Readings:

Review: K&F: *2.1*, 2.5, 2.6

K&F: 3.1

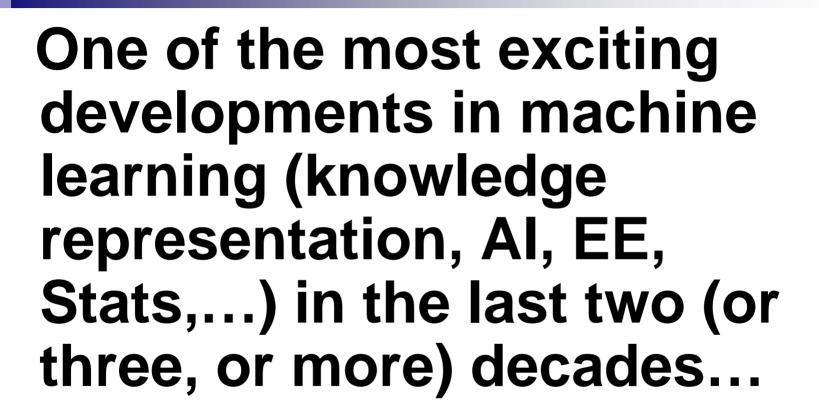
Introduction

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

September 13th, 2006

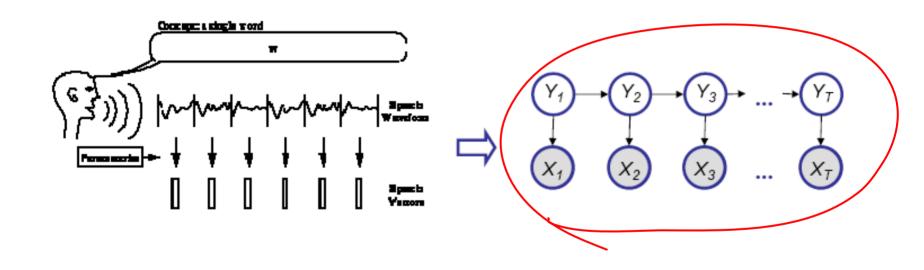


My expectations are already high... ©



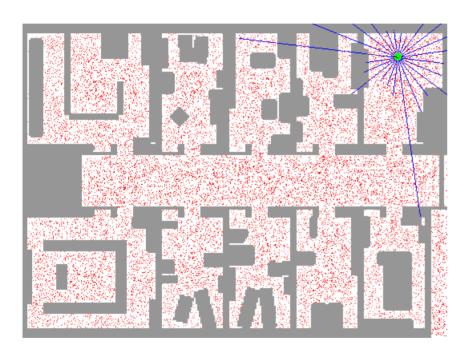
Speech recognition

Hidden Markov models and their generalizations

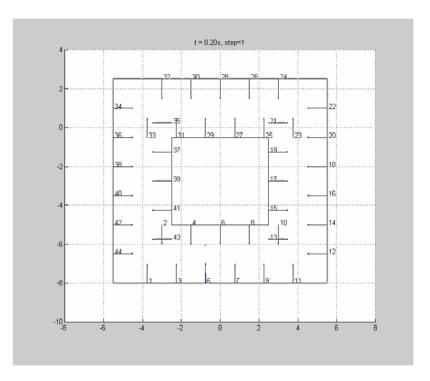


Tracking and robot localization

Kalman Filters



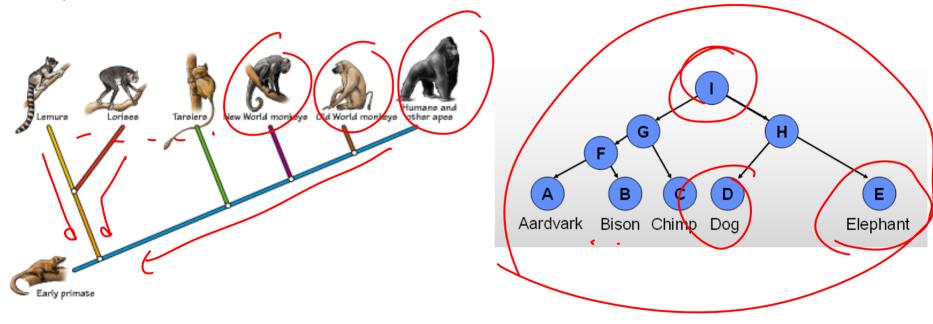
[Fox et al.]



[Funiak et al.]

Evolutionary biology

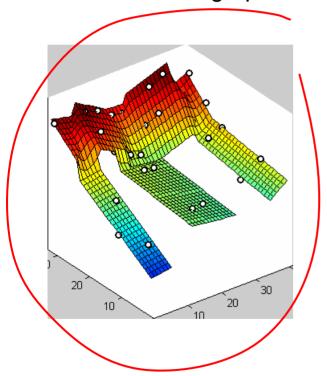
Bayesian networks



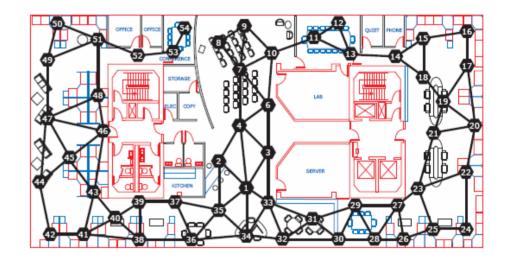
[Friedman et al.]

Modeling sensor data

Undirected graphical models



Markov Networks

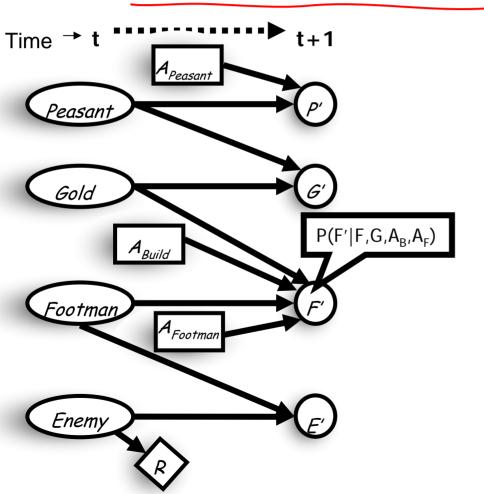


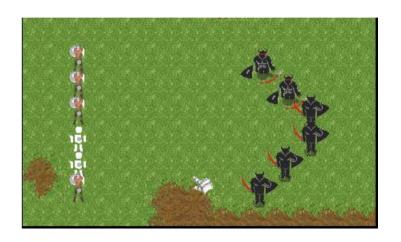


[Guestrin et al.]

Planning under uncertainty

Dynamic Bayesian networks Factored Markov decision problems

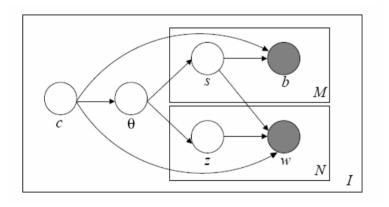


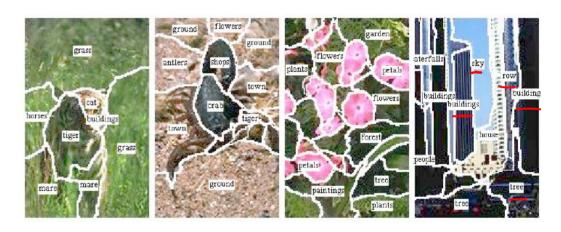


[Guestrin et al.]

Images and text data

Hierarchical Bayesian models

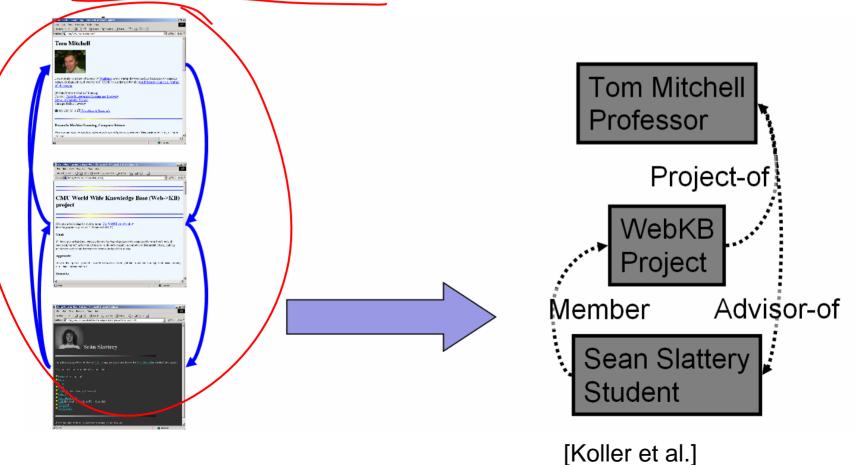




[Barnard et al.]

Structured data (text, webpages,...)

Probabilistic relational models



And many many many many many more...

Syllabus

- Covers a wide range of Probabilistic Graphical
 Models topics from basic to state-of-the-art
- You will learn about the methods you heard about:
 - Bayesian networks, Markov networks, factor graphs, decomposable models, junction trees, parameter learning, structure learning, semantics, exact inference, variable elimination, context-specific independence, approximate inference, sampling, importance sampling, MCMC, Gibbs, variational inference, loopy belief propagation, generalized belief propagation, Kikuchi, Bayesian learning, missing data, EM, Chow-Liu, structure search, IPF for tabular MRFs, Gaussian and hybrid models, discrete and continuous variables, temporal and template models, hidden Markov Models, Forwards-Backwards, Viterbi, Baum-Welch, Kalman filter, linearization, switching Kalman filter, assumed density filtering, DBNs, BK, Relational probabilistic models, Causality,...
- Covers algorithms, theory and applications
- It's going to be fun and hard work ©

Prerequisites

- 10-701 Machine Learning, especially:
 - Probabilities
 - Distributions, densities, marginalization...
 - Basic statistics
 - Moments, typical distributions, regression...
- Algorithms
 - Dynamic programming, basic data structures, complexity...
- Programming
 - Matlab will be very useful
- We provide some background, but the class will be fast paced
- Ability to deal with "abstract mathematical concepts"

Review Sessions

- Very useful!
 - Review material
 - □ Present background
 - □ Answer questions
- Thursdays, 5:00-6:30 in Wean Hall 4615A
- First recitation is tomorrow
 - Review of probabilities & statistics
- Sometimes this semester: Especial recitations on Mondays 5:30-7pm in Wean Hall 4615A
 - Cover special topics that we can't cover in class
 - □ These are optional, but you are here to learn... ☺
- Do we need a Matlab review session?

Staff

- Two Great TAs: Great resource for learning, interact with them!
 - □ Khalid El-Arini <kbe@cs.cmu.edu>

□ Ajit Paul Singh <ajit@cs.cmu.edu>



- Administrative Assistant
 - Monica Hopes, Wean 4619, x8-5527, meh@cs.cmu.edu



First Point of Contact for HWs

- To facilitate interaction, a TA will be assigned to each homework question – This will be your "first point of contact" for this question
 - □ But, you can always ask any of us
 - □ (Due to logistic reasons, we will only start this policy for HW2)
- For e-mailing instructors, always use:
 - □ 10708-instructors@cs.cmu.edu
- For announcements, subscribe to:
 - □ 10708-announce@cs
 - https://mailman.srv.cs.cmu.edu/mailman/listinfo/10708-announce

Text Books

from Monica Hopes.

Primary: Daphne Koller and Nir Friedman, Bayesian
 Networks and Beyond, in preparation. These chapters

 Secondary: M. I. Jordan, An Introduction to Probabilistic Graphical Models, in preparation. Copies of selected chapters will be made available.

are part of the course reader. You can purchase one

Grading

- 5 homeworks (50%)
 - □ First one goes today!
 - □ Homeworks are long and hard ☺
 - please, please, please, please, please, please start early!!!
- Final project (30%)
 - □ Done individually or in pairs
 - □ Details out October 4th
- Final (20%)
 - □ Take home, out Dec. 1st, due Dec. 15th

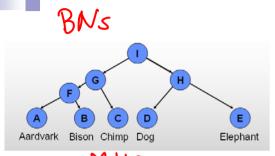
Homeworks

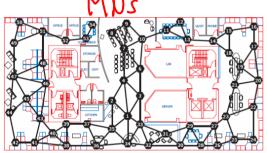
- Homeworks are hard, start early ©
- Due in the beginning of class
- 3 late days for the semester
- After late days are used up:
 - Half credit within 48 hours
 - Zero credit after 48 hours
- All homeworks must be handed in, even for zero credit
- Late homeworks handed in to Monica Hopes, WEH 4619
- Collaboration
 - You may discuss the questions
 - Each student writes their own answers
 - □ Write on your homework anyone with whom you collaborate
- IMPORTANT:
 - We may use some material from previous years or from papers for the homeworks. Unless otherwise specified, please only look at the readings when doing your homework → You are taking this advanced graduate class because you want to learn, so this rule is self-enforced ©

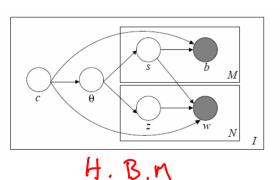
Enjoy!

- Probabilistic graphical models are having significant impact in science, engineering and beyond
- This class should give you the basic foundation for applying GMs and developing new methods
- The fun begins...

What are the fundamental questions of graphical models?







Representation:

- □ What are the types of models?
- □ What does the model mean/imply/assume? (Semantics)

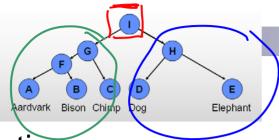
Inference:

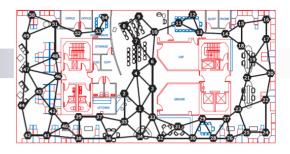
□ How do I answer questions/queries with my model?

Learning:

□ What model is the right for my data?

More details???





Representation:

- Graphical models represent exponentially large probability distributions compactly ~ Variables (binary) = P(X1, Xn) & represented by

Inference:

- What is the probability of X given some observations?
- What is the most likely explanation for what is happening?
- What decisions should I make?

Learning:

- What are the right/good parameters for the model?
- How do I obtain the structure of the model?

Where do we start?

- From Bayesian networks
- "Complete" BN presentation first
 - Representation
 - Exact inference
 - Learning
 - Only discrete variables for now
- Later in the semester
 - □ Undirected models
 - □ Approximate inference
 - Continuous
 - □ Temporal models
 - □ And more...
- Class focuses on fundamentals Understand the foundation and basic concepts

Today

- Probabilities
- Independence
- Two nodes make a BN
- Naïve Bayes
- Should be a review for everyone Setting up notation for the class

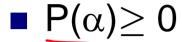
Event spaces

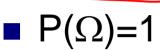
- Outcome space $\Omega = \{1, 2, 3, 7, 5\}$
- Measurable events S
 - $\alpha = \{1, 2\}$ \square Each $\alpha \in S$ is a subset of Ω

- Must contain
 - □ Empty event ∅
 - \square Trivial event Ω
- Closed under
 - □ Union: $\alpha \cup \beta \in S$
 - Complement: $\alpha \in S$, then Ω - α also in S



Probability distribution P over (Ω, S)









■ If $\alpha \cap \beta = \emptyset$, then $P(\alpha \cup \beta) = P(\alpha) + P(\beta)$

- From here, you can prove a lot, e.g.,
 - $\square P(\emptyset)=0$
 - $\Box P(\alpha \cup \beta) = P(\alpha) + P(\beta) P(\alpha \cap \beta)$



Interpretations of probability – A can of worms!

Frequentists

- flip HHTHT....lim # head = 0

 He limit

 P(H)
- \Box P(α) is the frequency of α in the limit
- Many arguments against this interpretation
 - What is the frequency of the event "it will rain tomorrow"?
- Subjective interpretation
 - \square P(α) is my degree of belief that α will happen
 - What the does "degree of belief mean?"
 - □ If I say $P(\alpha)=0.8$, then I am willing to bet!!!

For this class, we (mostly) don't care what camp you are in

Conditional probabilities

■ After learning that α is true, how do we feel about β ?

 \propto



Two of the most important rules of the semester: 1. The chain rule

$$P(\alpha \cap \beta) \neq P(\alpha)P(\beta | \alpha) = P(\beta) \cdot P(\alpha | \beta)$$

$$\begin{cases} 1,23 & 52,32 \end{cases}$$

More generally:

$$P(\alpha_{1}\cap...\cap\alpha_{k}) = P(\alpha_{1}) P(\alpha_{2}|\alpha_{1}) \cdots P(\alpha_{k}|\alpha_{1}\cap...\cap\alpha_{k-1})$$

$$= P(\lambda_{5}) \cdot P(\lambda_{7}|\lambda_{3}) \cdot P(\lambda_{4}|\lambda_{3},\lambda_{7}) \cdots$$

$$= P(\lambda_{5}) \cdot P(\lambda_{7}|\lambda_{3}) \cdot P(\lambda_{4}|\lambda_{3},\lambda_{7}) \cdots$$

$$= P(\lambda_{5}) \cdot P(\lambda_{7}|\lambda_{3}) \cdot P(\lambda_{4}|\lambda_{3},\lambda_{7}) \cdots$$

Two of the most important rules of the semester: 2. Bayes rule

$$P(\alpha \mid \beta) = \frac{P(\beta \mid \alpha)P(\alpha)}{P(\beta)}$$

More generally, external event γ:

$$P(\alpha \mid \beta \cap \gamma) = \frac{P(\beta \mid \alpha \cap \gamma)P(\alpha \mid \gamma)}{P(\beta \mid \gamma)}$$

Most important concept: a) Independence

and β independent, if $P(\beta|\alpha)=P(\beta)$ $P \models (\alpha \perp \beta)$ Ast. entails α indep. of β

■ **Proposition:** α and β *independent* if and only if $P(\alpha \cap \beta) = P(\alpha)P(\beta)$

Most important concept: b) Conditional independence

- Independence is rarely true, but conditionally...
- α and β conditionally independent given γ if $P(\beta|\alpha\cap\gamma)=P(\beta|\gamma)$

$$\Box P \models (\alpha \perp \beta \mid \gamma)$$
a indep β Jiven δ

Proposition: $P \models (\alpha \perp \beta \mid \gamma)$ if and only if

$$P(\alpha \cap \beta \mid \gamma) = P(\alpha \mid \gamma) P(\beta \mid \gamma)$$

Random variable

- Events are complicated we think about attributes
 - □ Age, Grade, HairColor
- Random variables formalize attributes:
 - \Box Grade \Rightarrow A+shorthand for event $\{\omega \in \Omega: f_{Grade}(\omega) = A+\}$
- Properties of random vars, X:

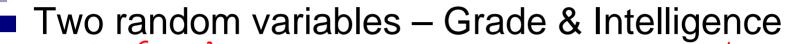
Marginal distribution

Probability P(X) of possible outcomes X

$$P(Grade > H) = \sum_{a} \sum_{h} P(G=A, Age=a, H=h)$$

$$P(G, A, H)$$

Joint distribution, Marginalization



$$G = \{A, B\}$$
 $I = \{H, VH\}$
 $VH = \{A, B\}$
 $VH = \{A, B\}$
 $VH = \{A, B\}$
 $VH = \{A, B\}$
 $VH = \{A, B\}$

Marginalization – Compute marginal over single var

$$\begin{array}{cccc}
(P(G=A) = 0.8 + 0.04 = P(G=A, I=H) + P(G=A, I=H) \\
= 0.84 \\
P(G): P(G=B) = 0.16
\end{array}$$

Marginalization – The general case

Compute marginal distribution P(X_i):

$$P(X_{1}, \dots, X_{n}) = \sum_{\substack{x_{i+1}, \dots, x_{n} \\ 2^{n-i} \text{ for each } X_{i} \dots X_{i}}} P(X_{1}, X_{2}, \dots, X_{i}, x_{i+1}, \dots, x_{n})$$

$$P(X_i) = \sum_{x_1, \dots, x_{i-1}} P(x_1, \dots, x_{i-1}, X_i)$$

$$\sum_{x_1, \dots, x_{i-1}} P(x_1, \dots, x_{i-1}, X_i)$$
for each x_i

Basic concepts for random variables

- Atomic outcome: assignment $x_1,...,x_n$ to $X_1,...,X_n$
- Conditional probability: P(X,Y)=P(X)P(Y|X)P(Y|X) : P(Y=g|X=x)
- Bayes rule: $P(X|Y) = P(Y|X) \cdot P(X)$ P(Y)
- Chain rule:

Conditionally independent random variables

- Sets of variables X, Y, Z ∈ \SATS
- X is independent of Y given Z if P(X=x|Y=y|Z=z)=P(x=z) $P \models (X=x\perp Y=y|Z=z), \forall x \in Val(X), y \in Val(Y), z \in Val(Z)$
- Shorthand:
 - □ Conditional independence: P ⊨ (X ⊥ Y | Z)
 - □ For $P \models (\mathbf{X} \perp \mathbf{Y} \mid \emptyset)$, write $P \models (\mathbf{X} \perp \mathbf{Y}) \times \text{indep.} \quad \forall \text{ given } \geq 0$
- Proposition: P statisfies (X ⊥ Y | Z) if and only if
 - $\square P(X,Y|Z) = P(X|Z) P(Y|Z)$

Properties of independence

Symmetry:

$$\square (X \perp Y \mid Z) \Rightarrow (Y \perp X \mid Z)$$

Decomposition:

$$\square$$
 (X \perp Y,W | Z) \Rightarrow (X \perp Y | Z)

Weak union:

$$\square$$
 (X \perp Y,W | Z) \Rightarrow (X \perp Y | Z,W)

Contraction:

$$\square$$
 (X \perp W | Y,Z) & (X \perp Y | Z) \Rightarrow (X \perp Y,W | Z)

Intersection:

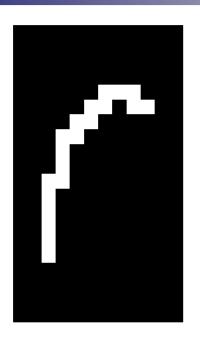
$$\square$$
 (X \perp Y | W,Z) & (X \perp W | Y,Z) \Rightarrow (X \perp Y,W | Z)

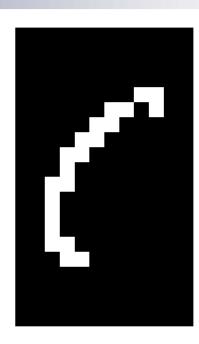
- □ Only for positive distributions!
- \square P(α)>0, $\forall \alpha, \alpha \neq \emptyset$
- **Notation** (1/P) independence properties entailed by P

Bayesian networks

- One of the most exciting recent advancements in statistical AI
- Compact representation for exponentially-large probability distributions
- Fast marginalization too
- Exploit conditional independencies

Handwriting recognition





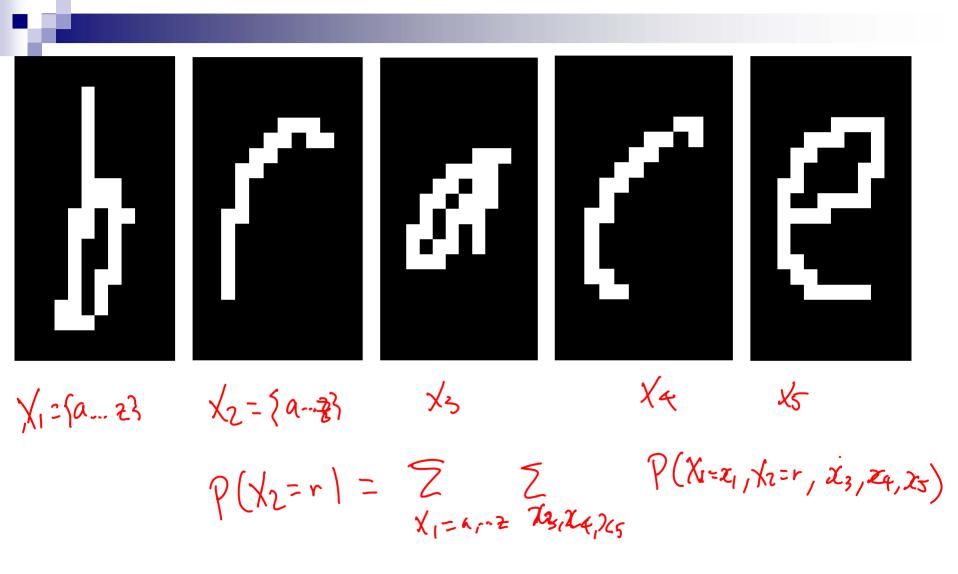
Webpage classification



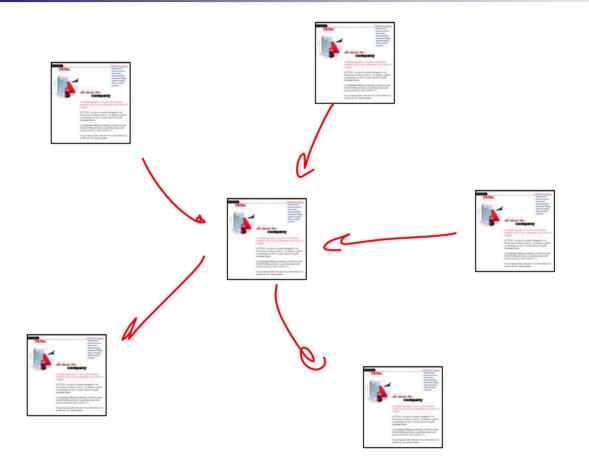
Company home page
vs
Personal home page
vs
Univeristy home page
vs

. . .

Handwriting recognition 2



Webpage classification 2



Let's start on BNs...

- Consider P(X_i)

 K-1 Parms.
 - \square Assign probability to each $x_i \in Val(X_i)$
 - Independent parameters
- Consider $P(X_1,...,X_n)$
 - □ How many independent parameters if |Val(X_i)|=k?

What if variables are independent?



- $\square (X_i \perp X_j), \forall i,j \qquad (\{x_1 x_3\} \perp \{x_7, x_7\})$
- □ Not enough!!! (See homework 1 ☺)
- \square Must assume that $(\mathbf{X} \perp \mathbf{Y}), \forall \mathbf{X}, \mathbf{Y} \text{ subsets of } \{X_1, \dots, X_n\}$

Can write

 $\square P(X_1, \dots, X_n) = \prod_{i=1\dots n} P(X_i)$

How many independent parameters now?

Conditional parameterization – two nodes



Conditional parameterization – three nodes

- Grade and SAT score are determined by Intelligence
- (G ⊥ S | I)

The naïve Bayes model – Your first real Bayes Net

- Class variable: C
- Evidence variables: X₁,...,X_n
- assume that $(\mathbf{X} \perp \mathbf{Y} \mid \mathbf{C})$, $\forall \mathbf{X}, \mathbf{Y}$ subsets of $\{X_1, ..., X_n\}$

What you need to know

- Basic definitions of probabilities
- Independence
- Conditional independence
- The chain rule
- Bayes rule
- Naïve Bayes

Next class

- We've heard of Bayes nets, we've played with Bayes nets, we've even used them in your research
- Next class, we'll learn the semantics of BNs, relate them to independence assumptions encoded by the graph