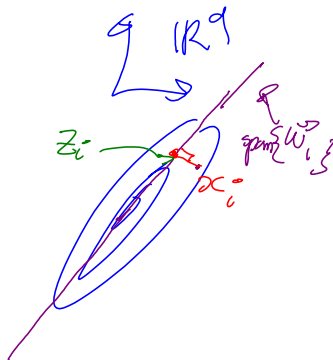today :  · PCA & factor analysis

· inference :  sum product alg.
HMM

## Factor analysis

$\vec{x}_i \in \mathbb{R}^d \rightsquigarrow \vec{z}_i \in \mathbb{R}^k \qquad k << d$

(dimensionality reduction)

## Synthesis view of PCA

get a orthonormal basis $\vec{w}_1, \dots, \vec{w}_k \in \mathbb{R}^d$

$$W = [ \vec{w}_1 \cdots \vec{w}_k ]$$
$d \times k$

$\rightarrow W^T W = I_k$ $\qquad \langle \vec{w}_i, \vec{w}_j \rangle = \delta_{ij}$
$k \times k$ identity

orthogonal projection on subspace $\mathrm{span}\{\vec{w}_1, \dots, \vec{w}_k\}$

$$P_W = W W^T \qquad (P_W^2 = P_W)$$

$$P_W \vec{x} = \sum_k \vec{w}_k \underbrace{\langle \vec{w}_k, \vec{x} \rangle}_{(\vec{z})_k}$$

$$\vec{z}_{\in \mathbb{R}^k} = W^T \vec{x}$$

PCA $\qquad \min_{\substack{W \in \mathbb{R}^{d\times k} \\ W^T W = I_k}} \sum_i \| \vec{x}_i - \underbrace{W \underbrace{W^T \vec{x}_i}_{\vec{z}}}_{W} \|^2 \qquad$ exercise... $\qquad X = \begin{pmatrix} -x_i^T- \\ \vdots \end{pmatrix}$
$n \times d$
data matrix

"analysis" view $\min \left( \text{cst.} - \sum_k \vec{w}_k^T X^T X \vec{w}_k \right)$

scaled "empirical covariance matrix" if data is centered

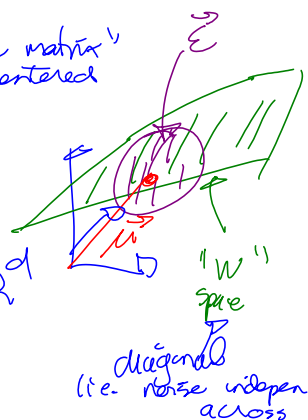## generative model (latent variable model)

$\vec{z} \sim N(0, I_k)$ $\qquad$ "somewhat uniform"

$\vec{x} = W\vec{z} + \vec{\mu} + \underset{\text{noise}}{\vec{\varepsilon}} \qquad \vec{\varepsilon} \sim N(0, \Psi) \in \mathbb{R}^d$

$\vec{\varepsilon} \perp\!\!\!\perp \vec{z}$

diagonal (i.e. noise indep. across)

$$\vec{x}|\vec{z} \sim N(W\vec{z} + \vec{\mu}, \Psi)$$

what is $p(\vec{x})$?  $\mathbb{E}[\vec{x}] = \mathbb{E}[\mathbb{E}[\vec{x}|\vec{z}]]$

$$\mathbb{E}[\underbrace{W\vec{z}}_{0} + \mu] = \vec{\mu}$$

$$cov(\vec{x}, \vec{x}) = cov(W\vec{z} + \vec{\mu} + \vec{\varepsilon}, \; W\vec{z} + \vec{\mu} + \vec{\varepsilon})$$

<span style="color:red">all independent</span>  (by independence)

$$= cov(W\vec{z}, W\vec{z}) + cov(\vec{\varepsilon}, \vec{\varepsilon})$$

$$= W\underbrace{\mathbb{E}[\vec{z}\vec{z}^T]}_{I_k} W^T \; + \; \Psi$$

$$= WW^T + \Psi$$

$$\vec{x} \sim N(\vec{\mu}, \underset{d \times k}{WW^T} + \Psi)$$

$d$ degrees of freedom

to get $\vec{z}$,  $p(\vec{z}|\vec{x})$  $\begin{pmatrix} \vec{x} \\ \vec{z} \end{pmatrix} \sim N\left( \begin{matrix} \mu_x \\ \mu_z \end{matrix}, \begin{pmatrix} \Sigma_{xx} & \Sigma_{zx} \\ \Sigma_{xz} & \Sigma_{zz} \end{pmatrix} \right)$

$$\mu_x = \vec{\mu} \quad \mu_z = 0 \quad \Sigma_{xx} = WW^T + \Psi$$

$$\Sigma_{zz} = I_k$$

$$cov(\vec{x}, \vec{z}) = cov(W\vec{z} + \mu + \vec{\varepsilon}, \vec{z}) \quad I_k$$

$$= cov(W\vec{z}, \vec{z}) = W\mathbb{E}[\vec{z}\vec{z}^T] = W$$

$$\mathbb{E}[\vec{z}|\vec{x}] = \mu_z + \Sigma_{zx} \Sigma_{xx}^{-1}(\vec{x} - \vec{\mu}_x)$$

$$= 0 + W^T(WW^T + \Psi)^{-1}(\vec{x} - \vec{\mu}_x)$$

$$\Psi = \sigma I_d \quad \rightarrow \quad \text{probabilistic PCA}$$

$$\lim_{\sigma \to 0} W^T(WW^T + \sigma I)^{-1} = \underset{\color{red}{\text{pseudo-inverse}}}{\color{red}{W^{+T}}}$$

$$= W^T \text{ here (exercise)}$$

PCA representation     $W^T(\vec{x} - \vec{\mu}_x)$

PCA as limit PPCA $\sigma \to 0$

to estimate $W, \Psi, \mu \rightarrow$ maximum likelihood

$$EM \qquad E \to p(\breve{z}|\breve{x})$$

model not identifiable: $\qquad WW^T + \Psi$

$$W' = WR$$

$R$ $k \times k$ rotation matrix
$R^T R = RR^T = I_k$

$$W'W'^T = WRR^TW^T = WW^T$$
$$\underbrace{\phantom{RR^T}}_{I_R}$$

why factor analysis?
- richer noise model $\Psi$
- easy to plug into bigger models

## Inference in graph model

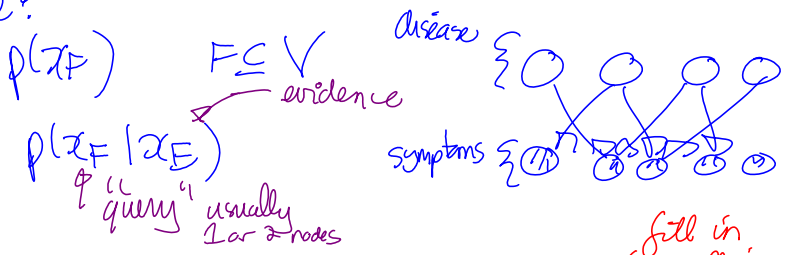say DGM: $\quad p(x) = \prod_i p(x_i | x_{\pi_i})$

or

UGM $\quad p(x) = \frac{1}{Z} \prod_{C \in \mathcal{C}} \psi_C(x_C)$

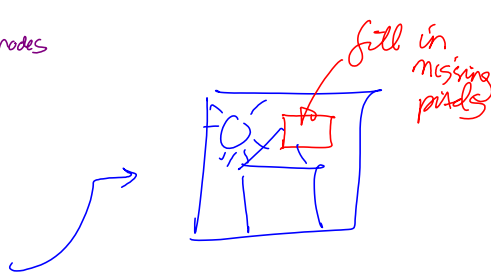$i$ & $\pi_i$ put in clique
moralization

want to compute:

a) marginal $\quad p(x_F)$ $\qquad F \subseteq V$
evidence

b) conditional $\quad p(x_F | x_E)$
marginal
"query" usually 1 or 2 nodes

disease $\{$ ... $\}$
symptoms $\{$ ... $\}$

c) for UGM compute $Z$

why? • prediction
• missing data
fill
fill in missing pixels

d) related: decoding $\quad \underset{x_F}{\text{argmax}} \ p(x_F | x_E)$

$F$ might be hide

need inference for estimation (ML) for UGM on discrete data

$$\log p(x | n) = \sum_{C \in \mathcal{C}} \sum_{y_C} \underbrace{\mathbb{1}\{x_C = y_C\}}_{[T(x)]_{y_C}} \underbrace{\log \psi_C(x_C)}_{n_{\{y_C\}}} - \log Z$$

$$= T(x)^T n - A(n)$$

$$(\partial \theta \ell) \quad \partial^c \quad [T(x)]_{y_c} \quad m_{4}(y_c)$$

$$= T(x)^T \eta - A(\eta)$$

$$\nabla_\eta [\quad\quad] = T(x) - \underbrace{\nabla_\eta A(\eta)}_{\mathbb{E}_\eta[T(X)]} \nearrow$$

$$\mathbb{E}[\mathbb{1}\{y_c = x_c\}] = p(y_c)$$

aside: $\quad |C| \leq 2 \quad \rightarrow$ for binary state $\quad$ "Ising model"

$\quad\quad\quad\quad\quad\quad\quad\quad\quad K$ states $\quad$ "pott's model"

# graph elimination
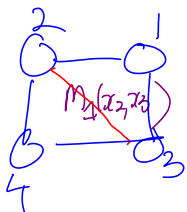
consider UGM $\quad p(x) = \frac{1}{Z} \prod_C \psi_C(x_C)$

goal: to compute $p(x_F)$

main trick: distributivity $\quad \sum_{x_1, x_2} f(x_1) g(x_2)$

convince yourself

$$\left(\sum_{x_1} f(x_1)\right) \left(\sum_{x_2} g(x_2)\right)$$

$$\sum_{x_{1:n}} \left(\prod_i f_i(x_i)\right) = \prod_i \left(\sum_{x_i} f_i(x_i)\right)$$

huge (exp) sum $\quad$ looks UGM factorization $\quad$ small sum (!)

$m_1(x_2, x_3)$

$$p(x_4) = \sum_{x_{1:3}} \frac{1}{Z} \psi_{12}(x_1, x_2) \, \psi(x_1, x_3) \, \psi(x_2, x_4) \, \psi(x_3, x_4)$$

$$p(x_4) = \frac{1}{Z} \sum_{x_3} \psi(x_3, x_4) \sum_{x_2} \psi(x_2, x_4) \sum_{x_1} \underbrace{\psi(x_1, x_3) \psi(x_1, x_2)}_{\psi_{123}(x_1, x_2, x_3)}$$

$$\underbrace{\phantom{\sum_{x_1} \psi(x_1, x_3) \psi(x_1, x_2)}}_{m_1(x_2, x_3)}$$

new interactions (neighbors of 1 were tied together)

$$p(x_4) = \frac{1}{Z} m_2(x_4)$$

$m_2(x_3, x_4)$

$p(x_4) \propto m_3(x_4)$

<u>graph eliminate algo.:</u>

$$Z = \sum_{x_4} M_3(x_4)$$

INIT: • choose ordering of elimination st. set F is at the end

• put all $\psi_c(x_c)$ on "active list"

UPDATE: repeat in order of elimination

suppose eliminating $x_i$

1) remove all factors in "active list" that $x_i$ as argument and take product; find all $\alpha$ s.t.

Let $S_i$ be variables appearing in these factors excluding $i$ $\quad f_\alpha(x_\alpha)$ was factor in list and $i \in \alpha$

2) sum out $x_i$ of product:

$$M_i(x_{S_i}) = \sum_{x_i} \left( \prod_{\substack{\alpha \\ i \in \alpha}} f_\alpha(x_\alpha) \right)$$

function ($x_{S_i \cup \{i\}}$)

3) add $\boxed{M_i(x_{S_i})}$ on active list

↳ new table to store $2^{|S_i|}$ ← new big clique
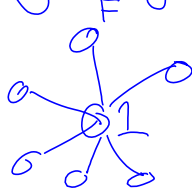
at end you're left with only factors of $x_F$

→ normalize to get $p(x_F)$

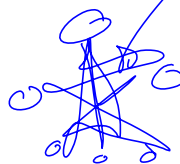<u>Cost of alg.:</u> store factors of size $2^{|S_i|}$

sum over $2^{|S_i|+1}$ values to compute message $M_i(x_{S_i})$

memory & time is exponential in size of <u>biggest clique</u> formed during elimination
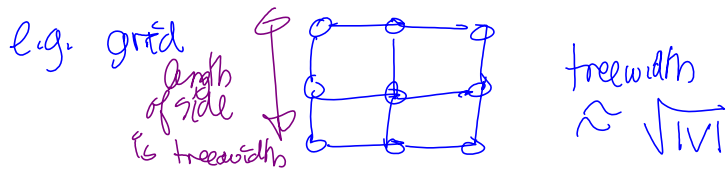
not all orderings are good



$M_1(x_{V \setminus \{i\}})$

big clique

convention → that 1 for trees

<u>treewidth of graph</u> = min over all orderings (size biggest clique appearing during elimination) $-1$

<u>bad news:</u>

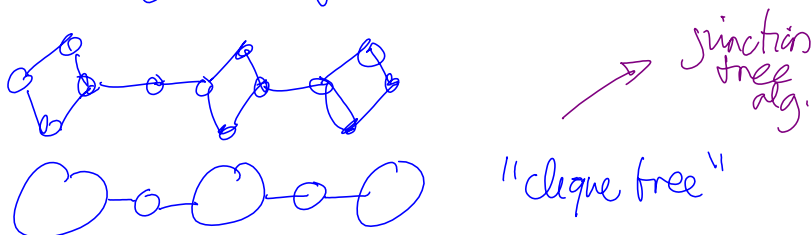• computing treewidth in general graph is NP-hard (and so also finding optimal order ...)

(and so also finding optimal order" " " )

- inference in general graph model is NP-hard
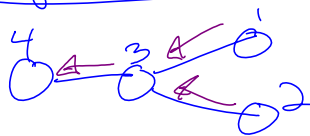
e.g. grid  treewidth $\sim \sqrt{|V|}$

$\rightsquigarrow$ need approximations

good news
- exact inference for trees in linear time in size of tree
- "tree like" graph of small treewidth

 $\longrightarrow$ junction tree alg.

"clique tree"

## Inference on a tree
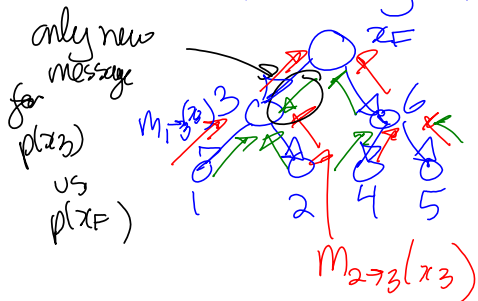


$$p(x_4) = \frac{\psi_4(x_4)}{Z} \sum_{x_3} \psi_3(x_3)\psi(x_3,x_4) \underbrace{\sum_{x_2}\psi(x_2)\psi(x_2,x_3)}_{m_{2\to3}(x_3)} \overbrace{\sum_{x_1}\psi(x_1)\psi(x_1,x_3)}^{m_{1\to3}(x_3)}$$

$\underbrace{\hspace{4cm}}_{m_{3\to4}(x_4)}$

$$p(x_4) = \frac{\psi_4(x_4)\, m_{3\to4}(x_4)}{Z}$$

compute $p(x_F)$ where $F$ is singleton

orient the tree by making $x_F$ the root

only new message for $p(x_3)$ vs $p(x_F)$



$m_{1\to3}(x_3)$

$m_{2\to3}(x_3)$

$$m_{i\to j} = \sum_{x_i} \psi_i(x_i)\psi_{ij}(x_i,x_j) \prod_{K\in ch(i)} m_{K\to i}(x_i)$$

child $\to$ parent

children of $i$

idea : use dynamic programming

idea: use dynamic programming
       to store messages
       and efficiently compute $p(x_i) \; \forall i$
    $\rightarrow$ sum-product alg.

at end
$$p(x_F) = \frac{\psi_F(x_F)}{Z} \prod_{k \in C(F)} \left( M_{k \to F}(x_F) \right)$$



$x_F = 7$

"collect phase" $\Rightarrow$ needed for $p(x_F)$

"distribute phase"

1    2    4    5

$$p(x_3) = M_{1 \to 3}(x_3) \, M_{2 \to 3}(x_3) \, M_{7 \to 3}(x_3) \, \frac{\psi_3(x_3)}{Z}$$

general rule:
     want to compute $M_{i \to j}(x_j)$ for all $\{i,j\} \in E$
               $M_{j \to i}(x_i)$

rule: $i$ send message to $j$ if it has received all
               messages from other neighbors

$k \to i \to j$

$$M_{i \to j}(x_j) = \sum_{x_i} \psi_i(x_i) \, \psi_{ij}(x_i, x_j) \cdot \prod_{k \in N(i) \setminus \{j\}} M_{k \to i}(x_i)$$

$$\left[ \text{loopy BP: } \quad M_{i \to j}^{new}(x_j) = \left( M_{i \to j}^{old}(x_j) \right)^{\alpha} \left( \triangleright \right)^{1-\alpha} \right]$$

$$0 \leq \alpha < 1 \quad \text{"damping"}$$
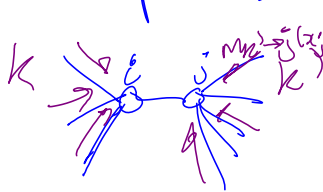
$\rightarrow \Rightarrow$ only leaves can start

parallel schedule:

1) initialize $M_{i \to j}(x_j)$ to uniform distribution

2) at every step (in parallel)

     compute $M_{i \to j}(x_j) = \sum \psi_i(x_i) \, \psi_{ij}(x_i, x_j)$

length of largest path $x_i$ $\prod_{k \in N(j) \setminus \{i\}} m_{k \to i}(z_i)$

$\longrightarrow$ can prove after $\underbrace{\text{diameter (tree)}}$ steps, all messages have converged to correct value

after sum-product

$$p(x_i) = \frac{1}{Z} \, \psi_i(x_i) \prod_{k \in N(i)} m_{k \to i}(x_i)$$



$$p(x_i, x_j) = \frac{1}{Z} \, \psi_i(x_i) \, \psi_j(x_j) \, \psi_{ij}(x_i, x_j) \left( \prod_{k \in N(i) \setminus \{j\}} m_{k \to i}(x_i) \right) \left( \prod_{k' \in N(j) \setminus \{i\}} m_{k' \to j}(x_j) \right)$$

conditionals:

$$p(x_i \mid \bar{x}_E) \propto p(x_i, \bar{x}_E)$$

just fix $\bar{x}_E$ during sum-product
(do not sum them out)

(formal trick)   redefine $\tilde{\psi}_j(x_j) = \psi_j(x_j) \, \delta(x_j, \bar{x}_j)$   $j \in E$

Kronecker delta

when have $\sum_{x_j} \tilde{\psi}_j(x_j) \cdot \text{stuff}(x_j, x_i)$

$$= \psi_j(\bar{x}_j) \cdot \text{stuff}(\bar{x}_j, x_i)$$

$$p(x_i \mid \bar{x}_E) = \frac{p(x_i, \bar{x}_E)}{\sum_{x_i'} p(x_i', \bar{x}_E)}$$

remarks: sum-product a.k.a. belief propagation (BP)
or message passing
is exact only on trees
[ approximation on graph with cycles ]
loopy BP

• only property used for sum product
is distributivity of $\oplus$ with respect $\odot$

$(\mathbb{R}, \oplus, \odot)$   as a semi-ring

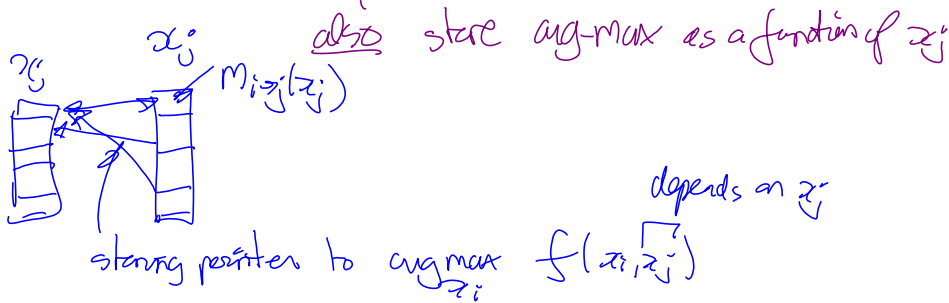$(\hookrightarrow$ "don't need universe to addition$)$

do sum-product on other semi-rings

$$(\mathbb{R}_+, \max, \odot) \quad \text{or} \quad (\mathbb{R}, \max, \oplus)$$

$\longrightarrow$ "max-product"  $\quad \max(a \cdot b, a \cdot c) = a \cdot \max(b, c)$
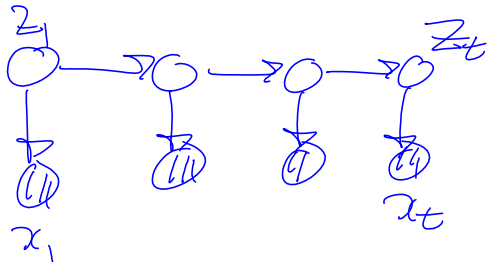
$\oplus ( \qquad )$

$$\max_{x_{1:n}} \prod_i (f_i(x_i)) = \prod_i (\max_{x_i} f_i(x_i))$$

$$m_{i \to j}(x_j) = \max_{x_i} \left[ \Psi_i(x_i) \, \Psi_{ij}(x_i, x_j) \prod_{k \in N(i) \setminus \{j\}} (m_{k \to i}(x_i)) \right]$$

$x_i$  $x_j$ ⟶ *also* store arg-max as a function of $x_j$

 $m_{i \to j}(x_j)$

storing pointers to $\underset{x_i}{\text{argmax}}$ $f(x_i, \boxed{x_j})$

depends on $x_j$

max-product a.k.a. viterbi algorithm $\Rightarrow$ decoding

$$\underset{x_{1:n}}{\text{argmax}} \; p(x_{1:n})$$

# HMM: Hidden Markov Model



$z_1$  $z_t$

$x_1$  $x_t$

directed G.M.

$z_t \in \{1, \ldots, K\}$

$\}$ phoneme

$\}$ speech signal $x_t \begin{cases} \text{cts.} \\ \text{discrete} \end{cases}$

emission prob.   transition prob.

$$p(x_{1:T}, z_{1:T}) = p(z_1) \prod_{t=1}^{T} (p(x_t \mid z_t)) \prod_{t=2}^{T} (p(z_t \mid z_{t-1}))$$

• *suppose* homogeneous ie $p(x_t \mid z_t) = f(x_t \mid z_t)$

$\qquad$ ↳ does not depend on $t$

$$p(z_t = i \mid z_{t-1} = i) = A_{ii} \qquad \int$$

$$p(z_t = i \mid z_{t-1} = j) = A_{ij}$$

$$\text{need} \quad \sum_i A_{ji} = 1$$

$$A \begin{pmatrix} \overset{j}{\prod} \end{pmatrix}$$

prob. dist over $z$

tasks:

   prediction: $\quad p(z_t \mid x_{1:t-1})$    "where next?"

   filtering: $\quad p(z_t \mid x_{1:\boxed{t}})$    "where now?"

   smoothing: $\quad p(z_t \mid x_{1:T})$    "what is the past?"
           $T > t$



$$p(z_t, \bar{x}_{1:t})$$

$$= \frac{1}{Z} m_{x_t \to z_t}(z_t) \, m_{z_{t-1} \to z_t}(z_t)$$
$$1 = Z$$

$$m_{x_t \to z_t}(z_t) = \sum_{x_t} \delta(x_t, \bar{x}_t) \, p(x_t \mid z_t)$$

$$= p(\bar{x}_t \mid z_t)$$

$$m_{z_{t-1} \to z_t}(z_t) = \sum_{z_{t-1}} p(z_t \mid z_{t-1}) \underbrace{m_{x_{t-1} \to z_{t-1}}(z_{t-1}) \, m_{z_{t-2} \to z_{t-1}}(z_{t-1})}_{p(z_{t-1}, \bar{x}_{1:t-1})}$$

$$\downarrow$$
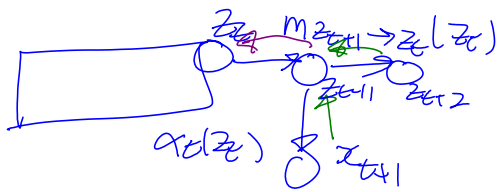$$\alpha_{t-1}(z_{t-1})$$

$$\boxed{\alpha_t(z_t) = p(z_t, \bar{x}_{1:t}) = p(\bar{x}_t \mid z_t) \sum_{z_{t-1}} p(z_t \mid z_{t-1}) \, \alpha_{t-1}(z_{t-1})}$$

$\alpha$-recursion     forward recursion

like the "collect phase" in sum-product

$$p(\bar{x}_{1:t}) = \sum_{z_t} \alpha_t(z_t)$$

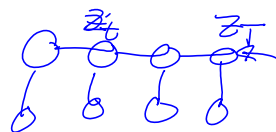for smoothing, also need backward message $\beta_t(z_t)$

$$M_{z_{t+1} \to z_t}(z_t)$$

$$M_{z_{t+1} \to z_t}(z_t) = \beta_t(z_{t+1})$$

$$\sum_{z_{t+1}} p(\bar{x}_{t+1} | z_{t+1}) \, M_{z_{t+2} \to z_{t+1}}(z_{t+1})$$
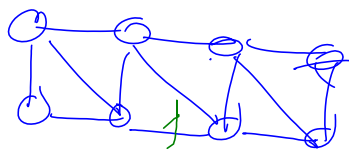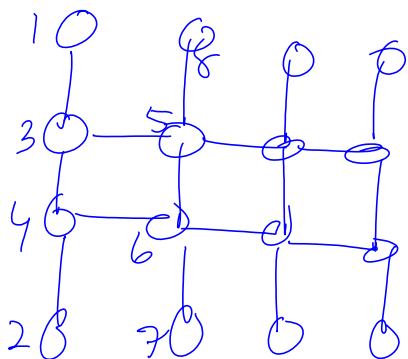
$\rightarrow$ $\beta$-recursion
(backward pass)

## initialization:

$$\beta_T(z_T) = 1$$



$$\alpha_1(z_1) = p(\bar{x}_1 | z_1) \, p(z_1)$$

$$p(z_t, \bar{x}_{1:T}) = \alpha_t(z_t) \, \beta_t(z_t)$$

$$\beta_t(z_t)$$

$$\alpha_t(z_t)$$



---



clique tree
running intersection property



$1,2,3,4$ — $3,4,5,6$ — $7,8,5,6$