

- Undirected Graphical Models

- record intuitions, questions - be selective

⊗ review of independence properties of DAGs.

- the properties of DAGs require more reading of Koller.

(*) A fully connected DAG G is an I-map for any distri $I(G) = \emptyset \subseteq I(P) \forall P$.

- Minimal I-map ⊗ (A1)

- other characteristics → readings

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

1. Moralised ancestral graph

Also:

- can be done algorithmically

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 independences.

P-Maps → see Koller ⊗ (A2)

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

Ⓚ: for any graph G and distri P , do we always expect a P-map / equivalence.

- Theorem - NO!

- By counterexample

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

- constant: you must draw with DAGs.

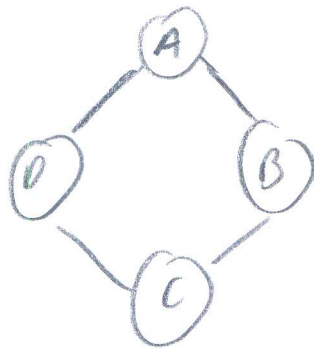
- Are 2 maps BN1 and BN2 equivalent to the $I = (A \perp C | \{B, D\}; B \perp D | \{A, C\})$

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

- BN2 - via v-structure - DIB.

⊗ 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.

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



motivates



- Graph separation in NGMS

NGMS

- more expressive, precise about conditional independencies.
- pairwise relationships, no causality (no parents/ancestors).
- Assign 'scores' to configurations
- not generative, i.e. using conditional probability distri.

CV problems \rightarrow pixel labelling

ex: seems like C.S. believe that all intuitions can be encoded algorithmically

- Grid model

- Ising model (magnetism)

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

\hookrightarrow contingency tables to express preferences over patterns

- DGM have semantics of directionality, causality, temporality

- NGMS do not

- maybe synchronic vs diachronic use?

- Canonical model for G0 \rightarrow PLM before AlphaGo.

- Information retrieval

$\psi(\cdot)$ is a table here

- Canonical definitions \rightarrow readings notes

- undirected graph H

- potential function ψ_c : mapping from

- potential functions config \rightarrow no.

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 | |

Ⓢ: But later normalise (!) using a partition function

Ⓢ: Partition function - sum of product of all potentials over all possible configurations

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

Ⓢ: 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

Ⓢ: formal defi

- max clique - maximal fully connected subset of graph

- sub-clique - pairs of nodes and singletons

Ⓢ: 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 - interp.

- $(X)-(Y)-(Z)$ - not max. clique
- 2 potentials over (X,Z) and (Y,Z) ; or over singletons X, Y, Z
- ^{ways.} A number of breaking down joint $P(X,Y,Z)$
- EX: Not sure what point about lack of correspondence between marginal, conditional probabilities; and potential functions
- (physics interp to of potential (e.g. singleton spin; pairwise magnetism))
- (useful for inference)
- (*) Forget the probabilistic semantics from PGMS for ψ_c .

UGM - using max cliques

- potential ψ_c is a mapping from a triple $\psi_c(z_{124})$
- $p(x_1, \dots, x_4)$ - 4 dim table, 2 states for each v.
- represent $\psi_c(z_{124}), \psi_c(z_{134}) \rightarrow 2$ 3D tables

UGMS - using subcliques

- using pairwise potentials
- pairwise Markov models (lazy)

UGM - canonical rep.

Q: Is canonical bad?

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

Q: Is overrep bad?

- EX: Not necessarily; think about NIV as over-represented which have everything possible as placeholder to anticipate complexity; even though spec. of sparse model.

EX Q: P_1 P_2
HW: max clique pairwise clique; $I(P_1) = I(P_2)$?
- same set of c.i.s. for 2 dist? $P(P_1) = P(P_2)$?
- Are no. ways dist. specified in the 2 ways the same?

(II) - Independence properties

- conditional independence properties via graph separation
- global Markov independencies + global Markov property (A5) } - pick up in notes.
- local Markov independencies (A6)
- in Markov networks, LCI 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:-

$$P(X_i | X_{i-1}) \equiv P(X_i | MB_i) \quad (\text{a lot smaller})$$

(A7) energy states in atomic bomb (!)

- statistical sampling algorithms used M.B. property to speed up comp.

- soundness & completeness \rightarrow (A7)

- Analogies

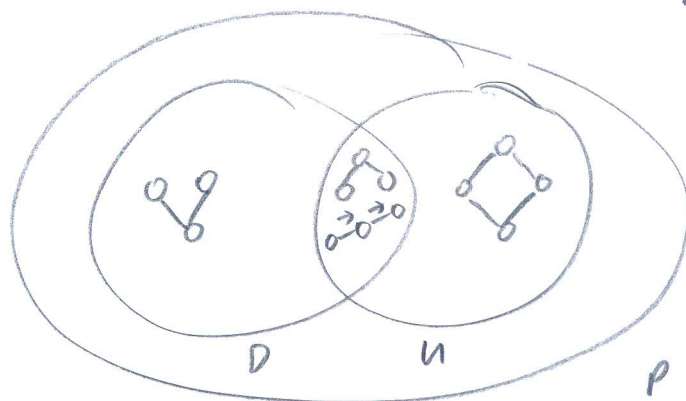
Hammersley-Clifford

- for arbitrary potentials in Gibbs distri.
- then the functional form of Gibbs distri is the ONLY way to write proper distri on graph i.e. $P(x_1, \dots, x_n) = \frac{1}{Z} \prod_{c \in C} \psi_c(x_c)$
- (A8) + theorem \rightarrow
- (~) to factorisation law for BNS
- can discover all cliques (max, sub) inside graph, write potential functions over cliques, multiply to get joint; captures all independencies in a graph

Perfect maps

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

- similarly, there are certain c.i.s. that could exist in a distri P that we cannot find/express in UGM.
- perfect-map is not an entitlement for arbitrary graphs
- see ven diagram (X) (AC)



- D cannot capture



- cannot capture v-structure



ex: specification of potential fns:

- when we specify c.p.d.s, we have probability constraints

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

- there may be symmetries over 0 e.g. -1, 1

- Gibbs distr. inadequate

→ define potential functions ψ_c as exponentiated energy functions ϕ_c

- free energy/Boltzmann distr (stat phys)

- log-linear (statistics)

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

Boltzmann Machines

- provide pairwise and singleton potentials

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

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

- Allows linear algebraic specification

- μ is an offset

preparative form for energy fn.

EX: Tabular expressions of pot. ψ_c

OR

EX: Physicists method

- directly define ψ_c
as simple linear function over state values x

$$\textcircled{\#} = \begin{bmatrix} \alpha_1 & \theta_{12} & \theta_{13} & \dots \\ \theta_{21} & & & \\ \vdots & & & \\ \alpha_N \end{bmatrix}$$

- C normalises the distri

- norm. const

- Ising model

- Quadratic and singleton states

- RBMs

- important in ML/DL

- bipartite connection of hidden, visible

- Globalist defn (*) - singleton, singleton, pairwise

- very rich, deep connection to PNN

- Inference in RBM; computational graph very related to backprop in PNN

- 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, input of algorithms

HGM semantics → constructive defn.

- RBM Text modelling (topic modelling in HGM 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 → very interesting interp.
- H-potentials → spell-checker e.g. gz co-occurrence frequency / prob.
- Potentials can capture interesting global effects
- CRFs great template for feature engineering, language modelling

ex: graphical models with c.i. properties - neighbours - Markov property

- 3 special UBMS: 1. Ising model

2. RBM

3. CRFs.

(*)

①: where does graph structure come from
 $\left\{ \begin{array}{l} \text{arbitrarily? OR} \\ \text{data-driven?} \end{array} \right.$

- graph structure itself

can be learned from data

- data-driven causality inference?

ex: algorithms for unique discovery of structure exist (GOOD)

- provably correct under certain conditions

ex: for most conditions, cannot uniquely discover a structure

- many structures will give you same score you want to optimise

- causality is a statistical effect,