


Probabilistic Graphical Models

Lecture 1 – Introduction

CS/CNS/EE 155

Andreas Krause



One of the **most exciting advances** in machine learning (AI, signal processing, coding, control, ...) in the last decades



How can we gain
global insight based on
local observations?

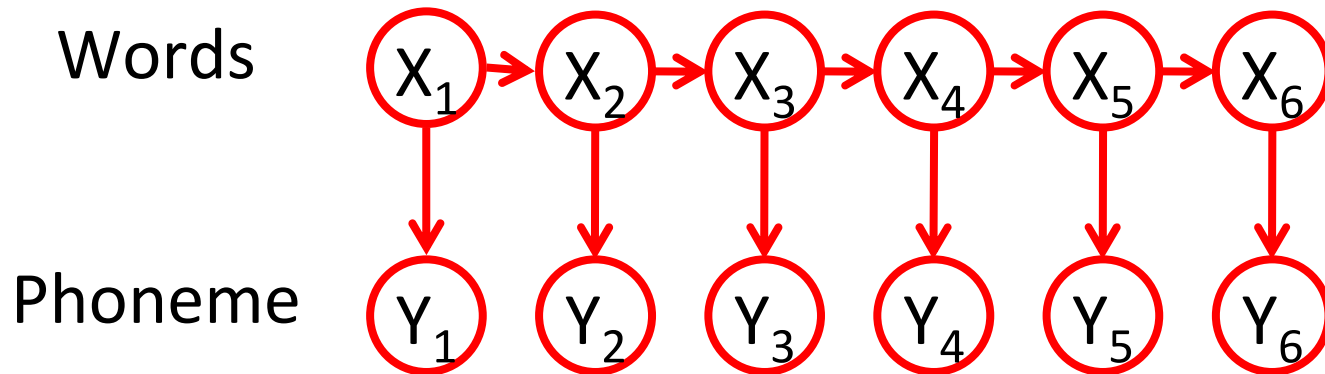
Key idea:

- **Represent** the world as a collection of random variables X_1, \dots, X_n with joint distribution $P(X_1, \dots, X_n)$
- **Learn** the distribution from data
- Perform “**inference**” (compute conditional distributions $P(X_i \mid X_1 = x_1, \dots, X_m = x_m)$)

Applications

Natural Language Processing

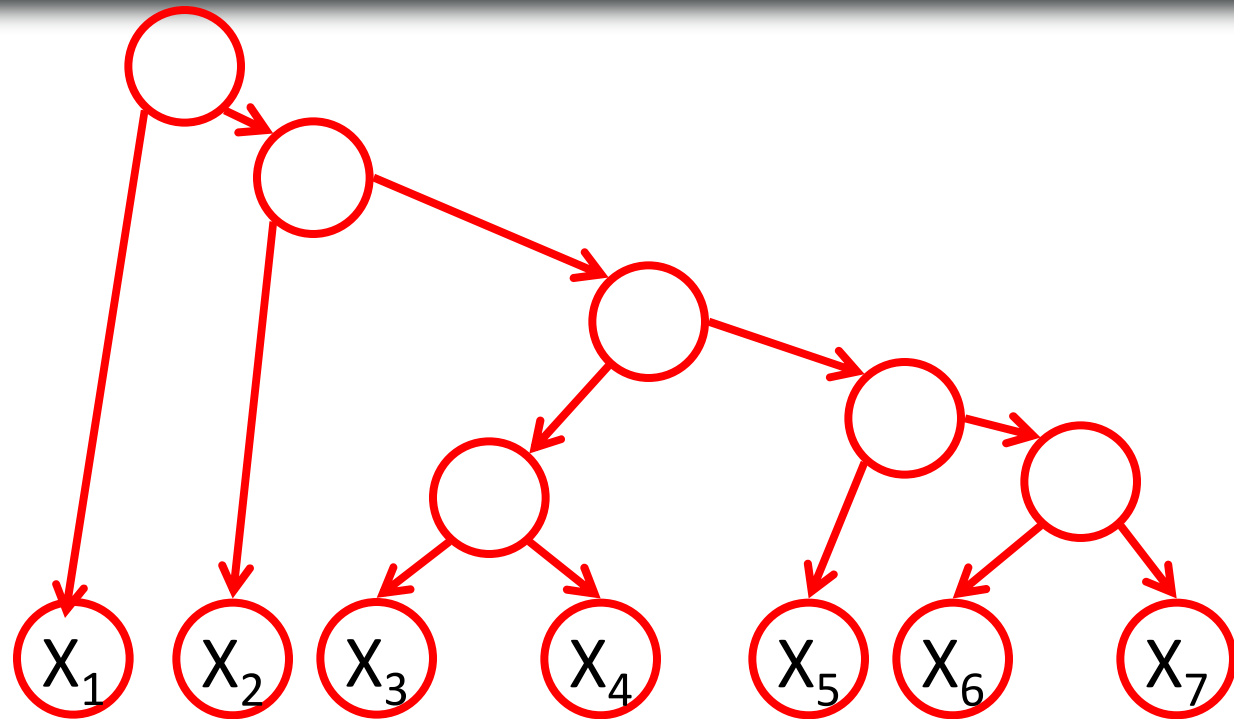
Speech recognition



“He ate the cookies on the couch”

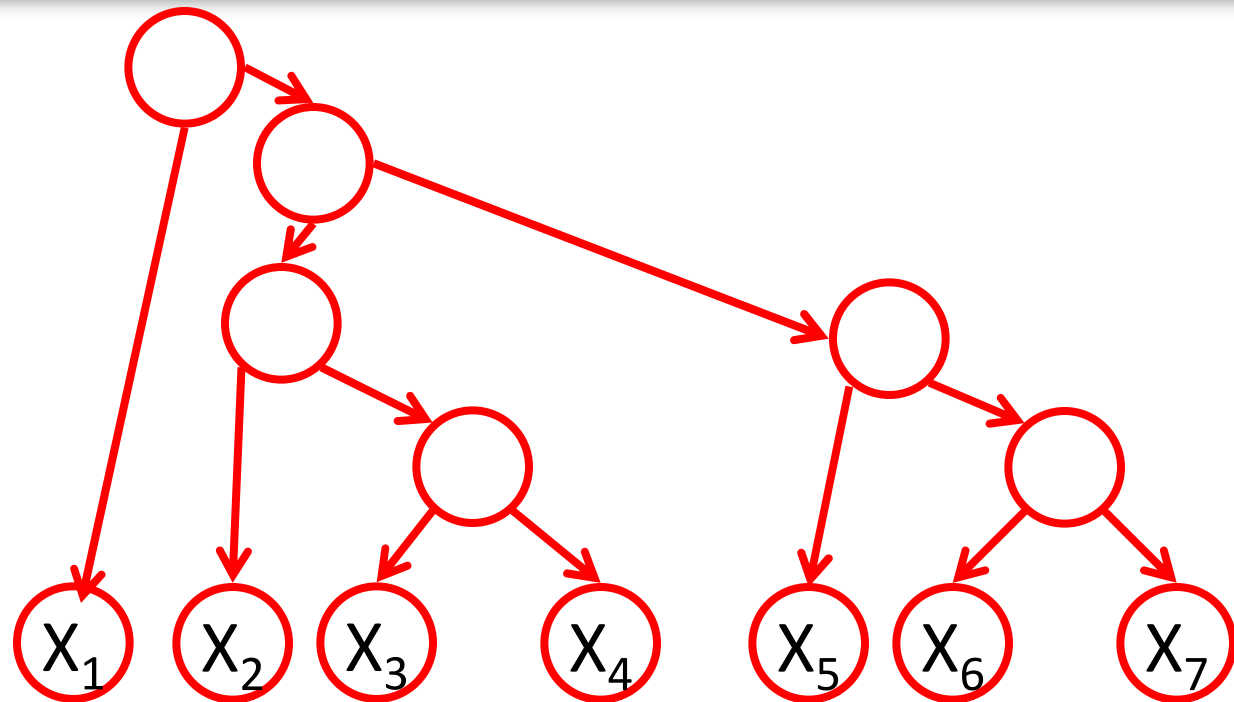
- Infer spoken words from audio signals
- “Hidden Markov Models”

Natural language processing



“He ate the cookies on the couch”

Natural language processing

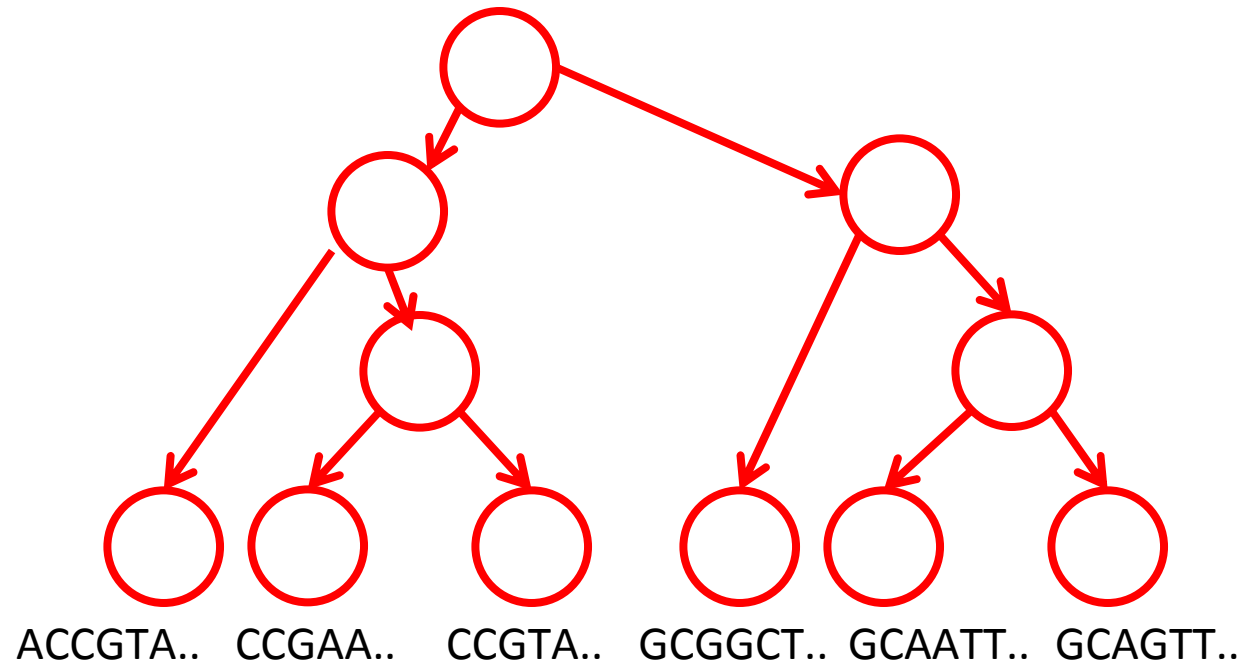


“He ate the cookies on the couch”

- Need to deal with ambiguity!
- Infer grammatical function from sentence structure
- “Probabilistic Grammars”

Evolutionary biology

[Friedman et al.]



- Reconstruct phylogenetic tree from current species (and their DNA samples)

Applications

Computer Vision

Image denoising

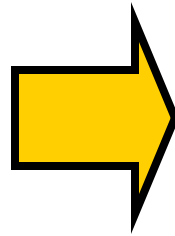
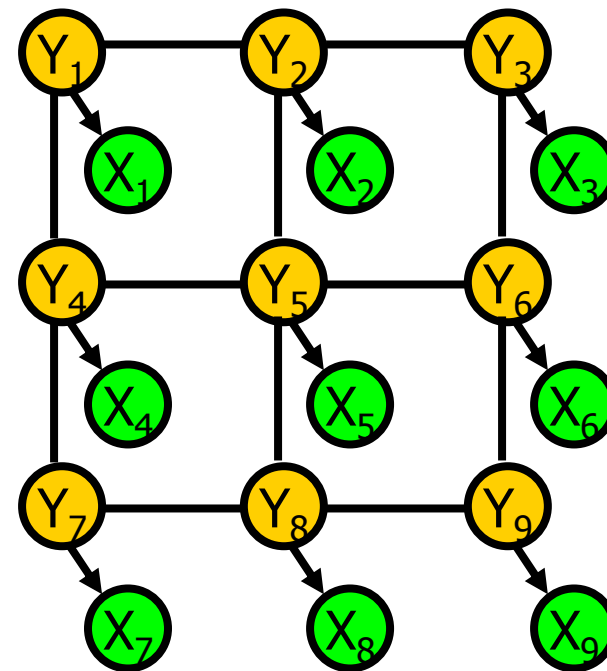


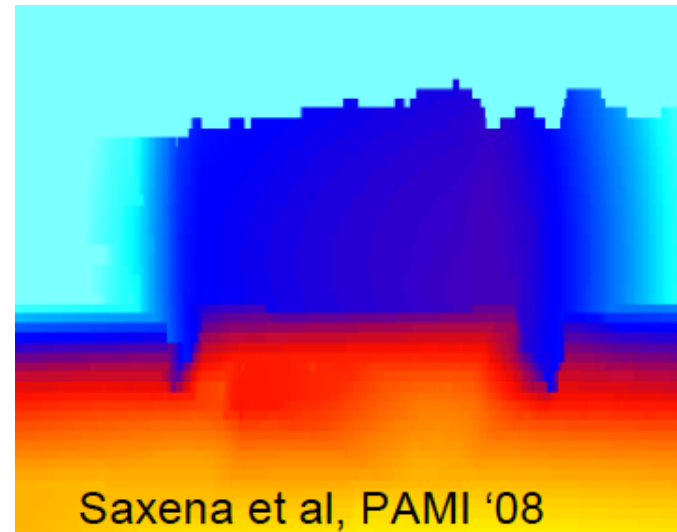
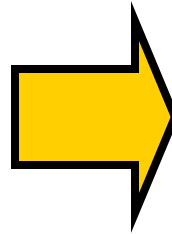
Image denoising

Markov Random Field



X_i : noisy pixels
 Y_i : "true" pixels

Make3D



- Infer depth from 2D images
- “Conditional random fields”

Applications

State estimation

Robot localization & mapping



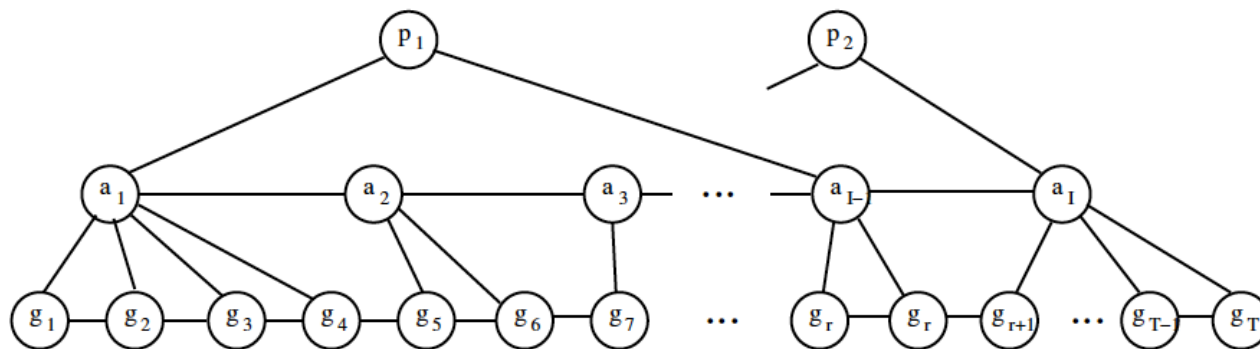
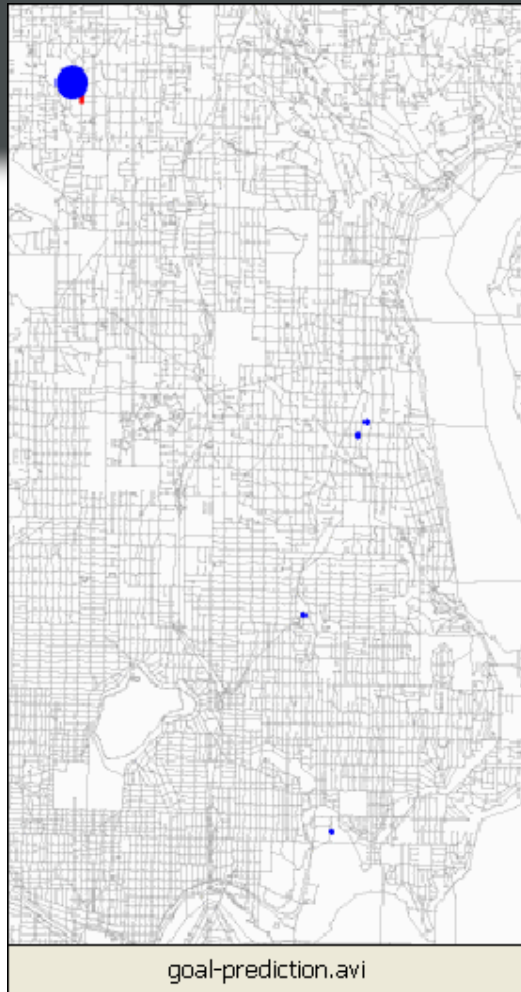
D. Haehnel,
W. Burgard,
D. Fox, and
S. Thrun.
IROS-03.

- Infer both location and map from noisy sensor data
- Particle filters

Activity recognition

L. Liao, D. Fox, and H. Kautz. *AAAI-04*

Predict “goals” from raw GPS data
“Hierarchical Dynamical
Bayesian networks”

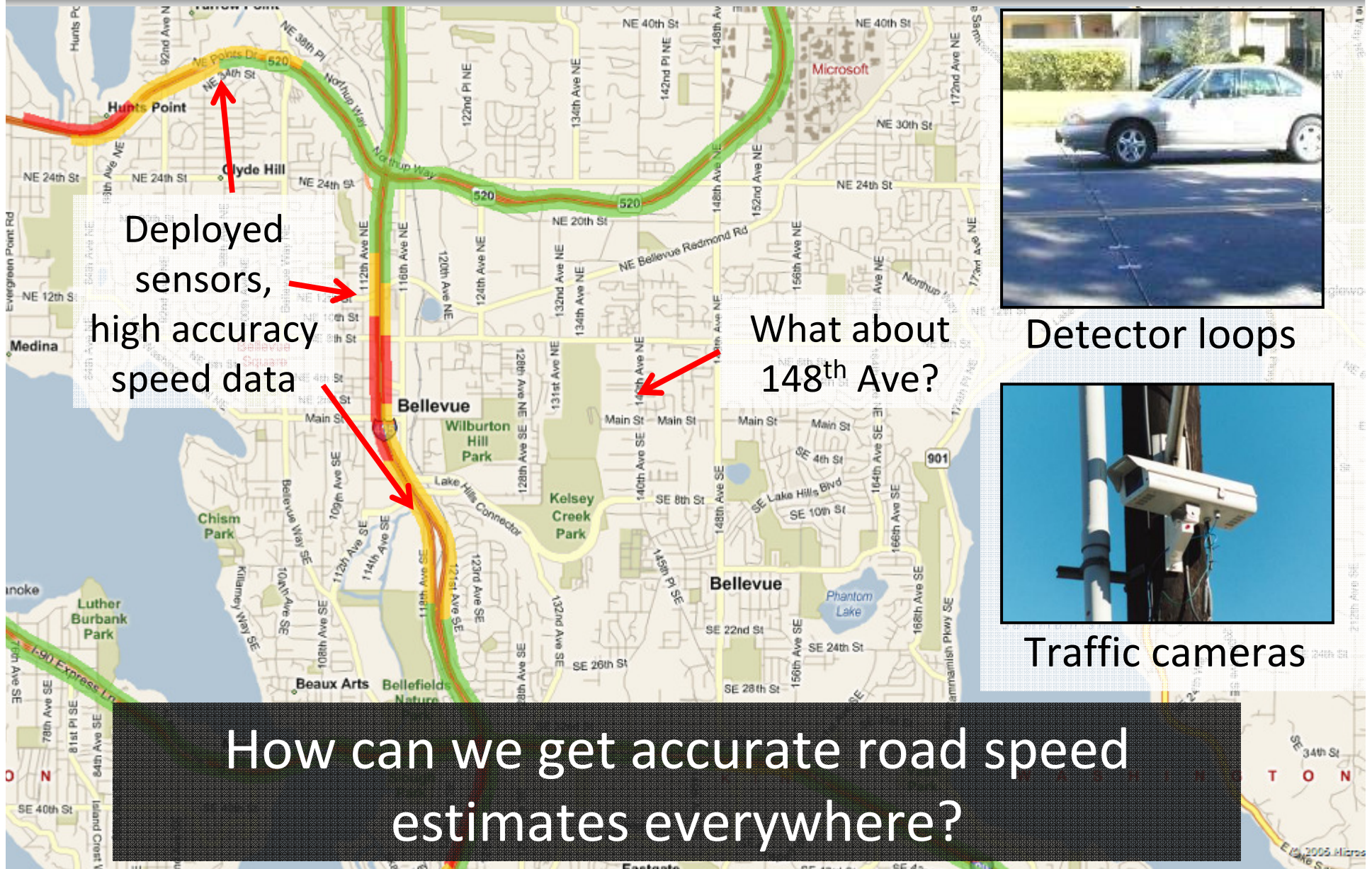


Significant places
home, work, bus stop, parking lot, friend

Activity sequence
walk, drive, visit, sleep, pickup, get on bus

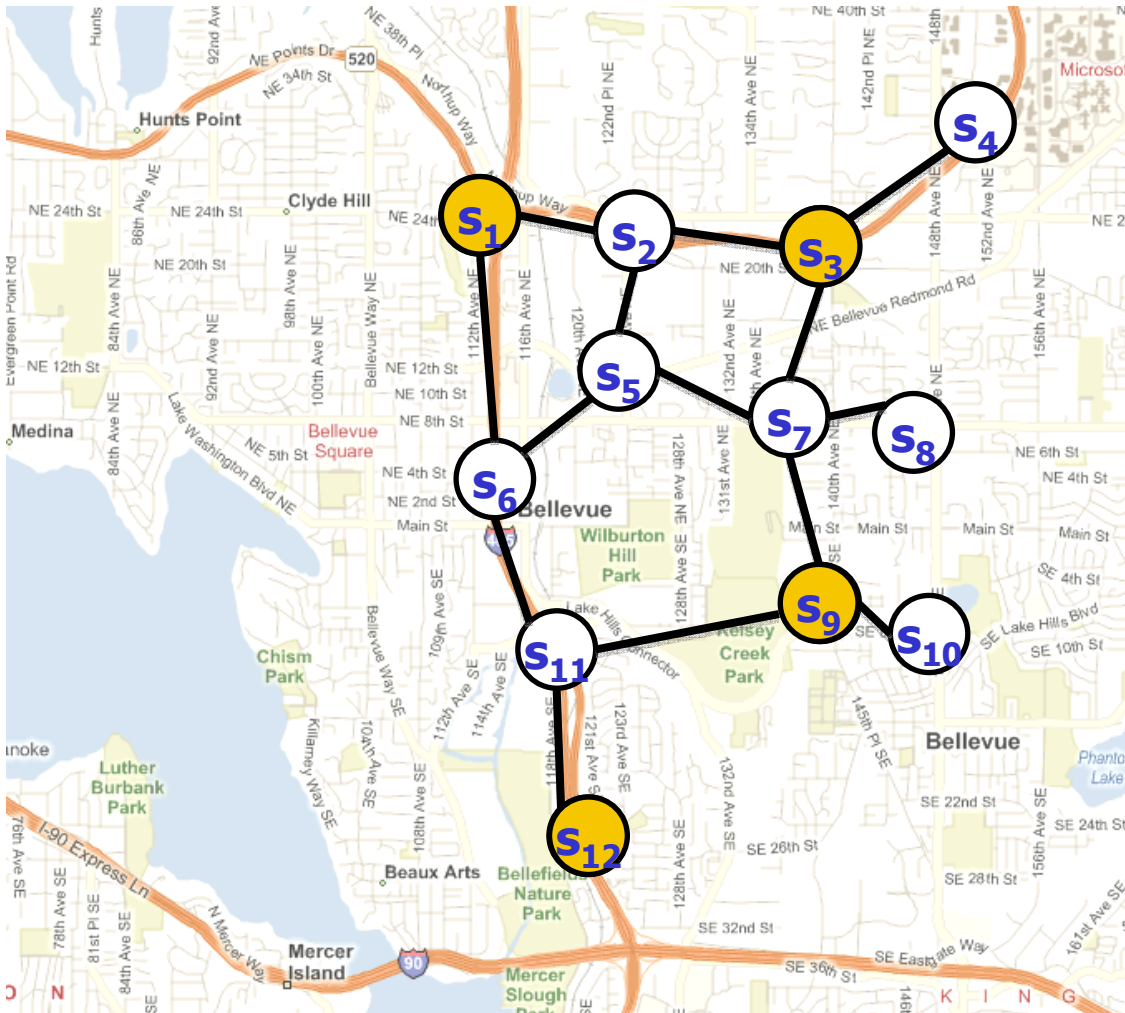
GPS trace
association to street map

Traffic monitoring



Cars as a sensor network

[Krause, Horvitz et al.]




- (Normalized) speeds as random variables
- Joint distribution allows modeling correlations
- Can **predict unmonitored** speeds from monitored speeds using $P(S_5 | S_1, S_9)$


Applications

Structure Prediction

Collaborative Filtering and Link Prediction

People you may know

 L. Brouwer [invite](#) | [x](#)

 T. Riley [invite](#) | [x](#)

[See more »](#)

NETFLIX

[Browse DVDs](#) [Watch Instantly](#) [Your Queue](#) [Movies You'll ♥](#)

[Suggestions \(731\)](#) [Rate Movies](#) [Taste Preferences](#)

Suggestions in [All Genres](#) ▼


Your Recent History [\(What's this?\)](#)

Recently Viewed Items



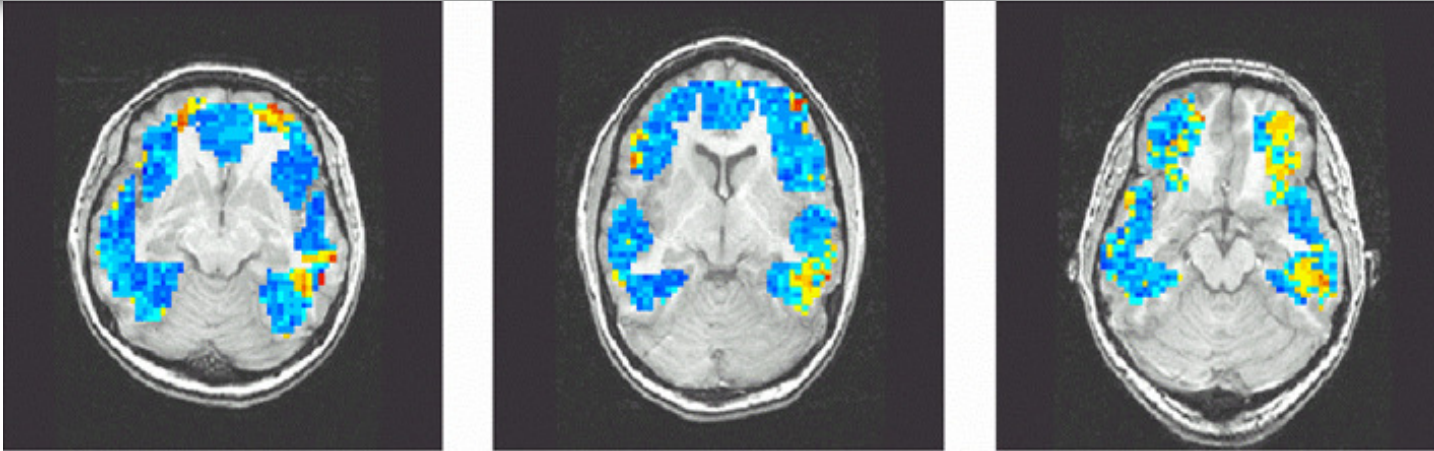
[Probabilistic Graphical Models: Principles and...](#)
by Daphne Koller

Continue shopping: Customers Who Bought Items in Your Recent History Also Bought



- Predict “missing links”, ratings...
- “Collective matrix factorization”, Relational models

Analyzing fMRI data



Mitchell et al.,
Science, 2008



- Predict activation patterns for nouns
- Predict connectivity (Pittsburgh Brain Competition)

Other applications

- Coding (LDPC codes, ...)
- Medical diagnosis
- Identifying gene regulatory networks
- Distributed control
- Computer music
- Probabilistic logic
- Graphical games
-

MANY MORE!!

Key challenges:

How do we

... **represent** such probabilistic models?

(distributions over vectors, maps, shapes,
trees, graphs, functions...)

... perform **inference** in such models?

... **learn** such models from data?

Syllabus overview

- We will study Representation, Inference & Learning
- First in the simplest case
 - Only discrete variables
 - Fully observed models
 - Exact inference & learning
- Then generalize
 - Continuous distributions
 - Partially observed models (hidden variables)
 - Approximate inference & learning
- Learn about algorithms, theory & applications

Overview

- Course webpage
 - <http://www.cs.caltech.edu/courses/cs155/>
- Teaching assistant: Pete Trautman
(trautman@cds.caltech.edu)
- Administrative assistant: Sheri Garcia
(sheri@cs.caltech.edu)

Background & Prerequisites

- Basic probability and statistics
- Algorithms
- CS 156a or permission by instructor
- Please fill out the questionnaire about background (not graded 😊)
- Programming assignments in MATLAB.
- Do we need a MATLAB review recitation? *No.*

Coursework

- Grading based on
 - 4 homework assignments (one per topic) (40%)
 - Course project (40%)
 - Final take home exam (20%)
- 3 ~~late days~~
- Discussing assignments allowed, but everybody must turn in their own solutions
- Start early! 😊 !!!

Course project

- “Get your hands dirty” with the course material
- Implement an algorithm from the course or a paper you read and apply it to some data set
- Ideas on the course website (soon)
- Application of techniques you learnt to your own research is encouraged
- Must be something new (e.g., not work done last term)

Project: Timeline and grading

- Small groups (2-3 students)
- October 19: Project proposals due (1-2 pages); feedback by instructor and TA
- November 9: Project milestone
- December 4: Project report due; poster session
- Grading based on quality of poster (20%), milestone report (20%) and final report (60%)



Review: Probability

- This should be familiar to you...
- Probability Space (Ω, \mathcal{F}, P)
 - Ω : set of “atomic events”
 - $\mathcal{F} \subseteq 2^\Omega$: set of all (non-atomic) events
 \mathcal{F} is a σ -Algebra
(closed under complements and countable unions)
 - $P: \mathcal{F} \rightarrow [0,1]$ probability measure
For $\omega \in \Omega$, $P(\{\omega\})$ is the probability that event ω happens

$$\omega \in \Omega$$
$$\Omega \setminus \{\omega\} \in \mathcal{F}$$

Interpretation of probabilities

- Philosophical debate..
- Frequentist interpretation
 - $P(\alpha)$ is relative frequency of α in repeated experiments
 - Often difficult to assess with limited data
- Bayesian interpretation
 - $P(\alpha)$ is “degree of belief” that α will occur
 - Where does this belief come from?
 - Many different flavors (subjective, pragmatic, ...)
- Most techniques in this class can be interpreted either way.

Independence of events

- Two events $\alpha, \beta \in F$ are independent if

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

- A collection S of events is independent, if for any subset $\alpha_1, \dots, \alpha_n \in S$ it holds that

$$P(\alpha_1, \dots, \alpha_n) = P(\alpha_1) P(\alpha_2) \cdot \dots \cdot P(\alpha_n)$$

~~Q: $\forall i, j: P(\alpha_i \cap \alpha_j) = P(\alpha_i) P(\alpha_j)$
 $\Rightarrow \alpha_1, \dots, \alpha_n$ independent?~~

No.

Conditional probability

- Let α, β be events, $P(\beta) > 0$
- Then:

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

Most important rule #1:

- Let $\alpha_1, \dots, \alpha_n$ be events, $P(\alpha_i) > 0$

- Then

$$P(\alpha_1 \cap \dots \cap \alpha_n) = P(\alpha_1) \cdot P(\alpha_2 | \alpha_1) \cdot \dots \cdot P(\alpha_n | \alpha_1 \dots \alpha_{n-1})$$

Chain rule

Most important rule #2:

- Let α, β be events with prob. $P(\alpha) > 0, P(\beta) > 0$
- Then

$$P(\alpha | \beta) = \frac{P(\alpha \cap \beta)}{P(\beta)} = \frac{P(\beta | \alpha) \cdot P(\alpha)}{P(\beta)}$$

$$P(\beta) = P(\beta | \alpha) \cdot P(\alpha) + P(\beta | \neg \alpha) \cdot P(\neg \alpha)$$

Bayes' rule



Random variables

- Events are cumbersome to work with.
- Let D be some set (e.g., the integers)
- A **random variable** X is a mapping $X: \Omega \rightarrow D$
- For some $x \in D$, we say
 $P(X = x) = P(\{\omega \in \Omega: X(\omega) = x\})$

“probability that variable X assumes state x ”

- Notation: $\text{Val}(X)$ = set D of all values assumed by X .

Examples

- Bernoulli distribution: “(biased) coin flips”

$$D = \{H, T\}$$

Specify $P(X = H) = p$. Then $P(X = T) = 1-p$.

Write: $X \sim \text{Ber}(p)$;

- Multinomial distribution: “(biased) m-sided dice”

$$D = \{1, \dots, m\}$$

Specify $P(X = i) = p_i$, s.t. $\sum_i p_i = 1$

Write: $X \sim \text{Mult}(p_1, \dots, p_m)$

Multivariate distributions

- Instead of random variable, have random vector

$$\mathbf{X}(\omega) = [X_1(\omega), \dots, X_n(\omega)]$$

- Specify $P(X_1=x_1, \dots, X_n=x_n)$

- Suppose all X_i are Bernoulli variables.

- How many parameters do we need to specify?

$$\begin{array}{cccc|c}
 x_1 & x_2 & \dots & x_n & P(x_1, \dots, x_n) \\
 0 & 0 & & 0 & ? \\
 0 & & & 0 & ? \\
 0 & & & 1 & \\
 & & 0 & 1 & \\
 & & 0 & 1 & \\
 & & 1 & 0 &
 \end{array}
 \quad 2^n - 1$$

Rules for random variables

- Chain rule

$$P(x_1 \dots x_n) = P(x_1) P(x_2 | x_1) \dots P(x_n | x_1 \dots x_{n-1})$$

- Bayes' rule

$$P(x|y) = \frac{P(y|x) P(x)}{P(y)}$$

↑
How do we get $P(y)$?

Marginal distributions

- Suppose, X and Y are RVs with distribution $P(X,Y)$

X : Intelligence
 Y : Grade

$X \backslash Y$	VH	H
A	0.7	0.15
B	0.1	0.05

$$P(\text{Grade} = A) = .85$$

Marginal distributions

- Suppose we have joint distribution $P(X_1, \dots, X_n)$
- Then

$$P(X_i = x_i) = \sum_{x_1} \sum_{x_2} \dots \sum_{x_{i-1}} \sum_{x_{i+1}} \dots \sum_{x_n} P(x_1 \dots x_n)$$

- If all X_i binary: How many terms? 2^{n-1}

Independent RVs

- What if RVs are independent?

RVs X_1, \dots, X_n are independent, if for any assignment

$$P(X_1=x_1, \dots, X_n=x_n) = P(x_1) P(x_2) \dots P(x_n)$$

$$\Leftrightarrow \{\omega: X_i(\omega) = x_i\} \quad \forall i, x_i \in \text{Val}(X_i) \quad \text{indep.}$$

- How many parameters are needed in this case?

$$n \ll 2^n$$

$$X_i, X_j \text{ indep} \Rightarrow P(X_i | X_j) = P(X_i)$$

- Independence too strong assumption... Is there something weaker?

Key concept: Conditional independence

- Events α, β conditionally independent given γ if

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

- Random variables X and Y cond. indep. given Z if for all $x \in \text{Val}(X)$, $y \in \text{Val}(Y)$, $z \in \text{Val}(Z)$

$$P(X = x, Y = y \mid Z = z) = P(X = x \mid Z = z) P(Y = y \mid Z = z)$$

- If $P(Y=y \mid Z=z) > 0$, that's equivalent to

$$P(X = x \mid Z = z, Y = y) = P(X = x \mid Z = z)$$

Similarly for sets of random variables $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$

We write: $P \models \mathbf{X} \perp \mathbf{Y} \mid \mathbf{Z}$

Why is conditional independence useful?

- $P(X_1, \dots, X_n) = P(X_1) P(X_2 | X_1) \dots P(X_n | X_1, \dots, X_{n-1})$

How many parameters?

$$2^0 + 2^1 + 2^2 + \dots + 2^{n-1} = 2^n - 1$$

- Now suppose $X_1 \dots X_{i-1} \perp X_{i+1} \dots X_n | X_i$ for all i

Then

$$P(X_1, \dots, X_n) = \underbrace{P(X_1)}_1 \cdot \underbrace{P(X_2 | X_1)}_2 \cdot \underbrace{P(X_3 | X_2)}_2 \cdot \dots \cdot \underbrace{P(X_n | X_{n-1})}_2$$

How many parameters?

$$2^{n-1} \ll 2^n$$

Exponential reduction in # params

- Can we compute $P(X_n)$ more efficiently? Yes (often)

Properties of Conditional Independence

- **Symmetry**

- $X \perp Y \mid Z \Rightarrow Y \perp X \mid Z$

- **Decomposition**

- $X \perp Y, W \mid Z \Rightarrow X \perp Y \mid Z$

- **Contraction**

"Inverse Decomposition"

- $(X \perp Y \mid Z) \wedge (X \perp W \mid Y, Z) \Rightarrow X \perp Y, W \mid Z$

- **Weak union**

- $X \perp Y, W \mid Z \Rightarrow X \perp Y \mid Z, W$

- **Intersection**

- $(X \perp Y \mid Z, W) \wedge (X \perp W \mid Y, Z) \Rightarrow X \perp Y, W \mid Z$

- Holds only if distribution is positive, i.e., $P > 0$

Key questions

- How do we specify distributions that satisfy particular independence properties?

→ **Representation**

- How can we exploit independence properties for efficient computation?

→ **Inference**

- How can we identify independence properties present in data?

→ **Learning**

Will now see examples: Bayesian Networks

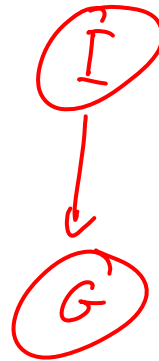
Bayesian networks

- A powerful class of probabilistic graphical models
- Compact parametrization of high-dimensional distributions
- In many cases, efficient exact inference possible
- Many applications
 - Natural language processing
 - State estimation
 - Link prediction
 - ...
- Demo..

Key idea

- Conditional parametrization
(instead of joint parametrization)
- For each RV, specify $P(X_i \mid \mathbf{X}_A)$ for set \mathbf{X}_A of RVs
- Then use chain rule to get joint parametrization
- Have to be careful to guarantee legal distribution...

Example: 2 variables

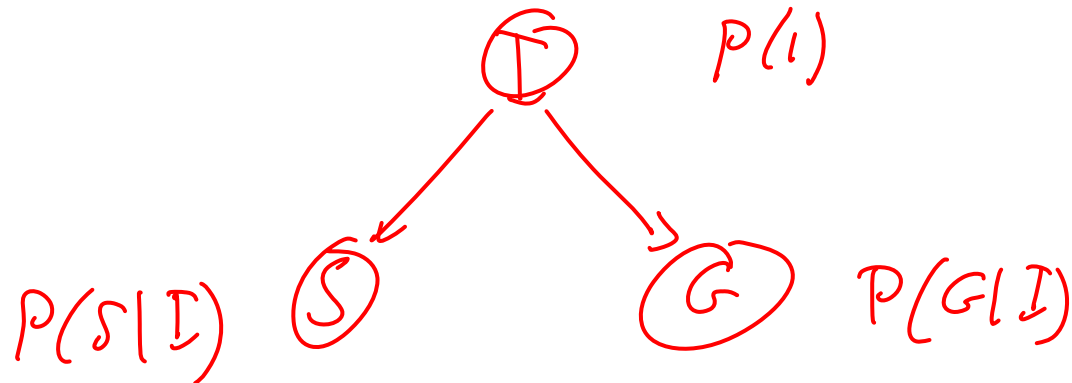


$$P(I = VH) = 0.8$$

$P(G|I)$

I \ G	A	B	
VH	0.8	0.2	$\Sigma_i = 1$
H	0.6	0.4	$\Sigma_i = 1$

Example: 3 variables



$$P(I, S, G) = P(I) P(G|I) P(S|I)$$

Example: Naïve Bayes models

- Class variable Y
- Evidence variables X_1, \dots, X_n
- Assume that $X_A \perp X_B \mid Y$
for all subsets X_A, X_B of $\{X_1, \dots, X_n\}$
- Conditional parametrization:
 - Specify $P(Y)$
 - Specify $P(X_i \mid Y)$
- Joint distribution

$$P(x_1, \dots, x_n, y) = P(y) \prod_i P(x_i \mid y)$$

What you need to know

- Basic probability
- Independence and conditional independence
- Chain rule & Bayes' rule
- Naïve Bayes models

Tasks

- By tomorrow (October 1, 4pm): hand in questionnaire about background to Sheri Garcia
- Read Chapter 2 in Koller & Friedman
- Start thinking about project teams and ideas (proposals due October 19)