

100f - OGMS

- undirected graphs $\textcircled{1} P(X_1, \dots, X_8) = \frac{1}{Z} \exp \{ \phi(X_1) + \dots + \phi(X_6, X_5, X_8) \}^2 \textcircled{2}$
- OGMS \rightarrow expert systems (complex causality webs)

HMM / dishonest casino

⑧ key areas to focus on \rightarrow 10708 learning objectives

- Areas to focus on:- - These should be included in 'lectures'
- BN definition
- formal specification of factorisation
- distinction between functional form of factorisation for OGMs, UGMS
- understand significance of conditional independence in PGMs
- know the 3 local structures visually, what C.I. relations they represent
- nuances of I-Maps; appl. to example (*)
- local Markovian assumptions
- d-separation: semantics, operation
- global Markovian assumptions - Bayes ball.
- Equivalence theorem; difference with I-maps (*)
- Soundness + completeness in context of (*).

Jordan (2003) Ch 2 (read; annotate/supplement)

- very clear formal setup
- in general, not making use of additional 'structure'; a joint probability distn of n random variables, each taking on r states requires full specification of r^n states in an n -dimensional probability table
- graphical models represent this more economically with ass. about structure via 'local relationships'
- formal notation:-

- directed graph $G(V, E)$ V -set of nodes E -edges G -acyclic
- 1-to-1 mapping r.v.s to nodes
- # For each node $i \in V \rightarrow$ r.v. X_i
- $V = \{1, \dots, n\}$ and r.v.s. = $\{X_1, \dots, X_n\}$

- each node has a set of parent nodes, may be empty
- for each node $i \in V$, we denote π_i as the set of parents of node i
- the set of random variables X_{π_i} are the parents of the r.v. X_i
- parent-child relations \rightarrow locality \rightarrow economical rep. of joint p.d.
- to each node i ; associate function $f_i(x_i, x_{\pi_i})$ with properties:-
 $f_i(x_i, x_{\pi_i}) \geq 0$ and $\sum_{x_i} f_i(x_i, x_{\pi_i}) = 1$
- the pre-conditional probability formalism is:
- let $V = \{1, \dots, n\}$; given a set of functions $\{f_i(x_i, x_{\pi_i}) : i \in V\}$, define joint p.d.:-
 $p(x_1, \dots, x_n) \stackrel{(2.1)}{=} \prod_{i=1}^n f_i(x_i, x_{\pi_i})$
- Have to verify definition obeys joint p.d. constants
- choice of numerical values for $f_i \Rightarrow$ generation of specific joint p.d.
- Ranging over all numerical values \Rightarrow family of joint prob distr. associated with graph G (*)
- This family can be characterised as products of local functions or graph-theoretically (edge patterns)
- (*) The relationship between different ways of characterising family of probability distributions associated with a graph that is key to PLM
- conditional probability formalism:

$$p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | x_{\pi_i}) \quad (2.2)$$

- $p(x_i | x_{\pi_i})$ - local condition prob. associated with graph G
- Building blocks merely joint p.d. synthesised associated with G .
- Example given: (see diag).

$$p(x_1, \dots, x_6) = p(x_1)p(x_2|x_1)p(x_3|x_1)p(x_4|x_2)p(x_5|x_3)p(x_6|x_2, x_5) \quad (2.3)$$

Product of LCD

representational economy:

- r^n joint prob. table $\rightarrow r^{m+1}$ table
 (no structure) (graph structure)

- m_i - no. of parent nodes of X_i
- local cond. dist. $X_i - (m+1)$ dim table
- each node/r.v. has r values

② exponential growth in n exchanged for exponential growth in m_i

(variables in domain)

(no. of parents of indiv. nodes X_i (fan-in))

- small fan-in \rightarrow enormous reduction in complexity

• ③ But there is more; not just data structure promise, but inferential machinery

2.1.1 conditional independence

independence:

x_A and x_B are independent, i.e. $x_A \perp\!\!\!\perp x_B$ if :-

$$p(x_A, x_B) = p(x_A)p(x_B) \quad (2.4)$$

conditional independence

x_A and x_C are conditionally independent, given x_B if:-

$$p(x_A, x_C | x_B) = p(x_A | x_B)p(x_C | x_B) \quad (2.5)$$

OR

$$p(x_A | x_B, x_C) = p(x_A | x_C) \quad (2.6)$$

} moving
from
here is
④

for all x_B : $p(x_B) > 0$

(2.5) - independence ; with addition of conditioning on $|x_B$
 formula

(2.6) - "probability of x_A given x_B and x_C is equal to
 probability of x_A given x_C " i.e. conditioning on x_B and x_C
 is the same as omitting x_B (conditioning variable) in prob.

- establish independence / conditional independence requires
 factorising the joint probability distribution

- PGM: representing a p.d. within gm formalism involves making independence assumptions, which are embedded in the structure of the graph
- This graphical structure, other indep. relations can be derived, reflecting that certain factorisations of joint p.d. imply other factorisations
- Factorisations - read off via graph search algorithms
- Graphical structure encodes conditional independence:-
- Chain rule of prob: PMF in a general factored form, given an ordering on $\{x_1, \dots, x_6\}$:-

$$p(x_1, \dots, x_6) = p(x_1)p(x_2|x_1)p(x_3|x_1, x_2)p(x_4|x_1, x_2, x_3)p(x_5|x_1, x_2, x_3, x_4)p(x_6|x_1, x_2, x_3, x_4, x_5)$$
- In an arbitrary node ordering
- In general: $p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i|x_1, \dots, x_{i-1}) \quad (2.7)$

- ① Compare (2.2) and (2.7) \Rightarrow conditioning variables dropped
- ② Missing variables in local c.p. functions \Rightarrow missing edges underlying graph
- ③ It is this transfer of interpretation from missing variables \rightarrow missing edges that underlies the probabilistic interpretation for missing edges in graph in terms of c.i.
- Formally:
 - Define an ordering I of nodes in a graph to be topological if for every node $i \in V$ the nodes in T_i appear before i in the ordering.
 - E.g. $I = \{1, 2, 3, 4, 5, 6\}$ (topological ordering for graph in 2.1)
- ④ Let v_i denote set of nodes that appear earlier than i in the ordering I , excluding the parent nodes T_i .
 - e.g. $v_5 = \{1, 2, 4\}$
 - v_i necessarily contains ancestors of node i (excluding parents T_i) and may contain other nondescendant nodes also

Given a topological ordering π for a graph G we associate to the graph the following set of c.i. statements:-

$$\{x_i \perp\!\!\!\perp x_{\pi_i} \mid X_{\pi_i}\} \quad \text{for } i \in V \quad (2.8)$$

"Given the parents of a node (X_{π_i}) , the node (X_i) is independent of all earlier nodes in the ordering (X_{π_i}) ".

The example encodes following mdp. ass:-

$$(x_i \perp\!\!\!\perp \emptyset \mid \emptyset) \quad (2.9)$$

$$(x_2 \perp\!\!\!\perp \emptyset \mid x_1) \quad (2.10)$$

$$x_3 \perp\!\!\!\perp x_2 \mid x_1 \quad (2.11)$$

$$x_4 \perp\!\!\!\perp \{x_1, x_3\} \perp x_2 \quad (2.12)$$

$$x_5 \perp\!\!\!\perp \{x_1, x_2, x_4\} \perp x_3 \quad (2.13)$$

$$x_6 \perp\!\!\!\perp \{x_1, x_3, x_4\} \perp \{x_2, x_5\} \quad (2.14)$$

Interpretation of missing edges as c.i. consistent with (2.2) - no variance later

verified an example of (2.12) by direct calculation from (2.3) :-

$$p(x_1, x_2, x_3, x_4) = \sum_{x_5} \sum_{x_6} p(x_1, \dots, x_6) \quad (\text{marginalise}) \quad (2.15)$$

m.g. $\{x_1, \dots, x_4\}$ (sub.s.)

$$p_{\text{prob}} = \sum_{x_5} \sum_{x_6} p(x_1) p(x_2|x_1) p(x_3|x_1) p(x_4|x_2) p(x_5|x_3) p(x_6|x_2, x_5) \quad (2.16)$$

$$= p(x_1) p(x_2|x_1) p(x_3|x_1) p(x_4|x_2) \underbrace{\sum_{x_5} p(x_5|x_3)}_{=1} \underbrace{\sum_{x_6} p(x_6|x_2, x_5)}_{=1} \quad (2.17)$$

$$= p(x_1) p(x_2|x_1) p(x_3|x_1) p(x_4|x_2) \quad (2.18)$$

m.g. prob $\{x_1, x_2, x_3\}$:-

$$p(x_1, x_2, x_3) = \sum_{x_4} p(x_1) p(x_2|x_1) p(x_3|x_1) p(x_4|x_2) \quad (2.19)$$

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

$$= p(x_1) p(x_2|x_1) p(x_3|x_1) \quad (2.20)$$

$$\text{m.u.: } p(x_4|x_1, x_2, x_3) = p(x_4|x_2) \quad (\text{recall all c.i. definition!}) \quad (2.21)$$

- (8) above shows we can interpret missing edges in graph in terms of conditional independencies
- (*) (Q): are there other conditional independence statements that are true of such joint probability distis; and do these have graphical interp.
 There are other conditional independencies e.g. $X_1 \perp\!\!\!\perp X_6 | \{X_2, X_3\}$, not in list (2.14), but implied by the list.
 (unifiable by algebra (tedious))
- (Q): we want to write down ALL conditional independencies implied by basic set.
- develop criterion for doing so without factorizing joint into all possible triplets of var. subsets
- (Q): That is, a general graph search algorithm to so find all implied independencies as well as explicit ones from graph structure.

2.17. conditional independence and Bayes ball

- Bayes ball - reachability and a def. of separation.
 - we can not only derive C.I. assertions from (2.2); but otherwise around.
- significance:
- to each graph we associate family of joint prob. distri.
 - family arises due to consideration of range over different choices of numerical values of local c.p.s. $p(x_i | x_{\pi(i)})$. list of
 - view c.i. statements gen by Bayes ball as constants on \mathcal{P} joint p.d.
 - those joint p.d.s that meet \rightarrow if not out .
- (Q): relationship between characterization of a family of p.d.s in terms of conditional independencies; and numerical character in terms of local cond. prob. (S.2.1.3.)

3 canonical graphs

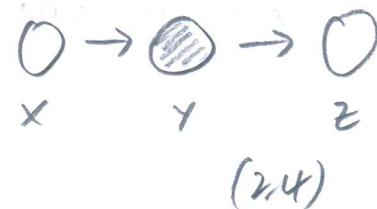
missing edges in PGM \rightarrow c.i.

e

Missing edge: $x \rightarrow z$

Hence C.I.:-

$\textcircled{1} X \perp\!\!\!\perp Z | Y$



(2.4)

② No other conditional independencies associated with this graph

Justifying ①: (DAG structure \Rightarrow C.I.)

$$\text{From DAG: } p(x, y, z) = p(x)p(y|x)p(z|y) \quad (2.23)$$

$$\text{From (2.23): } p(z|x, y) = \frac{p(x, y, z)}{p(x, y)} \quad (2.24)$$

$$= \frac{p(x)p(y|x)p(z|y)}{p(x)p(y|x)} \quad (2.25)$$

$$= p(z|y) \quad (2.26)$$

- Establishing C.I.

Justifying ②:- $\textcircled{6ii}$

"There are no further conditional independencies associated with this graph"

- does NOT mean that no further conditional independencies can arise in any of the distributions in the family associated with this graph i.e. distributions with factorised form (2.23) (?)

- There exist some distributions which exhibit conditional independencies

- e.g. free to choose any local c.p. $p(y|x) \rightarrow$ choose dist'n in which # prob. y same no matter value of x. then for this $p(y|x)$, $X \perp\!\!\!\perp Y$.

$\textcircled{6ii}$: Fig (2.4) $\not\Rightarrow$ X and Y are necessarily dependent (not independent).

$\textcircled{6ii}$: Edges do not necessarily imply dependence

BUT $\textcircled{6ii}$: Missing edges do imply independence

- universally vs existentially qualified statements with respect to family of distributions associated with a graph.

$\textcircled{6ii}$: asserted conditional independencies always hold for these distributions

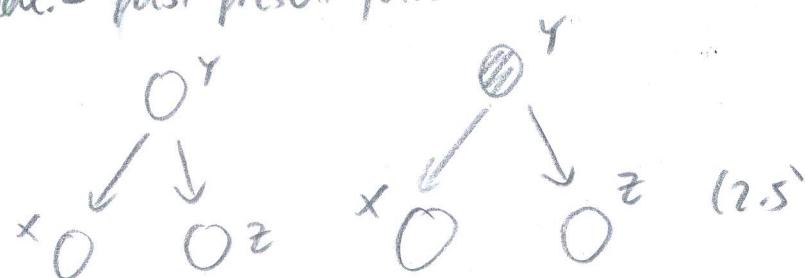
- implied/non-asserted conditional independencies sometimes fail to hold for distributions associated with a particular graph, but sometimes do hold.

- ⑤
- An algorithm based on conditional independencies will be correct for all distns associated with a graph
 - An algorithm based on absence of c.i.s. will sometimes be correct, sometimes not.

mental model for cascade:- past-present-future

common parent

- missing edge: $X \perp\!\!\!\perp Z | Y$



Justification ①:- (DAG \rightarrow C.I.)

$$p(x,y,z) = p(y)p(x|y)p(z|y)$$

(2.28)

$$p(x,z|y) = \frac{p(y)p(x|y)p(z|y)}{p(y)}$$

(2.29)

$$\Rightarrow p(x,z|y) = p(x|y)p(z|y) \Rightarrow X \perp\!\!\!\perp Z | Y$$

(2.30)

mental model: 'hidden variable'

- X -shoe size ; Z - 'gray hair' ; Y -age

- $X, Z \rightarrow$ m.pop., strong dep.

- But with Y -age, we might be willing to assert $X \perp\!\!\!\perp Z | Y$.

- Hidden Y explains all of observed $X-Z$ dep.

② "no other c.i.s. associated" with ^{this} graph

- No assertions of dependence, in particular we do not necessarily assume X and Z dependent because they depend on Y .

- But we assert the at least some distns. in which such dependence

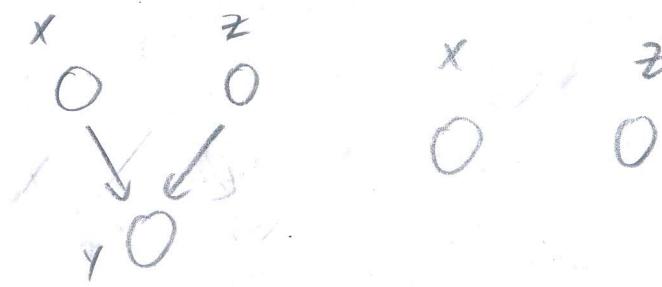
& for d.

V-structure

$X \perp\!\!\!\perp Z$

missing edge $X-Y$

(statement about subgraph)



- vo-06 MS readings
- Koller + Friedman (2009)
- (iii): use the intuitive examples here to supplement, reinforce understanding; they hold the key to applying the reasoning
 - (iv): the purposes of this reading → clarifying BN semantics, I-maps, equivalence, soundness, completeness
 - (v): recall 36705, 11D as cornerstone of statistical analysis?
 - (vi): 10.708 → complementary through consideration of dependence
 - reason for parametrisation through cpd rather than full joint → modularity
 - Adding node will not require new joint p.d.; modular in cpd representation.
 - local vs global

3.2.2.2. Bayesian Network Semantics

definition 3.1

- A Bayesian network structure G is a DAG whose nodes represent r.v.s. X_1, \dots, X_n .
- let Pa_i^G represent parents of X_i in G
- let NonDescendants $_i$ denote variables that are not descendants of X_i
- G encodes following set of conditional independence assumptions, called local independencies.
- These are denoted I_G :-
- for each variable X_i : $(X_i \perp\!\!\!\perp \text{NonDescendants}_i \mid \text{Pa}_i^G)$
- "The local independencies state that each node X_i is conditionally independent of its nondescendants given parents"

3.2.3. Graphs and distri

- (i) :- 1.) Bayesian network graph with conditional independence assertions
- 2.) Bayesian network graph with conditional probability distributions
- These definitions are equivalent

⑥⑦:

(*) distribution P satisfies local independencies associated with a graph G $\Leftrightarrow P$ is representable as a set of CPDs associated with graph G

• 3.2.3.1. - I-maps

- define set of independencies associated with distri P
definitions 3.2. independencies in P

- let P be a distri over X

(*) we define $I(P)$ to be the set of independence assertions of the form $(X \perp\!\!\!\perp Y \mid Z)$ that hold in P .

(*) - we can rewrite statement that

" P satisfies the local independencies associated with G " as $I_G(G) \subseteq I(P)$

(*) semantically, " G is an I-map (independency map) for P "

definition 3.3 - I-map

- let K be any graph object associated with a set of independencies $I(K)$

- we say " K is an I-map for a set of independencies I if $I(K) \subseteq I$ "

(*) semantically

" G is an I-map for P if G is an I-map for $I(P)$ "

⑥: do you understand these nuances?

- (i) direction of subset inclusion \Rightarrow for G to be an I-map of P ,

it is necessary that G does not mislead us regarding independencies in P .

- start example to help with this abstract definition

- note that $I_G(G) \subseteq I(P)$ implies that $I(P)$ may contain independence assertions not present in $I_G(G)$!

- these definitions are finicky \rightarrow do some scribblings

Intuition from I-map to factorisation

write statements
mathematically

- BN structure G encodes a set of c.i. assumptions
- Every distri for which G is an I-map must satisfy the assumptions
- This is key to understanding factorised representation
- Consider any distri P for which student BN G_{student} is an I-map.
- Decompose joint distribution via chain rule:-
- $P(I, D, G, L, S) = P(I)P(D|I)P(G|I, D)P(L|I, D, G)P(S|I, D, G, L)$ (3.15)
- No assumptions, but not helpful
- Apply c.i. assumptions, induced by B.N.
- Some G_{student} is an I-map for our distri. P
- (3.14) $\rightarrow (I \perp D) \in I(P) \Rightarrow P(D|I) = P(D)$
- (3.10) $\rightarrow (L \perp I, D|G) \in I(P) \Rightarrow P(L|I, D, G) = P(L|G)$
- (3.11) $\rightarrow (S \perp D, G, L|I) \in I(P) \Rightarrow P(S|D, G, L, I) = P(S|I)$
- Hence $P(I, D, G, L, S) = P(I)P(D)P(G|I, D)P(L|G)P(S|I)$ (3.16)

- Any entry in joint distri can be computed as a product of factors, one for each var.
- Each factor \rightarrow a conditional prob. of variable given network parents
- This factorisation applies to any distribution P for which G_{student} is an I-map

Formally:

Def 3.4 (Factorisation)

- Let G be a BN graph over X_1, \dots, X_n
- We say a distri P over state space factorises according to G if P can be expres:-

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | P_{X_i}^G)$$

- Chain rule for BN.
- Individual factors $P(X_i | P_{X_i}^G)$ are CPDs (conditional prob distri)

03.5. Bayesian Network

- A Bayesian network is a pair $B = (G, P)$ where P factorises over G , and where P is specified as a set of CPDs associated with G 's nodes.
- The distri P is often annotated P_B .
- Phenomenon holds for G without more generality.

Theorem (3.1)

- Let G be a BN structure over a set of r.v.s. X , and let P be a joint distribution over the same space.

(*) If G is an I-map for P , then P factorises according to G .

Proof: - see Koller & Friedman (many use for exercises)

Comments: - C.I. assumptions implied by BN structure G allow factorisation of a distri P for which G is an I-map into small CPDs.

(@) 1.3.1. Shows direction of fundamental connection between the c.i. independencies encoded by the BN structure and the factorisation of distri into local probabilistic models / CPDs.

- C.I. \Rightarrow factorisation

- Theorem 3.2 (converse: factorisation according to $G \Rightarrow$ C.I.)

- Let G be a BN structure over a set of r.v.s. X and let P be a joint distri over some space.

(*) If P factorises according to G , then G is an I-map for P .

3.3 Independencies in graphs

- Graph structure $G \rightarrow$ encodes a set of c.i. assumptions $I_G(G)$

- Knowing only that a distribution factorises over G ; we can conclude that distri satisfies $I_G(G)$. (i.e. $I_G(G) \subseteq I(P)$)

(@) recall Jordan :- Are there any other independencies that we can 'read off' directly from G .

- Are there other independencies that hold for every distribution P that factorises over G ?
(a distri family)

(@) A question about implied independencies i.e. $I(P) - I_G(G) (?)$

D-separation

- who can we guarantee that an independence $(X \perp\!\!\!\perp Y | Z)$ holds in a distribution associated with a BN structure G .
 - in interests of conciseness, I'm going to leave connects
 - the next section is on the canonical local structures / independencies;
 - ~ Jordan, cascade, common parent, v-structure.
 - Koller (2009) adds additional info/intuition and categories as:-
 - 1. Direct connection $X \rightarrow Y$
 - 2. Indirect connection, (3-node networks)
 - i) Indirect causal
 - ii) Indirect evidential
 - iii) Common cause
 - iv) Common effect
- } each contains intuition
- key is to note that the key event is whether a node is observed or unobserved (under which C.I statements, (de)separations) / decoupling is generated)
 - Koller (2009): (i) when probabilistic influence can flow from X to Y via Z ; then we have an active trail / blocking
 - 10/108 12 notes: Better summary linking Bayes Ball to Koller's idea of active trails.
 - ②: essentially: (Bayes Ball-reachability - Bell-blocked - Jordan)
 - probabilistic influence - active-trails-blocked (Koller)
 - ↳ both encode intuitions about what leads to d-separator
- general case:
- longer trail $X_1 \Rightarrow \dots \Rightarrow X_n$
 - for influence to flow from X_1 to X_n ; it needs to flow through every node on trail.
 - X_i can influence X_n if every two edge trail $X_{i-1} \Rightarrow X_i \Rightarrow X_{i+1}$ along trail allows influence to flow.

formally:

D.3.6.

Let G be a BN structure; and $X_1 \Rightarrow X_2 \Rightarrow \dots \Rightarrow X_n$ a trail in G .

Let Z be a subset of observed r.v.s.

The trail $X_1 \Rightarrow X_2 \Rightarrow \dots \Rightarrow X_n$ is active given Z if:-

- whenever we have a v-structure $X_{i-1} \rightarrow X_i \leftarrow X_{i+1}$ then X_i or one of its descendants are in \underline{Z} .
 - no other node along the trail is in \underline{Z}
 - If X_i or X_n are in \underline{Z} , the trail is not active.
 - need to account for graphs with more than one trail between nodes
 - d-separation (*): d-separation
 - i) Bayes ball algorithm
 - ii) separation on normalised ancestral graph
 - iii) Active trails (ii). Multiple algs for some distribution
- D.3.7 (d-separation)
- Let X, Y, Z be three sets of nodes in G .
- We say X and Y are d-separated given Z i.e. $d\text{-sep}_G(X; Y | Z)$, if there is no active trail between any nodes $X \in X$ and $Y \in Y$ given Z .
- $I(G)$ is the set of independencies that correspond to d-separation
- $$I(G) = \{ (X \perp\!\!\!\perp Y | Z) : d\text{-sep}_G(X; Y | Z) \}$$
- This set is called the set of global Markov independencies.
 The independencies in $I(G)$ are precisely those that are guaranteed to hold for every distribution over G .
- some formalities (mathematical) regarding dsep. (i.e. proof style)
desirable properties of d-separation
3. (Soundness of dsep)
- If a distri P factorises according to G , then $I(G) \subseteq I(P)$
- mark: Any independence reported by d-separation is satisfied by the underlying distri
- If two nodes X and Y are d-separated given some \underline{Z} , we are guaranteed they are c.i. given \underline{Z} .
- Completeness - d-separation detects all possible independencies
- If two variables X and Y are independent given \underline{Z} ; they are d-separated defined in this form \rightarrow no specification of distribution in which X and Y are independent
- c. actually a statement which involve a set of distri./distri family?

D3.8. (Faithful)

- A distri P is faithful to G if, whenever $(X \perp\!\!\!\perp Y | Z) \in I(P)$; then $dsep_G(X; Y | Z)$.
- Any independence in P is reflected in the d-sep. properties of graph (7)
- A candidate formalisation of completeness is in terms of converse/contrapositive (11)(c) (x)(y)
- (*) For any distri P that factorises over G , we have that P is faithful to G
- (*) i.e. if X and Y are not d-separated given Z in G ; then X and Y are dependent in all distrib. P that factors over G . (12)

(11) - converse to soundness (11)(c)

- If true, two together (which?) imply that for any P that factorises over G , we have $I(P) = I(G)$.

→ - Highly desirable property is false (?) which (?)

- (*) - even if a distri factorises over G , it can still contain additional independencies
- (*) not reflected in structure

completeness property does not hold for this candidate definition of completeness

Settle for a weaker definition:-

- If $(X \perp\!\!\!\perp Y | Z)$ in all distri P that factorise over G , then $dsep_G(X; Y | Z)$.
- contrapositive:- If X and Y are not d-separated given Z in G ; then X and Y are dependent in some distribution that factorises over G .

- formally construct theorem:

1.3.4

- let G be a BN structure.

- If X and Y are not d-separated given Z in G , then X and Y are dependent given Z in some distribution that factorises over G .

Proof:- See Molle (2009).

- Remark: completeness result tells us that our definition of $I(G)$ is maximal. for any independence assertion not a consequence of d-separation in G one can find a counterexample distri P that factors over G .

13.5 (measure theoretic qual.)

- for almost all distributions P that factorise over G , that is for all distributions except for a set of measure zero in the space of CPD parameterisations, we have $I(P) = I(G)$.

stronger than 13.4

- (*) - results state that for almost all parameterisations P of the graph G (that is for almost all possible choices of CPDs for the variables), the d-separation test precisely characterise all independences that hold for P .

3.3.4. I-equivalence

- $I(G)$ specifies a set of c.i. assertions associated with a graph.
- can abstract away details of graph structure and view as specification of indep. properties
- (*) - very different BN structures can be equivalent in that they encode precisely the same set of c.i. assertions (e.g. local structures/canonical maps).

0.3.9. (I-equivalence)

- two graph structures K_1 and K_2 over X are I-equivalent if $I(K_1) = I(K_2)$
- the set of all graphs over X is partitioned into a set of mutually exclusive and exhaustive I-equivalence classes.
- which is the set of equivalence classes induced by the I-equivalence relation

(*) I-equivalence of 2 graphs \Rightarrow any distn P that can be factorised over one of these graphs can be factorised over the other.

- (*) - no intrinsic property of P that would allow us to associate it with one graph rather than an equivalent one.
- (*) - important implications on our ability to determine dir. of influence

- Undirected Graphical Models

- record intuitions, questions - be selective

④ Review of independence properties of DAGs.

- the properties of DAGs require more reading of Koller.

(*) A fully connected DAG G has 1-map for any distri $I(G) = \phi \subseteq I(P) \forall P$.

- Minimal 1-map $\textcircled{A1}$

- other characteristics \rightarrow readings

⑤ How to find d-separation was a key question. Multiple ways.

1. moralised ancestral graph

And:

- combesore algorithmically - mechanical protocol for detecting c.i.

Ex: Bayes ball algorithm - mechanical protocol for detecting c.i.

- note that d-separation/Bayes ball is an encoding of obs. that simple appeals to 'blocking' do not hold for V-structures.
(conditioning on a variable)

- key exception is V-structure

- using Bayes ball algorithm; we build up d-separation on entire graph through consideration of each of these canonical graph structures / local independencies.

1-Maps \rightarrow see Koller $\textcircled{A1}$

- 2 sets $I(P), I(G)$ are equivalent?

⑥ for any graph G and distri P , above always expect a 1-map / equivalence.

- theorem - NO!

- by counterexample

- try and define CPD that satisfies these two C.I.s. $A \perp\!\!\!\perp C | \{B, D\}; B \perp\!\!\!\perp D | \{A, C\}$.

- constant: you must draw with DAGs.

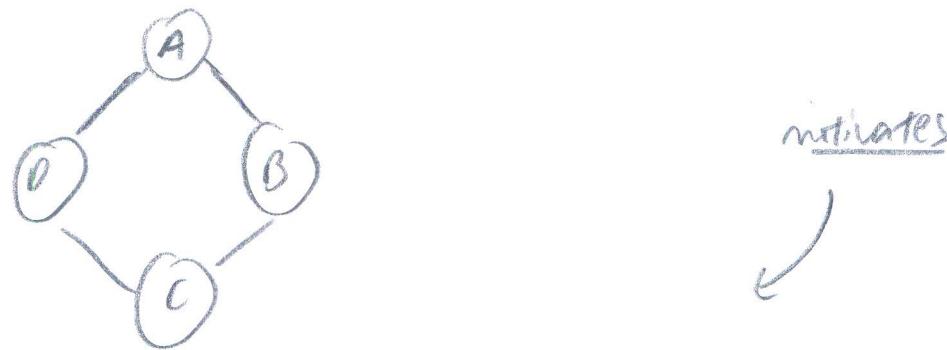
- Are 2 maps BN1 and BN2 equivalent to the $I = (A \perp\!\!\!\perp C | \{B, D\}; B \perp\!\!\!\perp D | \{A, C\})$

i.e. Are there C.I.s. in the graph structure not in the listed C.I.s.

- BN2 - via v-structure - DIB.

Q: There are no DAGs which encodes exactly the same set of conditional independencies in this set. \Rightarrow not every distri has a p-map in DAG.

Q: But $A \perp\!\!\!\perp C | \{B, D\}$ and $B \perp\!\!\!\perp D | \{A, C\}$ can be captured by:-



- graph separation in NGMS

NGMS

more expressive, precise about conditional independencies.

parallel relationships, no causality (no parents/ancestors).

assign 'scores' to configurations

not generative, i.e. using conditional probability distri.

w problems \rightarrow pixel labelling

Ex: sears like C.S. believe that all intuitions can be encoded algorithmically

- grid model

- Ising model (magnetism)

- probabilistic model with symmetrically connected r.v.s.

↳ contingency tables to express preferences over patterns

- NGMs have semantics of directionality, causality, temporality

- NGMS do not

- maybe synchronic vs diachronic use?

'Canonical model for Go' \rightarrow PGM before AlphaGo.

• information retrieval

(1) is a table here

- canonical definitions \rightarrow readings notes

- undirected graph H

- potential function ψ_c : mapping from config \rightarrow no.

- potential functions are pre-probabilistic (do not need to conform to rules of probability)

- gibbs distribution

x_1	x_2	$\psi_c(x)$	$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$
1	1	10	
1	0	1	
0	1	1	
0	0	2	

- (2): But later normalise (?) using a partition function

(3): partition function - sum of product of all potentials over all possible configurations

- i.e. across all subsets over which you have defined potential

- (4): define c appropriately i.e. not entire set of r.v.s.

- statistical physics \rightarrow probabilistic graphical models!

- quantitative specifications - cliques, potentials

ex: potential function defined over cliques

(5): formal defi

- max clique - maximal fully connected subset of graph

- sub-clique - pairs of nodes and singletons

(6): why are cliques important?

- clique - totally connected - no condit. independence within the clique

- every configuration is possible and has to be 'honoured'

- large configuration with certain disconnections \rightarrow there may be independences, certain configurations may not deserve attention.

- cliques - every configuration has to be associated with a characterisation (e.g. potential number or prob. mass)

ex: only define potential functions on cliques

Clique potentials - intro

- $\emptyset - \{Y\} - \{Z\}$ - not max clique
- 2 potentials over (X, Z) and (Y, Z) ; or on singletons X, Y, Z
- via a number of ways, breaking down joint $P(X, Y, Z)$
- ex: Not sure what point about lack of correspondence between marginal, conditional probabilities; and potential functions
- physics intro to of potential (e.g. singleton spin; pairwise magnetism)
- is useful for intuition
- ↓ - (*) forget the probabilistic semantics from PGMs for ψ_c .

PGM - using max cliques

- potential ψ_c is a mapping from a triple $\psi_c(\{1, 2, 4\})$
- $p(x_1, \dots, x_4) = 4$ dim table, 2 states for each r.v.
- represent $\psi_c(\{1, 2, 4\})$; $\psi_c(\{1, 2, 3, 4\}) \rightarrow 2$ 3D tables

PGMs - using subcliques

- using pairwise potentials
- pairwise Markov models (lazy)

PGM - canonical rep.

① is canonical bad?

- canonical repres. → potential functions for all cliques in the graph

② is overrep bad?

- (6) not necessarily; think about NIV as over-represented which have everything possible as placeholder to anticipate complexity; even though spec. of sparse model.

ex ③: HW: $\begin{matrix} p_1 & p_2 \\ \text{max clique} & \text{pairwise clique} \end{matrix}; I(p_1) = I(p_2) ?$

- same set of c.i.s. for 2 distn?
- Are no. ways distn specified in 1112 ways the same?

(II)- Independence properties

- conditional independence properties via graph separation
- global Markov independencies + global Markov property $\textcircled{A5}$
- weak Markov independencies $\textcircled{A6}$

3 - pick up in notes.

- in Markov networks, CCI are more intuitive - neighbour based
- Markov blankets

Ex: Graphical models unpopular nowadays due to deep learning; but certain properties are useful.

Markov blankets in DGMs \rightarrow drawing samples cond. on evidence (deep gen models)

- difficult - uses $P(X_i | X_{i-1})$

- gigantic, conditioning on many many things.

- If DGM can be modelled as a MN:-

$$P(x_i | x_{i-1}) \equiv P(x_i | MB_i) \quad (\text{a lot } \underline{\text{smaller}})$$

- soundness & completeness $\rightarrow \textcircled{A7}$

- Analogies

$\textcircled{A7}$ energy states in atomic bomb
- statistical sampling algorithms used M.B.
property to speed up comp.

Hammersley-Clifford

- for arbitrary potentials in Gibbs distri.

- then the functional form of Gibbs distri is the ONLY WAY to write proper distri's on graph i.e. $P(x_1, \dots, x_n) = \frac{1}{Z} \prod_{c \in C} \psi_c(x_c)$

- $\textcircled{A8}$ theorem \rightarrow

(a) to factorisation law for BNS

- can discover all cliques (max, sub) inside graph, write potential functions over cliques, multiply to get joint; captures all independences on graph

perfect maps

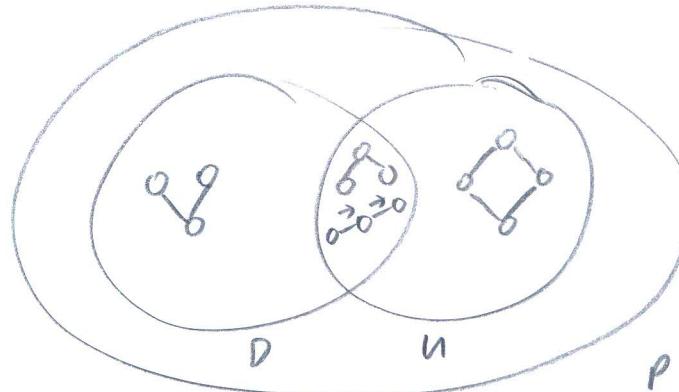
- high-level

- Ex: we know there are certain independences in a distri for which we cannot find a P-map in DGMs.

- similarly, there are certain c.s. that would exist in a distri P that we cannot find/express in UGM.

- perfect-map is not an entitlement for arbitrary graphs

- see venn diagram ~~(A)~~ ~~(B)~~



- D cannot capture



- U cannot capture
V-structure



ex: specification of potential fns:

^{in UGM}

- when we specify CPDs, we have probability constraints

- Specifying potential fns in Gibbs distribution → must be non-zero constraint

- there may be symmetries over O e.g. -1, 1

- Gibbs distri. inadequate

↳ refine potential functions ϕ_C as exponentiated negy functions ϕ_C

- free negy/Boltzmann distri (stat phys)

- log-linear (statistics)

- Allows specification of -ve nos (getting around constraint)

Boltzmann Machines

- provide pairwise and singleton potentials

$$\text{#}: p(x_1, \dots, x_N) = \frac{1}{Z} \exp \left\{ \sum_{ij} \phi_{ij}(x_i x_j) + \sum_i \phi_i \right\}$$

- turn state variables → nos → allow direct definition of pairwise potentials, singleton pot.

- allows incomparable specification

- μ is an offset

quadratic form
for negy fn.

ex: tabular expressions of pot.

ϕ_C

OR

ex: physicists method

- directly define ϕ_C
assume linear
function of
state values x

$$\text{#} = \begin{bmatrix} \alpha_1 & \alpha_2 & \alpha_3 & \dots \\ \alpha_1 & \vdots & & \\ \vdots & & & \\ \alpha_N & & & \end{bmatrix}$$

- C normalises
the distri
- norm. constat

Ising model

- quadratic and singleton states

RBMs

- important in ML/DL
- bipartite connection of hidden, visible
- globalist alfin (*) - singleton, singleton, pairwise
- very rich, deep connection to DNN
- Inference on RBM; computational graph very related to backprop on DNN
- computational and model analogies with DNN.

Properties of RBM

- ex: why no directionality?
- ex: ML scientists vs statisticians → after putting forward model, have to compute

(*) computational algorithm as important as model

- directionality has significant computational consequences.
- creates tons of V-structures → coupling → difficult
- hence Gibbs sampling possible due to this.
- ex: modelling mathematics tightly coupled with efficiency, comput. of algorithms

uGM semantics → constructive defn.

- RBM text modelling (topic modelling in uGM space)
- localist definitions give intuitive specifications to get globalist defn.

reference points:-

CRFs - lafferty, very famous (difficult paper)

- HMM model with counterpart in undirected space (no directionality)

- Potential functions \rightarrow very interesting interp.
 - H-potentials \rightarrow spell-checker e.g. g_2 co-occurrence frequency / prob.
 - Potentials can capture interesting global effects
 - Offers great template for feature engineering, language modelling
- ex: Graphical models with c.i. properties - neighbours - Markov property
- 3 special UGMs:
 1. Ising model
 2. RBM
 3. CRFs.
- (*)
- Q: where does graph structure come from
arbitrarily? or
data-driven?
 - graph structure itself
can be learnt from data
 - data-driven causality inference?
- (Ex): algorithms for unique discovery of structure exist (HOOD)
- Invariably correct under certain conditions

- (Ex): for most conditions, cannot uniquely discover a structure
- many structures will give you some score you want to optimise
 - causality is a statistical effect.