### CSE544: Principles of Database Systems

Review of MapReduce Datalog (cont'd)

#### Announcements

• Review 7 (Datalog) was due yesterday

- Project Milestone due next Friday
- HW3 is due the following Monday

### MapReduce Review

- What is the map function?
- What is the reduce function?
- What is a map task?
- What is a reduce task?
- What is a mapreduce job?

### MapReduce Review

- What is the map function?
  - Takes (k,v) returns {(k1,v1),(k2,v2),...}
- What is the reduce function?

– Takes (k,{v1,v2,...}) returns any result

• What is a map task?

- A set of (k,v) pairs that are scheduled as a unit

• What is a reduce task?

– A set of (k,{v1,...}) pairs scheduled as a unit

• What is a mapreduce job?

- The entire map reduce program

# Anatomy of a Query Execution

- Running Part B of HW2
- 20 nodes = 1 master + 19 workers

- Using PARALLEL 50
- Let's see what happened









3h 51min

#### 3h 52min

Kind	% Complete	Num Tasks	Pending	Running	Complete	Killed	Failed/Killed Task Attempts
map	100.00%	15816	0	0	<u>15816</u>	0	0 / <u>18</u>
reduce	37.72%	50	<u>19</u>	22	<u>9</u>	0	0/0

Kind	% Complete	Num Tasks	Pending	Running	Complete	Killed	Failed/Killed Task Attempts
<u>map</u>	100.00%	15816	0	0	<u>15816</u>	0	0 / <u>18</u>
reduce	42.35%	50	11	20	<u>19</u>	0	0/0









#### 4h 18min

Several servers failed: "fetch error". Their map tasks need to be rerun. All reducers are waiting....

Kind	% Complete	Num Tasks	Pending	Running	Complete	Killed	Failed/Killed Task Attempts
map	99.88%	15816	<u>2638</u>	<u>30</u>	<u>13148</u>	0	<u>15/3337</u>
reduce	48.42%	50	15	<u>16</u>	<u>19</u>	0	070



uce Completion Graph - close



#### 7h 10min Mappers finished, reducers resumed. Failed/Killed Pending Running Complete Kind % Complete Num Tasks Killed Task Attempts 100.00% 15816 0 <u>15816</u> 0 26 / 5968 0 map 94.15% 50 0 0 reduce 6 44 0/8





#### 7h 20min

#### Success! 7hrs, 20mins.

#### Hadoop job\_201203041905\_0001 on ip-10-203-30-146

User: hadoop Job Name: PigLatin:DefaultJobName Job File: hdfs://10.203.30.146:9000/mnt/var/lib/hadoop/tmp/mapred/staging/hadoop/.staging/job\_201203041905\_0001/job.xml Submit Host: ip-10-203-30-146.ec2/internal Submit Host Address: 10.203.30.146 Job-ACLs: All users are allowed Job Setup: Successful Status: Succeeded Started at: Sun Mar 04 19:08:29 UTC 2012 Finished at: Mon Mar 05 02:28:39 UTC 2012 Finished at: Successful Black-listed TaskTrackers:2

Kind	% Complete	Num Tasks	Pending	Running	Complete	Killed	Failed/Killed Task Attempts
<u>map</u>	100.00%	15816	0	0	<u>15816</u>	0	<u>26</u> / <u>5968</u>
reduce	100.00%	50	0	0	<u>50</u>	0	0 / <u>14</u>



#### Degree Sequence

### Famous Example of Big Data Analysis

Kumar et al., The Web as a Graph

Question 1: is the Web like a "random graph"?

Question 2: how does the Web graph look like?

# Graph as Databases

Many large databases are graphs

- The Web
- The Internet
- Social Networks
- Flights btw. Airports
- Etc,etc,etc



Source	Target
а	b
b	а
а	f
b	f
b	е
b	d
d	е
d	с
е	g
g	С
с	g

# Data Analytics on Big Graphs

Queries expressible in SQL:

- How many nodes (edges)?
- How many nodes have > 4 neighbors?
- Which are the "most connected nodes"?

Queries requiring recursion:

- Is the graph connected?
- What is the diameter of the graph?
- Compute <u>PageRank</u>
- Compute the <u>Centrality</u> of each node



Source	Target
а	b
b	а
а	f
b	f
b	е
b	d
d	е
d	с
е	g
g	С
С	g

### Histogram of a Graph a.k.a. Degree Sequence

- Outdegree of a node = number of outgoing edges
- For each d, let n(d) = number of nodes with oudegree d
- The outdegree histogram of a graph = the scatterplot (d, n(d))





#### Histograms Tell Us Something About the Graph



# **Exponential Distribution**

- $n(d) \approx c/2^d$  (generally,  $cx^d$ , for some x < 1)
- A random graph has exponential distribution
- Best seen when n is on a log scale



# **Zipf Distribution**

- $n(d) \approx 1/d^x$ , for some value x>0
- Human-generated data has Zipf distribution: letters in alphabet, words in vocabulary, etc.
- Best seen in a log-log scale



# The Histogram of the Web



Figure 2: In-degree distribution.

#### The Bowtie Structure of the Web



Figure 4: The web as a bowtie. SCC is a giant strongly connected component. IN consists of pages with paths to SCC, but no path from SCC. OUT consists of pages with paths from SCC, but no path to SCC. TENDRILS consists of pages that cannot surf to SCC, and which cannot be reached by surfing from SCC.

### Review

- What is datalog?
- What is the naïve evaluation algorithm?
- What is the seminaive algorithm?

### **Discussion in Class**

Assume a linear graph:

$$0 \longrightarrow 1 \longrightarrow 2 \longrightarrow 3 \longrightarrow ... \longrightarrow n$$

- How many iterations are needed by each algorithm?
- How many times is the tuple T(0,n) discovered by each algorithm?

**Right linear TC** 

Non-linear TC

T(x,y) := R(x,y)T(x,y) := T(x,z), T(z,y)

### **Discussion in Class**

The *Declarative Imperative* paper:

- What are the extensions to datalog in Dedalus?
- What is the main usage of Dedalus described in the paper? Discuss some applications, discuss what's missing.

#### Semantics of a Datalog Program

Three different, equivalent semantics:

• Minimal model semantics

Least fixpoint semantics

• Proof-theoretic semantics

To each rule r:  $P(x_1...x_k) := R_1(...), R_2(...), ...$ 

To each rule r:  $P(x_1...x_k) := R_1(...), R_2(...), ...$ All variables in the rule Associate the logical sentence  $\Sigma_r$ :  $\forall z_1...\forall z_n$ .  $[(R_1(...)\Lambda R_2(...)\Lambda ...) \rightarrow P(...)]$ 







 $\equiv \forall x. \forall y. (\exists z.R(x,z) \land T(z,y) \rightarrow T(x,y))$ 

**<u>Definition</u>**. A pair (I,J) where I is an EDB and J is an IDB is a *model* for P, if  $(I,J) \models \Sigma_P$ 

**Definition**. Given an EDB database instance I and a datalog program P, the minimal model, denoted J = P(I) is a minimal database instance J s.t.  $(I,J) \models \Sigma_P$ 

**Theorem**. The minimal model always exists, and is unique.

**<u>Definition</u>**. A pair (I,J) where I is an EDB and J is an IDB is a *model* for P, if  $(I,J) \models \Sigma_P$ 

**<u>Definition</u>**. Given an EDB database instance I and a datalog program P, the minimal model, denoted J = P(I) is a minimal database instance J s.t.  $(I,J) \models \Sigma_P$ 

**Theorem**. The minimal model always exists, and is unique.

#### Example:

$$1 \longrightarrow 2 \longrightarrow 3 \longrightarrow 4 \longrightarrow 5$$

Which of these IDBs are *models*? Which are *minimal models*?

R=	1	2	
	2	3	
	3	4	
	4	5	

T=	
1	2
2	3
3	4
4	5
1	3
2	4
3	5

**<u>Definition</u>**. A pair (I,J) where I is an EDB and J is an IDB is a *model* for P, if  $(I,J) \models \Sigma_P$ 

<u>**Definition**</u>. Given an EDB database instance I and a datalog program P, the minimal model, denoted J = P(I) is a minimal database instance J s.t.  $(I,J) \models \Sigma_P$ 

**Theorem**. The minimal model always exists, and is unique.

#### Example:

$$1 \longrightarrow 2 \longrightarrow 3 \longrightarrow 4 \longrightarrow 5$$

Which of these IDBs are *models*? Which are *minimal models*?

R=	1	2	
	2	3	
	3	4	
	4	5	

1=	
1	2
2	3
3	4
4	5
1	3
2	4
3	5

T=	-
1	2
2	3
3	4
4	5
1	3
2	4
3	5
1	4
2	5
1	5

**<u>Definition</u>**. A pair (I,J) where I is an EDB and J is an IDB is a *model* for P, if  $(I,J) \models \Sigma_P$ 

<u>**Definition**</u>. Given an EDB database instance I and a datalog program P, the minimal model, denoted J = P(I) is a minimal database instance J s.t.  $(I,J) \models \Sigma_P$ 

**Theorem**. The minimal model always exists, and is unique.

#### Example:

$$1 \longrightarrow 2 \longrightarrow 3 \longrightarrow 4 \longrightarrow 5$$

Which of these IDBs are *models*? Which are *minimal models*?

R=	1	2	
	2	3	
	3	4	
	4	5	

1=	
1	2
2	3
3	4
4	5
1	3
2	4
3	5

T=	
1	2
2	3
3	4
4	5
1	3
2	4
3	5
1	4
2	5
1	5

Т=

1	1
1	2
1	3
1	4
1	5
5	4
5	5

All 25 pairs of nodes

# **Minimal Fixpoint Semantics**

<u>Definition</u>. Fix an EDB I, and a datalog program **P**. The <u>immediate consequence</u> operator  $T_P$  is defined as follows. For any IDB J:  $T_P(J) = all IDB$  facts that are immediate consequences from I and J.

<u>**Fact</u></u>. For any datalog program P, the immediate consequence operator is monotone. In other words, if J\_1 \subseteq J\_2 then T\_P(J\_1) \subseteq T\_P(J\_2).</u>** 

# **Minimal Fixpoint Semantics**

<u>**Definition**</u>. Fix an EDB I, and a datalog program **P**. The <u>immediate consequence</u> operator  $T_P$  is defined as follows. For any IDB J:  $T_P(J) = all IDB$  facts that are immediate consequences from I and J.

<u>**Fact</u></u>. For any datalog program P, the immediate consequence operator is monotone. In other words, if J\_1 \subseteq J\_2 then T\_P(J\_1) \subseteq T\_P(J\_2).</u>** 

<u>**Theorem</u></u>. The immediate consequence operator has a unique, minimal fixpoint J: fix(T\_P) = J, where J is the minimal instance with the property T\_P(J) = J.</u>** 

Proof: using Knaster-Tarski's theorem for monotone functions. The fixpoint is given by:

fix  $(T_P) = J_0 \cup J_1 \cup J_2 \cup \dots$  where  $J_0 = \emptyset$ ,  $J_{k+1} = T_P(J_k)$ 

#### **Minimal Fixpoint Semantics**



T(x,y) :- R(x,y) T(x,y) :- R(x,z), T(z,y)

#### R=

1	2
2	3
3	4
4	5



ļ	J <sub>1</sub> = <sup>-</sup>	Γ <sub>Ρ</sub> (J <sub>0</sub>	))
	1	2	
	2	3	
	3	4	
	4	5	

_	J <sub>2</sub> =	T <sub>P</sub> (J	1)
	1	2	
	2	3	
	3	4	
	4	5	
	1	3	
	2	4	
	3	5	

•	$J_3 = -$	T <sub>P</sub> (J <sub>2</sub>	<u>,</u> )
	1	2	
	2	3	
	3	4	
	4	5	
	1	3	
	2	4	
	3	5	
	1	4	
	2	5	

J <sub>4</sub> =	$T_{P}(J_3)$
------------------	--------------

1	2
2	3
3	4
4	5
1	3
2	4
3	5
1	4
2	5
1	5

### **Proof Theoretic Semantics**

Every fact in the IDB has a *derivation tree*, or *proof tree* justifying its existence.



# Adding Negation: Datalog<sup>¬</sup>

**Example**: compute the complement of the transitive closure

What does this mean??

#### Recursion and Negation Don't Like Each Other

EDB:  $I = \{ R(a) \}$ 

Which IDBs are models of **P**?

$$J_1 = \{ \}$$
  $J_2 = \{S(a)\}$   $J_3 = \{T(a)\}$   $J_4 = \{S(a), T(a)\}$ 

Recursion and Negation Don't Like Each Other

EDB:  $I = \{ R(a) \}$ 

Which IDBs are models of **P**?



There is no *minimal* model!

Recursion and Negation Don't Like Each Other

EDB:  $I = \{ R(a) \}$ 

Which IDBs are models of **P**?



There is no *minimal* model!

There is no minimal fixpoint! (Why does Knaster-Tarski's theorem fail?)

# Adding Negation: datalog<sup>¬</sup>

- Solution 1: Stratified Datalog<sup>¬</sup>
  - Insist that the program be <u>stratified</u>: rules are partitioned into strata, and an IDB predicate that occurs only in strata ≤ k may be negated in strata ≥ k+1
- Solution 2: Inflationary-fixpoint Datalog<sup>¬</sup>
  - Compute the fixpoint of J  $\cup$  T<sub>P</sub>(J)
  - Always terminates (why ?)
- Solution 3: Partial-fixpoint Datalog<sup>-,\*</sup>
  - Compute the fixpoint of  $T_P(J)$
  - May not terminate

### Stratified datalog<sup>¬</sup>

A datalog<sup>¬</sup> program is <u>stratified</u> if its rules can be partitioned into k strata, such that:
If an IDB predicate P appears negated in a rule in stratum i, then it can only appear in the head of a rule in strata 1, 2, ..., i-1



Note: a datalog<sup>¬</sup> program either is stratified or it ain't!

Which programs are stratified?

 $\begin{array}{l} T(x,y) := R(x,y) \\ T(x,y) := T(x,z), \ R(z,y) \\ CT(x,y) := Node(x), \ Node(y), \ not \ T(x,y) \end{array}$ 

S(x) :- R(x), not T(x) T(x) :- R(x), not S(x)

# Stratified datalog<sup>¬</sup>

• Evaluation algorithm for stratified datalog<sup>-</sup>:

- For each stratum i = 1, 2, ..., do:
  - Treat all IDB's defined in prior strata as EBS
  - Evaluate the IDB's defined in stratum i, using either the naïve or the semi-naïve algorithm

Does this compute a minimal model?

T(x,y) := R(x,y)T(x,y) := T(x,z), R(z,y)

CT(x,y) :- Node(x), Node(y), not T(x,y)

# Stratified datalog<sup>¬</sup>

• Evaluation algorithm for stratified datalog<sup>-</sup>:

- For each stratum i = 1, 2, ..., do:
  - Treat all IDB's defined in prior strata as EBS
  - Evaluate the IDB's defined in stratum i, using either the naïve or the semi-naïve algorithm

Does this compute a minimal model?	T(x,y) :- R(x,y) T(x,y) :- T(x,z), R(z,y)
NO: J <sub>1</sub> = { T = transitive closure, CT = its complement} J <sub>2</sub> = { T = all pairs of nodes, CT = empty}	CT(x,y) :- Node(x), Node(y), not T(x,y)

# Inflationary-fixpoint datalog<sup>¬</sup>

Let **P** be any datalog<sup>¬</sup> program, and I an EDB. Let  $T_P(J)$  be the <u>immediate consequence</u> operator. Let  $F(J) = J \cup T_P(J)$  be the <u>inflationary immediate consequence</u> operator.

Define the sequence:  $J_0 = \emptyset$ ,  $J_{n+1} = F(J_n)$ , for  $n \ge 0$ .

<u>**Definition**</u>. The inflationary fixpoint semantics of **P** is  $J = J_n$  where n is such that  $J_{n+1} = J_n$ 

Why does there always exists an n such that  $J_n = F(J_n)$ ?

Find the inflationary semantics for:

T(x,y) := R(x,y) T(x,y) := T(x,z), R(z,y)CT(x,y) := Node(x), Node(y), not T(x,y)



# Inflationary-fixpoint datalog<sup>¬</sup>

- Evaluation for Inflationary-fixpoint datalog<sup>¬</sup>
- Use the naïve, of the semi-naïve algorithm

 Inhibit any optimization that rely on monotonicity (e.g. out of order execution)

# Partial-fixpoint datalog<sup>,\*</sup>

Let **P** be any datalog<sup>¬</sup> program, and I an EDB. Let  $T_P(J)$  be the *immediate consequence* operator.

Define the sequence:  $J_0 = \emptyset$ ,  $J_{n+1} = T_P(J_n)$ , for  $n \ge 0$ .

<u>**Definition**</u>. The partial fixpoint semantics of **P** is  $J = J_n$  where n is such that  $J_{n+1} = J_n$ , if such an n exists, undefined otherwise.

Find the partial fixpoint semantics for:

Note: there may not exists an n such that  $J_n = F(J_n)$ 

T(x,y) := R(x,y) T(x,y) := T(x,z), R(z,y)CT(x,y) := Node(x), Node(y), not T(x,y)



#### Discussion

• Which semantics does Daedalus adopt?

#### Discussion

Comparing datalog<sup>¬</sup>

- Compute the complement of the transitive closure in inflationary datalog<sup>¬</sup>
- Compare the expressive power of:
  - Stratified datalog<sup>¬</sup>
  - Inflationary fixpoint datalog<sup>¬</sup>
  - Partial fixpoint datalog