

Readings

- key equations; some intuition on threshold points
- Koller (2009) / Jordan (2003) (need to shorten & hence select + concise)
- BN cannot model a distri that satisfies $(A \perp\!\!\!\perp C | \{B, D\})$ and $(B \perp\!\!\!\perp A, C)$ only.
- $M, N \rightarrow$ undirected edges (prob. interaction)
- similar to BN; parametrisation of MN
 - local interactions
 - global model!
 - combine
- product of local factors, normalised
- normalising constant \rightarrow partition function (MRF in statistical physics)
- MN - connection between factorisation and model/choice properties

0.4.3 (Gibbs disti)

- an undirected graphical model represents a distri $p(x_1, \dots, x_n)$ defined by an undirected graph H , and a set of (positive) potential functions ψ_c associated with cliques of H s.t.

$$p(x_1, \dots, x_n) = \frac{1}{Z} \prod_{c \in C} \psi_c(x_c) \quad Z = \sum_{x_1, \dots, x_n} \prod_{c \in C} \psi_c(x_c)$$

0.4.4 (MN factors.)

- we say a distri $p(x_1, \dots, x_n)$ with $\Psi = \{\psi_1(x_1), \dots, \psi_c(x_c)\}$ factorises over a MN H if each x_c ($c=1, \dots, C$) is a complete subgraph of H
- Factors that parametrise MN are clique potentials

- (*) - informally, subs. with idea of a maximal clique
 (*) wlog, parametrisation using maximal cliques obscures structure in original
set of factors (of a graph)

- informally: clique \rightarrow fully connected subset of nodes
 (Jordan) maximal cliques \rightarrow cliques that cannot be extended to include additional nodes without losing property of being fully connected

formally :-

- for $G = \{E, V\}$, a complete subgraph (clique) is a subgraph $G' = \{V' \subseteq V, E' \subseteq E\}$ such that nodes V' are fully interconnected

- A maximal clique is a complete subgraph s.t. any superset $V'' \supset V'$ is not complete
- A sub-clique is not necessarily a maximal clique
- Sub-cliques - can be edges, singletons
- Koller (2009): see Box 4.B. - MN for CV

4.3. MN independencies

- Similar to Koller's pres. \rightarrow flows of prob. influence / active tails

04.8

- Let H be an MN structure
 - Let X_1, \dots, X_K be a path in H .
 - Let $Z \subseteq X$ be a set of observed variables
 - The path X_1, \dots, X_K is active given Z if none of the X_i 's $i=1, \dots, K$ is in Z
- allows a definition of separation

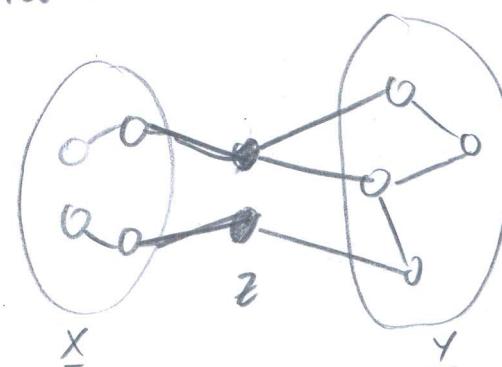
04.9 (separation/global independencies)

- We say a set of nodes Z separates X and Y in H , denoted $\text{sep}_H(X; Y | Z)$, if there is no active path between any node $X \in X$ and $Y \in Y$ given Z .
- We define the global independencies associated with H to be:-

$$I(H) = \{(X \perp Y | Z) : \text{sep}_H(X; Y | Z)\}$$

Remark:

Independencies $I(H)$ are precisely those guaranteed to hold for every distri P over H .



sep. criterion sound for detecting indep. properties over H .

(iii) BN \rightarrow connection with independence prop amplified by MN structure

+
factorising a distri over the graph

(iv) of both representation

(*)

Theorems for BN, equivalence of:- Gibbs fact. of
distri P over a graph H \iff H is an I-map for
 P (that is P
satisfies Markov ass. $I(H)$)

4.3.1.1. (soundness)

(*)

- (ii) to 4.3.3. i.e. Gibbs distri satisfies independencies associated with a graph.
 i.e. soundness of separation — (factorisation $\stackrel{4.1/3.2}{\Rightarrow} \stackrel{4.2/3.1}{\Leftarrow}$ C.I.)

T.4.1.

- let P be a distri over X , and H an MN structure over X
- If P is a Gibbs distri that factorises over H , then H is an I-map for P

other direction i.e. C.I. of distri \rightarrow factorisation: Hammersley-Clifford theorem

- unlike for BN; HCT does not hold in general.
- require additional assumption that P is a positive distri

T.4.2

- let P be a positive distri over X , and H a Markov network graph over X .
- If H is an I-map for P , then P is a Gibbs distri that factorises over H .

(*) For positive distris; the global independencies imply that distri factorises according to the network structure.

(*) For this class of distributions, we have that a distribution P factorises over a Markov Network if and only if H is an I-map of P .

(*) lecture defn: An H is an I-map for a distri P if $I(H) \subseteq I(P)$ i.e. entails $I(H)$

(*) If P is a Gibbs distri over H , then H is an I-map of P .

- soundness of separation as criterion for detecting independencies in MN.
- any distri that factorises over G satisfies the independence ass. implied by separation
- completeness - strong version of completeness does not hold
- It is not the case that every pair of nodes X and Y that are not separated in H are dependent in every distribution P which factorises over H
(use weak)

T.4.3

- let H be a MN structure.
- If X and Y are not separated given Z in H , then X and Y are dependent given Z in some distri P that factorises over H .

- same arguments as 1.3.5 to conclude:-
- ⑥ for almost all distri P that factorise over H (all distri except for a set of measure 0 in space of factor param.), we have $I(P) = I(H)$.
- our defn of $I(H)$ is maximal one.
- for any independence assertion that is not a consequence of separation on H , we can always find a counterexample distri P that factorises over H .

4.3.2. Indep. revisited

BN: local indep. (each node is indep. of nondesc. given parents)

Global indep. (induced by d-sep)

- showed that these are equiv. in the sense of one implies the other.

Q: (1) to BN; can we provide local indep. induced by MN, analogously to local indep. of BN

P(H): 3 diff poss. defin. of independencies associated with network structure ○
two local, one global in def 4.9.

4.3.2.1. - Local Markov ass.

D.4.10 → intuitively, two variables directly connected; potential for direct correlation via an mediated way.
conversely, two vars not directly linked, some way of reducing c.i.
- X and Y indep given all other nodes.

D.4.10 (Pairwise independencies)

Let H be a MN

we define pairwise independencies associated with H to be:-

$$I_p(H) = \{(X \perp Y \mid X - \{X, Y\} : \{X, Y\} \not\subseteq H\}$$

D.4.11 → analogue to local indep. associated with B.N.

D.4.11. (Markov blanket)

for a given graph H , we define the Markov blanket of X in H , $MB_H(X)$ to be the neighbors of X in H .

we define the local independencies associated with H to be:-

$$I_l(H) = \{(X \perp X - \{X\} - MB_H(X) \mid MB_H(X) : X \in X\}\}$$

- i.e. local independences state that X is indep. of nodes in graph given inmed. neighbours.
- we will show that these local indep. ass. hold for any distri that factorizes over H , so that X 's Markov blanket in H truly does sep. it from all other variables.

4.3.2.2. - Relationships between Markov properties

- 3 sets of indep. assertions assoc. with network structure H .
- For general distri $I_p(H)$ is weaker than $I_c(H)$ is weaker than $I(H)$.
- ~~All 3 are eq.~~ (*) All 3 are equivalent for positive distri

Prop 4.3

- For any MN H , and any distri P , we have that if $P \models I_c(H)$ then $P \models I_p(H)$

Prop 4.4

- For any MN H , and any distri P , we have that if $P \models I(H)$ then $P \models I_c(H)$

Th. 4.4

- Let P be a positive distri. If P satisfies $I_p(H)$, then P satisfies $I(H)$.

Corollary 4.1

- The following statements are equivalent for a positive distri P .

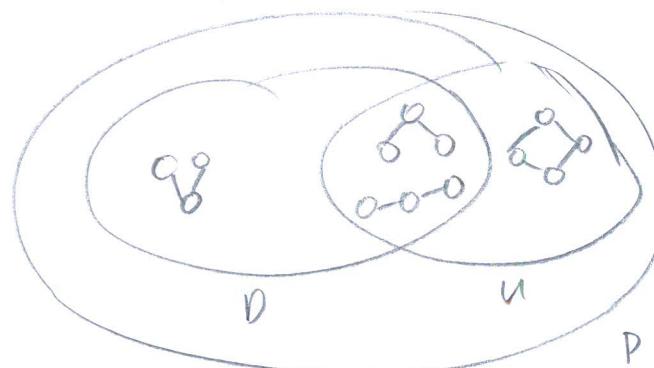
1. $P \models I_c(H)$
2. $P \models I_p(H)$
3. $P \models I(H)$

P & P-maps

- An MN H is a perfect-map for P if for any X, Y, Z , we have that

$$\text{sept}_H(X:Y|Z) \Leftrightarrow P \models (X|Y|Z)$$

Theorem: not every distri has a perfect map as NGM



From hereon, Jordan (2003)

Exponential form

- constraining clique potentials to be tree, unconstrained possibly

(Q1): what effect does this have on equivalence of local and global Markov properties?

- represent a clique potential as:- $\phi_c(x_c) = \exp\{-\psi_c(x_c)\}$ (*)

- $\phi_c(x_c)$ - a 'potential'

- Additive structure: $p(x) = \frac{1}{Z} \prod_{c \in C} \exp\{-\phi_c(x_c)\}$

$$= \frac{1}{Z} \exp\left\{-\sum_{c \in C} \phi_c(x_c)\right\} = \frac{1}{Z} \exp\left\{-H(x)\right\}$$

$$H(x) = \sum_{c \in C} \phi_c(x_c)$$

Boltzmann distri'

$H(x)$ - free energy

some notable forms / graph topologies:-

Ostheumann Machines

fully connected graph, pairwise edge pot., binary valued nodes (-1, 1)

$$p(x_1, x_2, x_3, x_4) = \frac{1}{Z} \exp\left\{\sum_{i,j} \phi_{ij}(x_i, x_j)\right\} = \frac{1}{Z} \exp\left\{\sum_{i,j} \theta_{ij} x_i x_j + \sum_i \alpha_i x_i + C\right\}$$

call negy fn:-

$$\psi(x) = \sum_{ij} (x_i - \mu) \theta_{ij} (x_j - \mu) = (x - \mu)^T \Theta (x - \mu)$$

w technique to learn this model, recover graph str. from data.

Models

model energy of a physical system (atom interaction)

regular grid topology

Mi-state \rightarrow Potts.

$$p(x) = \frac{1}{Z} \exp\left\{\sum_{(i,j) \in N_i} \theta_{ij} x_i x_j + \sum_i \alpha_i x_i\right\}$$

2. Restricted Boltzmann Machines

$$p(x, h | \theta) = \exp \left\{ \sum_i \theta_i \phi_i(x_i) + \sum_j \theta_j \phi_j(h_j) \sum_{i,j} \theta_{i,j} \phi_{i,j}(x_i, h_j) - A(\theta) \right\}$$

- See paper

CRFs
unlabeled graph rep.; encode cond. distri $p(y|x)$ y_i - target x - obs.

$$p_\theta(y|x) = \frac{1}{Z(\theta, x)} \exp \left\{ \sum_c \theta_c f_c(x, y_c) \right\}$$

- exact inference, variable elimination - you have notes already; read intuitions
- ex: 3 lectures on GM representation
- ex: focus on inference; then learning (which uses inference as a subroutine)

Query 1 - likelihood

- esimo - quantit. specification of probability of fishy sequence of all outcome - estimation/likelihood.

- marginal probability of evidence , likelihood

- marginalise over r.v.s. whose observations you do not have

out, 2 - C.P.

- use previous query as subproblem

- A posterior belief

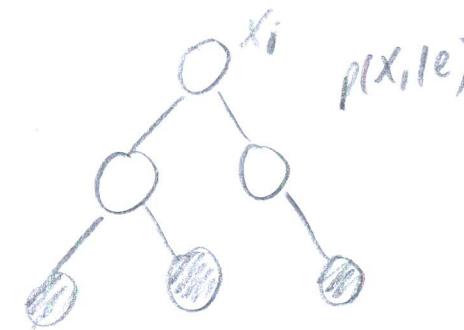
- don't query all; interested in subset of hidden r.v.s.

- marginalise out hidden variables which we are not interested in

applications of post. belief

prediction: $p(c|A,B) = p(c|B)$
as $C \perp\!\!\!\perp A | B$

agnosis: $p(A|B,C) = p(A|B)$
as $C \perp\!\!\!\perp A | B$.



systematic way of using PGMS for potentially large GMs.

- e.g. social network

- reduces necessity to deal with entire network due to C.I. properties

PBN

- optimisation semantics; representation learning/embedding is also an inference problem.

layers may correspond to different granularities of features. can be viewed as hidden r.v.s.

query 3-MPA (maximum a posterior config. given evidence)
or (most probable assignment)

- Again previous query is subsumed.

Q: recall the distinctions made by JP for HMM

- which y gives highest C.P. mass

- application - classification, explanation

- PGM may yield better results in situations where simple linear class not good.

(*) Mores

- MPA answer depends on framing

i) $\underset{y_1}{\operatorname{argmax}} p(y_1, y_2) \quad y_1^* = 1$ (expected value?)

ii) $\underset{y_1, y_2}{\operatorname{argmax}} p(y_1, y_2)$

$$y_1^* = 0 \quad y_2^* = 0$$

- MPA is different depending on whether there is context in the form of y_2

- PGM makes this explicit \rightarrow joint label or separate?

- Ex: y_1, y_2 connected / in Markov Blanket \rightarrow 'use context'

y_1, y_2 a-separated \rightarrow 'no need for context'

Complexity of Inference

Ex. Proof of NP results \rightarrow signpost of how to allocate your time

- computing $p(X=x|e)$ - NP hard
in GM

) no assumptions about graphical models \rightarrow no. of configuration increases exponentially.

*) And answer for any subset \rightarrow need to enumerate (without help of GM)

Ex: Hardness does not mean not soluble

- 10-708 - In many cases, certain graph structures offer polynomial time solution methods, or approximate with poly. complexity

Ex: class focuses on these generic (not special) cases; tradeoff between rich graphical models and comput. feasibility

Ex: deep learning paradigm \rightarrow no thinking about algorithm; not luxury in GM.

Approaches to inference

- exact inference - guaranteed theoretically to get exact answer dictated by model

- approx inference - approx 'the answer'; more heavily used practically

Ex: many deep learning methods / model architectures may be good; but inference algorithm not \rightarrow hence performance

marginalisation / Elimination

- trivial statistically

① likelihood of pattern E being active
(binary state)

• $P(E)$

• marginalisation

② How expensive is marginalisation; exponential $^{4+} (?)$ ③ clarity

• computation difficult with very long chain (too. summation, even; non PGM)

- use chain decomp (PGM) :-

- two strategies:

1) exponential cost - Hold node together, 2^k configurations; change only one record values of joint probability node (?) ④ clarity complexity: $O(k^n)$
(Naive) $P(E) =$

- 2) (*) enumeration savings \rightarrow not all c.p.s. function of a
- (*) Note $\sum_a p(a)p(b|a) = \phi(b) = p(b)$ (i.e. a function of b)
- more systematic method than 1)
 - repeat until we have $p(e) = \sum_a p(e|a)p(a)$
 - one off summation has cost(K^2)
cost of e.g. $\sum_a p(a)p(b|a)$ is $|b| \times |a|$ i.e. 4 (as $a=2$ states)
(quadratic)
 - # for n eliminations $\rightarrow O(nK^2)$ (quadratic complexity w.r.t. largest config of nodes K)
 - HMM: structure gives opportunities
HMM \rightarrow prob hidden state given entire seq.
 - C.P. $\therefore p(y_1|x_1, \dots, x_T) = \sum_{\{y_i\}_{t=1}^T \setminus y_1} \dots \sum_{y_1} p(y_1, \dots, y_T, x_1, \dots, x_T)$
 - via fact law: \rightarrow How does this affect complexity of inference?
Ex: illustrates inference complex. reduction for HMM (④ ⑤ ⑥ - key junction for full underst.)

(*) $\sum_{y_2} \sum_{y_3} \dots \sum_{y_T} \dots \sum_{y_1} p(y_1)p(x_1|y_1)p(y_2|y_1)$

$\underbrace{\quad}_{\text{excl. } y_1}$ $\underbrace{\quad}_{\text{(any)}}$ $m(x_1, y_2) = p(x_1, y_2)$ ④ ⑤ ⑥ - check

(*) $\sum_{y_3} \dots \sum_{y_T} \dots \sum_{y_2} p(x_2|y_2)p(y_3|y_2)f(x_1, y_2)$

$\underbrace{\quad}_{\text{excl. } y_1}$ $\underbrace{\quad}_{\text{(any)}}$ $m(x_1, x_2, y_3)$ ④ ⑤ ⑥

- repeat
- algorithm is incr. in no. of nodes; quadratic in no. of states
- This is the HMM forward algorithm (elimination machine)

HMM - different elimination (sequence)

$$(*) \sum_{y_1} \dots \sum_{y_T} p(x_1 | y_1) p(y_1 | y_{1-1}) \\ = n(x_1, y_{1-1}) = p(x_1 | y_{1-1})$$

- ultimately: - $n' = p(x_{1:T} | y_{1:T})$...

- Backward algorithm: eliminating nodes from the tail. (15)

- n' takes semantics: intern (prob. of 1st $\frac{1}{2}$ of sequence
given latent states)
of condit. prob. of partial sequence of observed given 1 hidden state.

- n' takes semantics: -

of joint prob of 1st half of sequence given n hidden states.

(2) (AS): Have to solidify understanding of HMM forward-backward

- major invention of the 70s (figuring it out algebraically not trivial)
- ordering the nodes for elimination \rightarrow highly specialised c.i. insight.
- PGMs allow determination of poss. / feasibility via simplified ag.

undirected chains

{ chain models
have recurring pattern

CRFs

(*) Sum-product operation

$$\sum_{\phi} \prod_{F \in \phi} \phi$$

F - set of factors

ϕ - factors

- ex: A key vehicle for understanding complexity for inference on GMS:-
(*) count no. of summations, multiplications

(*) Appl. of sum-product to GMS regional chains

write
(*) Query: $p(X_t, e) = \sum_{x_1} \dots \sum_{x_2} \prod_i p(x_i | p_{ai})$

(to define
an ordering
of summation signs)