

Compact Representation of Distributions
 Directed Acyclic Graphical Models
 Bayes-ball Algorithm

AG: - typo in Fenchel dual:

$$D(y) = -f^*(-y^*) - g^*(A^T y)$$

- extra hints now included
- Due Monday!

Midterm: - cancel class monday: 3pm in class!

- Midterm November 12, 17 TBD
- tutorial moved? TBD

Compact Representations of Joint Distributions

Notation for today:

- we'll use x_j for $(x_i)_j$

training example i ; vector, pick j

Consider $x \in \{0,1\}^d$ (binary vectors)

We want to model $p(x)$ (joint distribution)

- why? scientific discovery, outlier detection, $p(x_i | y_i)$ in generative model, $p(x_A | x_B)$
- previous: naive bayes, Gaussian, Mixture Models
- today: "graphical" models

Basic Idea behind directed acyclic graphical models:

- use product rule repeatedly:

$$p(x) = p(x_1) p(x_2 | x_1) p(x_3 | x_2, x_1) \dots p(x_d | x_{1:(d-1)})$$

$$= \prod_{j=1}^d p(x_j | x_{1:(j-1)})$$

2^{D-1} parameters

- too many parameters

Solutions:

- "parsimonious" parameterization

Can also do both \rightarrow

$$p(x_j | x_{1:(j-1)}) = f(w^T x_{1:(j-1)})$$

D parameters

- conditional independence

eg. naive bayes: $x_j \perp x_{1:(j-1)} | y_i \Rightarrow p(x_j | x_{1:(j-1)} | y_i) = p(x_j | y_i)$

useful for time: markov chain: $x_j \perp x_{1:(j-2)} | x_{j-1} \Rightarrow p(x_j | x_{1:(j-1)}) = p(x_j | x_{j-1})$

General: $x_j \perp x_{1:(j-1) \setminus \pi(j)} | \pi(j)$ "parents"

Directed Acyclic Graphical (DAG) Models

"Bayesian Networks"

"Belief Networks"

"Causal Networks" \rightarrow (you need a justification to interpret edges causally)

Graph $G = (V, E)$

\downarrow "vertices" \downarrow "edges"

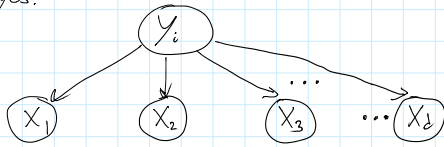
V : random variables x_j

E : $x_{\pi(j)} \rightarrow x_j$

$\pi(j)$: parents of j

$$p(x) = \prod_{j=1}^d p(x_j | x_{\pi(j)})$$

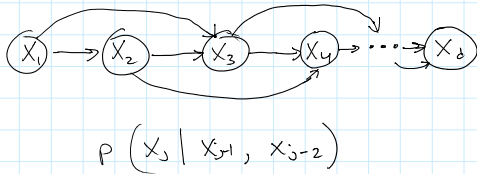
Naive Bayes:



Markov Chain:

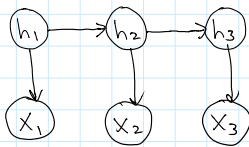


2nd Order Markov Chain:

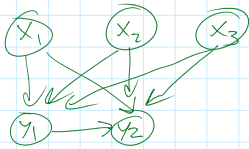


$$p(x_j | x_{j-1}, x_{j-2})$$

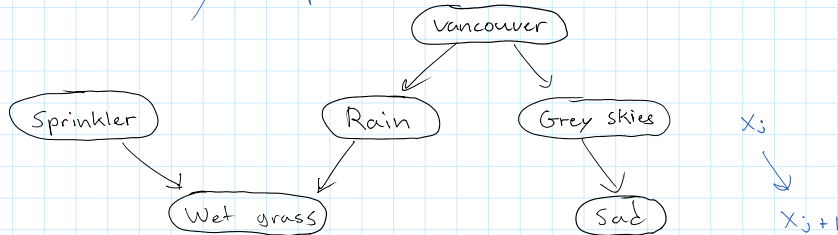
Hidden Markov Model: (Kalman filtering)



Conditional DAG:



Conditionality Example



Learning

- fit each $p(x_i | x_{\pi(i)})$ independently
- "inference"
 - computing $p(x)$
 - computing $p(x_i | x_{\pi(i)})$ } easy
 - computing $p(x_i | x_{j+1})$ } #P-Hard

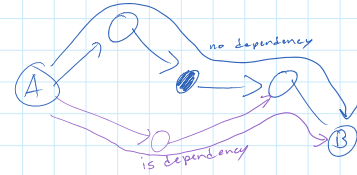
Conditional Independence Properties

- We can use the graph to check whether or not

$$X_A \perp\!\!\!\perp X_B \mid X_C$$

follows from conditional independence assumptions.

$$X_A \perp\!\!\!\perp X_B \mid X_E \iff A \text{ and } B \text{ are "d-separated" given } E$$



A and B are d-separated if for all paths 'P' between A and B, at least one of the following holds:

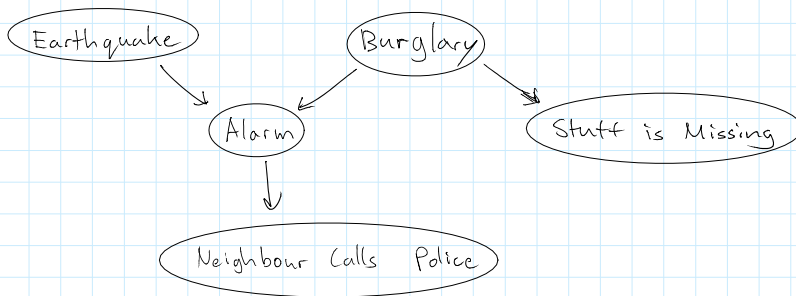
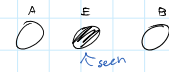
1. P includes a "chain": $A \rightarrow E \rightarrow B$

2. P includes a "fork": $A \leftarrow E \rightarrow B$

note: no E were

3. P contains a "collider": $A \rightarrow C \leftarrow B$

where C and all its descendants are unobserved



Earthquake ~~\perp~~ call

dependent: increase chances

Earthquake \perp call | Alarm

independent

Alarm ~~\perp~~ Stuff Missing

dependent

Alarm \perp stuff Missing | Burglary

independent

Earthquake \perp Burglary

independent

Earthquake ~~\perp~~ Burglary | Alarm

dependent
"explaining away"

Earthquake ~~\perp~~ Burglary | Call

dependent (descendant of alarm)

Burglary ~~\perp~~ Call

Call ~~\perp~~ stuff is missing

Gaussian DAGs

$$p(x_j | x_{\pi(j)}) \sim \mathcal{N}(\mu_j + \bar{w}_j^T x_{\pi(j)}, \sigma_j^2)$$

$$\Leftrightarrow x \sim \mathcal{N}(\mu, \Sigma)$$

$$\Sigma^{-1} = LL^T$$

non-zero pattern of L comes from G

Plate Notation

