

## - Review + readings

### Jordan Ch 3

- query node (unshaded)
- conditioning nodes - evidence nodes (dark shaded)
- marginalisation nodes - (lightly shaded)

(\*) Marginalisation of joint:-

$$p(x_1, \dots, x_n)$$

- naive  $\rightarrow$  via C.P. table

- complexity  $O(k^n)$

n - no. of states

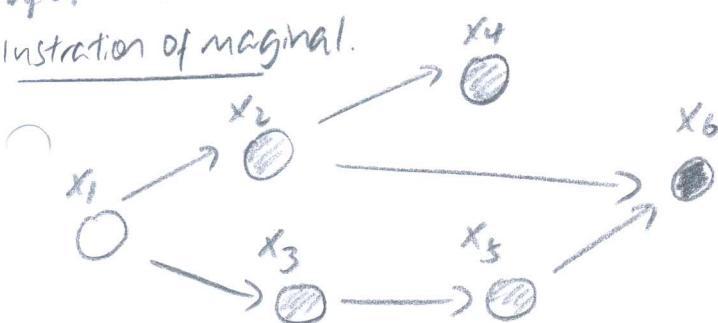
n - no. of r.v.s.

(\*) via factorisation law, use distrib.:-

- one-off summation has cost  $R'$  r-intermediate term  
(no. of variables)
- for n summations/elimin.
- overall cost  $O(nk')$

- graph theoretic methods for outlin. of r.

- illustration of marginal.



- (\*) Key idea: - 'pushing sums inside products'  
(6)
- sum and product op.  $\rightarrow$  commutative
  - products  $\rightarrow$  associative
  - scope, summation, product

- Fix  $x_6 = \bar{x}_6$

$$p(x_1, \bar{x}_6) = \sum_{x_2} \sum_{x_3} \sum_{x_4} \sum_{x_5} p(x_1) p(x_2|x_1) p(x_3|x_1) p(x_4|x_2) p(x_5|x_3) p(\bar{x}_6|x_2, x_5) \quad (3.8)$$

$$= p(x_1) \sum_{x_2} p(x_2|x_1) \sum_{x_3} p(x_3|x_1) \sum_{x_4} p(x_4|x_2) \sum_{x_5} p(x_5|x_3) p(\bar{x}_6|x_2, x_5) \quad (3.9)$$

$$= p(x_1) \sum_{x_2} p(x_2|x_1) \sum_{x_3} p(x_3|x_1) \sum_{x_4} p(x_4|x_2) m_5(x_2, x_3)$$

$$\text{i.e. } m_5(x_2, x_3) := \sum_{x_5} p(x_5|x_3) p(\bar{x}_6|x_2, x_5) \quad \begin{aligned} & \text{- omit dependence on } \bar{x}_6 \text{ for} \\ & \text{not prop.} \end{aligned}$$

(\*) Algebraic removal  
of  $x_5$  via elimination  $\rightarrow$  node removal.

$$p(x_1, \bar{x}_6) = p(x_1) \sum_{x_2} p(x_2|x_1) \sum_{x_3} p(x_3|x_1) m_5(x_2|x_3) \sum_{x_4} p(x_4|x_2)$$

$$\text{define } m_4(x_2) := \sum_{x_4} p(x_4|x_2)$$

$$= p(x_1) \sum_{x_2} p(x_2|x_1) m_4(x_2) \sum_{x_3} p(x_3|x_1) m_5(x_2, x_3)$$

$$\text{define } m_3(x_1, x_2) = \sum_{x_3} p(x_3|x_1) m_5(x_2, x_3)$$

$$= p(x_1) \sum_{x_2} p(x_2|x_1) m_4(x_2) m_3(x_1, x_2)$$

$$\text{define } m_2(x_1) = \sum_{x_2} p(x_2|x_1) m_4(x_2) m_3(x_1, x_2)$$

$$p(x_1, \bar{x}_6) = p(x_1) m_2(x_1) \quad (\text{reduced joint}) \quad (3.14)$$

$$p(\bar{x}_6) = \sum_{x_1} p(x_1) m_2(x_1) \quad (3.15)$$

$$p(x_1|\bar{x}_6) = \frac{p(x_1, \bar{x}_6)}{p(\bar{x}_6)} = \frac{p(x_1) m_2(x_1)}{\sum_{x_1} p(x_1) m_2(x_1)} \quad (3.16)$$

formal trick for including  
conditioning on a fixed observation.

$x_i$  - evidence node

$\bar{x}_i$  - observed value

evidence potential:  $\delta(x_i, \bar{x}_i) = \begin{cases} 1 & x_i = \bar{x}_i \\ 0 & \text{otherwise} \end{cases}$

evaluate  $g(x_i)$  at  $\bar{x}_i$  :-

$$g(\bar{x}_i) = \sum_{x_i} g(x_i) \delta(x_i, \bar{x}_i)$$

A way of putting condition into sum-product form

$$m_6(x_2, x_5) = p(\bar{x}_6 | x_2, x_5) = \sum_{x_6} p(x_6 | x_2, x_5) \delta(x_6, \bar{x}_6)$$

(\*) In general:-

For subset cond. nodes  $E$ , config.  $\bar{x}_E$ , compute  $p(x_E | \bar{x}_E)$

Total evidence pt:  $\delta(x_E, \bar{x}_E) = \prod_{i \in E} \delta(x_i, \bar{x}_i)$

(\*) Understanding the elimination algorithm in a formal, algorithmic sense

(it is a formalisation of the steps involved  
in the earlier example; by exploiting  
properties of summation, product operator;  
together with the 'scope' of a function.)

(\*) Do you understand elimination:-

- i) informal algebraic ✓
- ii) formal algorithmic (~) ✓
- iii) graph-theoretic ✓

(\*) Algorithmic aspects of elim.

(\*) Notes don't specify

- some supplements on elimination (\*) place all factors in active list → sum-product variable elim

(\*) Check you understand elimination or not-chain structured F

(\*) When eliminating a node  $z_i$ ; we split active list into  $F'$  and  $F''$

depending on whether the factor/potential/ICPD contains  $z_i$ .

- Take product of those factors  $\phi$  that contain the node we want to eliminate as arguments i.e.  $\phi(z_i, \dots)$  (i.e. those in  $F'$  on slides).

- sum over  $z_i$  over the product of factors

$$\left( \prod_{\phi \in F'} \phi(z_i, \dots) \right)$$

(\*) Then repeat the subroutine

$$\sum_{z_i}$$

- This resulting intermediate factor does NOT contain  $z_i$

and is placed in set  $\{T\}$

- The partition of factors with no  $z_i$  and  $\{T\}$  are combined;  $z_i$  having been eliminated → this forms new  $F$  (active list)

Non-chain structured elimination;  
and graph-theoretic elim.

- understood in Jordan and my notes ✓ ↴ (\*)
- Note that  $DGM \rightarrow UGM$  in this process
- (\*)- define elimination ordering
- remember conditioning/evidence potential.
- use scoping to take sum over a product of factors containing variable you wish to eliminate (delete on graph)
- This summation/marginalisation eliminates the node/variable.
- But creates an intermediate factor which you need to be a function of new nodes. (these need connecting/moralising)

(\*) Finally after all nodes other than query node eliminated;  
apply Bayes law over joint, conditional, normalised joint (mag.) ○  
to get query c.p.

understand reconstituted graph as graph which is same as original; but with addition of edges that are newly created in query elimination step.

elimination clique - the fully connected subset of nodes that include the neighbours of eliminated node (which get moralised and connected after elimination of relevant node); and eliminated node.

marginalisation by removing r.v. from joint:

↳ sum of products of all factors that contain r.v. we wish to eliminate.

↳ this couples all the other r.v.s. that appear in those factors

e.g.  $\sum_{x_5} \psi(x_3, x_5) M_6(x_2, x_5)$  creates an intermediate factor involving  $(x_2, x_3)$

i.e. inducing dependency between  $x_2$  and  $x_3$  (which get connected)

(\*) In graph elimination algorithm; the elimination clique are those nodes l/r.v.s. which in the variable elimination algorithm, are the scope of the generated intermediate term/r/factor.  
(excluding eliminated node)

- See slides → these are very intuitive and easy to understand

(\*) Overall complexity of variable elimination algo is determined by the no. of variables in the largest elimin. clique

- recall  $O(nR^r)$  - r-no. of intermediate term l.r.s.!

(\*) That is computational complexity question → graph theoretic

(\*) largest elimination clique → well studied ✓

(\*) tree width parameter  $R = \left(\min_i w_{G,i} - 1\right)$  minimum achievable value of cardinality of

(\*) ~~largest elim. clique that is small as possible~~ overall elim. orderings

(\*) BUT finding best elim ordering of a graph (i.e. an elim. ordering) that achieves tree width → NP-hard.

(\*) -NP-hardness → constraints on comput. efficiency of elimination algo.

→ comput. bottleneck of elimination algo

→ graph-theoretic (undirected graph eliminate) algo  
allows practically useful tool for assessing severity.

(\*) For a given elimination ordering, cheap estimate of running time of var. eliminate by running graph-elimination algo (NGE).

- If graph-theoretic algo yields elim cliques of small cardinality; elim is viable.

(\*) Limitations of variable elimination:

i) Single active list → inefficient traversal (S): buckets.

ii) Single query node e.g. cond. prob of all non-evidence nodes in graph.

## a way of reducing redundant computation

- Next:
  - i) sum-product algo. (marginals for trees)
  - ii) Junction-tree algo. (marginals for graphs)

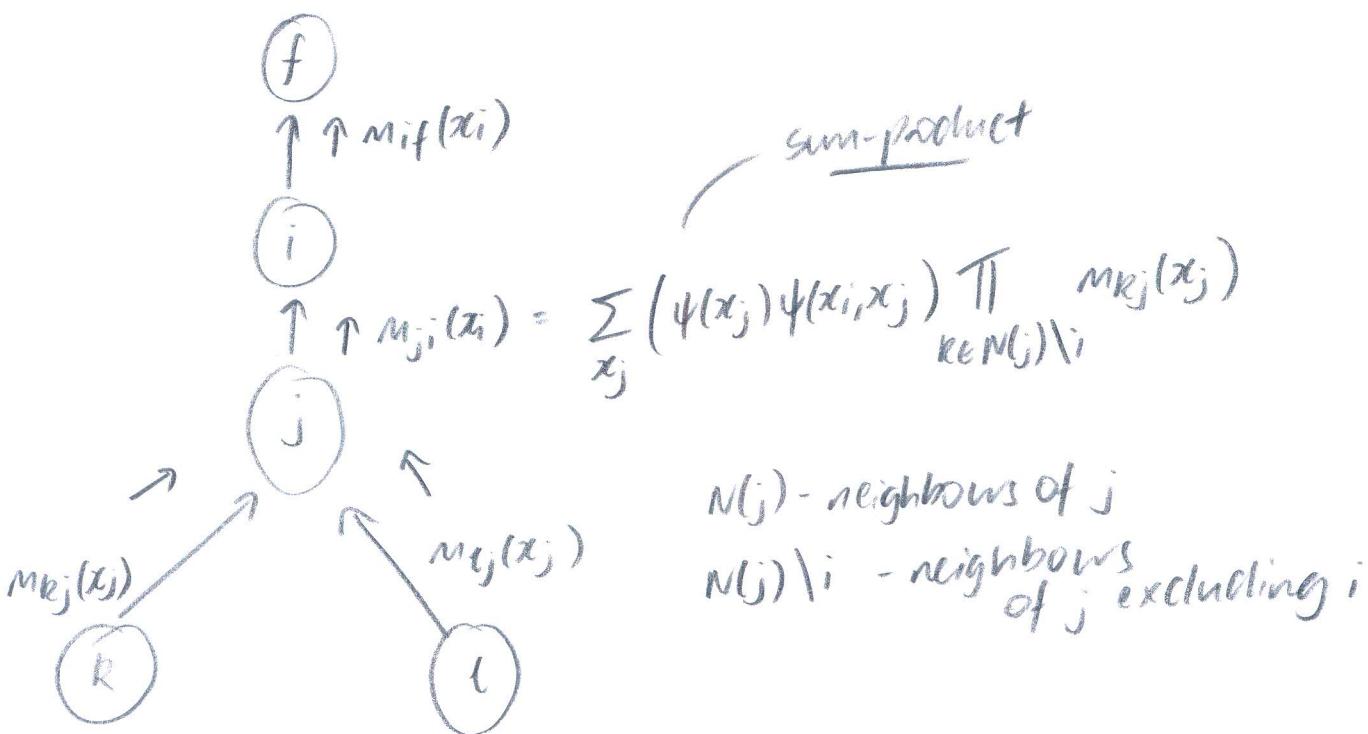
(\*) A clique tree: connects elimination cliques via undirected edges  
based on whether subsets of r.v.s. share.

(\*) Note correspondence between intermediate factors and elimination cliques

(\*) need to message passing on clique trees (Jordan)

(\*) review / familiarise eliminate algo

4.5.2020 10.708 (message passing)

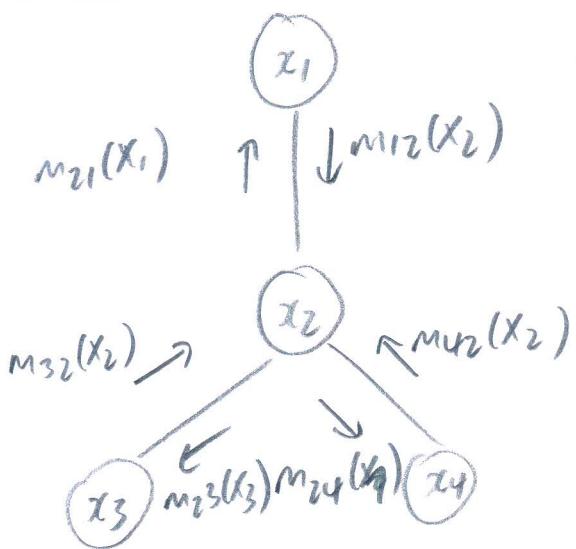


- message from  $j$  to  $i$  i.e.  $m_{ij}(x_i)$ ; collect messages flowing into it i.e.

$m_{kj}(x_j)$  and  $m_{lj}(x_j)$

- don't use  $\rho$ , use  $\psi$ ; singleton, pairwise potentials

## message passing protocol



## 2 pass algorithm

- pick a node, make it root
- create messages from the leaves
- collect them to the root node
- make contributions

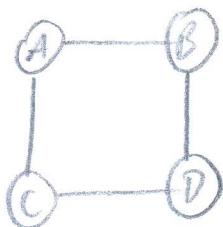
- 2 passes done
- All necessary terms from message passing produced
- Free transfer queries

## ① M message passing

- MPP - collect messages from leaves, pass from top. (schedule)
- Belief propagation - parallel synchronous implementation
- localist view of doing message passing → inspired approximations
- theorem: message passing guarantees all marginals on the tree:-

$$m_{ji}(x_i) = \sum_{x_j} \left( \psi(x_j) \psi(x_i, x_j) \prod_{k \in N(j) \setminus i} m_{kj}(x_j) \right)$$

## ② xy message passing on a non-tree



- message passing only on trees
- no guarantees on non-trees

Turn non tree → tree

↳ do graph relabelling, identify cliques, turn into clique tree; messages on c.t.

- If you can turn an arbitrary gm → junction tree
- can run message passing on tree cliques; but tree cliques may be huge



## 15. Parameter estimation (S.2020)

### - Ex: Parameter estimation in PGMs (at a clean reference point)

• Start with named/Indexed r.v.s.

• And training data (IID); fully observed - every r.v. has instantiation

#### 1.) Structure learning

- In principle, possible to learn structure from gdata
- often experts.

#### 2) Parameter estimation

PGMs - nos in CPT or potential function values

$\text{CPT}_M$

Ex: learning completely observed PGM - fairly trivial

PGMs (partially obs.)

- directed - focus on this

estimation principles

NLE  $\rightarrow$  classical setting statistics

Ex: ML has 'gone beyond this'; present a universal, standardised view that specifies these.

e.g. reinforcement learning - intrinsic/extrinsic

adversarial learning - adversarial score

- traditionally, characterise learning / param est in terms of statistical consistency etc.

- in more modern engineering ML  $\rightarrow$  may be compromised

#### Simplest case

- PGMs where structure is known

- Parameter learning for BN

(\*) Analytically write down loss function  $\rightarrow$  likelihood of data,  $g$  as a function of the parameters  
- Probability (likelihood)? as product of many local terms.

Ex: Point you to further reading

- Building blocks of GM:  $\text{W} \text{ A}$

- single node GM (e.g. root node in tree)  $\rightarrow$  supp slides

- instances of exponential family distri

(\*) ultimately; parameter (probability) is empirical frequency count

2 node graphical models  $\rightarrow$   $\text{D} \text{ A}$

$\hookrightarrow$  supp slides

Ex: Patterns underlying parameter learning for otherwise

exponential family

Ex: class interested in building blocks.

$$p(x|\eta) = h(x) \exp \{ \eta^T I(x) - A(\eta) \}$$

$I(\cdot)$  and  $\eta(\cdot)$  most important terms; note dot product

$A(\cdot)$  - log normaliser  
(canonical)

Ex: Estimation of parameters  $\eta$  only require  $I(\cdot)$  (i.e. sufficient statistic)

• examples - MVG

- canonical moments related

-  $\mu$  and  $\Sigma$  - moment parameters

-  $\text{D} \text{ A} \rightarrow$  Review exp.

exponential family representation

$\text{vec}(\cdot)$  - ~~measure~~ turn the entity into a vector

$\eta$  and  $I(\cdot)$  - of same dimensionality