

Lecture 10 - scribbles

Tuesday, October 4, 2016
14:17

today: directed graph model

general issues in this class

- A) representation \rightarrow DGM
 \searrow UGM
 parameterization \rightarrow exponential family } probabilities
- B) inference computing $p(x_Q | x_E)$
 \searrow sum-product junction alg.
 \rightarrow "query" \rightarrow "evidence"
- C) statistical estimation \rightarrow MLE
 maximum entropy
 method of moments

Notation:

n discrete R.V. X_1, \dots, X_n $V \leftarrow$ set of vertices; one R.V. per node

joint $p(X_1=x_1, X_2=x_2, \dots, X_n=x_n) = p(x_1, \dots, x_n)$
 $= p(x_V) = p(x)$

for any $A \subseteq V$ marginal on x_A
 $p(x_A) = P\{X_i = x_i : i \in A\} = \sum_{x_{A^c}} p(x_A, x_{A^c})$
subset of "subscripts" $x_{A^c} \leftarrow$ summing over all possible values for $\{x_i : i \in V \setminus A\}$

$1:n \triangleq \{1, \dots, n\}$

conditional independence: $X_A \perp\!\!\!\perp X_B \mid X_C \quad A, B, C \subseteq V$

(F) $\Leftrightarrow p(x_A, x_B | x_C) = p(x_A | x_C) p(x_B | x_C)$

(C) $\Leftrightarrow p(x_A | x_B, x_C) = p(x_A | x_C) \quad \forall x_A, x_B, x_C \text{ s.t. } p(x_C) > 0$
 $\leftarrow p(x_B, x_C) > 0$

"marginal independence": $X_A \perp\!\!\!\perp X_B \mid \emptyset$

3 facts about C.I.:

3 facts about C.I.:

1) can repeat variables e.g. $X \perp\!\!\!\perp Y, Z \mid Z, W$ is fine to say

2) decomposition: $X \perp\!\!\!\perp Y, Z \mid W \Rightarrow X \perp\!\!\!\perp Y \mid W$
and $X \perp\!\!\!\perp Z \mid W$

3) trick: extra conditioning on both side of equation
maintains true statements

e.g. $p(x, y) = p(x|y)p(y)$ (always true)

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

⊗ pairwise independence $\not\Rightarrow$ mutual independence

e.g. $X_3 = X_1 \text{ xor } X_2$ with $X_1 \perp\!\!\!\perp X_2$
 $\sim \text{Bernoulli}(\frac{1}{2})$

chain rule:
(always true)

$p(x_v) = \prod_{i=1}^n p(x_i | x_{1:i-1})$

last conditional is $p(x_n | x_{1:n-1})$ \rightarrow table with 2^n entries

graph model

$p(x_v) = \prod_{i=1}^n p(x_i | \pi_i)$

parents of i in graph

\rightarrow tables $2^{\max |\pi_i| + 1}$

Directed graph model:

let $G=(V, E)$ be a DAG

(DGM)

a directed graphical model (associated with G) is a family of distributions
aka Bayesian network
over X_v here $\mathcal{J}(G) \triangleq \{ p : \exists \text{ local factors } f_i \text{ s.t. } \}$
(not necessarily unique)

dist. over X_v $p(x_v) = \prod_{i=1}^n f_i(x_i, \pi_i)$

for potentials f_i s.t. 1) $f_i \geq 0$

$f_i: \Omega_{x_i} \times \Omega_{\pi_i} \rightarrow [0, 1]$

2) $\forall i, \sum_{x_i} f_i(x_i, \pi_i) = 1 \quad \forall \pi_i$

is like a CPT

(conditional proba table)

$\pi_i \triangleq \{ j \in V : (j, i) \in E \}$

terminology is: if $p(x_v) = \prod_{i=1}^n f_i(x_i, \pi_i)$ \leftarrow this is where G comes into play

terminology is: if $p(x) = \prod_{i=1}^n f_i(x_i, \pi_i)$ \leftarrow this is where G comes into play
 then say that " p factorizes according to G " (conditional proba table)

eg. $p(x,y) = p(x)p(y|x) \rightarrow \begin{matrix} G_x \\ G_y \end{matrix}$
 $= p(y)p(x|y) \rightarrow \begin{matrix} G_y \\ G_x \end{matrix}$

"leaf plucking" property Fundamental property of DGM

let n be a leaf (i.e. n is not parent of anything)

then $p(x_{1:n-1}) \in \mathcal{F}(G - \{n\})$
 $= \prod_{j \neq n} f_j(x_j, \pi_j) \in \mathcal{F}$

proof:

$p(x_n, x_{1:n-1}) = f_n(x_n, \pi_n) \prod_{j \neq n} f_j(x_j, \pi_j)$

$n \in \pi_j \forall j$
 \downarrow

$p(x_{1:n-1}) = \sum_{x_n} p(x_n, x_{1:n-1}) = \left(\sum_{x_n} f_n(x_n, \pi_n) \right) \prod_{j \neq n} f_j(x_j, \pi_j)$
 \downarrow by definition

proposition: if $p \in \mathcal{F}(G)$; let $\{f_i\}$ be a factorization w.r.t. G
 $\prod_{i=1}^n f_i$

then $\forall i \quad p(x_i | \pi_i) = f_i(x_i, \pi_i)$

proof: WLOG, let $\{1, \dots, n\}$ be top sort
 [let i be fixed, want to prove that $p(x_i | \pi_i) = f_i(x_i, \pi_i)$]

then $p(x_n, x_{1:n-1})$
 \downarrow
 leaf

factorization for DAG where x_{n-1} is now a leaf

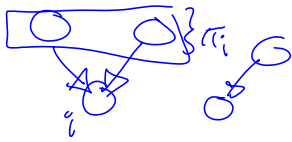
thus pluck it to get $p(x_{1:n-1}) = \prod_{j \neq n} f_j(x_j, \pi_j)$

keep doing this [in reverse top sort order $n-1, \dots, i+1$] to get:

$p(x_{1:i}) = f_i(x_i, \pi_i) \prod_{j < i} f_j(x_j, \pi_j)$

let $A \triangleq 1:i-1 \setminus \pi_i$ so that $1:i = \pi_i \cup A \triangleq \mathcal{A}(x_{1:i-1})$

$$p(x_i | x_{\pi_i}) = \frac{\sum_{x_A} p(x_i, x_{\pi_i}, x_A)}{\sum_{x_A} \sum_{x_i} p(x_i, x_{\pi_i}, x_A)} = \frac{f_i(x_i, x_{\pi_i}) \sum_{x_A} g(x_i, x_{\pi_i})}{\sum_{x_i} f_i(x_i, x_{\pi_i}) \sum_{x_A} g(x_i, x_{\pi_i})} = f_i(x_i, x_{\pi_i})$$



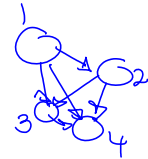
equivalently, we can write: $\mathcal{J}(G) = \{p: p(x_V) = \prod_{i=1}^n p(x_i | x_{\pi_i})\}$

Note: adding edges \Rightarrow more distributions

i.e. $G = (V, E) \hookrightarrow G' = (V, E')$ with $E' \supseteq E$
then $\mathcal{J}(G) \subseteq \mathcal{J}(G')$

Examples: $E = \emptyset \Rightarrow \mathcal{J}(G)$ contains only fully independent dist.
i.e. $\prod_i p(x_i)$
(trivial graph)

complete digraph: get $p(x_V) = \prod_i p(x_i | x_{1:i-1})$
for all $i: j$
either $(i \rightarrow j) \in E$ or $(j \rightarrow i) \in E$
 \Rightarrow all distributions are in $\mathcal{J}(G_{\text{complete}})$ (like chain rule)



removing edges impose more factorization constraints

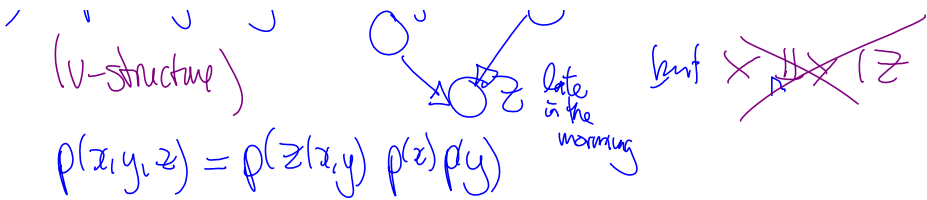
Basic 3 nodes graphs

Factorization \Leftrightarrow cond. indep. assumption

1) Markov chain: $X \rightarrow Z \rightarrow Y$
"future" "past" "present"
 $X \perp\!\!\!\perp Y \mid Z$
but $X \not\perp\!\!\!\perp Y$

2) Latent cause (hidden variable)
 $X \rightarrow Z \rightarrow Y$
"grey hair"
 $X \perp\!\!\!\perp Y \mid Z$
but $X \not\perp\!\!\!\perp Y$

3) explaining away (V-structure)
 $X \rightarrow Z \leftarrow Y$
"shoe size" "broken watch"
 $X \perp\!\!\!\perp Y$
but $X \not\perp\!\!\!\perp Z$



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

non-monotonic property of conditioning

$$p(\text{alien} | \text{tiny}) \quad p(\text{alien} | \text{late}) > p(\text{alien})$$

$$\text{but } p(\text{alien} | \text{late, broken watch}) < p(\text{alien} | \text{late})$$