

Factor Graphs

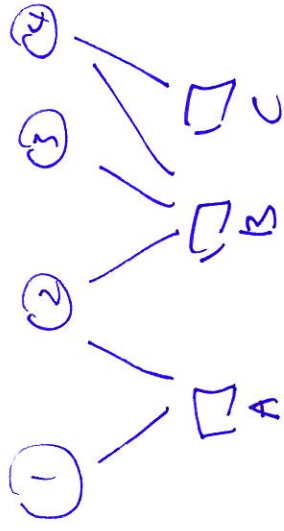
$X_1 \dots X_N$

$$P(X_1, \dots, X_N) = \frac{1}{Z} \prod_a f_a(x_a)$$

Interested in:

① estimating marginals $P_S(X_S) = \sum_{X \setminus X_S} P(X)$,

② estimating Z or $\log Z$



$$P(X_1, \dots, X_4) = \frac{1}{Z} f_A(x_1, x_2) f_B(x_2, x_3, x_4) f_C(x_3, x_4)$$

Belief Propagation

(Sum product, forward backward, Gallager decoding, turbo-decoding)

- Messages $M_{a \rightarrow i}(x_i)$ from factors to variables
- what values does factor a like variable X_i to take on?
- Messages $M_{i \rightarrow a}(x_i)$ from variables to factors
- what values X_i likes based on information from all but a .
- Beliefs

$$b_i(x_i) \propto \prod_{a \in N(i)} M_{a \rightarrow i}(x_i)$$

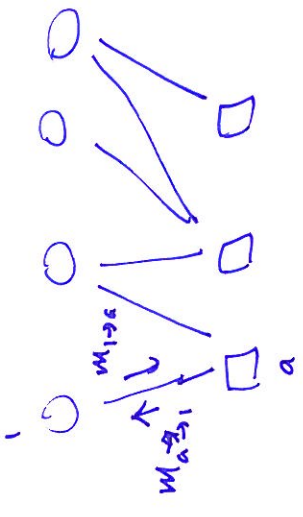
based on all pieces of information coming into X_i
 → note product is independent pieces of information

$$b_a(x_a) \propto f_a(x_a) \prod_{i \in N(a)} M_{i \rightarrow a}(x_i)$$

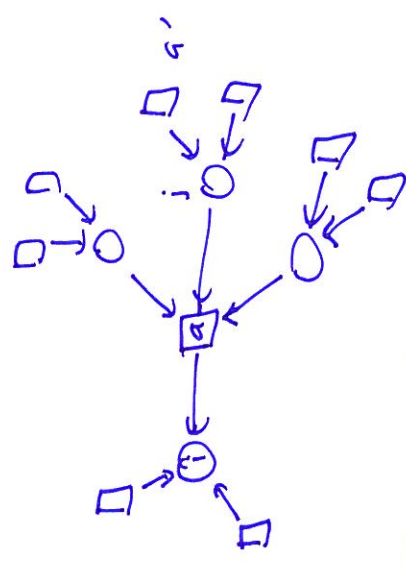
- Message updates

$$b_i(x_i) = \sum_{x_{a'}} b_a(x_{a'})$$

$$\begin{aligned} \prod_{a \in N(i)} M_{a \rightarrow i}(x_i) &= \sum_{x_{a'}} f_a(x_{a'}) \prod_{j \in N(a)} M_{j \rightarrow a}(x_j) \\ &= \sum_{x_{a'}} f_a(x_{a'}) \prod_{j \in N(a)} \prod_{a' \in N(j) \setminus a} M_{a' \rightarrow j}(x_j) \\ &= \prod_{a' \in N(i) \setminus a} M_{a' \rightarrow i}(x_i) \sum_{x_{a'}} f_a(x_{a'}) \prod_{j \in N(a) \setminus i} M_{a' \rightarrow j}(x_j) \\ M_{a \rightarrow i}(x_i) &= \sum_{x_{a'}} f_a(x_{a'}) \prod_{j \in N(a) \setminus i} M_{j \rightarrow a}(x_j) \end{aligned}$$



- Constraints:
 - nonnegative $b_i(x_i) \geq 0$. $b_a(x_a) \geq 0$
 - normalization $\sum_{x_i} b_i(x_i) = 1$ $\sum_{x_a} b_a(x_a) = 1$
 - marginalization $b_i(x_i) = \sum_{x_{a'}} b_a(x_{a'})$



- replace \sum by \max
 → max product for MAP state

Free Energies

- energy $E(x) = - \sum_a \log f_a(x_a)$ ~~$f_a(x_a)$~~
- entropy $H(b) = - \sum_x b(x) \log b(x)$
- Gibbs free energy

$$F(b) = - \sum_x b(x) E(x) - H(b) \\ = U(b) - H(b)$$

- system prefers low energy, high entropy
- minimize $F(b)$

$$\Rightarrow b(x) = \frac{1}{Z(\tau)} e^{-E(x)/T}$$

attaining

$$F(b_{\text{min}}) = - \log Z(\tau) = F_H$$

Z Helmholtz free energy

$$F(b) = F_H + \text{KL}(b \| P)$$

Two computational bottlenecks:

- ~~$H(b)$~~ b very large
- $H(b)$ expensive to compute

\Rightarrow

Solutions

- assume b comes from tractable family
 - \rightarrow variational approximation
- give up on b being a distribution
 - \rightarrow replaced by a family of marginal beliefs $b_a(x_a) \dots$
 - \rightarrow introduce constraints among them
 - \rightarrow usually set of constraints ensuring global consistency is very large
 - \rightarrow just use a subset of them
 - \rightarrow local consistency
- approximate $H(b)$ as a sum of marginal beliefs

Bethe free energy

- family of beliefs

$$b_i(x_i), b_a(x_a)$$

- nonnegative, normalization constraints

$$- U_{\text{Bethe}} = - \sum_a \sum_{x_a} b_a(x_a) \log f_a(x_a)$$

$$H_{\text{Bethe}} = - \sum_a \sum_{x_a} b_a(x_a) \log b_a(x_a)$$

$$+ \sum_i (d_i - 1) \sum_{x_i} b_i(x_i) \log b_i(x_i)$$

$\mathcal{L}(\text{INCL})$

maximize

$$\mathcal{L}(b, \lambda, \pi) = - \sum_a \sum_{x_a} b_a(x_a) \log f_a(x_a) + \sum_a \sum_{x_a} b_a(x_a) \log (b_a(x_a)) - \sum_i (d_i - 1) \sum_{x_i} b_i(x_i) \log b_i(x_i)$$

$$= \sum_a \sum_{x_a} \gamma_a \left(\sum_{x_a} b_a(x_a) - 1 \right) + \sum_i \gamma_i \left(\sum_{x_i} b_i(x_i) - 1 \right) + \sum_{a \in \text{NCL}} \lambda_{a_i}(x_i) \left(b_i(x_i) - \sum_{x_{a_i}} b_a(x_a) \right)$$

$$= - \log f_a(x_a) + \log b_a(x_a) + 1 + \gamma_a - \sum_{i \in \text{NCL}} \lambda_{a_i}(x_i) \Rightarrow b_i(x_a) = f_a(x_a) e^{\gamma_a - 1 + \sum_i \lambda_{a_i}(x_i)}$$

$$= \pi (1 - d_i) (\log b_i(x_i) + 1) - \gamma_i + \sum_a \lambda_{a_i}(x_i) \Rightarrow b_i(x_i) = e^{-1 + \frac{1}{d_i - 1} (-\gamma_i + \sum_a \lambda_{a_i}(x_i))}$$

$$\lambda_{a_i}(x_i) = \log \pi_{i \rightarrow a}(x_i) = \log \prod_{a' \in \text{NCL}(i)} \pi_{a' \rightarrow i}(x_i)$$

$$b_a(x_a) \propto f_a(x_a) \prod_{i \in \text{NCL}(a)} \pi_{i \rightarrow a}(x_i)$$

$$b_i(x_i) \propto \left(\prod_{a \in \text{NCL}(i)} \prod_{a' \in \text{NCL}(i)} \pi_{a' \rightarrow i}(x_i) \right)^{\frac{1}{d_i - 1}} = \prod_{a \in \text{NCL}(i)} \pi_{a \rightarrow i}(x_i)$$

- BP updates from marginal beliefs constraints again

- ① can be not globally consistent (unrealizable)

- ② can obtain negative Bethe entropy -

+ ③ exact on trees

+ ④ maxent-normal

entropy of X_i over-counted by d_i times.

Region Graphs

Region R is a set of variables V_R and factors A_R such that if $a \in A_R$ then $N(a) \subseteq V_R$.

- true marginal $P_R(x_R)$
- beliefs $b_R(x_R)$
- energy $-\sum_{a \in A_R} \log f_a(x_a) = E_R(x_R)$
- free energy $F_R(b_R) = \sum_{x_R} b_R(x_R) E_R(x_R) - H_R(b_R)$

Region-based free energy $U_R(b_R)$ counting numbers.

$$F(\{b_R\}) = \sum_R c_R U_R(b_R) - \sum_{R \in \mathcal{R}} c_R H_R(b_R)$$

- ↳ treat each region exactly
- ↳ approximate free energy (in fact entropy)
- ↳ treat interactions among regions using marginalization constraints.

① Valid if $\sum_R c_R \mathbb{1}(a \in A_R) = \sum_R c_R \mathbb{1}(i \in V_R) = 1 \quad \forall a, i$

exact energy if $b_R = P_R$ exact entropy if uniform distribution under constraints

- ② maxent-normal if $H(\{b_R\}) = \sum_R c_R H_R(b_R)$ achieves maximum when all b_R uniform
- ③ short loops should be regions (loop-based region graphs)
- ④ perfect-correlation if $\sum_R c_R = 1$ exact when all variables fully coupled.
- ⑤ non-singular if message passing algorithms converge to unique exact uniform fixed point. loop graphs singular if $\sum_R c_R > 1$.

① Tree EP
② planar graphs
 $\sum_R c_R = L - E + V = 1$ for faces, planar.

Region Graphs

- ~~constraints~~ organize regions into a

directed acyclic graph

$R_1 \rightarrow R_2$ only if $R_2 \subset R_1$ / all factors g
all variables i .

enforce marginalization constraint

$$\sum_{X_{R_1} \setminus X_{R_2}} b_{R_1}(X_{R_1}) = b_{R_2}(X_{R_2})$$

↳ require: subgraph of all regions containing

extended to: $\#$ regions containing any
subset of a 's, i 's are connected.

each a or i is connected

so that region graph gives consistent beliefs

about a, i .

↳ require $C_R = 1 - \sum_{R' \in \text{Ancestors}(R)} C_{R'}$ degree of freedom computed properly.

$$C_R = 1 - \sum_{R' \in \text{Ancestors}(R)} C_{R'}$$

- exact if region graph forms a tree (and satisfies all constraints above)

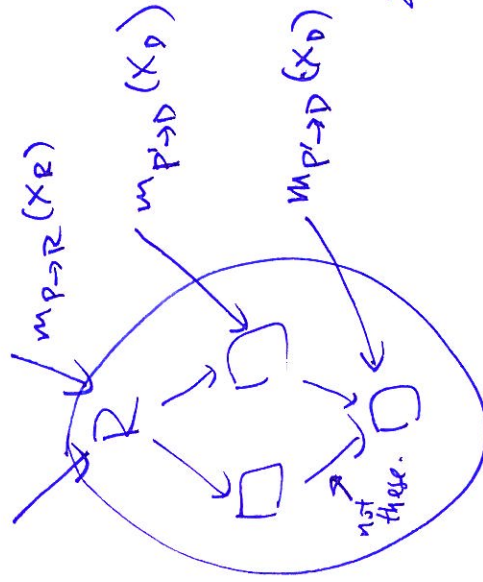
Parent-to-Child algorithm

For every $P \rightarrow R$ in region graph

message $M_{P \rightarrow R}(X_R)$

beliefs

$$b_R(X_R) \propto \prod_{a \in A_R} f_a(x_a) \prod_{P \in \mathcal{P}(R)} M_{P \rightarrow R}(X_R) \prod_{D \in \mathcal{D}(R)} M_{P \rightarrow D}(X_D)$$

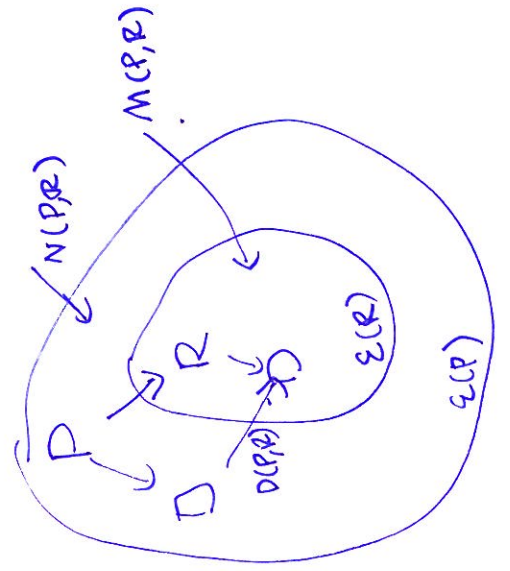


all messages from outside $\mathcal{E}(R)$

$$b_P(X_P) \propto \prod_{a \in A_P} f_a(x_a) \prod_{(I,J) \in N(P,R) \cup M(P,R)} M_{I \rightarrow J}(X_J)$$

$$\neq b_R(X_R) \propto \prod_{a \in A_R} f_a(x_a) \prod_{(I,J) \in \mathcal{D}(P,R) \cup M(P,R) \cup \{C(P,R)\}} M_{I \rightarrow J}(X_J)$$

$$\sum_{X_P \setminus X_R} b_P(X_P) = b_R(X_R)$$



$$\sum_{X_P \setminus X_R} \prod_{a \in A_P \setminus A_R} f_a(x_a) \prod_{N(P,R)} M_{I \rightarrow J}(X_J) = M_{P \rightarrow R}(X_R) \prod_{\mathcal{D}(P,R)} M_{I \rightarrow J}(X_J)$$

$\mathcal{D}(R)$ = descendants
 $\mathcal{E}(R)$ = $R \cup \mathcal{D}(R)$
 $\mathcal{P}(R)$ = parents

Further Details of message-passing

① initialization

- random
- uniform

② termination

- convergence criteria
 - message change
 - belief change

③ damping

$$M_{P \rightarrow R}^{\text{new}} := \alpha M_{P \rightarrow R}^{\text{old}} + (1-\alpha) M_{P \rightarrow R}^{\text{update}}$$

$$\text{OR } := (w^{\text{old}})^{\alpha} (M_{P \rightarrow R}^{\text{update}})^{1-\alpha}$$

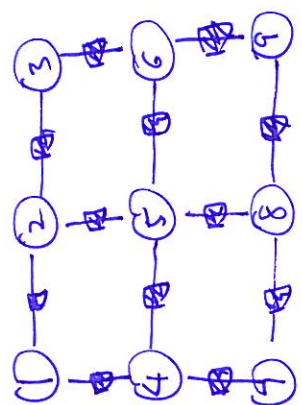
$$M_{P \rightarrow R}^{\text{update}} = \frac{\sum_{X \in \mathcal{X}_R} \prod_{A \in \mathcal{A}_{P \rightarrow R}} f_A(X_A) \prod_{N \in \mathcal{N}(P,R)} M_{I \rightarrow J}^{\text{old}}(X_J)}{\prod_{D \in \mathcal{D}(P,R)} M_{I \rightarrow J}^{\text{update}}(X_J)}$$

start from bottom of region graph, work your way up.

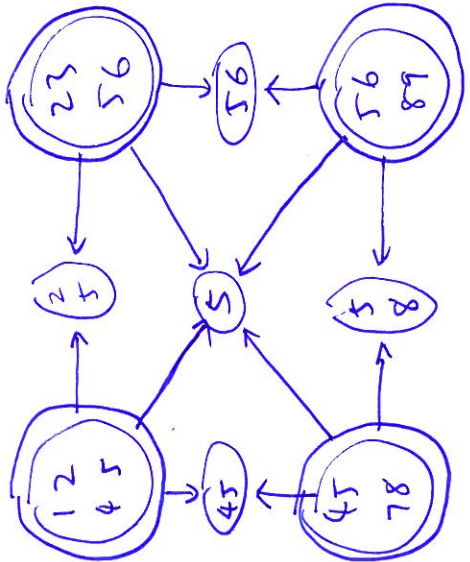
④ normalization

- normalize messages
- or work in log domain.

Construction Method 1

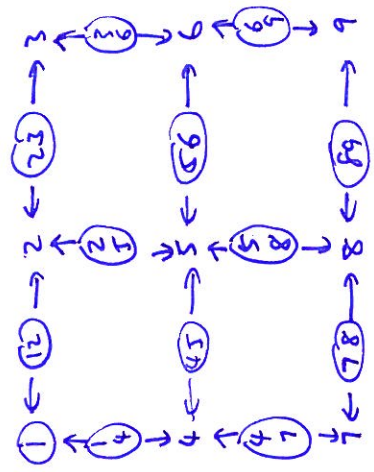
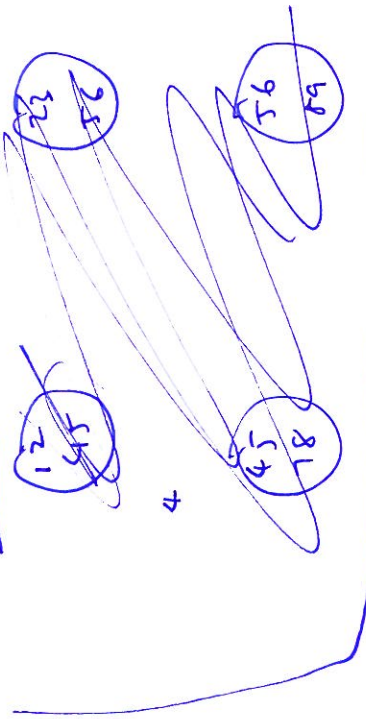


Junction graph

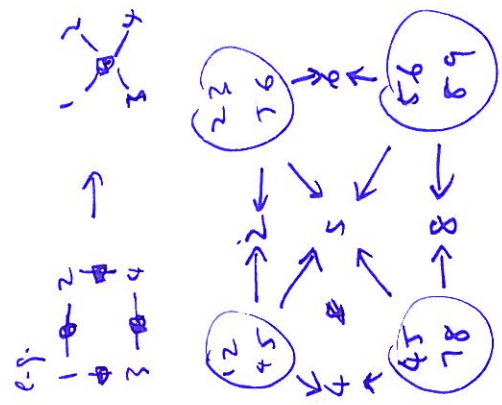


for each i remove it from regions until induced region graph forms a tree

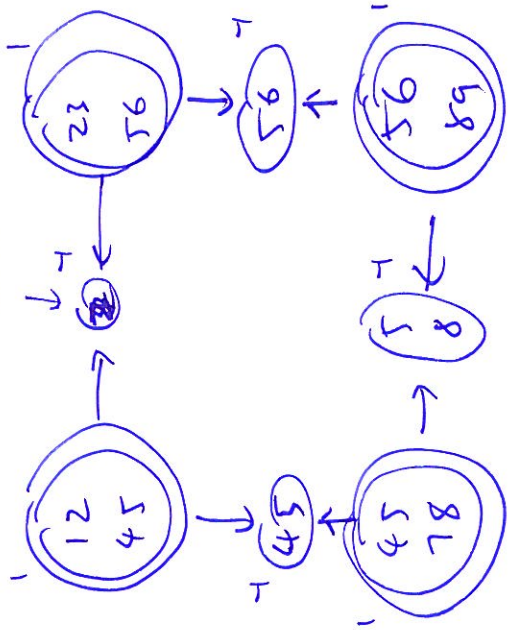
Bethe



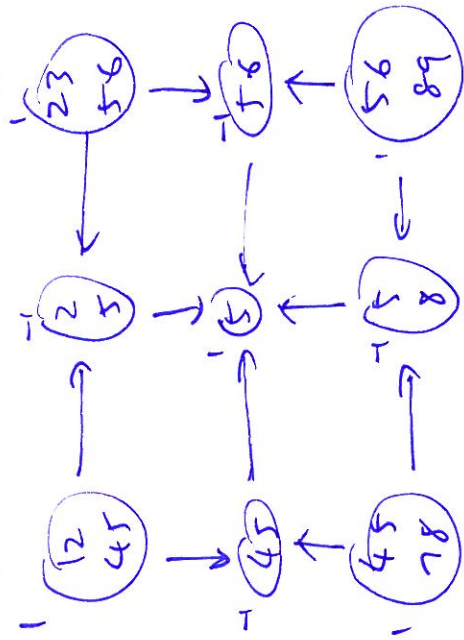
or: merge factors first.



remove 5.



Cluster Variation method



form all intersects, intersects of intersects

Note: need not be maxent-normal

Loop-based approximation,

- start with Bethe
- add strongest coupled loops first,

Maintaining non-singularity.

(~~temp~~ cycle space)

loops are non-singular (independent in cycle space))

- stop at maximal.

Tree EP

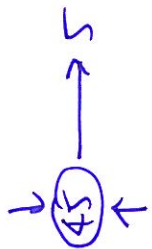
- take spanning tree
- loops are those in spanning tree & individual edges

Planar Graphs

- use faces of planar graph

produce non-singular loop-based approximations

Parent-child updates — Cluster Variation method 3x3 grid

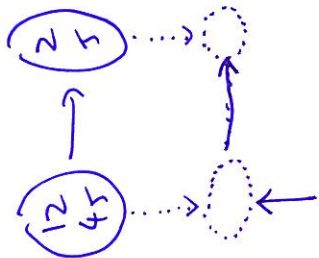


$$M_{45 \rightarrow 5}(X_5) = \sum_{X_4} f_{45}(X_4 \times X_5) M_{1245 \rightarrow 45} M_{4578 \rightarrow 45} M_{4578 \rightarrow 45} (X_4 \times X_5) M_{4578 \rightarrow 45} (X_4 \times X_5)$$

$$N(45, 5) = \left\{ \begin{matrix} (1245, 45) \\ (4578, 45) \end{matrix} \right\}$$

$$D(45, 5) = \{ \}$$

Similarly $M_{56 \rightarrow 5}, M_{58 \rightarrow 5}, M_{58 \rightarrow 5}$.

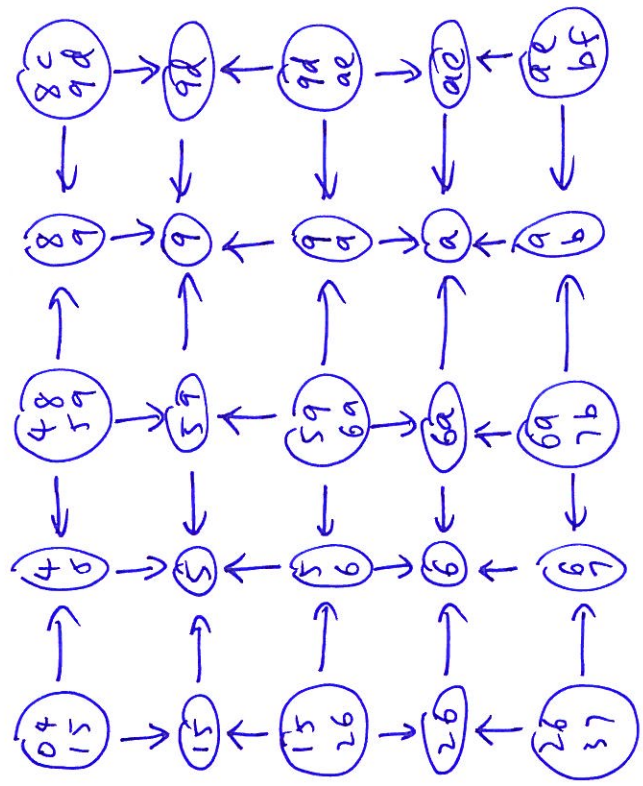
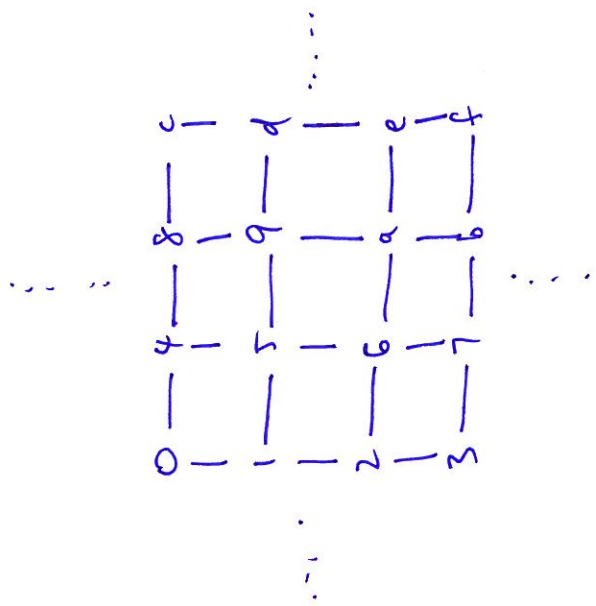


$$M_{1245 \rightarrow 25}(X_2 \times X_5) M_{45 \rightarrow 5}(X_5) = \sum_{X_1 \times X_4} f_{12} f_{45} f_{14} f_{45} M_{1245 \rightarrow 45} (X_4 \times X_5)$$

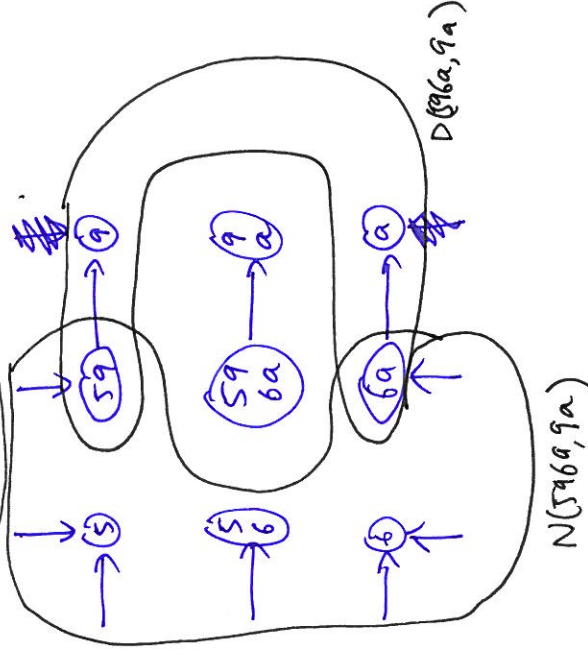
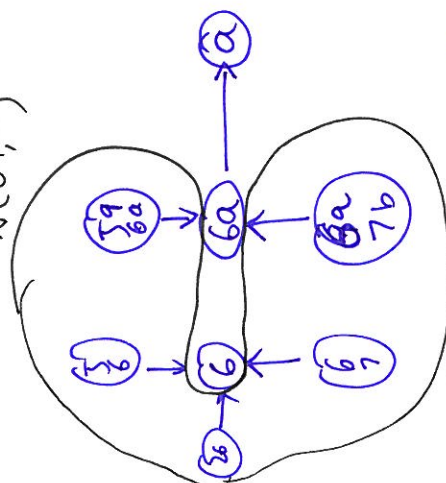
$$N(1245, 45) = \left\{ (4578, 45) \right\}$$

$$D(1245, 45) = \left\{ (45, 5) \right\}$$

Parent-child updates - Cluster Variation Method nxm grid



$N(6a, a)$



$$N(6a, a) = \{ (596a, 6a), (6a7b, 6a), (56, 6), (26, 6), (67b) \}$$

$$D(6a, a) = \{ \}$$

$$M(6a, a) = \{ (9a, a), (ae, a), (ab, a) \}$$

$$\sum_{x_5 \times x_6} f_{56} f_{59} f_{6a} f_{7b} \quad M_{4589 \rightarrow 59} \quad M_{45 \rightarrow 5} \quad M_{1256 \rightarrow 58} \quad M_{26 \rightarrow 6} \quad M_{67 \rightarrow 6} \quad M_{6a7b \rightarrow 6a}$$

$$M_{596a \rightarrow 9a} (X_9 X_a) =$$

$$M_{59 \rightarrow 9} \quad M_{6a \rightarrow a}$$