

- today:
- finish DGM
 - undirected GM
 - exponential family

3 facts C.I.

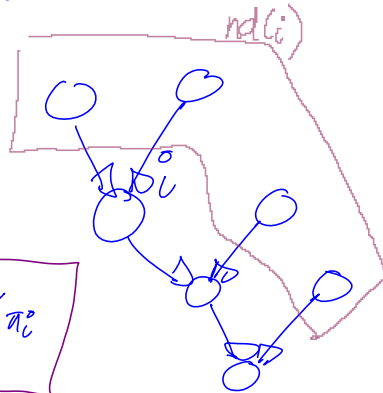
1) can repeat variables $X \perp\!\!\!\perp Y \mid Z, W$

2) decomposition: $X \perp\!\!\!\perp Y, Z \mid W$
 $\Rightarrow \begin{cases} X \perp\!\!\!\perp Y \mid W \\ \text{and} \\ X \perp\!\!\!\perp Z \mid W \end{cases}$

3) trick: extra-conditioning on both sides of eqn
 doesn't change anything
 e.g. $p(x, y) = p(x|y) p(y)$ [always true]

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

let $nd(i) \triangleq \{j : \text{no path from } i \text{ to } j\}$
 "non-descendants"

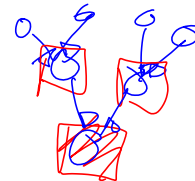


prop: $p \in \mathcal{P}(G) \Leftrightarrow X_i \perp\!\!\!\perp X_{nd(i)} \mid X_{\pi_i}$
 $\forall i \in V$

proof:

\Rightarrow key pt.: let i be fixed

\exists a top. sort. s.t. $nd(i)$ are just before i
 i.e. $(nd(i), i, V \setminus (\{i\} \cup nd(i)))$



let $A \triangleq nd(i) \setminus \pi_i$

(plucking leaves)

[marginalizing out $X_{V \setminus (\{i\} \cup nd(i))}$]

$p(x_i, x_{\pi_i}, x_A) = p(x_i | x_{\pi_i}) \prod_{j \in A} p(x_j | x_{\pi_j})$

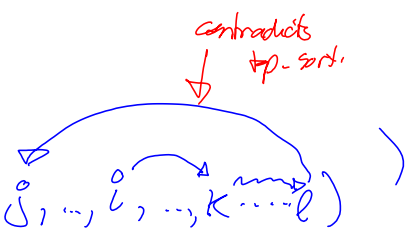
$p(x_i | x_{nd(i)}) = p(x_i, x_{\pi_i}, x_A) = p(x_i | x_{\pi_i}) \prod_{j \in A} p(x_j | x_{\pi_j})$

$$p(x_{\pi_i} | x_A) = \sum_{x_i} [p(x_i | x_{\pi_i})] \prod_{j \in \text{nd}(i)} p(x_j | x_i)$$

= $p(x_i | x_{\pi_i})$

$$\Rightarrow X_i \perp\!\!\!\perp X_{\text{nd}(i)} | X_{\pi_i}$$

\Leftrightarrow Let $1:n$ is a top sort.
 then $\{1, \dots, i-1\} \subseteq \text{nd}(i)$
 (suppose $j \in \{1, \dots, i-1\}$ and $j \notin \text{nd}(i)$)
 $\Rightarrow \exists$ path from i to j i.e. I have $(\dots, j, \dots, i, \dots, k, \dots, \ell)$



$$p(x_V) = \prod_{i=1}^n p(x_i | x_{\{j: j < i\}})$$

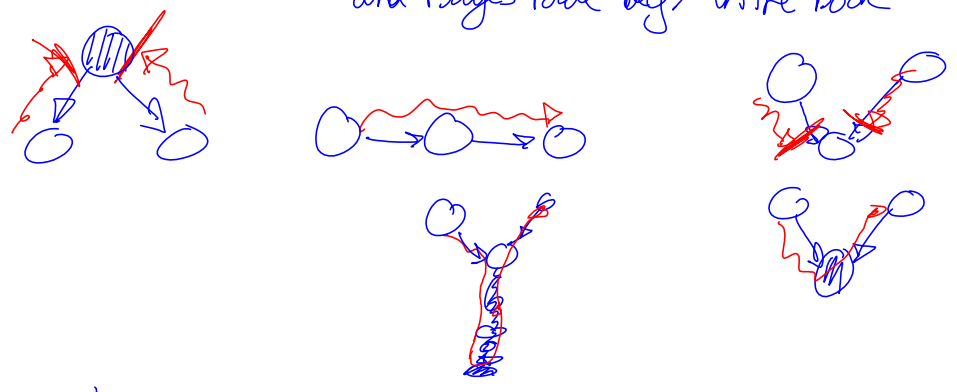
(chain rule)

$$= \prod_{\substack{i=1 \\ i \in \mathcal{I}(G)}}^n p(x_i | x_{\pi_i})$$

[by cond. indep.]

\otimes any other cond. indep properties beyond $X_i \perp\!\!\!\perp X_{\text{nd}(i)} | X_{\pi_i}$ for DAG G ?

yes? \rightarrow see d-separation and Bayes Ball alg. in the book "v-structure"



other properties

inclusion: $E \subseteq E'$ then $\mathcal{I}(G) \subseteq \mathcal{I}(G')$

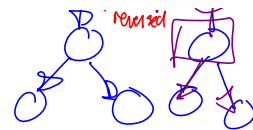
reversal: if G is a directed tree [or forest] ($|\pi_i| \leq 1$, i.e. no v-structure)

then consider G' any other orientation of the tree [take undirected tree version of G]





and pick unique root to orient



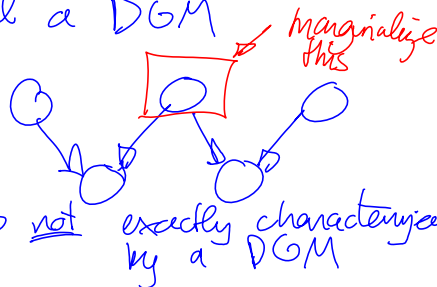
$$I(G) = I(G')$$

→ why direction of edge is not causal in a chain

marginalizing

• marginalize leaves → still a DGM

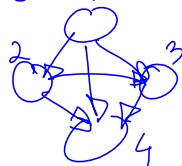
• it's not true in general



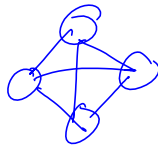
set of distributions obtained is not exactly characterized by a DGM

complete graph

Digraph:



undirected



(clique)

↳ set of nodes with edge between every pair

i.e. C clique $\Leftrightarrow \forall i, j \in C, \exists \{i, j\} \in E$

undirected graphical model

→ aka Markov random field or Markov networks

let $G = (V, E)$ be undirected graph

let \mathcal{C} the set of cliques of G

$$\text{then } I(G) \triangleq \left\{ p : p(x) = \frac{1}{Z} \prod_{C \in \mathcal{C}} \psi_C(x_C) \right\}$$

where $\psi_C(x_C) \geq 0$ "potential"

$$Z \triangleq \sum_x \left(\prod_{C \in \mathcal{C}} \psi_C(x_C) \right) \quad \text{"partition or normalizer"}$$

notes: • $\psi_C(x_C)$ is not related directly to $p(x)$ unlike in DGM

• $\psi_C(x_C) = \text{constant} \cdot \psi_C(x_C)$ doesn't change anything

• sufficient to only consider \mathcal{C}_{\max} = set of maximal cliques

i.e. $C' \subseteq C$ can redefine

$$\psi_{C'}(x_{C'}) = \psi_C(x_C) \cdot \psi_C(x_C)$$

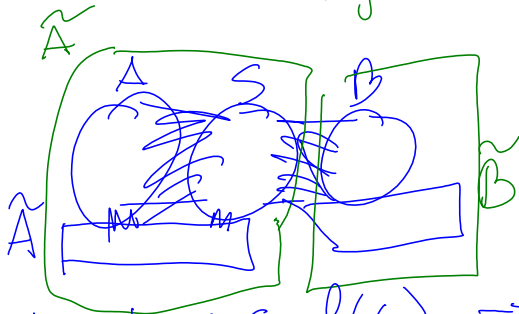
as before: $E \subseteq E' \Rightarrow \mathcal{J}(G) \subseteq \mathcal{J}(G')$
 $E = \emptyset \Rightarrow \mathcal{J}(G) = \{p : \text{all } X_i \text{ are independent}\}$
 $E = \text{all pairs (i.e. complete graph)} \Rightarrow \mathcal{J}(G) = \{\text{all distributions}\}$

if $\psi_c(x_c) > 0 \forall x_c \in C$ write $p(x) = \exp(\underbrace{\sum_{c \in C} \log \psi_c(x_c)}_{\text{negative energy function}} - \log Z)$
 "exp. family"

conditional indep.

def: p satisfies global Markov property \rightarrow any path from A to B passes through S
 iff $\forall A, B, S \subseteq V$ s.t. S separates A from B in G
 disjoint

then $X_A \perp\!\!\!\perp X_B \mid X_S$



prop: $p \in \mathcal{J}(G) \Rightarrow p$ satisfies global Markov property

proof: WLOG $A \cup B \cup S = V$ $\& \triangleq \& \triangleq$ "and"

$\tilde{A} \triangleq A \cup \{a \in V : a \notin A \text{ are not separated by } S\}$
 $\tilde{B} \triangleq V \setminus (S \cup \tilde{A}) \Rightarrow \tilde{A} \& \tilde{B} \text{ are } S\text{-separated}$

* let $C \in \mathcal{C}$; then we can't have both $C \cap A \neq \emptyset$ and $C \cap B \neq \emptyset$

$$p(x) = \frac{1}{Z} \prod_{\substack{C \in \mathcal{C} \\ C \cap A \neq \emptyset}} \psi_C(x_C) \prod_{\substack{C \in \mathcal{C} \\ C \cap B \neq \emptyset}} \psi_C(x_C)$$

$$= f(x_A, x_S) \underbrace{g(x_B, x_S)}_{\text{constant w.r. } x_A}$$

$$p(x_A | x_S) \propto \sum_{x_B} f(x_{A \cup S}) g(x_{B \cup S}) = f(x_{A \cup S}) \sum_{x_B} g(x_{B \cup S})$$

$$p(x_A | x_S) = \frac{f(x_A, x_S)}{\sum_{x_A} f(x_A, x_S)} \quad \text{similarly: } p(x_B | x_S) = \frac{g(x_B, x_S)}{\sum_{x_B} g(x_B, x_S)}$$

$$p(x_A | x_S) p(x_B | x_S) = \frac{f(x_A, x_S) g(x_B, x_S)}{\sum_{x_A} \sum_{x_B} f(x_A, x_S) g(x_B, x_S)} p(x_A, x_B, x_S)$$

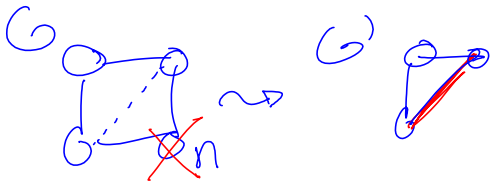
$$= p(x_A, x_B | x_S) \Rightarrow X_A \perp\!\!\!\perp X_B | X_S$$

thm.: Hammersley Clifford

if $\forall x p(x) > 0$, then $p \in \mathcal{P}(G) \Leftrightarrow p$ satisfies global Markov property

marginalization:

let $V' = V \setminus \{n\}$ $E' =$ edges of G connecting all the neighbors of n and removing n



$$p \in \mathcal{P}(G) \Rightarrow p(x_{1:n-1}) \in \mathcal{P}(G')$$

(true for any $a \Rightarrow$ we get closure)

directed vs. undirected GM

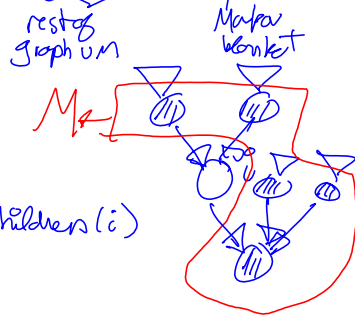
def: Markov blanket for i is the smallest set of nodes $M \subseteq V$ s.t. $X_i \perp\!\!\!\perp X_{V \setminus \{i\}} | X_M$

for UGM:

$$M = \{j : \{i, j\} \in E\} \text{ (neighbors of } i)$$

for DGM

$$M = \pi_i \cup \text{children}(i) \cup \pi_j | j \in \text{children}(i)$$



| | DGM | UGM |
|-----------------|---|--|
| factorization | $p(x) = \prod_i p(x_i x_{\pi_i})$ | $p(x) = \prod_{c \in \mathcal{C}} \psi_c(x_c)$ |
| cond. indep | d-separation [more than $X_i \perp\!\!\!\perp X_{nd(i)} X_{\pi_i}$] | separation $X_A \perp\!\!\!\perp X_B X_S$ |
| marginalization | not closed in general [fine for a leaf] | closed |
| 1.1.1.1.1.1 | " | " |

my guess is:



question: G is a DAG; when can we transform to equivalent UGM?

def: for G a DAG, call \bar{G} the moralized graph of G where \bar{G} is undirected graph with same V

$$E = \left\{ \{i, j\} : (i, j) \in E \right\} \cup \left\{ \{k, l\} : k, l \in \pi_i \text{ for some } i \right\}$$

undirected version of E
"moralization"

no additions (\emptyset) when no "v-structure"

prop: for DAG G with no v-structure [forest]

$$\text{then } I(G) = I(\bar{G}) = I(\text{undirected}(G))$$

but in general, only have $I(G) \subseteq I(\bar{G})$

note: not all subsets of cond. indep. statements can be precisely described by a graph model

[e.g. XOR example where pairwise indep. but not mutual indep.]

$$X_1 \perp\!\!\!\perp X_2$$

$$X_2 \perp\!\!\!\perp X_3 \quad X_1 \perp\!\!\!\perp X_3 \quad \text{but not } X_1 \perp\!\!\!\perp X_2, X_3$$

[$X_3 = X_1 \text{ XOR } X_2$]

general issues in this class

A) representation / parametrization
↳ DGM / UGM
↳ exp. family

B) inference (computing $p(x_A | z_B)$)
↳ sum-product alg.

C) statistical estimation
→ maximum likelihood
→ " " entropy

KL & MLE

let X be discrete set

$$KL(p \parallel q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)}$$

n observations x_1, \dots, x_n

define empirical distribution \hat{p}_n
with p.m.f. (probab. mass fct.)

dist. over X

$$\hat{p}_n(x) \triangleq \frac{1}{n} \sum_{i=1}^n \delta(x-x_i)$$

Kronecker-delta fct.

$$\delta(z) = \begin{cases} 1 & \text{if } z=0 \\ 0 & \text{o.w.} \end{cases}$$

prop.: $\{p_\theta\}_{\theta \in \Theta}$ parametric family on X

then ML for $p_\theta \Leftrightarrow \min KL(\hat{p}_n \parallel p_\theta)$

proof:

$$\begin{aligned} KL(\hat{p}_n \parallel p_\theta) &= \sum_{x \in X} \hat{p}_n(x) \log \frac{\hat{p}_n(x)}{p_\theta(x)} \\ &= -H(\hat{p}_n) - \sum_{x \in X} \hat{p}_n(x) \log p_\theta(x) \\ &= -H(\hat{p}_n) - \frac{1}{n} \sum_{i=1}^n \sum_{x \in X} \delta(x-x_i) \log p_\theta(x) \\ &= -H(\hat{p}_n) - \frac{1}{n} \sum_{i=1}^n \underbrace{\left(\sum_{x \in X} \delta(x-x_i) \log p_\theta(x) \right)}_{\log p_\theta(x_i)} \\ &= \text{const.} - l(\theta) \end{aligned}$$

max. entropy principle : (different than ML) //
in general

idea: consider $\mathcal{P}(X)$ a subset of distributions on X
which satisfies some constraints (usually from data)

• pick $p \in \mathcal{P}(X)$ which maximize entropy

choose \hat{p} by solving $\arg \max_{p \in \mathcal{P}(X)} H(p) = \arg \min_{p \in \mathcal{P}(X)} KL(p \parallel \text{uniform distribution on } X)$

Exponential family

a (flat/canonical) exponential family on X

is a set of distributions defined by two quantities

I) $h(x) d\mu(x) \rightarrow$ reference measure on \mathcal{X}
 reference density \leftarrow base measure \leftarrow counting measure

II) $T: \mathcal{X} \rightarrow \mathbb{R}^p$ called "sufficient statistics" vector
 (aka feature vector)
 members of family have dist.

$$p(x; \eta) d\mu(x) = \exp(\eta^T T(x) - A(\eta)) h(x) d\mu(x)$$

"canonical" parameter $\rightarrow \eta$
 "log-normalizer" defining pieces aka, log-partition fct, cumulant generating fct. $\rightarrow A(\eta)$

* want $1 = \int_{\mathcal{X}} p(x; \eta) d\mu(x)$
 $= \int_{\mathcal{X}} \exp(\eta^T T(x)) e^{-A(\eta)} h(x) d\mu(x)$
 $\Rightarrow A(\eta) \triangleq \log \left(\int_{\mathcal{X}} \exp(\eta^T T(x)) h(x) d\mu(x) \right)$

domain $\Omega \triangleq \{ \eta \in \mathbb{R}^p \mid A(\eta) < \infty \}$

* more generally, consider reparameterization as a subset of the family by defining mapping

$$\eta: \Theta \rightarrow \Omega$$

consider $p(x; \theta) \triangleq p(x; \eta(\theta))$ for $\theta \in \Theta$
 parameter $\rightarrow \theta$

("curved exp-family" if $\eta(\Theta)$ is a curved manifold in Ω)

example: (multinomial family)

$$X \sim \text{Mult}(1, \pi); \quad \mathcal{X} = \{0, 1\}^K$$

$$\pi \in \Delta_K = \left\{ \pi : \sum_{l=1}^K \pi_l = 1, \pi_l \geq 0 \right\}$$

then for valid x (i.e. indicator vector)

Suppose $\pi_c > 0$

$$p(x; \pi) = \prod_{c=1}^k \pi_c^{x_c} = \exp\left(\sum_{c=1}^k (\log \pi_c) x_c\right)$$

$$= \exp(\langle \eta(\pi), x \rangle)$$

where $\eta_c(\pi) = \log \pi_c$

$T(x) = x$

$d\mu(x) =$ counting measure on \mathcal{X}

$h(x) = \mathbb{1}\left\{\sum_c x_c \text{ has exactly one entry equal to } 1\right\}$

here $A(\eta(\pi)) = \mathbb{O}$ here, $\Theta = \text{int}(\Delta_K)$
 $A(\eta(\pi)) = \mathbb{O} \forall \pi \in \Theta$

$$A(\eta) = \log\left(\sum_{x \in \mathcal{X}} h(x) \exp(\langle \eta, x \rangle)\right)$$

$$= \log\left(\sum_{c=1}^k \exp(\eta_c)\right) < \infty \forall \eta \in \mathbb{R}^k$$

$\Omega = \mathbb{R}^k$

$(\eta(\pi))_c = \log \pi_c$

thus $A(\eta(\pi)) = \log\left(\sum_{c=1}^k \frac{1}{\pi_c}\right) = 0$

remark: $\Theta \rightarrow$ dimension $k-1$

$\Omega \rightarrow \mathbb{R}^k$

for any x s.t. $h(x) \neq 0$

$$\sum_{c=1}^k x_c = 1$$

ie. $\sum_{c=1}^k T_c(x) - 1 = 0$

affine linear dep. between components of T

multiple η 's map to same distribution

(no identifiability) "overparameterized"

[not minimal]

for multinomial, minimal exp-family $T(x) = \begin{pmatrix} x_1 \\ \vdots \\ x_{k-1} \end{pmatrix}$