## Graphical models and their Role in Databases VLDB 2007 Tutorial

## Amol Deshpande<sup>1</sup> Sunita Sarawagi<sup>2</sup>

<sup>1</sup>University of Maryland

<sup>2</sup>IIT Bombay

1

# Why a tutorial on graphical models at VLDB?

## • VLDB

- Many DB tasks use probabilistic modeling
  - ★ Core: Selectivity estimation, Imprecise databases,
  - ★ Appication: Information extraction, Duplicate elimination, sensor networks.
  - ★ Data mining: Classification (naive Bayes, logistic), clustering (EM)
- Probabilistic modeling is simultaneously
  - intuitive (low barrier to entry)
  - subtle (important to understand well for correctness & efficiency)
- Graphical models
  - Fundamental tools for intuitively and efficiently modeling probabilities
  - Distilled body of knowledge from many fields (let us build upon them, instead of reinveting)

VLDB wants to broaden, GM a fun and useful candidate for broadening

## Probabilistic modeling

- Given: several variables:  $x_1, \ldots x_n$ , *n* is large.
- Task: build a joint distribution function  $Pr(x_1, \ldots x_n)$
- Goal: Answer several kind of projection queries on the distribution
- Basic premise
  - Explicit joint distribution is dauntingly large
  - Queries are simple aggregates over the joint distribution.

## Example: Selectivity estimation in databases

• Variables are columns of a table

Age	Income	Experience	Degree	Location
10 ranges	7 scales	7 scales	3 scales	30 places

- An explicit joint distribution over all columns not tractable: number of combinations:  $10 \times 7 \times 7 \times 3 \times 30 = 44100$ .
- Queries: Estimate number of people with
  - Income > 200K and Degree="Bachelors",
  - Income < 200K, Degree="PhD" and experience > 10 years.
  - Many, many more.

## Alternatives to an explicit joint distribution

- Assume all columns are independent of each other: bad assumption
- Use data to detect pairs of highly correlated column pairs and estimate their pairwise frequencies

  - Ad hoc methods of combining these into a single estimate
- Go beyond pairwise correlations: understand finer dependencies
  - ▶ income ⊥⊥ age, but income ⊥⊥ age | experience
  - experience  $\perp \perp$  degree, but experience  $\perp \perp$  degree | income

Graphical models make explicit an efficient joint distribution from these independencies

# Graphical models

Model joint distribution over **several** variables as a product of smaller factors that is

- Intuitive to represent and visualize
  - Graph: represent structure of dependencies
  - Potentials over subsets: quantify the dependencies
- 2 *Efficient* to query
  - given values of any variable subset, reason about probability distribution of others.
  - many efficient exact and approximate inference algorithms

## Graphical models = graph theory + probability theory.

# Graphical models in use

- Roots in statistical physics for modeling interacting atoms in gas and solids [1900]
- Early usage in genetics for modeling properties of species [ 1920]
- Al: expert systems (1970s-80s)
- Now many new applications:
  - Error Correcting Codes: Turbo codes, impressive success story (1990s)
  - Robotics and Vision: image denoising, robot navigation.
  - Text mining: information extraction, duplicate elimination, hypertext classification, help systems
  - Bio-informatics: Secondary structure prediction, Gene discovery
  - Data mining: probabilistic classification and clustering.

# **Overall plan**

- Fundamentals
  - Representation
  - Exact inference
- Applications
  - Selectivity estimation
  - Probabilistic databases
- Applications: Sensor data management
- Fundamentals
  - Learning a graphical model
  - 2 Conditional Random Fields
- Applications
  - Information extraction
  - Duplicate elimination

# Part I: Fundamentals of Graphical Models

## Part I: Outline

#### Representation

- Directed graphical models: Bayesian networks
- Undirected graphical models

## 2 Inference Queries

- Exact inference on chains
- Exact inference on general graphs

## 3 Constructing a graphical model

- Graph Structure
- Parameters in Potentials

## Approximate inference

- Generalized belief propagation
- Sampling: Gibbs, Particle filters

## Representation

## Structure of a graphical model: Graph + Potential

## Graph

- Nodes: variables  $\mathbf{x} = x_1, \dots x_n$ 
  - Continuous: Sensor temperatures, income
  - Discrete: Degree (one of Bachelors, Masters, PhD), Levels of age
- Edges: direct interaction
  - Directed edges: Bayesian networks
  - Undirected edges: Markov Random fields



## Representation

Potentials:  $\psi_c(\mathbf{x}_c)$ 

- Scores for assignment of values to subsets *c* of directly interacting variables.
- Which subsets? What do the potentials mean?
  - Different for directed and undirected graphs

## Probability

Factorizes as product of potentials

$$\Pr(\mathbf{x} = x_1, \dots, x_n) \propto \prod \psi_S(\mathbf{x}_S)$$

## Directed graphical models: Bayesian networks

- Graph G: directed acyclic
  - Parents of a node:  $Pa(x_i) = set$  of nodes in G pointing to  $x_i$
- Potentials: defined at each node in terms of its parents.

$$\psi_i(x_i, \mathsf{Pa}(x_i)) = \mathsf{Pr}(x_i | \mathsf{Pa}(x_i))$$

Probability distribution

$$\Pr(x_1 \dots x_n) = \prod_{i=1}^n \Pr(x_i | pa(x_i))$$

## Example of a directed graph



$$\psi_1(L) = \Pr(L)$$
**NY CA London Other**
0.2 0.3 0.1 0.4

$$\psi_2(A) = \Pr(A)$$
  
20-30 30-45 > 45  
0.3 0.4 0.3  
or, a Guassian distribution  
 $(\mu, \sigma) = (35, 10)$ 

 $\psi_2(E,A) = \Pr(E|A)$ 10–15 0–10 > 15 20-30 0.9 0.1 0 30-45 0.4 0.5 0.1 **> 45** 0.8 0.1 0.1

 $\psi_2(I, E, D) = \Pr(I|D, A)$ 

3 dimensional table, or a histogram approximation.

Probability distribution  $Pa(\mathbf{x} = L, D, I, A, E) = Pr(L) Pr(D) Pr(A) Pr(E|A) Pr(I|D, E)$ 

Fundamentals of graphical model

Sunita Sarawagi

## Popular Bayesian networks

• Hidden Markov Models: speech recognition, information extraction



- State variables: discrete phoneme, entity tag
- Observation variables: continuous (speech waveform), discrete (Word)
- Kalman Filters: State variables: continuous
  - Discussed later
- PRMs: Probabilistic relational networks:
  - An important relevant class for relational data
  - Discussed later
- QMR (Quick Medical Reference) system

# Undirected graphical models

- Graph G: arbitrary undirected graph
- Useful when variables interact symmetrically, no natural parent-child relationship
- Example: labeling pixels of an image.
- Potentials defined on arbitrary subcliques *C* of *G*. Popular choices:
  - Node potentials
- Edge potentials
   Probability distribution

$$\Pr(\mathbf{x} = y_1 \dots y_n) = \frac{1}{Z} \prod_{C \in G} \psi_C(\mathbf{y}_C)$$

where  $Z = \sum_{\mathbf{y}'} \prod_{C \in G} \psi_C(\mathbf{y}'_C)$ 



## Example

 $y_7 - y_8 - y_9 y_i = 1$  (part of foreground), 0 otherwise.

Node potentials

• 
$$\psi_1(0) = 4$$
,  $\psi_1(1) = 1$ 

• 
$$\psi_2(0) = 2, \ \psi_2(1) = 3$$

▶ ....

• 
$$\psi_9(0) = 1, \ \psi_9(1) = 1$$

- Edge potentials: Same for all edges
  - $\psi(0,0) = 5$ ,  $\psi(1,1) = 5$ ,  $\psi(1,0) = 1$ ,  $\psi(0,1) = 1$
- Probability:  $\Pr(y_1 \dots y_9) \propto \prod_{k=1}^9 \psi_k(y_k) \prod_{(i,j) \in E(G)} \psi(y_i, y_j)$

## Popular undirected graphical models

- Interacting atoms in gas and solids [ 1900]
- Markov Random Fields in vision for image segmentation
- Conditional Random Fields for information extraction

# Comparing directed and undirected graphs

• Some distributions can only be expressed in one and not the other.



- Potentials
  - Directed: conditional probabilities, more intuitive
  - Undirected: arbitrary scores, easy to set.
- Dependence structure
  - Directed: Complicated d-separation test
  - ► Undirected: Graph separation: A ⊥⊥ B | C iff C separates A and B in G.
- Often application makes the choice clear.
  - Directed: Causality
  - Undirected: Symmetric interactions.

# Part I: Outline

#### Representation

- Directed graphical models: Bayesian networks
- Undirected graphical models

#### 2 Inference Queries

- Exact inference on chains
- Exact inference on general graphs

#### 3 Constructing a graphical model

- Graph Structure
- Parameters in Potentials

#### 4 Approximate inference

- Generalized belief propagation
- Sampling: Gibbs, Particle filters

## Inference queries

Marginal probability queries over a small subset of variables:

- Find Pr(Income='High & Degree='PhD')
- Find  $Pr(pixel y_9 = 1)$

$$\Pr(x_1) = \sum_{x_2...x_n} \Pr(x_1...x_n)$$

Most likely labels of remaining variables: (MAP queries)

- Find most likely entity labels of all words in a sentence
- Find likely temperature at sensors in a room

$$\mathbf{x}^* = \operatorname{argmax}_{x_1...x_n} \Pr(x_1...x_n)$$

## Exact inference on chains

• Given,

$$y_1 \longrightarrow y_2 \longrightarrow y_3 \longrightarrow y_4 \longrightarrow y_5$$

- Graph
- Potentials:  $\psi_i(y_i, y_{i+1})$
- $Pr(y_1,\ldots,y_n) = \prod_i \psi_i(y_i,y_{i+1})$
- Find,  $Pr(y_i)$  for any *i*, say  $Pr(y_5 = 1)$ 
  - Exact method:  $Pr(y_5 = 1) = \sum_{y_1,\dots,y_4} Pr(y_1,\dots,y_4,1)$  requires exponential number of summations.
  - ► A more efficient alternative...

## Exact inference on chains

$$\begin{aligned} \mathsf{Pr}(y_5 = 1) &= \sum_{y_1, \dots, y_4} \mathsf{Pr}(y_1, \dots, y_4, 1) \\ &= \sum_{y_1} \sum_{y_2} \sum_{y_2} \sum_{y_3} \sum_{y_4} \psi_1(y_1, y_2) \psi_2(y_2, y_3) \psi_3(y_3, y_4) \psi_4(y_4, 1) \\ &= \sum_{y_1} \sum_{y_2} \psi_1(y_1, y_2) \sum_{y_3} \psi_2(y_2, y_3) \sum_{y_4} \psi_3(y_3, y_4) \psi_4(y_4, 1) \\ &= \sum_{y_1} \sum_{y_2} \sum_{y_2} \psi_1(y_1, y_2) \sum_{y_3} \psi_2(y_2, y_3) B_3(y_3) \\ &= \sum_{y_1} \sum_{y_2} \sum_{y_2} \psi_1(y_1, y_2) B_2(y_2) \\ &= \sum_{y_1} B_1(y_1) \end{aligned}$$

An alternative view: flow of beliefs  $B_i(.)$  from node i + 1 to node i

$$y_1 \longrightarrow y_2 \longrightarrow y_3 \longrightarrow y_4 \longrightarrow y_5$$

Fundamentals of graphical model

Sunita Sarawagi

## Adding evidence

Given fixed values of a subset of variables **x**<sub>e</sub> (evidence), find the *Marginal probability queries over a small subset of variables:* 

Find Pr(Income='High | Degree='PhD')

$$\Pr(x_1) = \sum_{x_2...x_m} \Pr(x_1...x_n | \mathbf{x}_e)$$

- Most likely labels of remaining variables: (MAP queries)
  - Find likely temperature at sensors in a room given readings from a subset of them

$$\mathbf{x}^* = \operatorname{argmax}_{x_1...x_m} \mathsf{Pr}(x_1 \dots x_n | \mathbf{x}_e)$$

Easy to add evidence, just change the potential.

## Inference in HMMs

• Given,

**y**<sub>1</sub>



- Evidence variables:  $\mathbf{x} = x_1 \dots x_n = o_1 \dots o_n$ .
- Find most likely values of the hidden state variables.

$$\mathbf{y} - y_1 \dots y_n$$
  
argmax<sub>y</sub>  $Pr(\mathbf{y} | \mathbf{x} = \mathbf{0})$   
• Define  $\psi_i(y_{i-1}, y_i) = Pr(y_i | y_{i-1}) Pr(x_i = o_i | y_i)$   
• Reduced graph only a single chain of y nodes.

• Algorithm same as earlier, just replace "Sum" with "Max"

 $y_3 \longrightarrow y_4 \longrightarrow y_5 \longrightarrow y_6 \longrightarrow y_7$ 

## This is the well-known Viterbi algorithm

 $y_2$ 

## Exact inference on trees

- Basic steps for marginal and MAP queries.
  - Perform sum/max over leaf node potential and send resulting "belief" to parent.
  - Each internal node, on getting beliefs from its children
    - Multiplies incoming beliefs with its own potentials
    - Performs sum/max on the result
    - Sends resulting "belief" factor to parent.
- Root has the answer.

Linear in the number of nodes in the graph

## Junction tree algorithms

- An optimal general-purpose algorithm for exact marginal/MAP queries
- Simultaneous computation of many queries
- Efficient data structures
- Complexity: O(m<sup>w</sup>N) w= size of the largest clique in (triangulated) graph, m = number of values of each discrete variable in the clique. → linear for trees.
- Basis for many approximate algorithms.
- Many popular inference algorithms special cases of junction trees
  - Viterbi algorithm of HMMs
  - Forward-backward algorithm of Kalman filters

## Creating a junction tree from a graphical model

1. Starting graph



2. Triangulate graph



3. Create clique nodes



4. Create tree edges such that variables connected.



5) Assign potentials to exactly one subsumed clique node.

 $x_1 x_2$ 



## Belief propagation on junction trees

- Each node *c* 
  - sends *belief*  $B_{c \rightarrow c'}(.)$  to each of its neighbors c'
    - ★ once it has beliefs from every other neighbor  $N(c) \{c'\}$ .
  - B<sub>c→c'</sub>(.) = belief that clique c has about the distribution of labels to common variables s = c ∩ c'

$$B_{c\to c'}(\mathbf{x}_s) = \sum_{\mathbf{x}_{c-s}} \psi_c(\mathbf{x}_c) \prod_{d \in N(c) - \{c'\}} B_{d\to c}(\mathbf{x}_{d\cap c})$$

Replace "sum" with "max" for MAP queries.

Compute marginal probability of any variable  $x_i$  as

•  $c = clique in JT containing x_i$ 

**2** 
$$\Pr(x_i) \propto \sum_{\mathbf{x}_{c-x_i}} \psi_c(\mathbf{x}_c) \prod_{d \in N(c)} B_{d \to c}(\mathbf{x}_{d \cap c})$$

## Example



 $\psi_{234}(\mathbf{y}_{234}) = \psi_{23}(\mathbf{y}_{23})\psi_{34}(\mathbf{y}_{34})$  $\psi_{345}(\mathbf{y}_{345}) = \psi_{35}(\mathbf{y}_{35})\psi_{45}(\mathbf{y}_{45})$  $\psi_{234}(\mathbf{y}_{12}) = \psi_{12}(\mathbf{y}_{12})$ 

- Clique "12" sends belief  $B_{12\rightarrow234}(y_2) = \sum_{y_1} \psi_{12}(\mathbf{y}_{12})$  to its only neighbor.
- ② Clique "345" sends belief  $B_{345→234}(\mathbf{y}_{34}) = \sum_{y_5} \psi_{234}(\mathbf{y}_{345})$  to "234"

# Part I: Outline

#### Representation

- Directed graphical models: Bayesian networks
- Undirected graphical models

#### 2 Inference Queries

- Exact inference on chains
- Exact inference on general graphs

## 3 Constructing a graphical model

- Graph Structure
- Parameters in Potentials

#### 4 Approximate inference

- Generalized belief propagation
- Sampling: Gibbs, Particle filters

# Graph Structure

- Manual: Designed by domain expert
  - Used in applications where dependency structure is well-understood
  - Example: QMR systems, Kalman filters, Vision (Grids), HMM for speech recognition and IE.
- 2 Learnt: from examples
  - NP hard to find the optimal structure.
  - Widely researched, mostly posed as a branch and bound search problem.
  - Useful in dyanmic situations
  - Example: Selectivity estimation over attributes of arbitrary tables.

## Parameters in Potentials

- Manual: Provided by domain expert
  - Used in infrequently constructured graphs, example QMR systems
  - Also where potentials are an easy function of the attributes of connected graphs, example: vision networks.
- 2 Learnt: from examples
  - More popular since difficult for humans to assign numeric values
  - Many variants of parameterizing potentials.
    - Each potential entry a parameter, example, HMMs
    - Potentials: combination of shared parameters and data attributes: example, CRFs. (Discussed in later with extraction)

# Part I: Outline

#### Representation

- Directed graphical models: Bayesian networks
- Undirected graphical models

#### 2 Inference Queries

- Exact inference on chains
- Exact inference on general graphs

## 3 Constructing a graphical model

- Graph Structure
- Parameters in Potentials

#### Approximate inference

- Generalized belief propagation
- Sampling: Gibbs, Particle filters

## Why approximate inference

- Exact inference is NP hard. Complexity:  $O(w^m)$ 
  - w= tree width = size of the largest clique in (triangulated) graph-1,
  - m = number of values of each discrete variable in the clique.
- Many real-life graphs produce large cliques on triangulation
  - A  $n \times n$  grid has a tree width of n
  - A Kalman filter on K parallel state variables influencing a common observation variable, has a tree width of size K + 1

# Generalized belief propagation

- Approximate junction tree with a cluster graph where
  - Nodes = arbitrary clusters, not cliques in triangulated graph. Only ensure all potentials subsumed.
  - 2 Separator nodes on edges = subset of intersecting variables.


## Belief propagation in cluster graphs

- Graph can have loops, tree-based two-phase method not applicable.
- Many variants on scheduling order of propagating beliefs.
  - Simple loopy belief propagation [Pea88]
  - Tree-reweighted message passing [Kol04]
  - Residual belief probagation [EMK06]
- Most have no guarantees of convergence
- Works well in practice, default method of choice.
  - Success story: Error correction using Turbo code

## MCMC (Gibbs) sampling

- Useful when all else failes, guaranteed to converge to the optimal over infinite number of samples.
- Basic premise: easy to compute conditional probability Pr(x<sub>i</sub>|fixed values of remaining variables)

### Algorithm

- Start with some initial assignment, say  $\mathbf{x}^1 = [x_1, \dots, x_n] = [0, \dots, 0]$
- For several iterations
  - For each variable  $x_i$

Get a new sample  $\mathbf{x}^{t+1}$  by replacing value of  $x_i$  with a new value sampled according to probability  $Pr(x_i|x_1^t, \dots, x_{i-1}^t, x_{i+1}^t, \dots, x_n^t)$ 

### Others

- Combinatorial algorithms for MAP [BVZ01, DTEK07, GDS07]
- Greedy algorithms: relaxation labeling
- Variational methods
- LP and QP based approaches

### Inference Task in DBNs



• Simplied representation of a dynamic Bayesian network

- Hidden state variables: x; Observed variables: o
- Assumed to be vector valued
- Given:
  - Prior on the initial state:  $p(\mathbf{x}_0)$
  - How state evolves:  $p(\mathbf{x}_t | \mathbf{x}_{t-1})$
  - How obsevations depend on state:  $p(\mathbf{o}_t | \mathbf{x}_t)$
- Estimate the *state at time t* given *observations till time t* 
  - The posterior distribution:  $p(\mathbf{x}_t | \mathbf{o}_{1:t})$

### Alternative Inference Tasks in DBNs



- Estimate the most likely sequence of states (for discrete x)
   argmax<sub>x1:t</sub> p(o<sub>1:t</sub>|x<sub>1:t</sub>) (*Cf. Viterbi Algorithm*)
- Estimate the distribution of all states till time t
  - $p(\mathbf{x}_{1:t}|\mathbf{o}_{1:t})$
- Estimate the state at time t given measurements till time t + I (fixed-lag smoothing)
  - $p(\mathbf{x}_t | \mathbf{o}_{1:t+l})$
  - Why ? Belief about the state at time t may change drastically given future observations

### Exact Inference in DBNs



- Easy to write down
  - Using Bayes rule and Chain rule, we get:

$$p(\mathbf{x}_t | \mathbf{o}_{1:t}) = \frac{p(\mathbf{o}_t | \mathbf{x}_t) \int p(\mathbf{x}_t | \mathbf{x}_{t-1}) p(\mathbf{x}_{t-1} | \mathbf{o}_{1:t-1}) d\mathbf{x}_{t-1}}{p(\mathbf{o}_t | \mathbf{o}_{1:t-1})}$$

- Where:
  - \*  $p(\mathbf{o}_t | \mathbf{x}_t)$  and  $p(\mathbf{x}_t | \mathbf{x}_{t-1})$  are known model parameters
  - ★  $p(\mathbf{x}_{t-1}|\mathbf{o}_{1:t-1})$  is available from the previous time
  - ★  $p(\mathbf{o}_t | \mathbf{o}_{1:t-1}) = \int p(\mathbf{o}_t | \mathbf{x}_t) p(\mathbf{x}_t | \mathbf{o}_{1:t-1}) d\mathbf{x}_t$  is a normalization constant (so may not need to be evaluated)

### Exact Inference in DBNs



- However, can solve exactly in very few cases:
  - Kalman filters: if the system is linear Gaussian
    - ★ If  $p(\mathbf{x}_{t-1}|\mathbf{o}_{1:t-1})$  is Gaussian and the system is linear Gaussian,  $p(\mathbf{x}_t|\mathbf{o}_{1:t})$  is Gaussian
    - ★ Very efficient
    - ★ Backward smoothing also easily doable
  - Grid-based method: if the state space is discrete and finite
    - ★ Can compute the integral as a sum exactly

### Approximate Inference in DBNs



- Extended Kalman Filter
  - Approximate the process as a linear Gaussian system
  - Will fail if the posterior density not close to a Gaussian (e.g. if it is bimodal or heavily skewed)
- Approximate Grid-based methods
  - Discretize the continuous state space using a grid
  - Need sufficiently dense grid for good approximation
  - Suffers from "curse of dimensionality"



- Approximate the state using a set of weighted samples, called particles
- At time t 1, approximate  $p(x_{t-1}|o_{1:t-1})$  using n particles:

• {
$$\mathbf{x}_{t-1}^1, w_{t-1}^1$$
}, { $\mathbf{x}_{t-1}^2, w_{t-1}^2$ },  $\cdots$ , { $\mathbf{x}_{t-1}^n, w_{t-1}^n$ }

• Can estimate any statistic using these particles

• e.g. 
$$E(\mathbf{x}_{t-1}|\mathbf{o}_{1:t-1}) \approx \sum_{i=1}^{n} w_{t-1}^{i} \mathbf{x}_{t-1}^{i}$$

• Inference Task: Generate a set of particles corresponding to  $p(\mathbf{x}_t | \mathbf{o}_{1:t})$  given  $\mathbf{o}_t$ 



- Generate one sample each from:  $p(\mathbf{x}_t | \mathbf{x}_{t-1}^i, \mathbf{o}_t)$
- Assign weights as:

$$w_t^i \propto w_{t-1}^i 
ho(\mathbf{o}_t | \mathbf{x}_{t-1}^i) = w_{t-1}^i \int 
ho(\mathbf{o}_t | \mathbf{x}_t') 
ho(\mathbf{x}_t' | \mathbf{x}_{t-1}^i) d\mathbf{x}_t'$$

- Problems:
  - Requires sampling from  $p(\mathbf{x}_t|...)$  and computing  $p(\mathbf{o}_t|...)$
  - Requires evaluating complex integrals
- Can solve in very few cases:
  - x<sub>t</sub> is discrete, or
  - $p(\mathbf{x}_t | \mathbf{x}_{t-1}^i, \mathbf{o}_t)$  is Gaussian (evolution can still be non-linear)



- Must use *importance sampling* 
  - Use an *importance density* q() to generate samples from
  - ...that closely approximates the true density p()
  - No magic bullet for choosing q()
- Degeneracy issues
  - After a while, a single particle has all the weight
  - Need to resample periodically



- Many other extensions/variations have been considered
  - A lot more art than science at this point
- For an approachable introduction, see "A Tutorial on Particle Filters for On-line Nonlinear/Non-Gaussian Bayesian Tracking"; Arulampalam et al.; IEEE Trans. Signal Processing; 2002
  - Our discussion heavily borrows from it

## More on graphical models

- Koller and Friedman book (Structured Probabilistic Models) not published yet but you could request authors for a draft.
- Kevin Murphy's brief online introduction (http://www.cs.ubc.ca/~murphyk/Bayes/bnintro.html)
- Graphical models. M. I. Jordan. Statistical Science (Special Issue on Bayesian Statistics), 19, 140-155, 2004. (http: //www.cs.berkeley.edu/~jordan/papers/statsci.ps.gz)
- Other text books:
  - R. G. Cowell, A. P. Dawid, S. L. Lauritzen and D. J. Spiegelhalter. "Probabilistic Networks and Expert Systems". Springer-Verlag. 1999.
  - J. Pearl. "Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference." Morgan Kaufmann. 1988.
  - Graphical models by Lauritzen, Oxford science publications F.
     V. Jensen. "Bayesian Networks and Decision Graphs". Springer.
     2001.



#### Yuri Boykov, Olga Veksler, and Ramin Zabih.

Fast approximate energy minimization via graph cuts. *IEEE Trans. Pattern Anal. Mach. Intell.*, 23(11):1222–1239, 2001.

#### J. Duchi, D. Tarlow, G. Elidan, and D. Koller.

Using combinatorial optimization within max-product belief propagation. In Advances in Neural Information Processing Systems (NIPS 2006), 2007.

#### G. Elidan, I. McGraw, and D. Koller.

Residual belief propagation: Informed scheduling for asynchronous message passing. In *Proceedings of the Twenty-second Conference on Uncertainty in AI (UAI)*, Boston, Massachussetts, July 2006.



Rahul Gupta, Ajit A. Diwan, and Sunita Sarawagi.

Efficient inference with cardinality-based clique potentials.

In Proceedings of the 24<sup>th</sup> International Conference on Machine Learning (ICML), USA, 2007.

Vladimir Kolmogorov.

Convergent tree-reweighted message passing for energy minimization. Technical Report MSR-TR-2004-90, Microsoft Research (MSR), September 2004.



#### Judea Pearl.

Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference. Morgan Kaufmann, 1988.

# Part II: Applications



# Part II: Outline

# Selectivity Estimation and Query Optimization

# > Probabilistic Relational Models

## > Probabilistic Databases

# Sensor/Stream Data Management

## >References

- Estimating the intermediate result sizes that may be generated during query processing
  - > Equivalently, selectivities of predicates over tables
- Key to obtaining good plans during optimization

SSN	 Income	 Homeowner?
	 100000	 Yes
	 11000	 Yes

### Customer

### **Purchases**

SSN	Store	 Amount

### Single-table predicates:

income > 90000 and homeowner = yes (on customer)

### Multi-table predicates:

p.ssn = c.ssn and c.homeowner = "no" and p.amount > 10000 (over Customer c and Purchases p)

- Optimizers make several independence assumptions
- Attribute value independence assumption
  - Attributes assumed to be independently distributed
  - Rarely true in practice

#### **Customer**

SSN	 Income	:	Homeowner?
	 100000	:	Yes
	 11000		Yes
	 50000		No
	 30000		No
	 200000		Yes

### Estimate

*p(income > 90000 and homeowner = yes)* 

#### as

p(income > 900000) \* p(homeowner = yes)

### Can result in severe underestimation

#### In reality:

p(income > 900000, homeowner = yes) ≈ p(homeowner = yes)

### Join uniformity assumption

Tuples from one relation assumed equally likely to join with tuples from other relation

Real datasets exhibit large skews

#### Customer

#### **Purchases**

SSN	 Income		Homeowner?	SSN	Store	 Amount
	 100000		Yes			
	 11000		Yes			
	 50000		No			
	 30000	:	No			
	 200000		Yes			

- Errors propagate exponentially [IC'91]
- > Optimizers highly sensitive to underestimation
  - May choose nested-loop joins
- Proposed solutions:
  - Multi-dimensional histograms, wavelets [PI'97,MVW'98, GKTD'00]
    - Expensive to build and maintain
    - > Suffer from "curse of dimensionality" in high dimensions
  - Random sampling [CDN'07]
    - Not as storage efficient
    - > Few matching tuples for high dimensional queries
    - > Need different sampling techniques for joins [AGPR'99]

Eliminating attribute value independence assumption [GTK'01,DGR'01,LWV'03,PMW'03]

#### Customer

SSN	age	Income	zipcode	Home owner?
		100000		Yes
		11000		Yes
		50000		No
		30000		No
		200000		Yes



Eliminating attribute value independence assumption [GTK'01,DGR'01,LWV'03,PMW'03]

#### Customer

SSN	age	Income	zipcode	Home
				owner?
		100000		Yes
		11000		Yes
		50000		No
		30000		No
		200000		Yes

**Selectivity** 

**Estimates** 



Eliminating join uniformity assumption ??

# Part II: Outline

# Selectivity Estimation and Query Optimization

# Probabilistic Relational Models

## > Probabilistic Databases

# Sensor/Stream Data Management

## >References

# Probabilistic Relational Models

- Real-world data often has highly relational structure
  - > There are *entities* and *relationships* between them etc
  - Bayesian networks treat each one individually
  - Will need a huge Bayesian network if we want to represent the uncertainties in such data
- PRMs: Generalization of PGMs to relational framework [FGKP'99]
  - Allows dependence over attributes in different relations through joins
- Significantly enrich both Bayesian networks and relational model

# **Relational Schema**



## Describes the types of objects and relations in the database

# Probabilistic Relational Model



# **PRM: Semantics**



Fixed relational skeleton

Objects and links between them

Non-key (descriptive) attributes uncertain

Amol Deshpande

Part II: Applications



### PRM defines distribution over instantiations of attributes

Amol Deshpande

Part II: Applications



Amol Deshpande

University of Maryland

# PRMs: Inference/Generalizations

### > Inference

> Option 1: Construct and use the *ground* Bayesian network

> Allows exact inference

> Too large for any reasonable dataset

> Option 2: Approximate inference

E.g. using *loopy belief propagation* 

### Generalizations

- Link uncertainty [GGFKT'02]
- Finer granularity dependencies using class hierarchies [dGK'00]
- > Undirected dependencies

Relational Markov networks [TAK'02]

Relational dependency networks [ND'04]

Exciting research area with huge potential impact in databases !!

# Part II: Outline

# Selectivity Estimation and Query Optimization (continued)

# > Probabilistic Relational Models

## > Probabilistic Databases

# Sensor/Stream Data Management

## >References

- Eliminating join uniformity assumption [GTK'01]
- Using a Probabilistic Relational Model



Can estimate selectivities of joint predicates across relations

- Eliminating join uniformity assumption [GTK'01]
- Using a Probabilistic Relational Model
- Caveat:
  - > Should not use them blindly
  - Need to add and reason about a new join indicator variable

- Called a Statistical Relational Model
- Details in Getoor, Tasker, Koller; SIGMOD 2001.

# **Discussion and Open Problems**

- Approximate query processing ?
   Can use the proposed techniques as they are
   However, no guarantees on the accuracy of results
   Optimize accuracy for a given storage
   To obtain guarantees, optimize for accuracy alone
   May result in large CPDs
- Using learned PGMs during optimization
   Optimizers get better selectivity estimates, but otherwise unaware of the modeling
   May be beneficial to explore tighter integration

# **Discussion and Open Problems**

- Can exploit new types of query plans
  - Based on horizontal partitioning of the relations [BBDW'05,DGHM'05,P'05]
  - Use different plans for different partitions of relations based on attribute values

- > Adaptive query processing
  - >PGMs ideal for learning the distribution properties
  - > Significantly fewer parameters  $\rightarrow$  easier to learn
  - > Many research challenges
#### Part II: Outline

#### Selectivity Estimation and Query Optimization

#### > Probabilistic Relational Models

#### Probabilistic Databases

#### Sensor/Stream Data Management

#### >References

#### **Probabilistic Databases**

- > Motivation: Increasing amounts of *uncertain* data
  - From sensor networks
    - > Imprecise data, data with confidence/accuracy bounds
    - Human-observed data
  - Statistical modeling/machine learning
    - Many models provide a distribution over a set of labels (e.g. classification models, HMMs)
  - > Approximate/vague queries
  - Information extraction
- Probability theory provides a strong foundation to reason about this
  - Caveat: It is not always clear if the underlying uncertainty measure follows probability theory semantics

### Probabilistic Databases

- Goal: Managing and querying data annotated with probabilities using databases
- > Types of uncertainties
  - > Existence uncertainty

> Don't know if a tuple exists in the database for sure

E.g. a sensor may detect a bird, but not 100% sure

> Attribute-value uncertainty

> The value of an attribute is not known for sure

Instead a distribution over the possible values is provided

- E.g. a sensor detects a bird for sure, but it may be a sparrow or a dove or something else
- Much work in recent years on both [DS'07]

#### Correlations in Probabilistic Databases

- Much of the probabilistic data is naturally correlated
  - > E.g. sensor data, data integration [AFM'06]
- Even if not..
  - Correlations get introduced during query processing

#### Example Probabilistic Database

Example from Dalvi and Suciu [2004]



#### Possible worlds

instance	probability		
{s1, s2, t1}	0.12		
{s1, s2}	0.18		
{s1, t1}	0.12		
{s1}	0.18		
{s2, t1}	0.08		
{s2}	0.12		
{t1}	0.08		
{}	0.12		

#### Correlations during query processing

Example from Dalvi and Suciu [2004]



Part II: Applications

#### Correlations in Probabilistic Databases

Much of the probabilistic data is naturally correlated

> E.g. sensor data, data integration

Even if not..

Correlations get introduced during query processing

Can use PGMs to capture such correlations

# **Example: Mutual Exclusivity**

#### Possible worlds

					Worrao			640
S				instance	nrohahility	<u>X<sub>s1</sub></u>	<b>X</b> <sub>t1</sub>	<b>†1()</b>
			1		probability	0	0	0
	Α	B	<u>prob</u>	{\$1, \$2, t1}	0		1	01
s1	m	1	0.6	{s1, s2}	0.3			0.4
s2	n	1	0.5	{s1, t1}	0		0	0.6
				{s1}	0.3		1	0
Т				{s2, t1}	0.2			
	С	D	<u>prob</u>	{s2}	0	,	/   s	٦/١
t1	1	р	0.4	{t1}	0.2		s2 14	<u>()</u>
				{}	0	_		.5

#### Possible worlds (if desired) computed using inference

Amol Deshpande

Part II: Applications

0.5

1

Introduce new factors as new tuples generated



Introduce new factors as new tuples generated

X <sub>s1</sub>	<b>X</b> <sub>t1</sub>	<b>X</b> <sub>i1</sub>	<b>f</b> <sup>AND</sup>
0	0	0	1
0	) 1 (		1
1	0	0	1
1	1	0	0
0	0	1	0
0	1	1	0
1	0	1	0
1	1	1	1



Introduce new factors as new tuples generated



- > Query evaluation ≡ Inference !!
- Can use variable elimination or junction tree..
- Can also use approximate inference algorithms



#### Discussion

- Similar to intensional semantics [FR'97,DS'04]
  - > Except this exposes the structure of the problem
  - > Can exploit for more efficient execution

- Safe plans on independent tuples generate tree-structured models
  - > Highly efficient inference

#### Part II: Outline

#### Selectivity Estimation and Query Optimization

#### > Probabilistic Relational Models

#### Probabilistic Databases

#### Sensor/Stream Data Management

#### >References

#### Motivation

#### Unprecedented, and rapidly increasing, instrumentation of our every-day world



Distributed measurement networks (e.g. GPS)



**RFID** 



Network Monitoring



Wireless sensor networks



Industrial Monitoring Part II: Applications Sensor Data Management: Challenges

- Data streams generated at very high rates
- > Strong spatio-temporal correlations in the data
- In-network, distributed processing tasks
  - Global inference needed to achieve consistency
- Need for higher-level modeling over the data
  - Typically imprecise, unreliable and incomplete data
     Measurement noises, failures, biases ...
  - > Application often need higher-level, *hidden* variables
  - > Pattern recognition, forming stochastic descriptions...

#### Sensor Data Management

- > A statistical/probabilistic model of the data must be incorporated in the sensor data processing
- Probabilistic graphical models are a natural
  - Can capture and exploit the spatial and temporal nature of the underlying process
  - > Minimize the number of parameters
  - >Amenable to distributed processing

#### Outline

#### A generic temporal model for sensor stream data

#### > Applications

- Online estimation and filtering
- Inferring hidden variables
- Model-based query processing
- In-network inference
- Miscellaneous









#### **Markov Property**

Interpretation:  ${X_{i,t+1}}$  independent of  ${X_{i,t-1}}$  given  ${X_{i,t}}$ 





#### State evolution can be modeled as a Dynamic Bayesian Network

Part II: Applications



#### **Parameters** ?

(1) System model Prior:  $p(X_{1,0}, X_{2,0}, X_{3,0})$ Evolution:  $p(X_{1,t}, X_{2,t}, X_{3,t} | X_{1,t-1}, X_{2,t-1}, X_{3,t-1})$ 

Part II: Applications



#### **Parameters** ?

(2) Measurement model  $p(O_{1,t}, O_{2,t}, O_{3,t} | X_{1,t}, X_{2,t}, X_{3,t})$ 

#### Outline

#### A generic temporal model for sensor stream data

- > Applications
  - Online estimation and filtering
  - Inferring hidden variables
  - Model-based query processing
  - >In-network inference
  - Miscellaneous

#### Application: Online Estimation and Filtering

- Using linear Gaussian dynamical systems
  - E.g. Kalman Filters
- Task: Estimating velocity and location from noisy GPS readings



#### Application: Online Estimation and Filtering

- Using linear Gaussian dynamical systems
  - E.g. Kalman Filters
- Task: Estimating velocity and location from noisy GPS readings

$$p(v_{t}|v_{t-1}) = N(v_{t-1}, \sigma_{v})$$

$$p(x_{t}|x_{t-1}, v_{t}) = N(x_{t-1} + v_{t}, \sigma_{x})$$

$$p(o_{t}|x_{t}) = N(x_{t}, \sigma_{o})$$

$$Prior: p(v_{0}), p(x_{0})$$

$$(v_{t-1})$$

X<sub>f</sub>

#### Application: Online Estimation and Filtering

- > Using linear Gaussian dynamical systems
  - E.g. Kalman Filters
- Closed-form equations for state estimation [Kalman'60]
  - > Because of the *linear Gaussian* assumption
- LDS Applications:
  - > Autopilot
  - Inertial guidance systems
  - Radar tracker
  - Economics...
- In databases:
  - > Adaptive stream resource management [JCW'04]
  - > Approximate querying in sensor networks [DGHHM'04]

#### Outline

#### A generic temporal model for sensor stream data

- > Applications
  - Online estimation and filtering
  - Inferring hidden variables
  - Model-based query processing
  - In-network inference
  - Miscellaneous

- Inferring "transportation mode"/ "activities" [P+04]
  - Using easily obtainable sensor data (GPS, RFID proximity data)
  - > Can do much if we can infer these automatically



- Inferring "transportation mode"/ "activities" [P+04]
  - Using easily obtainable sensor data (GPS, RFID proximity data)
  - > Can do much if we can infer these automatically



Desired data: Clean path annotated with transportation mode Online, in real-time

Use a dynamic Bayesian network to model the system state

 $Time = t \qquad Time = t+1$ 

*Transportation Mode: Walking, Running, Car, Bus* 

True velocity and location

**Observed** location



Amol Deshpande

Part II: Applications

University of Maryland

- Given a sequence of observations  $(O_t)$ , infer most likely  $M_t$ 's that explain it.
- Alternatively, could provide a probability distribution on the possible  $M_t$ 's.

Time = t

Time = t+1



University of Maryland

Amol Deshpande

#### Outline

#### A generic temporal model for sensor stream data

- > Applications
  - Online estimation and filtering
  - Inferring hidden variables
  - Model-based query processing
  - In-network inference
  - Miscellaneous

#### **Application: Model-based Query Processing** [DGMHH'04, SBEMY'06]



Amol Deshpande

University of Maryland

# Application: Model-based Query Processing [DGMHH'04,SBEMY'06]







#### Advantages:

Exploit correlations for efficient approximate query processing Handle noise, biases in the data Predict missing or future values




#### Outline

#### A generic temporal model for sensor stream data

- > Applications
  - Online estimation and filtering
  - Inferring hidden variables
  - Model-based query processing
  - In-network inference
  - Miscellaneous

- Often need to do in-network, distributed inference
  - Target tracking through information fusion
  - Optimal control (for actuation)
  - Distributed sensor calibration (using neighboring sensors)
  - In-network regression or function fitting



- Often need to do in-network, distributed inference
  - > Target tracking through information fusion
  - Optimal control (for actuation)
  - Distributed sensor calibration (using neighboring sensors)
  - > In-network regression or function fitting
- Obey a common structure:
  - Each sensor has/observes some *local* information
  - Information across sensors is correlated
  - The information must be combined together to form a global picture
  - The global picture (or relevant part thereof) should be sent to each sensor

#### Naïve option:

- Collect all data at the centralized base station too expensive
- > Using graphical models
  - Form a junction tree on the nodes directly
  - Use message passing (or loopy propagation [CP'03]) to form a global consistent view



Amol Deshpande

#### Naïve option:

- Collect all data at the centralized base station too expensive
- > Using graphical models:
  - Form a junction tree on the nodes directly
  - Use message passing (or loopy propagation [CP'03]) to form a global consistent view



Amol Deshpande

### Outline

#### A generic temporal model for sensor stream data

- > Applications
  - Online estimation and filtering
  - Inferring hidden variables
  - Model-based query processing
  - In-network inference
  - Miscellaneous

### Applications: Miscellaneous

- Data compression [CDHH'06]
  - Central task in sensor networks
    - Collect all observed data at the base station at specified frequency
  - > Challenge: How to exploit the correlations
  - > Probabilistic graphical models ideally suited:
    - > Can capture the correlations/pattern
    - > Allow for local checking of constraints/correlations
- Fault/anomaly detection
- Distributed regression
- Sensor calibration

### Part II: Outline

### Selectivity Estimation and Query Optimization

### > Probabilistic Relational Models

#### > Probabilistic Databases

#### Sensor/Stream Data Management

#### > References

- [AFM'06] Periklis Andritsos and Ariel Fuxman and Renee J. Miller; Clean Answers over Dirty Databases; ICDE 2006
- [BBDW'05] Pedro Bizarro, Shivnath Babu, David J. DeWitt, Jennifer Widom: Content-Based Routing: Different Plans for Different Data. VLDB 2005: 757-768
- [BDHW'06] O. Benjelloun, A. Das Sarma, A. Halevy, and J. Widom. ULDBs: Databases with uncertainty and lineage. In VLDB, 2006.
- [BGP'92] D. Barbara, H. Garcia-Molina, and D. Porter; The management of probabilistic data; In IEEE Trans. of Knowledge Data Eng.; 1992.
- [BGR'01] Shivnath Babu, Minos N. Garofalakis, Rajeev Rastogi: SPARTAN: A Model-Based Semantic Compression System for Massive Data Tables. SIGMOD Conference 2001: 283-294
- [CDN'07] Surajit Chaudhuri and Gautam Das and Vivek Narasayya; "Optimized stratified sampling for approximate query processing"; TODS 2007.
- [CP'03] Christopher Crick and Avi Pfeffer; Loopy Belief Propagation as a Basis for Communication in Sensor Networks; UAI 2003.
- [dGK'00] Marie desJardins, Lise Getoor, Daphne Koller: Using Feature Hierarchies in Bayesian Network Learning. SARA 2000: 260-270
- [DFG'05] Arnaud Doucet, Nando de Freitas, and Neil Gordon. Sequential Monte Carlo methods in practice. Springer, 2005.
- [DGHM'05] Amol Deshpande, Carlos Guestrin, Wei Hong, Samuel Madden: Exploiting Correlated Attributes in Acquisitional Query Processing. ICDE 2005: 143-154
- [DGMHH'04] Amol Deshpande, Carlos Guestrin, Samuel Madden, Joseph M. Hellerstein, Wei Hong: Model-Driven Data Acquisition in Sensor Networks. VLDB 2004: 588-599

- [DGMHH'04] Amol Deshpande, Carlos Guestrin, Samuel Madden, Joseph M. Hellerstein, Wei Hong: Model-Driven Data Acquisition in Sensor Networks. VLDB 2004: 588-599
- [DGR'01] Amol Deshpande, Minos N. Garofalakis, Rajeev Rastogi: Independence is Good: Dependency-Based Histogram Synopses for High-Dimensional Data. SIGMOD Conference 2001: 199-210
- [DS'07] Nilesh Dalvi and Dan Suciu; Management of Probabilistic Data: Foundations and Challenges; PODS 2007.
- [FGKP'99] Nir Friedman, Lise Getoor, Daphne Koller, Avi Pfeffer; Learning Probabilistic Relational Models; IJCAI 1999: 1300-1309.
- [FT'97] N. Fuhr and T. Rolleke; A probabilistic relational algebra for the integration of information retrieval and database systems; ACM Trans. on Info. Syst.; 1997
- [G'06] Lise Getoor; An Introduction to Probabilistic Graphical Models for Relational Data; IEEE Data Engineering Bulletin; March 2006.
- [G'98] Zoubin Ghahramani. Learning dynamic Bayesian networks. Lecture Notes in Computer Science, 1387, 1998.
- [GGFKT'02] Lise Getoor, Nir Friedman, Daphne Koller, Benjamin Taskar: Learning Probabilistic Models of Link Structure. Journal of Machine Learning Research 3: 679-707 (2002)
- [GKTD'00] D.Gunopulos, G.Kollios, V.J.Tsotras and C.Domeniconi. Approximating Multi-Dimensional Aggregate Range Queries Over Real Attributes. SIGMOD 2000.
- [GTK'01] Lise Getoor, Benjamin Taskar, Daphne Koller: Selectivity Estimation using Probabilistic Models. SIGMOD Conference 2001: 461-472
- [IC'93] Y. E.Ioannidis and S.Christodoulakis. "Optimal Histograms for Limiting Worst-Case Error Propagation in the Size of Join Results"; TODS 1993.

- [IL'84] T. Imielinski and W. Lipski, Jr; Incomplete information in relational databases; Journal of the ACM; 1984.
- [JCW'04] Jain, E. Change, and Y. Wang. Adaptive stream resource management using kalman filters. In SIGMOD, 2004.
- [KG'06] Daniel Kifer and Johannes Gehrke; Injecting utility into anonymized datasets; SIGMOD 2006.
- [LWV'03] Lipyeow Lim, Min Wang, Jeffrey Scott Vitter: SASH: A Self-Adaptive Histogram Set for Dynamically Changing Workloads. VLDB 2003: 369-380.
- [M'01] Kevin Murphy. The Bayes net toolbox for matlab. Computing Science and Statistics, 33, 2001.
- [M'02] Kevin Murphy. Dynamic Bayesian Networks: Representation, Inference and Learnig. PhD thesis, UC Berkeley, 2002.
- [MP'01] V. Mihajlovic and M. Petkovic. Dynamic bayesian networks: A state of the art. University of Twente Document Repository 2001.
- [MVW'98] Y. Matias, J. S. Vitter, and M. Wang. Wavelet-Based Histograms for Selectivity Estimation. SIGMOD 1998.
- > [ND'04] J. Neville and D. Jensen; Dependency Networks for Relational Data; ICDM 2004.
- P+'04] Matthai Philipose et al; Inferring Activities from Interactions with Objects; IEEE Pervasive Computing, October 2004
- PGM'05] Mark A. Paskin, Carlos Guestrin, Jim McFadden: A robust architecture for distributed inference in sensor networks. IPSN 2005.

- PLFK'03] D. Patterson, L. Liao, D. Fox, and H. Kautz. Inferring high level behavior from low level sensors. In UBICOMP, 2003.
- [R'89] L Rabiner. A tutorial on hidden Markov models and selected applications in speech recognition. 1989.
- [SBEMY'06] Adam Silberstein, Rebecca Braynard, Carla Ellis, Kamesh Munagala and Jun Yang. A Sampling-Based Approach to Optimizing Top-k Queries in Sensor Networks. ICDE 2006.
- [SD'07] Prithviraj Sen, Amol Deshpande: Representing and Querying Correlated Tuples in Probabilistic Databases. ICDE 2007: 596-605
- SKGGM'05] Vipul Singhvi, Andreas Krause, Carlos Guestrin, James H. Garrett Jr., H. Scott Matthews: Intelligent light control using sensor networks. SenSys 2005: 218-229
- [TAK'02] Benjamin Taskar, Pieter Abbeel, Daphne Koller; Discriminative Probabilistic Models for Relational Data; UAI 2002: 485-492
- [W'05] J. Widom. Trio: A system for integrated management of data, accuracy, and lineage. In CIDR, 2005.
- [Y'00] Jie Ying. A hidden markov model-based algorithm for fault diagnosis with partial and imperfect tests. IEEE Trans. on Systems, Man, and Cybernetics, Part C, 2000.

# Part III: Graphical models for Information extraction and data integration



# Graphical models for Information extraction and data integration

Sunita Sarawagi IIT Bombay http://www.cse.iitb.ac.in/~sunita

Part III: Information Extraction and Data Integration

#### Information Extraction (IE) & Integration

The Extraction task: Given,

- E: a set of structured elements
- S: unstructured source S

extract all instances of E from S

The integration task: Given

- database of existing inter-linked entities

Resolve which entities are the same.

- Many versions involving many source types
- Actively researched in varied communities
- Several tools and techniques
- Several commercial applications

### IE from free format text

- Classical Named Entity Recognition
  - Extract person, location, organization names

According to Robert Callahan, president of Eastern's flight attendants union, the past practice of Eastern's parent, Houston-based Texas Air Corp., has involved ultimatums to unions to accept the carrier's terms

- Several applications
  - -News tracking
    - Monitor events
  - -Bio-informatics
    - Protein and Gene names from publications
  - -Customer care
    - •Part number, problem description from emails in help centers

### **Text segmentation**

House

numberBuildingRoadCityStateZip4089Whispering PinesNobel DriveSan DiegoCA92122



# Information Extraction (IE)

- Many different uses
  - Disease outbreaks from news articles
  - Addresses/Qualifications from resumes for HR DBs
  - Titles/Authors/Venue/Year from citations
  - Room attributes from hotel websites
- Many approaches
  - Rules-based ----- Statistical learners
- Varying levels of difficulty
  - Wrappers for machine generated pages
  - **.**..
  - Fact extraction from speech transcripts

### Graphical models in Extraction & Dedup

- State of the art: Conditional Random Fields
- IE Models
  - Basic IE model (Chain)
  - IE with collective labeling of repeated words
- De-duplication models
  - Basic pair-wise model
  - Collective de-duplication of relational data
  - Collective de-duplication of multiple networked entities

### **Conditional Random Fields**

Special undirected graphical model

- 1. Conditional distribution Pr (**y**|**x**) where **y** =  $y_1y_2...y_n$
- 2. Graph: over the interdependent components of **y**
- 3. Potentials: weighted sum of features over **x**



#### [Lafferty et al 2001]

### Chain model

My review of Fermat's last theorem by S. Singh

y	Other	Other	Other	Title	Title	Title	other	Author	Author
x	Му	review	of	Fermat's	last	theorem	by	S.	Singh
t	1	2	3	4	5	6	7	8	9

$$y_1 - y_2 - y_3 - y_4 - y_5 - y_6 - y_7 - y_8 - y_9$$
  
 $\mathbf{f}(y_i, y_{i-1}, i, \mathbf{x})$ 



 $f_2(y_i, \mathbf{x}, i, y_{i-1}) = 1$  if  $y_i$  is Person &  $x_i$  is Douglas

 $f_3(y_i, \mathbf{x}, i, y_{i-1}) = 1$  if  $y_i$  is Person &  $y_{i-1}$  is Other

Parameters: weight for each feature (vector)  $\mathbf{W} = W_1 W_2 \dots W_{|\mathbf{f}|} \text{ Machine learnt}$ 

### Features in typical extraction tasks

- Words
- Orthographic word properties
  - Capitalized? Digit? Ends-with-dot?
- Part of speech
  - Noun?
- Match in a dictionary
  - Appears in a dictionary of people names?
  - Appears in a list of stop-words?
- Fire these for each label and
  - The token,
  - W tokens to the left or right, or
  - Concatenation of tokens.

# Examples: features with weights (publications).

#	Name	Person	Location	Other
1	x <sub>i</sub> is noun	1.2	1.2	-0.5
4	"at" in {x <sub>i-1</sub> , x <sub>i-2</sub> }	-0.3	3	0.2
7	x <sub>i-1</sub> x <sub>i</sub> in people names dictionary	3	-0.4	0
10	x <sub>i-1</sub> is single caps & dot.	2.1	-1.0	-0.1
13	y <sub>i-1</sub> is Location	-1.5	0.3	1.0
•				
•				
100000				

#### A large number

Part III: Information Extraction and Data Integration

# **Typical numbers**

Seminars announcements (CMU):

- speaker, location, timings
- SVMs for start-end boundaries
- 250 training examples
- F1: 85% speaker, location, 92% timings (Finn & Kushmerick '04)
- Jobs postings in news groups
  - 17 fields: title, location, company, language, etc
  - 150 training examples
  - F1: 84% overall (LP2) (Lavelli et al 04)

### Graphical models in Extraction & Dedup

- State of the art: Conditional Random Fields
- IE Models
  - Basic IE model (Chain)
  - IE with collective labeling of repeated words
- De-duplication models
  - Basic pair-wise model
  - Collective de-duplication of relational data
  - Collective de-duplication of multiple networked entities

### **Collective labeling**

- Y has character.
- Mr. X lives in Y.
- X buys Y Times daily.



Other applications of associative potentials Social network analysis: "friends of smokers are smokers" Image segmentation: "nearby pixels get the same label" Spam detection: "spam pages are pointed to by spams"

### Starting graphs (.. of an extraction task from addresses)

9	0	•	0	0	0	•	•	0	0	0	ø	•	0	Ø	Ø	Ø	0	0	0	•	•	•	0	Ø	•	•	0	Ø	0	0	0	۵	0
•	0	•	0	ø	0	\$	\$	•	6	ø	•	•	6	ø	6	6	6	6	0	ø	\$	\$	ø	\$	•	•	ø	6	ø	0	0	ø	ø
	•	•	•	•	•	\$	\$	\$	•	0	\$	\$	•	•	•	•	•	\$	0	\$	\$	\$	0	\$	\$	•	\$	\$	\$	0	0	•	•
Þ	•	•	•	0	0	•	\$	•	0	0	\$	•	0	0	0	0	0	0	0	\$	•	•	0	0	•	•	0	0	0	0	0	•	0
Þ	•	•	•	•	0	\$	\$	•	•	0	•	\$	•	0	•	0	•	•	0	\$	\$	\$	0	0	•	•	\$	\$	0	0	0	0	0
Þ	•	•	0	0	0	•	•	0	0	0	\$	•	0	0	0		0	0	0	\$	0	6	0		•	6	0	0	0	0	0	0	0
Þ	•	0	•	0	0	\$	\$	•	•	0	•	\$	0	0	•		0	•	0	\$	0		0		•		0	\$	0	0	0	0	0
Þ	•		0	0	0	•	•	•	0	0	\$	•	0	0	0			0	0	•			0		6		0	0	0	0	0	0	0
Þ	•		0	0		•	•	•		0	•	•	0	0	0			•		0			0				0		0	0	0	0	0
Þ	0		0			0	0	0			0		0	0				0		0			0				0		0		0	0	0
Þ	0					0	0	0						0									0				0		0		0	0	0
														T									Ι										

### Graph after collective edges



### Algorithms for collective inference

- Exact: intractable
- Approximate
  - Loopy Belief Propagation
    - Message passing (MP) on edges of the graph
    - [Bunescu & Mooney '04], [Sutton & McCallum '04]
  - Gibbs Sampling
    - [Finkel & Manning. '05]
  - Greedy local search (ICM)
    - [Lu & Getoor'03]
    - A special two-pass variant: [Krishnan & Manning '06]

#### Generic techniques, no guarantees

### **Associative Markov Networks**

• Graph with only associative edge potentials and node potentials  $y_{11}$   $y_{21}$   $y_{21}$   $y_{21}$   $y_{21}$   $y_{22}$ 



- Optimal for m=2. ½ approximation for m > 2
  - Min-cut with  $\alpha$ -expansion (Boykov '99)
  - LP-based metric labeling algorithms (Klienberg & Tardos '02)
  - BP with TRW-S message schedules (Kolmogorov & Wainwright, '05)

Not directly usable Slow Worse guarantees

Part III: Information Extraction and Data Integration

### **Generalized Belief propagation**

BP on clusters of cliques and chains with single node separators Clique X



- Basic MP step: Compute max-marginals for a separator node → MAP for each label of the node.
- MAP algorithms for chains  $\rightarrow$  easy and efficient.
- MAP algorithms for cliques → combinatorial algorithms can be used for this. (Gupta et al 2007)

### Graphical models in Extraction & Dedup

- State of the art: Conditional Random Fields
- IE Models
  - Basic IE model (Chain)
  - IE with collective labeling of repeated words
- De-duplication models
  - Basic pair-wise model
  - Collective de-duplication of relational data
  - Collective de-duplication of multiple networked entities

### Basic dedup problem

Given a pair of records x<sub>1</sub>, x<sub>2</sub>, predict "y" to denote if they are the same or not.

 $X_1$  Johnson Laird, Philip N. (1983). Mental models. Cambridge, Mass.: Harvard University Press.

X<sub>2</sub> P. N. Johnson-Laird. Mental Models: Towards a Cognitive Science of Language, Inference, and Consciousness. Cambridge University Press, 1983

#### • CRF: $Pr(y | x_1, x_2)$ where

- Features: list of similarity functions between record pairs.
- Graph: trivial single node graph

### Multi Attribute Similarity

		I				1			All-Ngrams*0.4 + AuthorTitleNgram*0.2
			$f_1$	$f_2$ .	f <sub>n</sub>				$\pm 0.2*$ PageMatch $3 > 0$
	Record 1 [ Record 2	D	1.0	0.4	0.2	1		No	- 0.2 ragemator - 3 > 0
	Record 1 N Record 3	N	0.0	0.1	0.3	0			Learners:
	Record 4 [ Record 5	D	0.3	0.4	0.4	1			Support Vector Machines (SVM) Logistic regression,
Unla	beled list	Ν	Map	ped	exam	ple	s	N	Linear regression, Perceptron
	Record 6		0.0	0.1	0.3	?		$\neg$	
	Record 7		1.0	0.4	0.2	?		$\langle \rangle$	
	Record 8		0.6	0.2	0.5	?			0.7 0.1 0.6 0
	Record 9		0.7	0.1	0.6	?			0.3 0.4 0.4 1
	Record 10		0.3	0.4	0.4	?			$0.0 \ 0.1 \ \dots \ 0.1 \ 0$
I	Record 11		0.0	0.1	0.1	?			0.3 0.8 0.1 1
	P	1	0.3	0.8	0.1	?			0.6 0.1 0.5   1
			0.6	0.1	0.5	?			

Part III: Information Extraction and Data Integration

### Graphical models in Extraction & Dedup

- State of the art: Conditional Random Fields
- IE Models
  - Basic IE model (Chain)
  - IE with collective labeling of repeated words
- De-duplication models
  - Basic pair-wise model
  - Collective de-duplication of relational data
  - Collective de-duplication of set-oriented data
# **De-duplication of relational records**

Collectively de-duplicate entities and its many attributes

	a <sup>1</sup>	a <sup>2</sup>	<b>— a</b> <sup>3</sup> <b>—</b>
Record	Title	Author	Venue
$b_1$	"Record Linkage using CRFs"	"Linda Stewart"	"KDD-2003"
$b_2$	"Record Linkage using CRFs"	"Linda Stewart"	"9th SIGKDD"
$b_3$	"Learning Boolean Formulas"	"Bill Johnson"	"KDD-2003"
$b_4$	"Learning of Boolean Expressions"	"William Johnson"	"9th SIGKDD"

Associate variables for predictions for each attribute k each record pair (i,j)  $A_{ij}^{k}$ 

for each record pair

from Parag & Domingos 2005

R<sub>ii</sub>

## **Graphical model**



### **Potentials**

- Independent scores
  - s<sub>k</sub>(A<sup>k</sup>,a<sub>i</sub>,a<sub>j</sub>) Attribute-level
    - Any classifier on various text similarities of attribute pairs
  - s(R,b<sub>i</sub>,b<sub>j</sub>) Record-level
    - Any classifier on various similarities of all k attribute pairs
- Dependency scores
  - d<sub>k</sub>(A<sup>k</sup>, R): record pair, attribute pair

	0	1
0	4	2
1	1	7

## Joint de-duplication steps

- Jointly pick 0/1 labels for all record pairs Rij and all K attribute pairs A<sup>k</sup><sub>ij</sub> to maximize sum of potentials
- Typical graphical model inference problem
- Efficient algorithm possible because of special forms of potentials
  - dependency scores associative
  - $dk(1,1) + dk(0,0) \ge dk(1,0) + dk(0,1)$

### Graphical models in Extraction & Dedup

- State of the art: Conditional Random Fields
- IE Models
  - Basic IE model (Chain)
  - IE with collective labeling of repeated words
- De-duplication models
  - Basic pair-wise model
  - Collective de-duplication of relational data
  - Collective de-duplication of set-oriented data

## Collective linkage: set-oriented data

P1	D White, J Liu, A Gupta	A Gupta J Cupta
P2	Liu, Jane & J Gupta & White, Don	D White White, Don
P3	Anup Gupta	A Gupta
P4	David White	David White D White D White

#### **Scoring functions**

- S(A<sub>ii</sub>) Attribute-level
  - Text similarity
- S(A<sub>ij</sub>, N<sub>ij</sub>) Dependency with labels of co-author set
  - Fraction of co-author set assigned label 1.
- Score: α s(A<sub>ij</sub>) + (1-α) s(A<sub>ij</sub>, N<sub>ij</sub>)

#### Inference Algorithm

- Exact inference hard
  - MCMC algorithm in (Bhattacharya and Getoor, 2007)

# **Concluding remarks**

- Graphical models provide a unified and flexible modeling of many extraction and integration tasks
- Much work is still needed in converting these to methods of choice in commercial systems
  - Scalable algorithms
  - Skillful integration of manual rules with statistical methods
  - Feedback on when the statistical method failes
  - Robust feature design so as to not overfit on the training data.