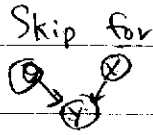


Matlab Tutorial See online handouts...

Discriminative Learning

ML & OML (hard)



Skip for now...



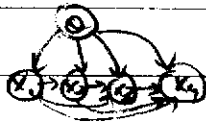
pdf is suboptimal

with Bayesian Deriv... $p(y|x, \theta)$

Graphical Models

- How to do ML & inference & efficiently using what we've learned?

"multinomial" learning with ML



just estimate counts

$$\begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$



but, what if structured?

- Outline: Chapter 2: Directed Graphs, Undirected, Separation
- Maximization, Bayes Ball, Chapter 3: Node Elimination
- Chapter 8: Complete Graphs, Directed Learning, Undirected Learning, IPF

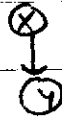
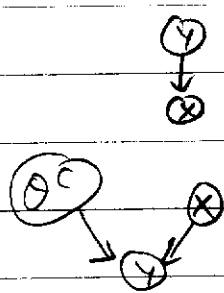
Next Week

ch. 3 & 8

Assignment #2
avail tomorrow

due in 2
weeks at noon

Discriminative Learning



pdf is suboptimal
 $p(y|x, \theta)$ is used
 as marginal

- ML & CML
- Bayesian Derivation

$p(y|x, \theta) = \text{bernoulli}$ (two-class)

$p(y=1|x, \theta) = \mu(x)$

$p(y=0|x, \theta) = 1 - \mu(x)$

$p(y|x) = \mu(x)^y (1 - \mu(x))^{1-y}$

$p'(y|x) = \mu^y (1 - \mu)^{1-y}$

$\mu(x) := \frac{1}{1 + e^{-\eta(x)}} = \frac{1}{1 + e^{-\theta^T x}}$

$\mu = \frac{1}{1 + e^{-\eta}}$

$\eta = \log\left(\frac{\mu}{1-\mu}\right)$

$\frac{d\eta}{d\mu} = \frac{d \log\left(\frac{\mu}{1-\mu}\right)}{d\mu} = \frac{1}{\mu(1-\mu)}$

$\frac{d\mu}{d\eta} = \mu(1-\mu)$

$L(\theta) = \prod_{i=1}^N p(y_i|x_i; \theta) = \prod_{i=1}^N \mu_i^{y_i} (1 - \mu_i)^{1-y_i}$

$l(\theta) = \sum_{i=1}^N y_i \log \mu_i + (1 - y_i) \log (1 - \mu_i)$

$\nabla_{\theta} l = \sum_{i=1}^N \frac{y_i}{\mu_i} - \frac{1-y_i}{1-\mu_i} \frac{d\mu_i}{d\eta} \frac{d\eta}{d\theta} = \sum_{i=1}^N \frac{y_i - \mu_i}{\mu_i(1-\mu_i)} X_i$

$= \sum_{i=1}^N (y_i - \mu_i) X_i = 0$
 $\sum_{i=1}^N y_i X_i = \sum_{i=1}^N \mu_i X_i = \sum_{i=1}^N \theta^T X_i X_i$ $X^T X \theta = X^T y$
 $\theta = (X^T X)^{-1} X^T y$

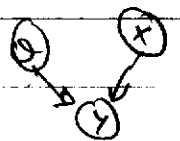
- Conditional Bayesian Inference
& Conditional Max Likelihood

- previous def'n of $p(y|x)$ was implicit, did not specify $p(y|x)$ model as a conditioning of $p(x,y)$, i.e. $p(y|x) = \frac{p(x,y)}{\sum_y p(x,y)}$
- explicit definition of $p(x,y)$ requires different treatment here $p(x,y|\theta)$ is computable yet we say θ only parametrizes the distrib. conditionally

joint inference

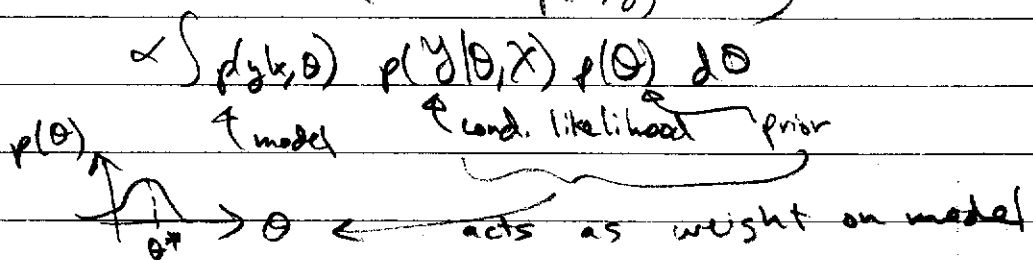


conditional



want conditional from data:

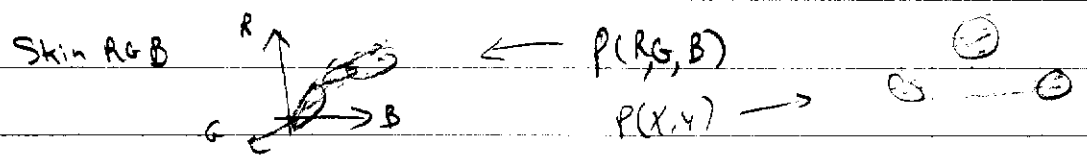
$$\begin{aligned}
 p(y|x) &:= p(y|x, X, Y) = \int p(y, \theta | x, X, Y) d\theta \\
 &= \int p(y|x, \theta, X, Y) p(\theta | x, X, Y) d\theta \\
 &= \int p(y|x, \theta) \{ p(\theta | x, Y) \} d\theta \\
 &= \int p(y|x, \theta) \left\{ \frac{p(y|\theta, X) p(\theta, X)}{p(X, Y)} \right\} d\theta \\
 &= \int p(y|x, \theta) \left\{ \frac{p(y|\theta, X) p(X|\theta) p(\theta)}{p(X, Y)} \right\} d\theta \\
 &= \int p(y|x, \theta) \left\{ \frac{p(y|\theta, X) p(X) p(\theta)}{p(X, Y)} \right\} d\theta
 \end{aligned}$$



$$\begin{aligned}
 p(y|x) &\approx p(y|x, \theta^*) \quad \theta^* = \underset{\theta}{\operatorname{argmax}} p(y|\theta, X) p(\theta) \leftarrow \text{conditional MAP} \\
 p(y|x) &\approx p(y|x, \theta^*) \quad \theta^* = \underset{\theta}{\operatorname{argmax}} p(y|\theta, X) \leftarrow \text{conditional ML}
 \end{aligned}$$

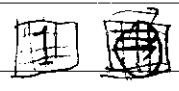
- Assignment #2, EM, due date, 4 datasets, Matlab, C, ...

- EM motivation:



- Matlab Tutorial - 4 functions

- 4 datasets



digits, how represented
lexicographically
as vectors
→ x

help fn
load filename
print -depsec filename

- Review ML & CML — implicit vs. explicit

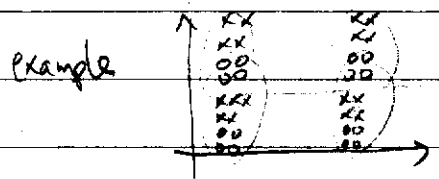
$$l(\theta) = \sum_i \log p(x_i, y_i | \theta) \leftarrow \text{splits data set}$$

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

$$Q(\theta) = \sum_i \log p(y_i | x_i, \theta)$$

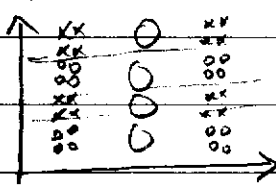
$$= \sum_i \log p(x_i, y_i | \theta) - \sum_i \log \sum_y p(x_i, y | \theta)$$

log-sum again is EM-like



EM
l = -8.0
Q = -1.7

but negated, even worse
repulsion term



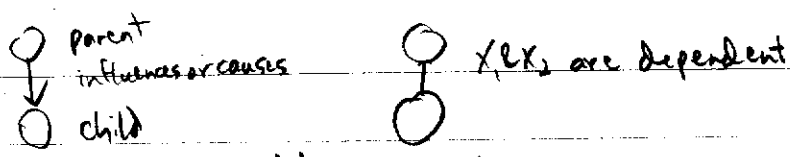
l = -54.7
Q = 0.4

← used in speech rec.

- seen ML, etc. for simple models (learning, inference, using pdf)
have graphical models...

discrete but easily extendable to continuous vars

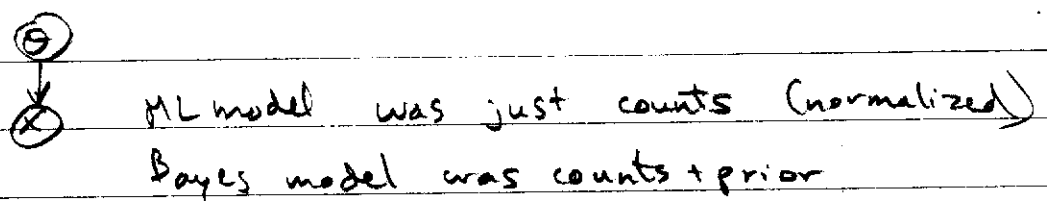
- Next week: ch. 5 & 16



Graphical Models - Directed vs, Undirected

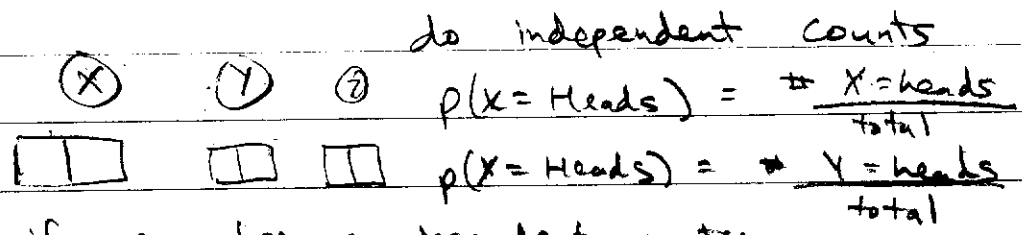
- simplest case: Bernoulli: $X \begin{matrix} \text{Heads} & \text{Tails} \\ 0.7 & 0.3 \end{matrix} \rightarrow p(X)$

multi dim: multinomial $X \in \{A, B, C, \dots\}$

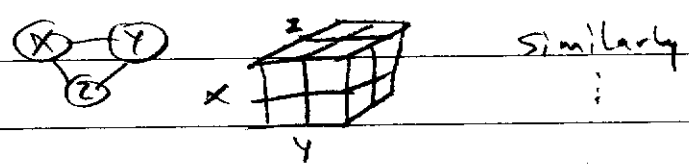
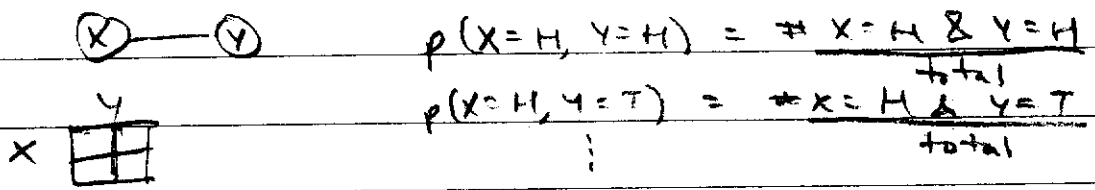


- if we observe independent events

i.e. rolling 2 or 3 dice at same time independent

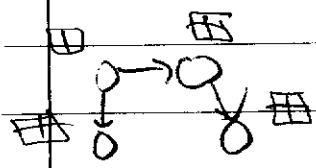


if we observe dependant counts:



get tables of size C^M for M vars.

- Graphical Models: how to learn parameters?
how to do inference?



how to fill tables from data?

EFFICIENCY \rightarrow storage / memory
 \rightarrow computation

already saw this

- Graph is called $G = (X, E)$ or (V, E)

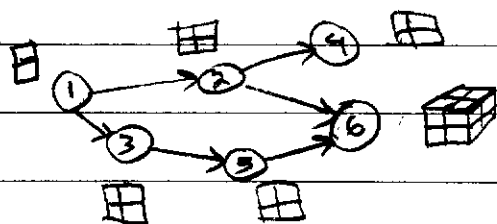
\uparrow nodes (random vars) \uparrow edges \uparrow vertices (random vars) \uparrow edges

$$X = \{X_1, \dots, X_n\} \quad E = \{(X_i, X_j) : i \neq j\}$$

$$X_c = \{X_1, X_3\} \quad C \text{ is a subset}$$

if G is Acyclic then $p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | x_{\pi_i})$

graph specifies pdf qualitatively / structurally
 still need to specify tables quantitatively with numerical
 ∴ graph \equiv family of possible pdfs.



- 1: flu
- 2: fever
- 3: sinus infection
- 4: temperature
- 5: sinuses swell
- 6: headache

Topological Graph: Order the nodes $1, \dots, n$ such that all nodes in X_{π_i} appear before X_i

If topological $\{X_i \perp\!\!\!\perp X_{V_i} \mid X_{\pi_i} \quad \forall i\}$

\uparrow all vars appearing before X_i

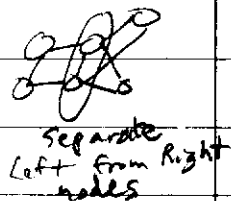
in other words $p(x_i | X_{V_i}, X_{\pi_i}) = p(x_i | X_{\pi_i})$

Example above: $p(x_1, \dots, x_6) = p(x_1) p(x_2 | x_1) p(x_3 | x_1) p(x_4 | x_2) p(x_5 | x_3) p(x_6 | x_2, x_5)$

$$p(x_4 | x_1, x_2, x_3) = \frac{p(x_1, x_2, x_3, x_4)}{p(x_1, x_2, x_3)} = \frac{\sum_{x_5} \sum_{x_6} p(x)}{\sum_{x_4} \sum_{x_5} \sum_{x_6} p(x)} = p(x_4 | x_2)$$

∴ $x_4 \perp\!\!\!\perp x_1, x_3 \mid x_2$
 read off the graph

- Trickier case: $X_1 \perp\!\!\!\perp X_6 \mid X_2, X_3$? harder to read off graph
- Cumbersome to compute $p(X_1 \mid X_6, X_2, X_3)$
- Want fast algorithm to check cond. independence
- Intuition: "separation" or "blocking" could imply cond. indep?
i.e. X_2, X_3 separate path from X_1 to X_6 so indep?



this definition is true for undirected graphs
for directed, need more sophisticated algorithm which uses arrows to determine "d-separation".

- Bayes Ball Algorithm: efficiently tests cond. indep. statements for directed graphs.

"side note": an edge doesn't necessarily imply dependence

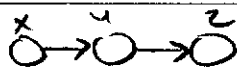
i.e. $X \rightarrow Y$ could have a table with numerical entries that are independent.

i.e. $p(x,y) = x \begin{matrix} y \\ \hline \end{matrix}$ where $x \perp\!\!\!\perp y$

So conditional indep. tests for a property that must hold for this graph. But, due to numerical entries some extra independence properties could exist we won't detect.

3 types of subgraphs:

causal chain

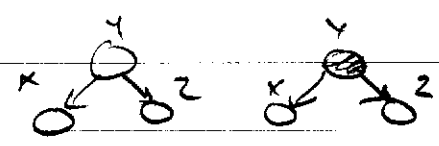


$$p(x,y,z) = p(x)p(y|x)p(z|y) \quad x \perp\!\!\!\perp z \mid y$$

$$p(z|x,y) = p(z|y) \quad \text{conditioning on } y \text{ blocks path from } x \text{ to } z$$

$x = \text{trip}$
 $y = \text{fall}$
 $z = \text{bruise}$

one cause multi effects

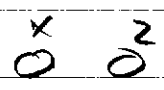
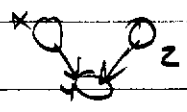


x = sore throat
y = flu
z = temperature

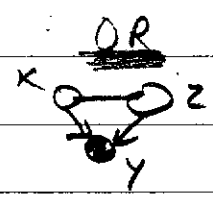
cond. on y
"blocks" path from x to z

$$p(x, y, z) = p(y) p(x|y) p(z|y)$$
$$p(x, z|y) = p(x|y) p(z|y)$$
$$x \perp\!\!\!\perp z | y$$

multi cause one effect



x = rain
y = wet driveway
z = car oil leak



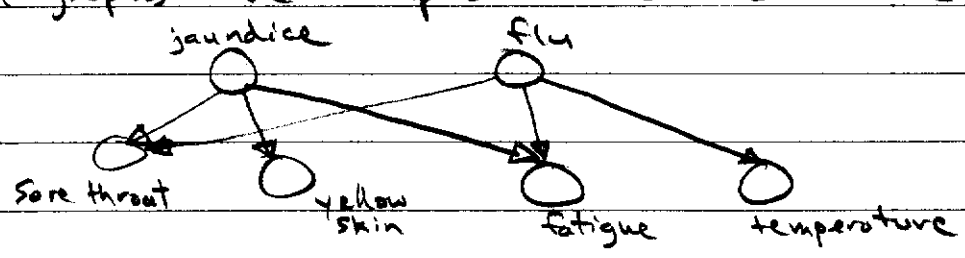
$$p(x, y, z) = p(x) p(z) p(y|x, z)$$
$$x \perp\!\!\!\perp z$$

cond. on y "unblocks" path from x to z

"V" structure, "explaining away"

- General graphs are composed of combos of the above

i.e.
multicause
multieffect



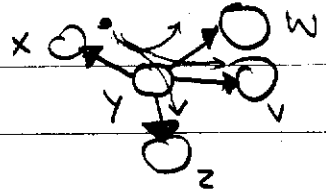
- need to compute d-separation for cond. indep.

- Bayes Ball Alg asks if $X_A \perp\!\!\!\perp X_B | X_C$

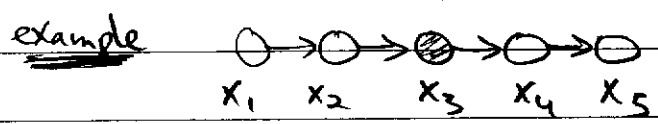
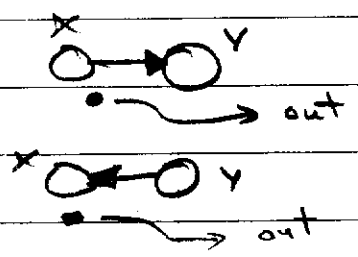
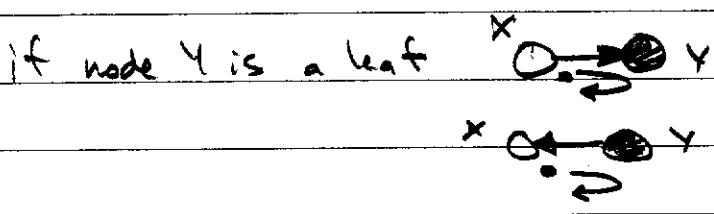
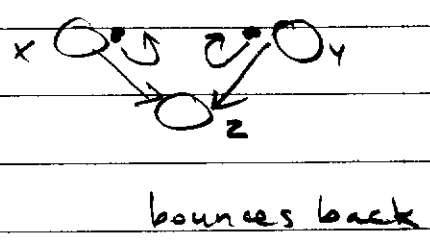
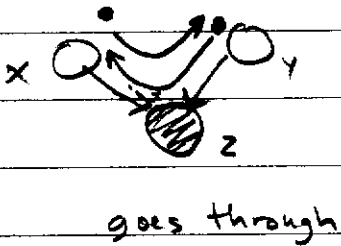
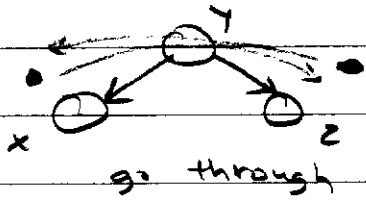
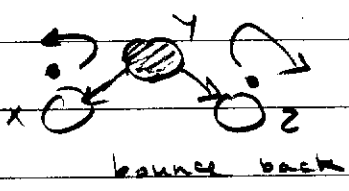
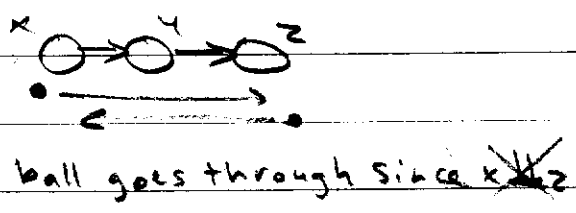
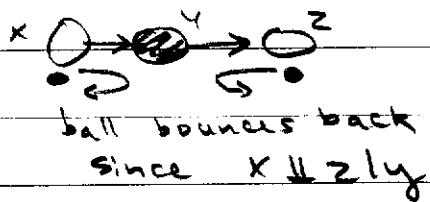
i.e. for a directed graph do nodes X_C d-separate nodes X_A & X_B

- 1) shade nodes X_C (i.e. observed)
- 2) place a ball at each node in X_A
- 3) bounce around the graph according to some rules to see if any of the balls reaches X_B
- 4) if no balls reach X_B , then $X_A \perp\!\!\!\perp X_B | X_C$ IS TRUE else FALSE

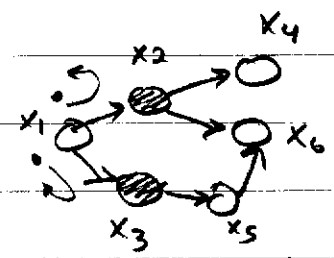
- Rules :
- Balls can travel along/against arrows
 - pick any path & any outgoing path & test each at a time.



look at 3 canonical subgraphs



example

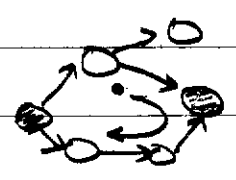


X_1, X_2, X_4 fails
 X_1, X_2, X_6 fails
 X_1, X_3, X_5 fails

X_1 to X_5 blocked by X_3
 X_4 to X_2 blocked

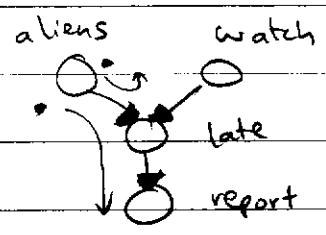
$X_1 \perp\!\!\!\perp X_6 \mid \{X_2, X_3\}$

But $X_2 \perp\!\!\!\perp X_3 \mid \{X_1, X_6\}$
 false

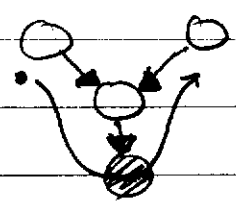


because of "V" structure

example



$aliens \perp\!\!\!\perp watch$



$aliens \not\perp\!\!\!\perp watch \mid report$

- Bob is waiting for alice but can't know if she is late instead a security guard says if she is
- she can be late if aliens abduct or Bob's watch is ahead (day light saving)
- report of her being late connects probs of aliens & watch, If watch is ahead, $p(alien=true) \downarrow$

Final Note: - graph implies $\prod p(x_i \mid x_{\pi_i})$

- numerical can imply more
- set of pdfs implied by graph
- set of pdfs that satisfy all cond. indep. properties given by Bayes Ball

i.e. $X_A \perp\!\!\!\perp X_C \mid X_B \Rightarrow$ all pdfs in graph satisfy this
 $X_A \not\perp\!\!\!\perp X_C \mid X_B \Rightarrow$ some pdfs in graph don't satisfy this