



PROBKB Web-Scale Probabilistic Knowledge Base

Yang Chen, Xing Liu

{yang,xinliu}@cise.ufl.edu

Computer and Information Science and Engineering
University of Florida

Mar 12, 2013

Outline

Introduction

Introduction

The PROKB System

PROKB Architecture

Grounding

Inference

Discussion

Discussion

Knowledge bases–Introduction

- A *knowledge base* is a collection of entities, facts, and relationships that conforms with a certain data model.
- A knowledge base helps machines understand humans, languages, and the world.

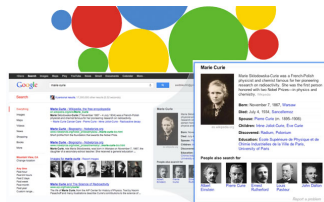
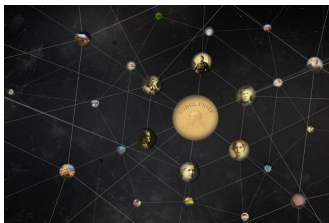
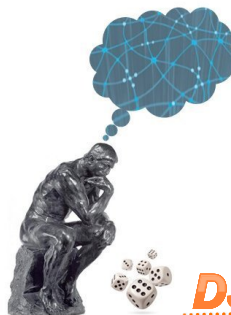


Figure: Google knowledge graph

Challenges & Motivation

Knowledge Acquisition

- **Statistical Inference**
 - **Markov logic**
- Information extraction
 - NELL (CMU), OpenIE (UW)
 - Entities, relations, rules
- Human collaboration
 - Wikipedia
 - Freebase



Challenges & Motivation

Uncertainty Management

- **Statistical Inference**
 - **Probabilistic graphical models**
 - **Markov chain Monte Carlo**
- Data integration
 - Merging multiple data sources
 - Crowdsourcing/user feedback
- Data cleaning
 - Conflict, incomplete, outdated data



Challenges & Motivation

Scalability

- **Scalable data management systems**
 - **Relational DBMS**
 - Hadoop
 - Spark, GraphLab, **Datapath**, etc
- Scalable algorithms
 - Incremental inference
 - Query-driven inference



Outline

Introduction

Introduction

The PROKB System

PROKB Architecture

Grounding

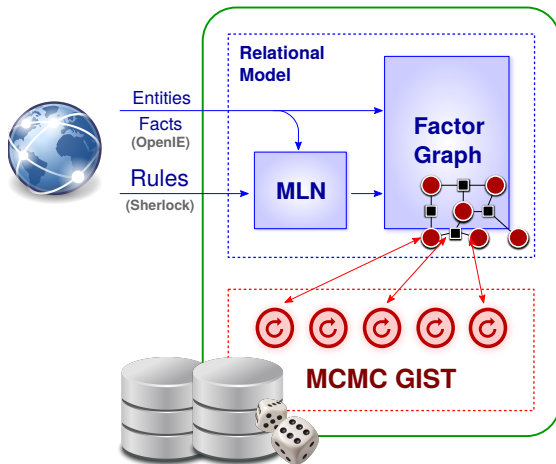
Inference

Discussion

Discussion

Architecture

PROBabilistic Knowledge Base



Markov Logic–Probabilistic Inference Framework

A *Markov logic network* (MLN) is a set of formulae with weights. Together with a finite set of constants $C = \{c_1, \dots, c_{|C|}\}$, it defines a Markov network.

Weight	First-Order Logic
0.7	$\text{Fr}(x, y) \wedge \text{Fr}(y, z) \rightarrow \text{Fr}(x, z)$
1.5	$\text{Sm}(x) \rightarrow \text{Ca}(x)$
1.1	$\text{Fr}(x, y) \wedge \text{Sm}(x) \rightarrow \text{Sm}(y)$

A set of constants (entities, or objects)

Table: Example Markov logic network.

$$C = \{A, B\}.$$

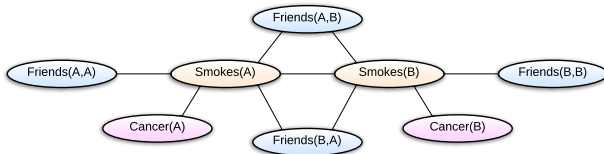


Figure: Grounded Markov network.

Outline

Introduction

Introduction

The PROKB System

PROKB Architecture

Grounding

Inference

Discussion

Discussion

Grounding

Grounding is the process of substituting constants into MLN clauses.

The result of grounding is a *factor graph* (or *Markov network*) from which we can infer marginal probabilities for individual facts.

Key Challenges

- Time-consuming, especially if the numbers of rules and entites are large.
- Grounded network has an intractably large size, making inference tasks slow.

Markov Logic: A Relational Point of View

- State-of-the-art (TUFFY, NELL): one table for each relation
- By considering only Horn clauses, we store the rules and relationships in a few tables:

Table: MLN (M)

head	body1	body2
p_1	q_1	r_1
p_2	q_2	r_2
p_3	q_3	r_3
p_4	q_4	r_4
	...	

Table: Relationships (R)

pred	ent1	ent2
p_1	x_1	y_1
p_1	x_2	y_2
p_2	x_1	y_1
p_2	x_2	y_2
	...	

Markov Logic: A Relational Pointer of View

The grounding of *ALL* rules of form

$$p(x, y) \leftarrow q(x, z), r(z, y)$$

is then expressed as a relational operation:

$$\begin{aligned} R &\leftarrow \rho_R(\text{pred}, \text{ent1}, \text{ent2})(\pi_{M.\text{head}, R_2.\text{ent1}, R_3.\text{ent2}} \quad (1) \\ &\quad ((M \bowtie_{M.\text{body1}=R_2.\text{pred}} R_2) \\ &\quad \bowtie_{M.\text{body2}=R_3.\text{pred} \text{ AND } R_2.\text{ent2}=R_3.\text{ent1}} R_3)) \\ G &\leftarrow R_1 \bowtie R \end{aligned}$$

Grounding Results

The relational Markov logic model saves us from managing thousands of tables as in previous approaches. As a result, We grounded the whole SHERLOCK-HOLMES dataset in about 1.5 hours using Greenplum, while the state-of-the-art implementation MLN crashes during its grouding phase.

#relations	10,672
#rules	31,000
#constants	1.1M
#evidence	250,000
#queries	10,672
TUFFY	Crash
PROKB	1.5 h

Table: Dataset statistics and performance.

Outline

Introduction

Introduction

The PROBKB System

PROBKB Architecture

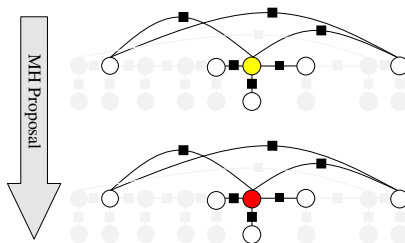
Grounding

Inference

Discussion

Discussion

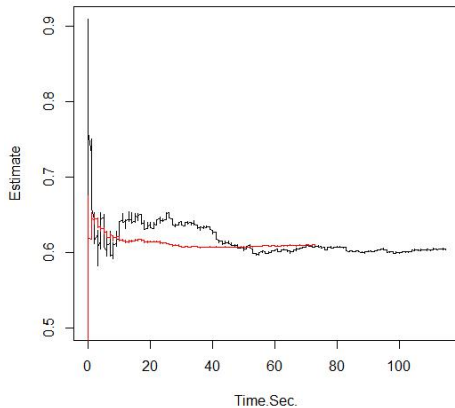
Inference: MCMC-MH



Markov locality property allows for parallel computing.

Mar 12, 2013

Preliminary Results





Preliminary Results

# Vertices	5000	10,000	25,000	250,000	2,500,000
Single Thread	7 sec	28 sec	228 sec	Hours	N/A
Datapath	7 sec	13 sec	30 sec	2661 sec	N/A

Table: Time to generate 5000 joint samples for different graphs.



Outline

Introduction

Introduction

The PROKB System

PROKB Architecture

Grounding

Inference

Discussion

Discussion



Next Steps

- Grounding/inference connection
 - Construct graph structures from relational format.
 - Partition and Merge: opportunity of GLA parallelism.
- DB-memory synchronization
 - Determine an update and write-back policy to synchronize DB and in-memory data.
 - Buffer manager; IO efficiency
- Incremental inference
 - Schedule Datapath execution so that computation is focused on the least convergent portion of the factor graph.
- Evaluation
 - MCMC evaluation: multi-chain convergence test
 - Result evaluation: Manual check (AMT)



Responsibilities

Yang Grounding

Xing Inference



Questions?

Thank you!