

$$\psi_{ij} = \psi_{ij}(x_i, x_j)$$

$$p(\underline{x}) = \frac{1}{Z} \psi_{12} \psi_{23} \dots \psi_{N-1,N}$$

full joint distr.

each x_i has K states

\Rightarrow the potential function ~~the~~ $\psi_{n-1,n}(x_{n-1}, x_n)$ is a $K \times K$ table.

\Rightarrow joint distr. $p(\underline{x})$ has $(N-1)K^2$ params

[?] How do we compute the marginal distr $p(x_n)$?

• Marginalize over joint:

$$\begin{aligned} p(x_n) &= \sum_{x_1} \dots \sum_{x_{n-1}} \sum_{x_{n+1}} \dots \sum_{x_N} p(\underline{x}) \\ &= \sum_{\underline{x} \in G \setminus n} p(\underline{x}) \end{aligned}$$

• Each x_n has K states $\Rightarrow K^N$ values for \underline{x}

\Rightarrow naive calculation of ~~p(x_n)~~ $p(x_n)$ is exponential in N .

• Can do better, just like variable elimination in prev. lecture for DAGs

By inspecting sum $\psi_{N-1,N}(x_{N-1}, x_N)$ is the only ~~one~~ term that depends on x_N

$$\Rightarrow \text{sum this first} \quad \sum_{x_N} \psi_{N-1,N}(x_{N-1}, x_N) = f(x_{N-1})$$

• Now this is the only term that depends of x_{N-1} , so

$$\text{compute} \quad \sum_{x_{N-1}} \psi_{N-2,N-1}(x_{N-2}, x_{N-1}) = f(x_{N-2})$$

• and so on...

② 491-L7 Int. GM2

- Similarly: $\psi_{1,2}(x_1, x_2)$ is the only one that depends on x_1
 - so perform this: $\sum_{x_1} \psi_{1,2}(x_1, x_2) = f(x_2)$
 - then compute: $\sum_{x_2} \psi_{2,3}(x_2, x_3) = f(x_3)$
 - and so on...
- Each sum removes a variable, which is equiv. to removing a node from the graph
- Can group the potentials and summations this way as:

$$p(x_n) = \frac{1}{Z} \underbrace{\left[\sum_{x_{n-1}} \psi_{n-1,n}(x_{n-1}, x_n) \cdots \left[\sum_{x_2} \psi_{2,3}(x_2, x_3) \left[\sum_{x_1} \psi_{1,2}(x_1, x_2) \right] \right] \cdots \right]}_{\mu_\alpha(x_n)} \quad \text{only var. left}$$

$$\times \underbrace{\left[\sum_{x_{n+1}} \psi_{n,n+1}(x_n, x_{n+1}) \cdots \left[\sum_{x_{N-1}} \psi_{N-2,N-1}(x_{N-1}, x_N) \left[\sum_{x_N} \psi_{N-1,N}(x_{N-1}, x_N) \right] \right] \right]}_{\mu_\beta(x_n)}$$

$$ab + ac = a(b+c) \quad \mu_\beta(x_n)$$

$$\frac{1}{Z} \left[\sum_{x_{n-1}} \psi_{n-1,n} \cdots \sum_{x_2} \psi_{2,3} \left[\sum_{x_1} \psi_{1,2} \right] \cdots \right] \times \left[\sum_{x_{n+1}} \psi_{n,n+1} \cdots \left[\sum_{x_{N-1}} \psi_{N-2,N-1} \left[\sum_{x_N} \psi_{N-1,N} \right] \right] \cdots \right]$$

[?] What's the cost of computing $\mu_\alpha(x_n) \mu_\beta(x_n)$?

- How many summations? Over everything but $x_n \Rightarrow N-1$
- How ~~many~~ many terms in each? Each var has K states.
Each ~~many~~ Ψ is a function of two vars, e.g.

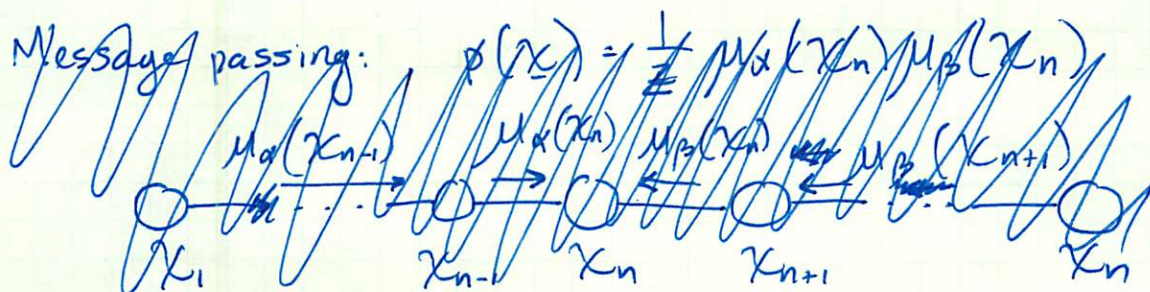
$$\Psi_{1,2}(x_1, x_2) \Rightarrow K \times K \text{ table}$$

$$\sum_{x_1} \Psi_{1,2}(x_1, x_2) \quad \text{sum over } K \text{ \#s for each val} \\ \Rightarrow \underline{K^2 \text{ cost}}$$

$$\Rightarrow \text{Computing } p(x_n) \text{ is } (N-1)K^2 = O(NK^2)$$

Naïve approach was exponential in N .

- This only works because the graph is much less than fully connected.
- Fully connected \Rightarrow need to use full joint ~~representation~~



④ 491-L7 Inf GM2

Message passing:

$$P(X) = \frac{1}{Z} \mu_\alpha(X_n) \mu_\beta(X_n)$$

Can be evaluated recursively:

$$\mu_\alpha(X_n) = \sum_{X_{n-1}} \Psi_{n-1,n}(X_{n-1}, X_n) \left[\sum_{X_{n-2}} \dots \right]$$

$$= \sum_{X_{n-1}} \Psi_{n-1,n}(X_{n-1}, X_n) \mu_\alpha(X_{n-1})$$

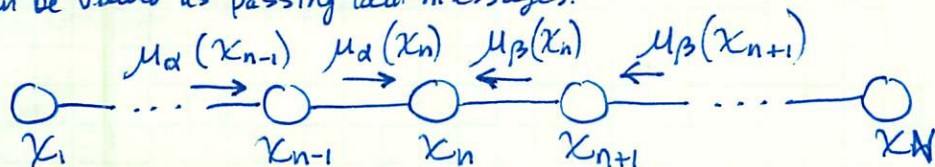
base case: $\mu_\alpha(X_2) = \sum_{X_1} \Psi_{1,2}(X_1, X_2)$

$$\mu_\beta(X_n) = \sum_{X_{n+1}} \Psi_{n,n+1}(X_n, X_{n+1}) \left[\sum_{X_{n+2}} \dots \right]$$

$$= \sum_{X_{n+1}} \Psi_{n,n+1}(X_n, X_{n+1}) \mu_\beta(X_{n+1})$$

base case: $\mu_\alpha(X_{N-1}) = \sum_{X_N} \Psi_{N-1,N}(X_{N-1}, X_N)$

Can be viewed as passing local messages:



Each has this form.
Outgoing message is
a point wise multiplication

over of the incoming
message and the
local potential Ψ_n
and summing over
the different vals of
the node var.

$$\mu_\alpha(X_n) = \sum_{X_{n-1}} \Psi_{n-1,n}(X_{n-1}, X_n) \overbrace{\mu_\alpha(X_{n-1})}^{\text{incoming}}$$

$$\mu_\beta(X_n) = \sum_{X_{n+1}} \Psi_{n,n+1}(X_n, X_{n+1}) \mu_\beta(X_{n+1})$$

⑤ 491 - ~~U18~~ L7 Inf. GM2

- What about Z ? $p(x_n) = \frac{1}{Z} \mu_\alpha(x_n) \mu_\beta(x_n)$

$$Z = \sum_x \prod_c \psi_c(\underline{x}_c)$$

in general, here cliques are just pairwise.

- Now we have:

$$p(x_n) = \frac{1}{Z} \mu_\alpha(x_n) \mu_\beta(x_n)$$

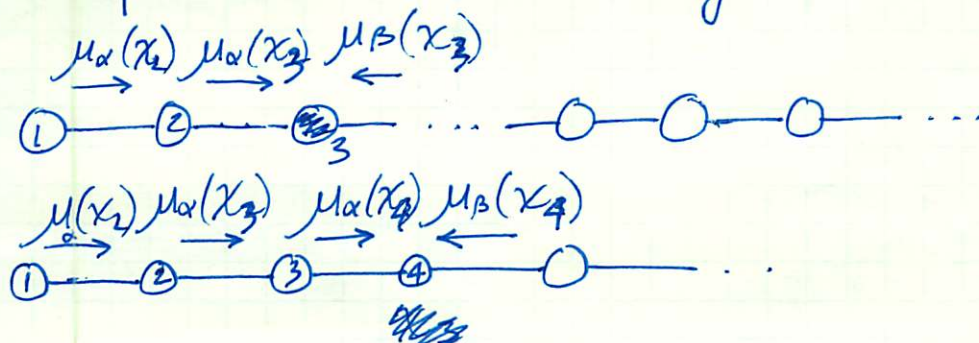
$$\text{So } Z = \sum_{x_n} \mu_\alpha(x_n) \mu_\beta(x_n)$$

- Only depends of x_n , so computation is $O(K)$.

What about computing $p(x_n)$ for other nodes?

If we do the procedure above to every node, it's $O(N \times N K^2)$
 $= O(N^2 K^2)$

□ Is this optimal in terms of efficiency?



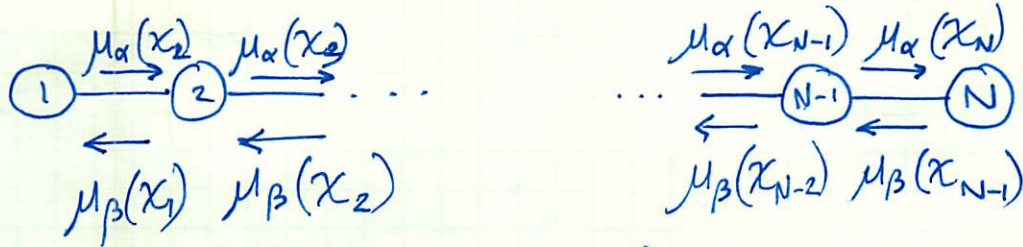
□ A: No it's very redundant.

~~Idea:~~ Store messages so they don't have to be recomputed.

The diagram shows a chain of nodes with nodes x_2 and x_{N-1} highlighted. Messages μ_α and μ_β are shown passing between nodes. There are some scribbles and corrections in the diagram.

⑥ 49-L7 Inf GM2

Idea: Store messages for efficient computation of marginals.



[?] What's the cost now? was $O(N \times NK^2)$

Messages only pass once
in each direction.

Now: $O(2 \times NK^2) = O(NK^2)$

Only twice as much computation
to compute every marginal!

[?] What about Z , the normalization constant?

A: Only needs to be computed once, it's the same for everything.

[?] What if some of the nodes are observed?

~~E.g. $\mu_\alpha(x_{n+1}) = \sum_{x_n} \Psi_{n,n+1}(x_n, x_{n+1}) \mu_\alpha(x_n)$~~

If $x_n = a$ is given, then this
is equivalent to setting the joint distri. $p(\underline{x}) = a$ for x_n .

~~$\mu_\alpha(x_{n+1}) = \sum_{x_n} \Psi_{n,n+1}(x_n, x_{n+1}) \mu_\alpha(x_n)$~~ This means all values ~~that~~ in the definition
of $p(\underline{x})$ where x_n could vary are now fixed to a .
Same for computation of Z .

$$\mu_\alpha(x_{n+1}) = \sum_{x_n=a} \Psi_{n,n+1}(x_n=a, x_{n+1}) \mu_\alpha(x_n=a)$$

⑦ 491-L7 Int GM2

[?] What about ~~about~~ computing other quantities?

E.g. $p(x_{n-1}, x_n)$? (two neighboring nodes in a chain)

$$p(x_{n-1}, x_n) = \frac{1}{Z} \sum_{\underline{x} \setminus x_n \setminus x_{n-1}} p(\underline{x})$$

$$= \sum_{x_1} \cdots \sum_{x_{n-2}} \sum_{x_{n+1}} \cdots \sum_{x_N} p(\underline{x})$$

i.e. sum over all but x_{n-1} & x_n

Now it proceeds as before, the only difference, however, is that we're not summing over x_{n-1} . Therefore we get:

$$p(x_{n-1}, x_n) = \frac{1}{Z} \mu_\alpha(x_{n-1}) \psi_{n-1,n}(x_{n-1}, x_n) \mu_\beta(x_n)$$

(refer to eqn that defines $\mu_\alpha(x_n) \propto \mu_\beta(x_n)$)

→ This means it's easy to compute joint distributions over sets of vars, once we've computed $\mu_\alpha \propto \mu_\beta$ for every node.

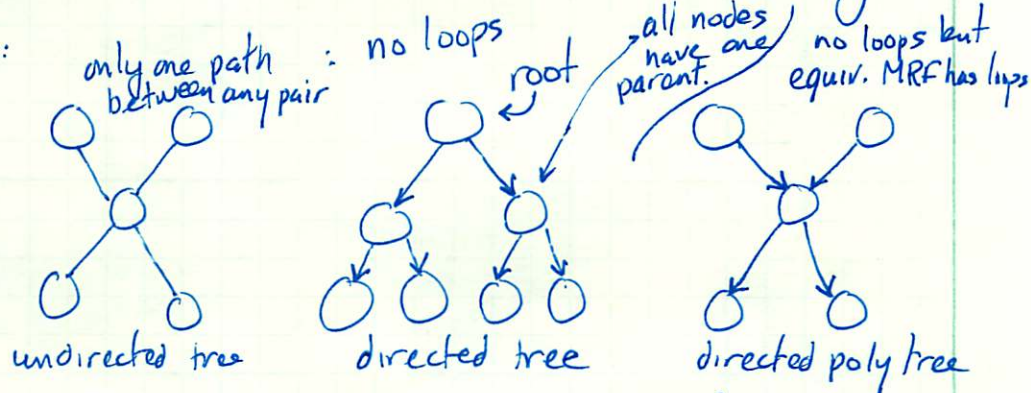
This gives us parametric forms for clique potentials,
→ or the conditional distributions, if we ~~we~~ start with a directed graph.

Trees

~~Why~~ We can do ^{exact} inference on a chain in linear time.

[?] Can we do it on other types of graphs?

We will show that inference can be done efficiently on trees:



Sum-product algorithm: ~~is~~ a generalization of ~~the~~ the message passing alg. which provides an efficient framework for exact inference on tree-structured graphs.

Sum product alg applies to:

- undirected trees
- directed trees
- poly trees

First: introduce a new graphical construction:
factor graphs

Factor graphs

Both ~~directed~~ directed and undirected graphs decompose joint pdfs into products:

$$p(\underline{x}) = \prod_i p(x_i | pa(\underline{x}_i))$$

cond. prob of node ~~and~~ given parents

$$p(\underline{x}) = \frac{1}{Z} \prod_c \psi_c(\underline{x}_c)$$

clique potential functions

More generally write pdf as a product of factors:

$$p(\underline{x}) = \prod_s f_s(\underline{x}_s) \quad \underline{x}_s = \begin{matrix} \text{some} \\ \text{subset of} \\ \text{vars} \end{matrix}$$

for DGs f_s = local conditional distr.

UGs f_s = potential fns over maximal cliques

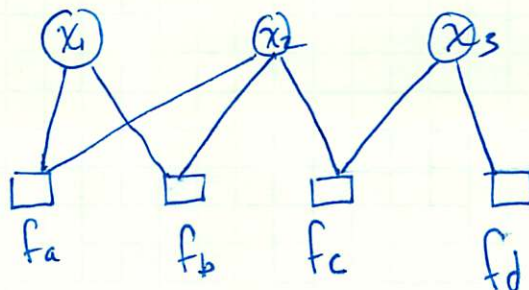
Note: Z isn't represented explicitly but could be defined as a factor that

doesn't depend on the vars.

Examples:

note: these are two factors for same vars

$$p(\underline{x}) = f_a(x_1, x_2) f_b(x_1, x_2) f_c(x_2, x_3) f_d(x_3)$$



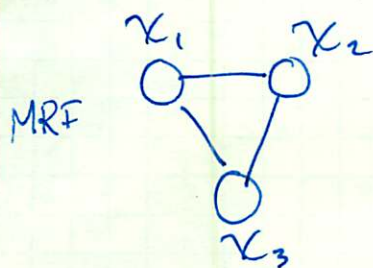
nodes

this is only over one var.

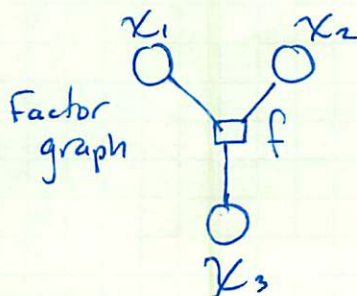
factors

Factor graphs keeps all factors explicit.

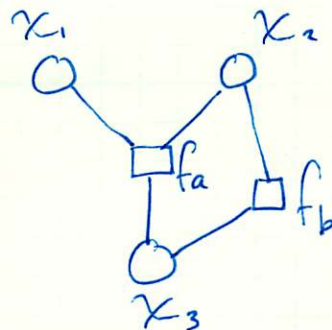
⑩ 491-L7 Int GM2



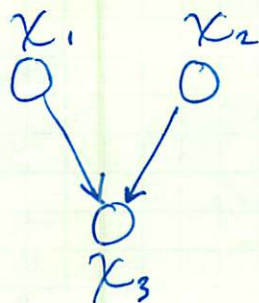
$$p(\underline{x}) = \frac{1}{Z} \Psi(x_1, x_2, x_3)$$



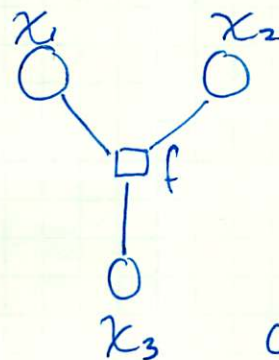
could add:



~~p(x)~~ $f_a(x_1, x_2, x_3) f_b(x_2, x_3)$
 $= \Psi(x_1, x_2, x_3)$



~~p(x)~~ $p(\underline{x}) = p(x_1) p(x_2) p(x_3 | x_1, x_2)$



$$f(x_1, x_2, x_3) = \frac{p(x_1)}{f_a} \frac{p(x_2)}{f_b} \frac{p(x_3 | x_1, x_2)}{f_c}$$

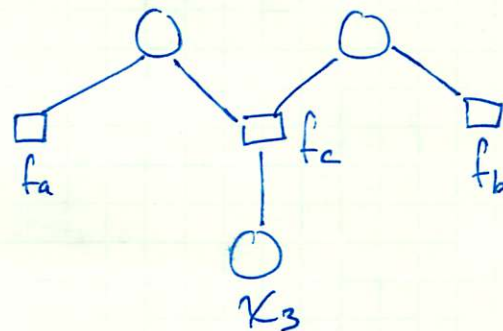
Factor graphs are bipartite:

- two kinds of nodes
- all links are between nodes of opposite types

could make all factor graphs like this



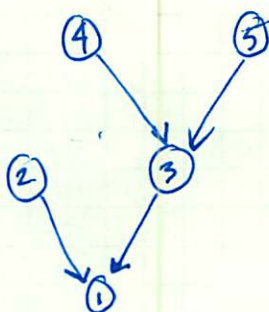
OR:



Factor graphs (cont'd):

For UGs: $\begin{cases} 1) \text{ factor nodes are maximal cliques } X_s. \\ 2) \text{ var nodes are same} \end{cases}$ $\begin{matrix} 3) \text{ fs}(X_s) \text{ clique potentials } \Psi_c(\underline{x}_c) \end{matrix}$

For a polytree:

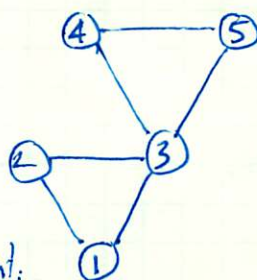


Could convert to UG:

~~$p(\underline{x}) = p(x_1) p(x_2) p(x_3) p(x_4) p(x_5)$~~

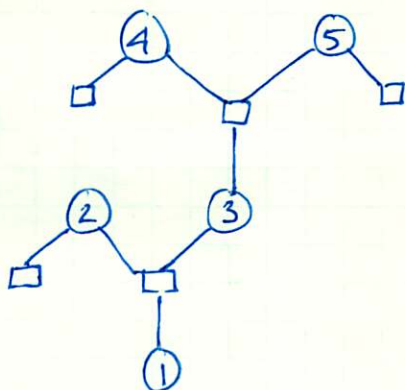
$$p(\underline{x}) = p(x_4) p(x_5) p(x_3 | x_4, x_5) = \Psi_2 \\ \times p(x_2) p(x_1 | x_2, x_3) = \Psi_1$$

The equiv. UG:
has loops



$$p(\underline{x}) = \frac{1}{Z} \Psi_1(x_1, x_2, x_3) \times \Psi_2(x_3, x_4, x_5)$$

But the equiv. factorgraph doesn't:



Note: It should be clear that the factorization fs does not correspond to any cond. indep. properties.