

today: inference → graph eliminate

general themes in this class

A) representation → DGM
 parameterization → UGM
 parameterization → exponential family } model high. dim. distribution

B) inference $p(x_Q | x_E)$
 "query" "evidence"
 → today: elimination algorithm
 → next: sum-product / belief propagation (for trees)
 later: approximate inference e.g. MCMC variational

C) statistical estimation / learning → MLE
 maximum entropy
 method of moments

Inference

want to compute

a) marginal: $p(x_F)$ for some $F \subseteq V$

b) conditional: $p(x_Q | x_E)$
 "query" "evidence"

c) for UGM = partition function $Z = \sum_{x \in V} (\prod_C \psi_C(x_C))$

why?

• missing data $p(x_{\text{unobs}} | x_{\text{obs}})$

example →



• prediction $p(x_{\text{future}} | x_{\text{past}})$

• "latent cause"

$p(x_{\text{cause}} | x_{\text{obs}})$

eg



* also related to inference:

$$\underset{x_F}{\text{argmax}} p(x_F | x_E)$$

↓
F could be big here (e.g. speech recognition)

* inference is also needed during estimation (parameter fitting MLE)
[e.g. during E-step $p(z|x)$]

* present inference alg. for UGM [for simplicity and more generality]
but note that sometimes it's more efficient to work directly on DGM

make DGM \rightarrow \subseteq UGM via moralization

$$\begin{aligned} \text{i.e. } p \in \text{DGM}; \quad p(x) &= \prod_i p(x_i | x_{\pi_i}) & \text{moralize: } C_i &\triangleq \{i\} \cup \pi_i \\ &= \frac{1}{Z} \prod_i \psi_{C_i}(x_{C_i}) & \psi_{C_i}(x_{C_i}) &\triangleq p(x_i | x_{\pi_i}) \\ & & Z &= 1 \end{aligned}$$

graph elimination alg. (for inference)

• consider $p \in \mathcal{J}(\mathcal{G})$
undirected $p(x) = \prod_{C \in \mathcal{C}} \psi_C(x_C)$

Say want to compute $p(x_F)$ for $F \subseteq V$ "query nodes"

main trick: use distributivity of \oplus over $\odot \rightarrow \odot(a \oplus b) = \odot a \oplus \odot b$

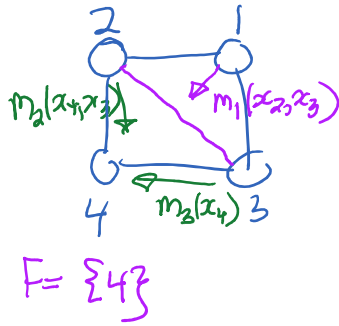
$$\sum_{x_1, x_2} f(x_1) g(x_2) = \left(\sum_{x_1} f(x_1) \right) \left(\sum_{x_2} g(x_2) \right) \quad [\text{convince yourself}]$$

↙ first term of product
e.g. $\sum_{x_2} f(x_1) g(x_2)$

more generally

$$\sum_{x_1, \dots, x_n} \prod_i f_i(x_i) = \prod_{i=1}^n \left(\sum_{x_i} f_i(x_i) \right)$$

$$O(k^n \cdot n) \quad O(k \cdot n)$$



$$\begin{aligned}
 p(x_4) &= \frac{1}{Z} \sum_{x_1, x_2, x_3} \psi(x_1, x_2) \psi(x_3, x_4) \psi(x_1, x_3) \psi(x_2, x_4) \\
 &= \frac{1}{Z} \left(\sum_{x_3} \psi(x_3, x_4) \left(\sum_{x_2} \psi(x_2, x_4) \underbrace{\left(\sum_{x_1} \psi(x_1, x_2) \psi(x_1, x_3) \right)}_{m_1(x_2, x_3)} \right) \right) \\
 &= \frac{1}{Z} m_3(x_4)
 \end{aligned}$$

↑ stored as table
↑ message

→ last message is proportional to marginal $p(x_4)$

$$\sum_{x_4} m_3(x_4) = Z$$

general dg: graph Eliminate

general alg: graph Eliminate

init.
 a) choose an elimination ordering s.t. F are the last nodes
 b) put all $\psi_c(x_c)$ on "active list"

"update" c) repeat in order of variables to eliminate
 (say x_i is variable to eliminate)

1) remove all factors from active list with x_i in it & take their product

$$\text{i.e. } \prod_{\substack{\psi \text{ st. } \\ i \in \alpha}} \psi_\alpha(x_\alpha)$$

2) sum ^{over} x_i to get a new factor $m_i(x_{S_i})$ (think as $\psi_{S_i}(x_{S_i})$)

i.e. S_i are all variables in those factors except i

$$\text{get } m_i(x_{S_i}) \triangleq \sum_{x_i} \prod_{\substack{\psi \text{ st. } \\ i \in \alpha}} \psi_\alpha(x_\alpha)$$

new degree to sum over $S_i \cup \{i\}$

$$S_i \triangleq \left(\bigcup_{\substack{\alpha \text{ st. } \\ i \in \alpha}} \alpha \right) \setminus \{i\}$$

3) put back $m_i(x_{S_i})$ in active list ($\psi_{S_i}(x_{S_i})$)

"normalize" d) last product of factors left has only $x_F \Rightarrow$ proportional $p(x_F)$

memory needed? $\approx 2^{\max_i |S_i|} \cdot (\# \text{ of factors})$

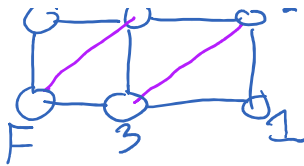
computational cost $\approx 2^{\max_i |S_i| + 1} \cdot n$

later, related "treewidth" of a graph

14h34

* "augmented graph" \rightarrow graph obtained by running graph Eliminate + keeping track of all edges added (for a fixed ordering)

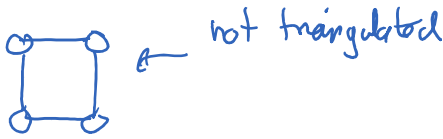




"augmented graph" after graphEliminate is always a triangulated graph

def: graph with no cycle of size 4 or more that cannot be broken by a "chord"

↳ on edge between two non-neighboring nodes in the cycle

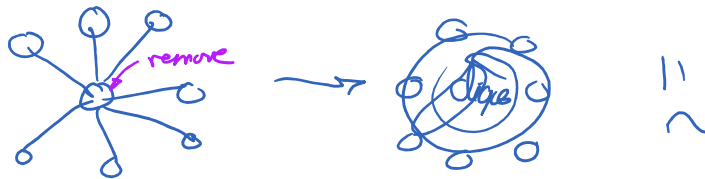


treewidth of a graph $\triangleq \min \left\{ \text{size of biggest clique in possible elimination orderings} - 1 \right\}$

convention treewidth(tree) = 1

both memory & running time of graphEliminate is dominated by $2^{\text{size of biggest clique}}$
best ordering give $\approx 2^{\text{treewidth}+1}$

not all orderings are good



bad news:

a) NP hard to compute treewidth of general graph (or find best ordering)

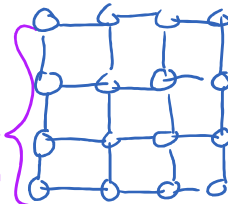
b) NP hard to do (exact) inference in general UGM

⇒ approximate methods

example:

treewidth of a grid $\approx \sqrt{|V|}$

side treewidth



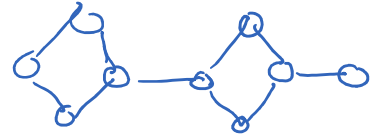
good news:

* inference is linear time for trees (treewidth=1) ("sum-product alg.")
 $\rightarrow |V| + |E|$

(HMM, Markov chain, etc.)

* efficient for "small treewidth graph"

→ use junction tree alg.



inference on trees

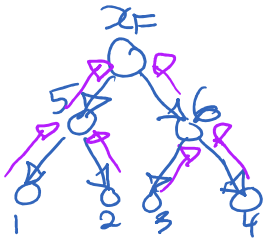
graphEliminate on a tree

good order is to eliminate leaves first

$$p(x) = \frac{1}{Z} \prod_i \psi_i(x_i) \prod_{i,j \in E} \psi_{ij}(x_i, x_j)$$



order: make a directed tree by using x_F as a root
singleton



$$m_{i \rightarrow j}(x_j) = \sum_{x_i} \psi_i(x_i) \psi_{ij}(x_i, x_j) \prod_{k \in \text{children}(i)} m_{k \rightarrow i}(x_i)$$

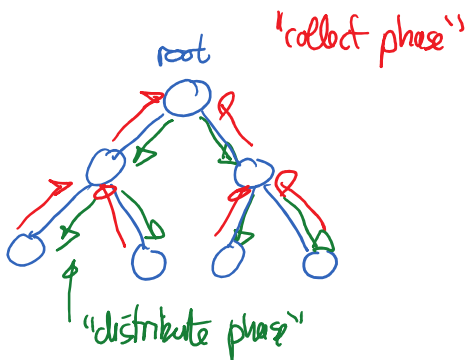
new factors containing i on or true list

Sum-product alg. (for trees)

alg. to get all node/edge marginals cheaply by storing (caching)

↳ re-using messages

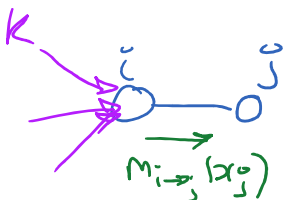
(dynamic programming)



goal: $\forall i, j \in E$, compute $m_{i \rightarrow j}(x_j)$

$m_{j \rightarrow i}(x_i)$

rule: i can only send message to neighbor j when it has received all messages from other neighbors



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



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

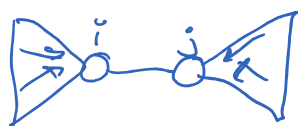
↑
neighbors

at end
(node marginal)

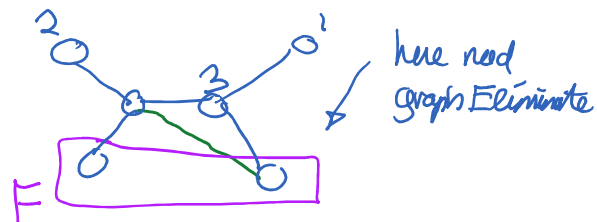
$$p(x) \propto \prod_{j \in N(i)} m_{j \to i}(x_i) \psi_i(x_i)$$

normalize $Z = \sum_{x_i} (\dots)$

(edge marginal)



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



Sum-product schedules:

a) above, collect/distribute schedule

b) (flooding) parallel schedule:

- 1) initialize all $m_{i \to j}(x_j)$ messages to uniform dist. $\forall (i,j) \neq (j,i) \in E$
- 2) at every step (in parallel) compute $m_{i \to j}^{new}(x_j)$

as if the neighbor messages were correctly computed

→ one can prove that "diameter of tree" # of steps

all messages are correctly computed (for a tree) (and are exact pt.)

