



PROKB Web-Scale Probabilistic Knowledge Base

Yang Chen, Xing Liu
{yang,xinliu}@cise.ufl.edu

Computer and Information Science and Engineering
University of Florida

Apr 9, 2013

Outline

Introduction

Introduction

The PROBKB System

PROBKB Architecture

Grounding

Inference

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.

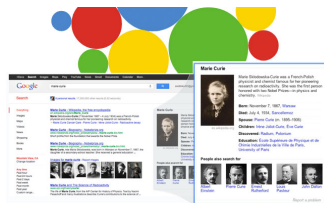
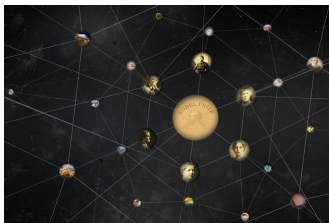
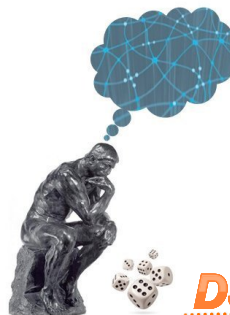


Figure: Google knowledge graph

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 - Conflict, incomplete, outdated data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 PROBKB System

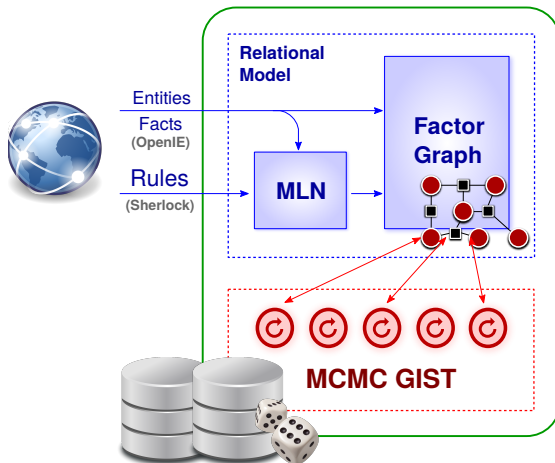
PROBKB Architecture

Grounding

Inference

Architecture

PROBabilistic Knowledge Base



Outline

Introduction

Introduction

The PROBKB System

PROBKB Architecture

Grounding

Inference

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 SHERLOCK-HOLMES dataset in 85 seconds using Greenplum, while the state-of-the-art implementation MLN crashes during its grouding phase.

#entities	480K
#relations	10,672
#relationships	5K
#rules	31,000
TUFFY	Crash
PROKB	85s

Table: Dataset statistics and performance.

Evaluation

- We learned 100K facts from the original 5K.
- Not all results are correct; errors propagate.



- Errors come from word ambiguity, incorrect extractions and rules.
- Need pruning and re-evaluating.

Outline

Introduction

Introduction

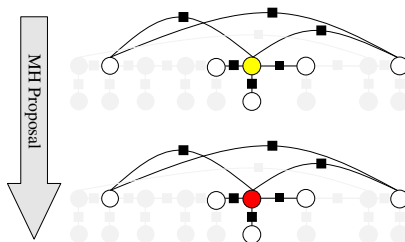
The PROBKB System

PROBKB Architecture

Grounding

Inference

Inference: MCMC-MH



Markov locality property allows for parallel computing.

Datapath GIST

Generalized Iterable State Transforms (GIST)

- GIST Performs *transitions* upon a *state* until that state has converged to the desired result.
 - Transition* MCMC Proposal function.
 - State* Factor graph with its samples.
- A user-defined local scheduler allows general MCMC proposal implementation.
- The GIST state keeps track of the inference result.

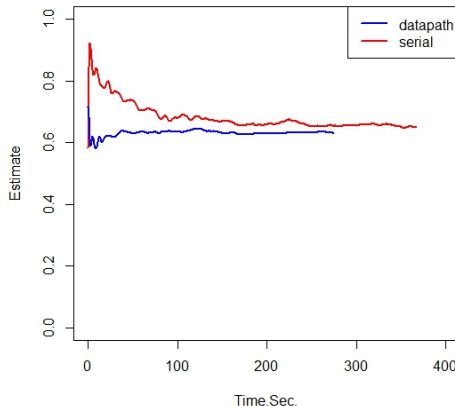
Datapath GIST

Parallel Mcmc implementation modification on datapath

- Previous approach: cut the graph to smaller graphs, run mcmc on each partition.
 - Lost information when a factor is accross multiple subgraphs. Inaccurate, but the faster performance.
- Modified approach: Do not cut the graph. Add write lock on each variables.
 - More accurate, but lower performance.

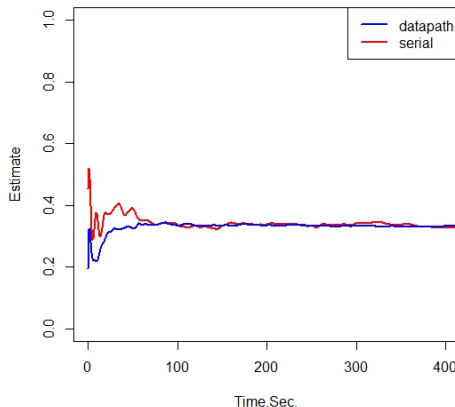
Parallel mcmc with write lock in datapath

50,000 vertices graph, parallel mcmc with write lock



Parallel mcmc with write lock in datapath

50,000 vertices graph, parallel mcmc with write lock



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.



Responsibilities

Yang Grounding

Xing Inference



Questions?

Thank you!