

LECTURE 3: DIRECTED GRAPHICAL MODELS (a.k.a. Bayesian Networks)

1. Statistical independence

→ From Lecture 1:

Most general model for $p(x_1, \dots, x_d)$ is intractable

$$\text{e.g. } x_i \in \{0, 1\}, \quad \begin{cases} p(x_1, \dots, x_d) = \theta_{x_1, \dots, x_d} > 0 \\ \sum_{x_1, \dots, x_d} \theta_{x_1, \dots, x_d} = 1 \end{cases}$$

representation
inference
learning } are cursed by dimension

→ Possible solutions:

- (Lecture 2) Forget everything except mean/variance
 \Rightarrow Gaussian model, tractable

- Independence:

- assume $p(x_1, \dots, x_d) = p(x_1) \cdots p(x_d)$

for $x_i \in \{0, 1\}$, only need d parameters

$$p(x_i = 1) = \theta_i \in \{0, 1\}$$

But: very strong assumption!

- in-between: e.g.

$$p(x_1, \dots, x_d) = p(x_1, \dots, x_s) \cdots p(x_{d-s+1}, \dots, x_d)$$

$$\Rightarrow \sum_{s=1}^d \binom{d}{s} \text{ parameters}$$

More expressive, but still limited ...

- Conditional independence:

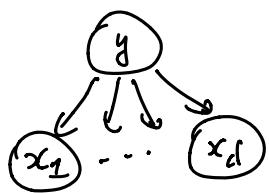
Def: $X \perp Y | Z$ if $p(x, y | z) = p(x | z) p(y | z)$
 equivalently, $p(y | z, x) = p(y | z)$

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

Note: \neq then $X \perp Y$!!

Ex: - Naive Bayes

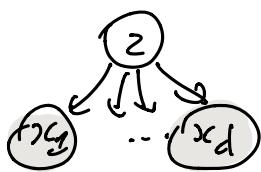
$$p(x_1, \dots, x_d, y) = p(y) \prod_i^{d+1} p(x_i | y)$$



e.g. $y = \text{label}$, $\{x_i\}$ C.I. given label

- Latent variable / cause

$$p(x_1, \dots, x_d) = \int p(z) \prod_i p(x_i | z) dz$$



Note: always true if (x_1, \dots, x_d) are "exchangeable"
 i.e. $(x_1, \dots, x_d) \stackrel{d}{=} (x_{\pi(1)}, \dots, x_{\pi(d)})$ for any permutation π

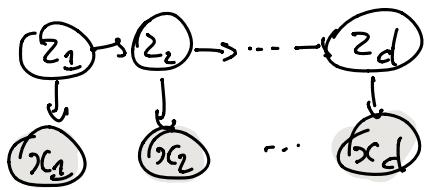
(De Finetti's theorem)

but: $p(x_i | z)$ may be complicated ...

- Markov / Hidden-Markov Models



$$p(x_1, \dots, x_d) = p(x_1) p(x_2 | x_1) \cdots p(x_d | x_{d-1})$$



$$p(x_1, \dots, x_d) = \int \dots \int p(z_1) p(z_2 | z_1) \cdots p(z_d | z_{d-1}) \times \\ \prod_i p(x_i | z_i) \cdot d z_1 \dots d z_d$$

e.g. • $z_i \in \{1, \dots, K\}$

$$x_i | z_i = k \sim \text{Cat}(\theta_k)$$

$$x_i | z_i = k \stackrel{\text{or}}{\sim} \mathcal{N}(\mu_k, \Sigma_k)$$

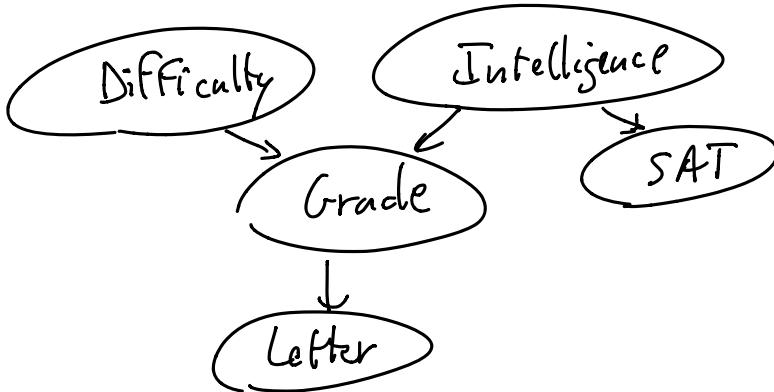
$$\bullet z_i | z_{i-1} \sim \mathcal{N}(A z_i + b, \Sigma_z)$$

$$x_i | z_i \sim \mathcal{N}(C z_i + d, \Sigma_x)$$

(state-space model
Kalman Filtering)

2. Directed Graphical Models

Ex:



(From Koller & Friedman)

- . $D \perp I ?$ yes
- . $D \perp I | G ?$ no e.g. $P(I|G, D=\text{hard}) > P(I|G)$
("explaining away")
- . $G \perp S ?$ no
- . $G \perp S | I ?$ yes
- . $D \perp I | L ?$ no

Factorize: $P(D, I, G, S, L) = P(D) P(I) P(G|D, I) P(S|I) P(L|G)$

More generally, consider a DAG $G = (V, E)$

$$\text{pa}(i) = \{j \text{ s.t. } (j, i) \in E\} \quad (\text{parents of node } i)$$

We want: $P(x_1, \dots, x_d) = \prod_{i=1}^d P(x_i | x_{\text{pa}(i)}) \quad (*)$

[Def: If this holds, we say P factorizes according to G]

Def: A **Directed graphical model** structure $G = (V, E)$ encodes the following (local) independencies, denoted $I_P(G)$:

$$\forall i, (X_i \perp X_{\text{non-descendants}(i)} | X_{\text{pa}(i)})$$

Fact:

If p obeys the C.I.s in $I_p(G)$, it factorizes as \star

Proof: use the chain rule with any topological ordering of the variables relative to G
(i.e. $(i,j) \in E \Rightarrow i \leq j$) 

Independencies and I-maps

Def: For a prob. distribution $p(x)$, we denote by $I(p)$ the set of C.I. of the form $(X \perp Y | Z)$ that hold in p .

Def: (I-map)

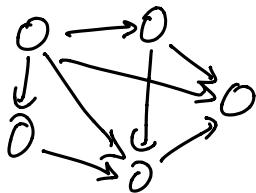
We say the K (e.g. graph) is an I-map for $I(p)$ if K "encodes" a set of C.I.s $I(K)$ with $I(K) \subseteq I(p)$

Fact:

If G is an I-map for p , then p factorizes according to G

indeed, G encodes at least the C.I.s in $I_p(G)$, which is sufficient for factorization -

Ex:



Fully-connected graph

- G : $\begin{matrix} & \text{---} \\ & | \\ \text{---} & \text{---} \end{matrix}$ G is always an I-map

- G : $\begin{matrix} & \text{---} \\ & | \\ \text{---} & \text{---} \end{matrix}$: $I(G)$ is large, requires independence!

→ Do we have the converse? Yes!

[Fact]: If p factorizes according to G , then G is an I-map for p .

exercise: check on Student example!

D-separation

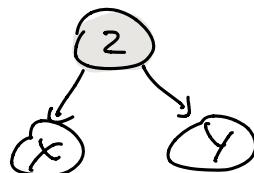
→ three node graphs: for (X, Y, Z) , when do we have $X \perp\!\!\!\perp Y \mid Z$?

- chain: $\begin{matrix} & X & \text{---} & Z & \text{---} & Y \\ & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \end{matrix}$ ("causal trail")

$$p(y|x,z) = p(y|z) \quad \checkmark$$

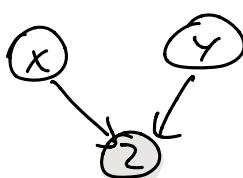


- "common cause":



✓

- "v-structure":



~~$X \perp\!\!\!\perp Y \mid Z$~~ !

(explaining away)

→ d-separation (directed separation)

Def: A path $X \rightsquigarrow Y$ is active in the following scenarios:

- $X \rightarrow Z \rightarrow Y$ and Z is not observed
- $X \leftarrow Z \leftarrow Y$ " "
- $X \leftarrow Z \rightarrow Y$ " "
- $X \rightarrow Z \leftarrow Y$ and Z or any descendant of Z is observed

L

Def: We say X and Y are d-separated given Z , denoted $d\text{-sep}_G(X, Y|Z)$ if there is no active path path between any $X_i \in X$ and $Y_i \in Y$ given Z .
We then write $\mathcal{I}(G) = \{(X \perp Y | Z) : d\text{-sep}_G(X, Y | Z)\}$

Note: we have $\mathcal{I}(G) \subseteq \mathcal{I}(P)$ iff P factorizes over G .
(Thm)

$\stackrel{\text{def}}{\iff}$
any d-sep statement is satisfied by such P .

Q: What about $\mathcal{I}(G) = \mathcal{I}(P)$?

Not true in general:

| | | |
|-------|-------|-------|
| | $y=0$ | $y=1$ |
| $x=0$ | 0.4 | 0.6 |
| $x=1$ | 0.4 | 0.6 |

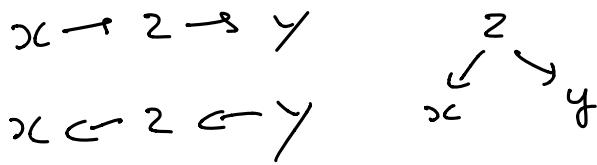
$\rightarrow X \perp Y$, yet we have that

p factorizes over $(X \rightarrow Y)$

Thm: For almost all distributions p that factorize over G , we have $I(G) = I(p)$
(i.e. except a set of measure 0)

Remark: generally, no uniqueness in the choice
of graph

e.g.



Note: this is different for causal models!