x_N)$$

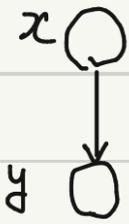


How to convert directed graph into undirected graph.

- 1) add additional undirected links between all pairs of parents
- 2) drop the arrows
- 3) initialize all of the clique potentials to 1
- 4) multiply each conditional distribution factor into one of the clique potentials.

8.4 Inference in Graphical Models

Bayes' theorem



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



y is observed



direction change

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

8.4.1 Inference on a chain

Consider chain of nodes, which will lay the foundation for general graphs.

$$p(x) = \frac{1}{Z} \psi_{1,2}(x_1, x_2) \psi_{2,3}(x_2, x_3) \dots \psi_{N-1,N}(x_{N-1}, x_N) \dots \textcircled{1}$$

Consider the inference problem of finding marginal distribution $p(x_n)$

$$p(x_n) = \sum_{x_1} \dots \sum_{x_{n-1}} \sum_{x_{n+1}} \dots \sum_{x_N} p(x) \dots \textcircled{2}$$

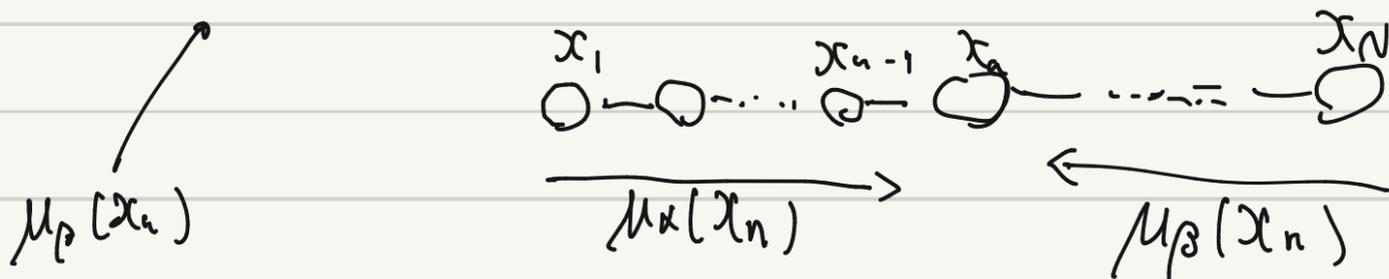
↳ naive implementation : hard to do

↳ graphical model gives us efficient algorithm.

from $\textcircled{1}$ and $\textcircled{2}$,

$$p(x_n) = \frac{1}{Z} \left[\sum_{x_{n+1}} \psi_{n,n+1}(x_n, x_{n+1}) \dots \left[\sum_{x_2} \psi_{2,3}(x_2, x_3) \left[\sum_{x_1} \psi_{1,2}(x_1, x_2) \right] \dots \right] \right]$$

$$\mu_\alpha(x_n) \left[\sum_{x_{n+1}} \psi_{n,n+1}(x_n, x_{n+1}) \dots \left[\sum_{x_N} \psi_{N-1,N}(x_{N-1}, x_N) \right] \dots \right]$$



interpret $\mu_\alpha(x_n)$ as a message x_{n-1} to x_n .

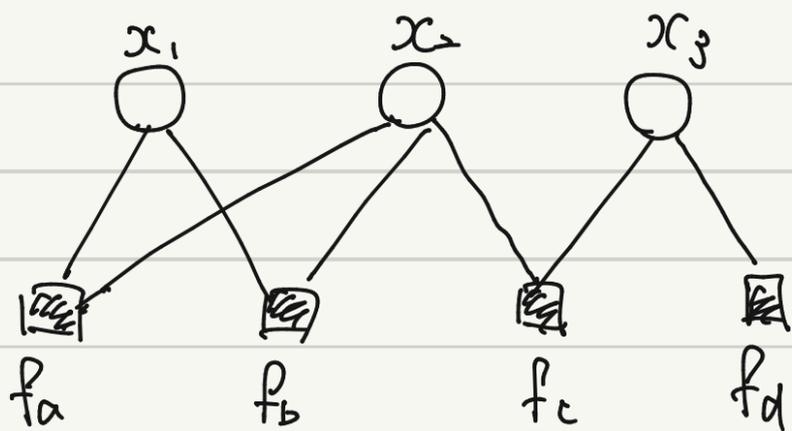
$$\mu_\alpha(x_n) = \sum_{x_{n-1}} \psi_{n-1,n}(x_{n-1}, x_n) \cdot \mu_\alpha(x_{n-1})$$

↑ inductive!! → O(k)

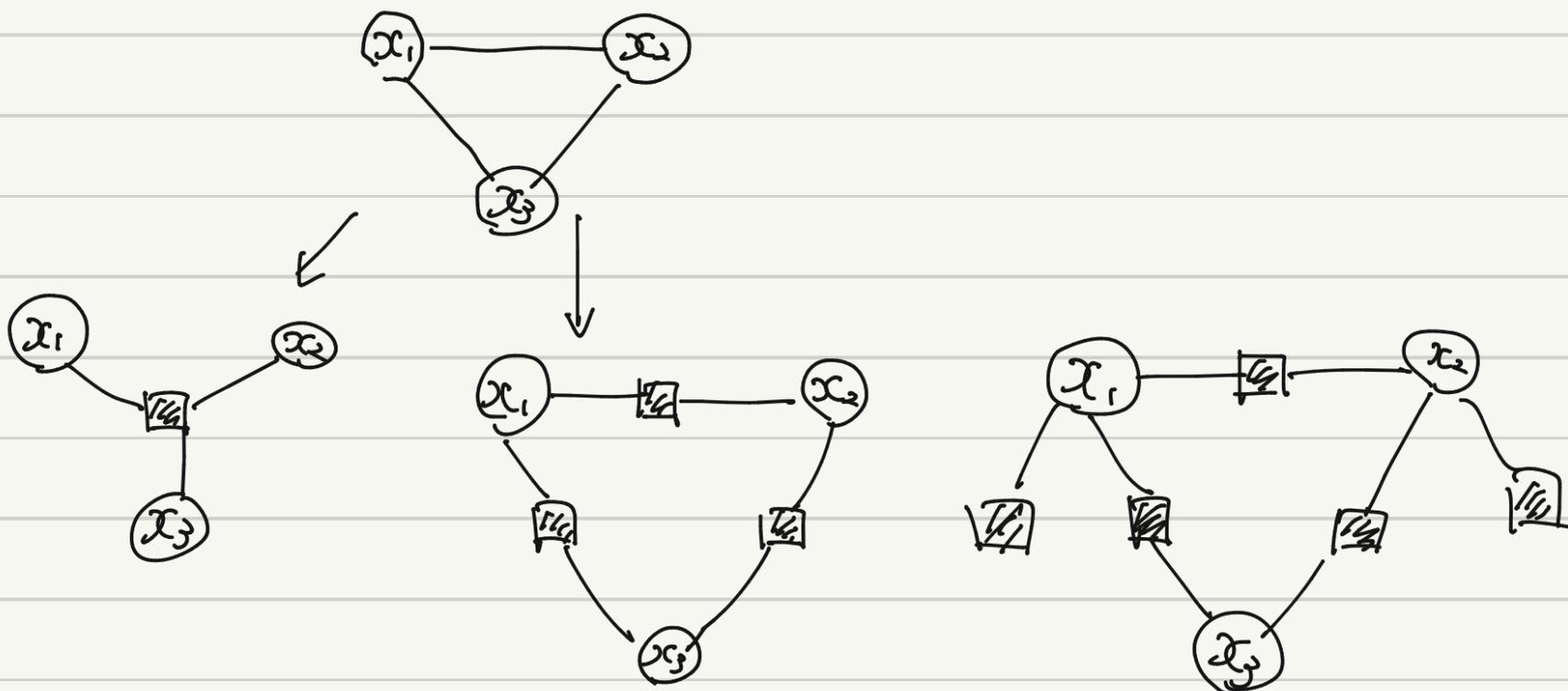
message passing can also be applied to tree structures.

8.4.3 Factor graphs

$$p(x) = f_a(x_1, x_2) f_b(x_1, x_2) f_c(x_2, x_3) f_d(x_3)$$



• factor graph is not unique



• but there is an algorithm that convert undirected graph into factor graph.

8.4.4 The sum-product algorithm

Assume that original graph is an undirected tree or directed tree or polytree.

corresponding factor graph has tree structure.

Goal: (1) to obtain an efficient, exact inference algorithm for finding marginals

(2) in situations where several marginals are required to allow computations to be shared efficiently

• How to evaluate marginal $p(x)$

1. view the variable node x as the root of the factor graph.

2. initiate message at the leaves of graph using

$$\mu_{x \rightarrow f(x)} = 1, \quad \mu_{f \rightarrow x}(x) = f(x)$$

3. the messages are passed through

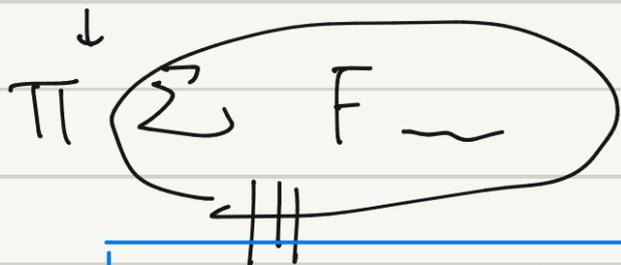
$$\mu_{x_m \rightarrow f_s}(x_m) = \prod_{d \in \text{ne}(x_m) \setminus f_s} \mu_{f_d \rightarrow x_m}(x_m)$$

$$\mu_{f_s \rightarrow x}(x) = \sum_{x_1} \dots \sum_{x_m} f_s(x, x_1, \dots, x_m) \prod_{m \in \text{ne}(f_s) \setminus x} \mu_{x_m \rightarrow f_s}(x_m)$$

applied recursively until messages have been propagated along every link, and the root node has received messages from all of their neighbors.

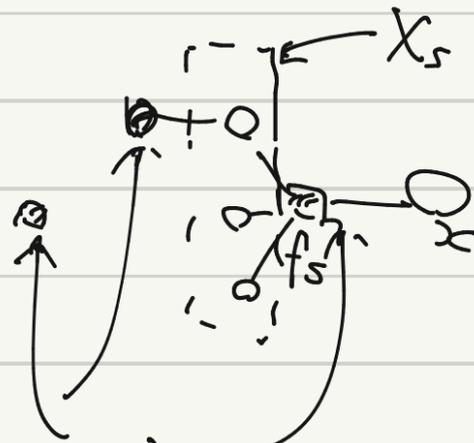
4. $p(x) = \prod_{s \in \text{ne}(x)} \mu_{f_s \rightarrow x}(x)$

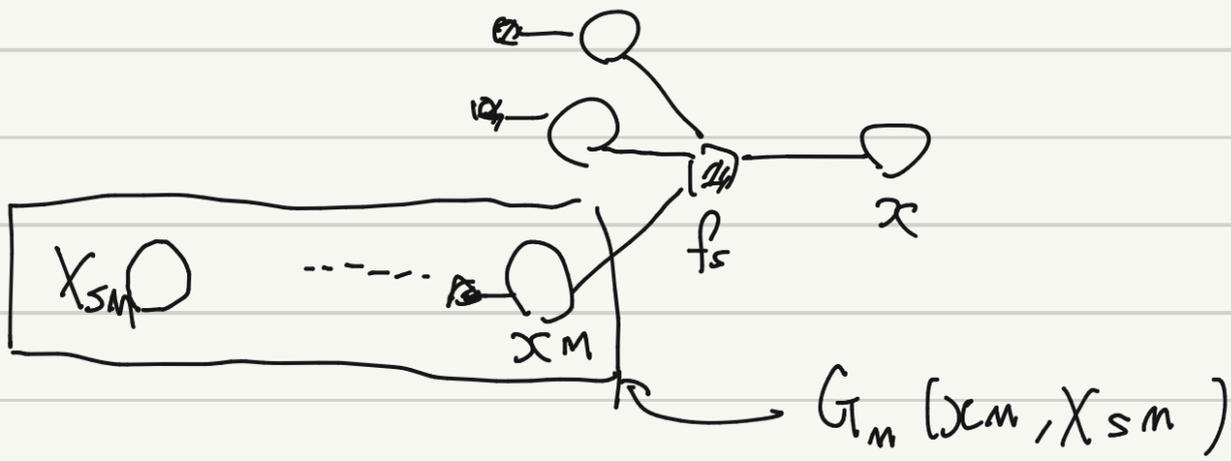
$$p(x_n) = \sum_{x \setminus x_n} \prod_{s \in \text{ne}(x)} F_s(x, X_s)$$



$$\mu_{f_s \rightarrow x}(x) \equiv \sum_{X_s} F_s(x, X_s)$$

F_s : multiply all factor that has relation with f_s





1.

< base case > | < step case >

$$\mu_{x \rightarrow f}(x) = 1$$

$$\mu_{f_s \rightarrow x}(x)$$

$$\Rightarrow F_s(x, X_s)$$

$$\mu_{f \rightarrow x}(x) = f(x)$$

$$\mu_{x_m \rightarrow f_s}(x_m)$$

$$= f_s(x, x_1, \dots, x_m) G_1(x_1, X_s)$$

$$\dots \dots G_m(x_m, X_{sm})$$

3.

$$P(X) = \sum_{\text{state}(x)} F_s(x, X_s)$$

$$P(x) = \sum_{X \setminus x} P(X)$$

8.4.5 max-sum algorithm

in order $\left\{ \begin{array}{l} - \text{ to find a setting of variables that has the largest prob} \\ - \text{ to find the value of that probability.} \end{array} \right.$

$$\max_{\mathbf{x}} p(\mathbf{x}) = \max_{x_1} \dots \max_{x_n} p(\mathbf{x})$$

consider the chain nodes.

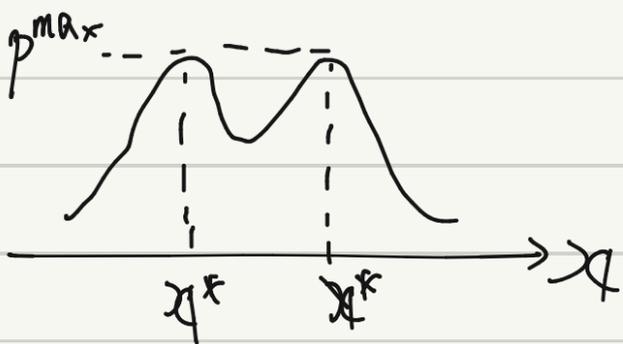
$$\max_{\mathbf{x}} p(\mathbf{x}) = \frac{1}{Z} \max_{x_1} \dots \max_{x_n} [\psi_{1,2}(x_1, x_2) \dots \psi_{N-1,N}(x_{N-1}, x_N)]$$

$$= \frac{1}{Z} \max_{x_1} \left[\max_{x_2} \left[\psi_{1,2}(x_1, x_2) \left[\dots \max_{x_{N-1}} \psi_{N-1,N}(x_{N-1}, x_N) \right] \dots \right] \right]$$

→ take same process as 8.4.4

$$p^{\max} = \max_x \left[\sum_{S \in \mathcal{H}(x)} \mu_{S \rightarrow x}(x) \right]$$

$$x^{\max} = \operatorname{argmax}_x \left[\sum_{S \in \mathcal{H}(x)} \mu_{S \rightarrow x}(x) \right]$$



p^{\max} is max value of $p(x)$

$$x^* (x_1^*, x_2^*, \dots, x_N^*)$$

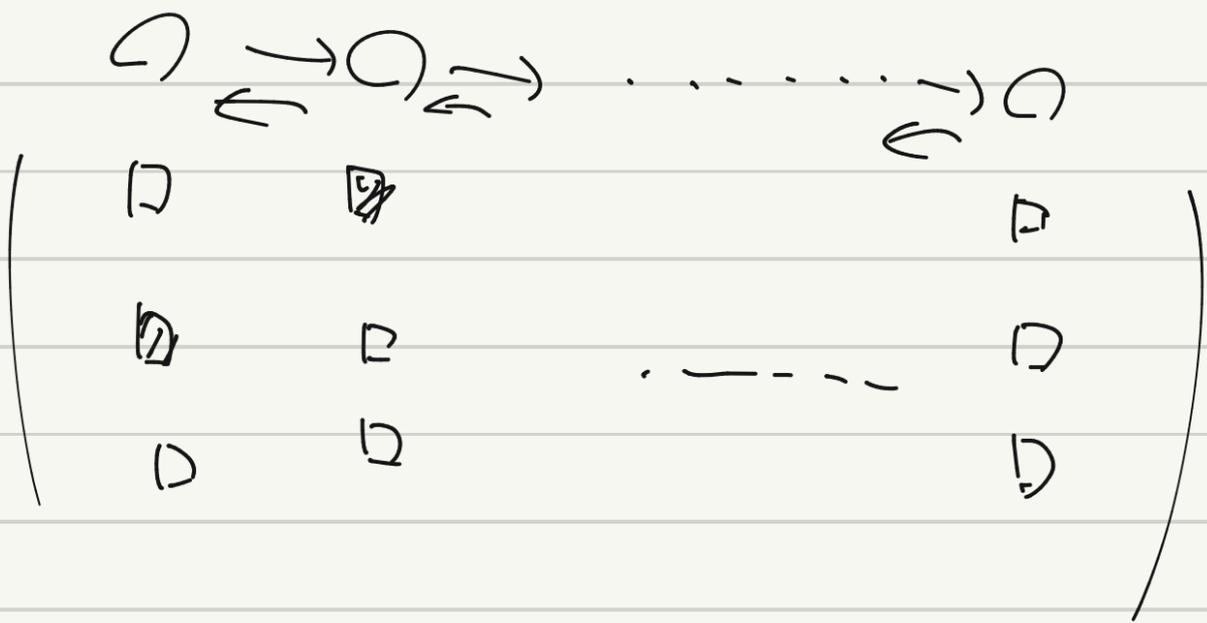
↑
root

↳ memorize (trellis diagram)

every node except root is not sure.

&

x^* may not be unique.



$$x_{n-1}^{\max} = \phi(x_n^{\max}) : \text{backtrack.}$$

$$x_N^{\max} = \arg \max [\mu_{f_{N-1}, N \rightarrow x_N}(x_N)]$$

memorize ϕ :

$$\phi(x_n) \equiv \arg \max_{x_{n-1}} [\ln P_{n-1, n}(x_{n-1}, x_n) + \mu_{x_{n-1} \rightarrow f_{n-1, n}}(x_{n-1})].$$

\uparrow
 x_{n-1} such that maximize a message