CSE544

Data Provenance

Project Presentations

Friday, June 7th, 9:30-2:30, in CSE 405

What to include:

- Describe the problem:
 - why is it important, why is it non-trivial
- Overview prior approaches,
 - related work
- Your approach
- Your results
 - theoretical, empirical, experimental
- Discuss their significance
 - do they work ? do they solve the problem you set out to do ? do they improve over existing work ?
- Conclusions

Rule of thumb: 1 slide / minute, less slack. 15' → 12 slides.

Outline

Sources:

- Karnouvarakis et al., *Provenance Semirings*, PODS 2007
- Cheney, Chiticariu, Tan, *Provenance in Databases: Why, How, and Where*, 2007
- Tannen, Tutorial on Provenance in EDBT 2010

Data Provenance

Cheney, Chiticariu, Tan, *Provenance in Databases: Why, How, and Where*, 2007

 Provenance information describes the origins and the history of data in its life cycle. Such information (also called lineage) is important to many data management tasks.

Data Provenance

Provenance inside the DBMS
– Will discuss today

- Provenance outside of the DBMS
 - Much more messy; there is a standard, OPM (Open Provenance Model)

Provenance Annotations

- Some query produces an output table T(A,B,C)
- We store it over some period of time
- Later we ask: "where did this tuple come from?"
- The "provenance annotation" answers this.

Α	В	С
a1	b1	c1
a2	b1	c1
a2	b2	c2
a2	b2	c3

provenance1 provenance2 provenance3 provenance4

Provenance Annotations

- Start by annotating each tuple in the original database with a unique identifier; can be the Tuple Id (TID)
- Next, compute the provenance expression inductively, based on the query plan



Join Operator



Projection Operator



Union Operator



Α	В	
a1	b1	X1
a2	b2	X2+Y1
a3	b3	X3

Selection Operator



Α	В	
a1	b1	X1
a1	b2	X2

We could simply remove the tuples filtered out. But it's better to keep them around (we'll see why). What is their annotation?

Selection Operator



Α	В	
a1	b1	X1·1
a1	b2	X2 · 1
a2	b1	X3·0
a2	b2	X4 · 0
a2	b3	X5·0

We could simply remove the tuples filtered out. But it's better to keep them around (we'll see why). What is their annotation?

Complex Example

$\sigma_{\mathsf{C=e}} \Pi_{\mathsf{AC}}(\ \Pi_{\mathsf{AC}}(\mathsf{R}) \bowtie \Pi_{\mathsf{BC}}(\mathsf{R}) \cup \Pi_{\mathsf{AB}}(\mathsf{R}) \bowtie \Pi_{\mathsf{BC}}(\mathsf{R})) =$

R =

Α	B	С	
а	b	С	X
d	b	е	Y
f	g	е	Z

Α	С	
а	С	$(X \cdot X + X \cdot X) \cdot 0 = 2 \cdot X^2$
а	е	$X \cdot Y \cdot 1 = X \cdot Y$
d	С	$Y \cdot X \cdot 0 = 0$
d	е	$(\mathbf{Y} \cdot \mathbf{Y} + \mathbf{Y} \cdot \mathbf{Z} + \mathbf{Y} \cdot \mathbf{Y}) \cdot 1 = 2 \cdot \mathbf{Y}^2 + \mathbf{Y} \cdot \mathbf{Z}$
f	е	$(Z \cdot Z + Z \cdot Y + Z \cdot Z) \cdot 1 = 2 \cdot Z^2 + Y \cdot Z$

Discuss in class what these annotations mean

K-Relations

Definition. A K-relation is a relation where each tuple is annotated with an element from the set K.

What we have described so far is an extension of the positive operations of the relational algebra to K-relations

We assumed that K has the operators +, ·

The problem:

- We have defined provenance for a query plan P
- Given a query Q, we want the provenance to be independent of the plan
- Needed: if P1=P2, then provenance(P1) = Provenance(P2)

Definition. A structure $(K, +, \cdot, 0, 1)$ is called a commutative semiring if:

- 1. (K,+,0) is a commutative monoid:
 - a. + is associative: (x+y)+z=x+(y+z)
 - b. + is commutative: x+y=y+x
 - c. 0 is the identity for +: x+0=0+x=x
- 2. $(K, \cdot, 1)$ is a commutative monoid:
 - a. ... (similar identities)
- 3. distributes over +: $x \cdot (y+z) = x \cdot y + x \cdot z$

4. For all x: $x \cdot 0 = 0 \cdot x = 0$

Definition. A structure (K, +, ·, 0, 1) is called a commutative semiring if:

- 1. (K,+,0) is a commutative monoid:
 - a. + is associative: (x+y)+z=x+(y+z)
 - b. + is commutative: x+y=y+x
 - c. 0 is the identity for +: x+0=0+x=x
- 2. $(K, \cdot, 1)$ is a commutative monoid:
 - a. ... (similar identities)
- 3. distributes over +: $x \cdot (y+z) = x \cdot y + x \cdot z$
- 4. For all x: $x \cdot 0 = 0 \cdot x = 0$

<u>Theorem</u>. The standard identities of the Bag algebra hold for K-relations iff $(K, +, \cdot, 0, 1)$ is a commutative semiring.

Discuss in class:

$$q(x,u) := R(x,y), S(y,z), T(z,u)$$

Given two plans, why are the annotations equal?

$\sigma_{\mathsf{C=e}} \Pi_{\mathsf{AC}}(\ \Pi_{\mathsf{AC}}(\mathsf{R}) \bowtie \Pi_{\mathsf{BC}}(\mathsf{R}) \cup \Pi_{\mathsf{AB}}(\mathsf{R}) \bowtie \Pi_{\mathsf{BC}}(\mathsf{R})) =$



ACac
$$2 \cdot X^2$$
ae $X \cdot Y$ de $2 \cdot Y^2 + Y \cdot Z$ fe $2 \cdot Z^2 + Y \cdot Z$

Q: Suppose we delete the tuple (d,b,e) from R. Which tuple(s) disappear from the result?

$\sigma_{\mathsf{C=e}} \Pi_{\mathsf{AC}}(\ \Pi_{\mathsf{AC}}(\mathsf{R}) \bowtie \Pi_{\mathsf{BC}}(\mathsf{R}) \cup \Pi_{\mathsf{AB}}(\mathsf{R}) \bowtie \Pi_{\mathsf{BC}}(\mathsf{R})) =$



	D		
а	b	С	Х
d	b	е	Y
f	g	е	Z

ACac
$$2 \cdot X^2$$
ae $X \cdot Y$ de $2 \cdot Y^2 + Y \cdot Z$ fe $2 \cdot Z^2 + Y \cdot Z$



A: Set Y=0

Q: Suppose we delete the tuple (d,b,e) from R. Which tuple(s) disappear from the result?

$\sigma_{\mathsf{C=e}} \Pi_{\mathsf{AC}}(\ \Pi_{\mathsf{AC}}(\mathsf{R}) \bowtie \Pi_{\mathsf{BC}}(\mathsf{R}) \cup \Pi_{\mathsf{AB}}(\mathsf{R}) \bowtie \Pi_{\mathsf{BC}}(\mathsf{R})) =$



ACac
$$2 \cdot X^2$$
ae $X \cdot Y$ de $2 \cdot Y^2 + Y \cdot Z$ fe $2 \cdot Z^2 + Y \cdot Z$

Q: Suppose each tuple in R occurs 3 times (bag semantics). How many times occurs each tuple in the answer?

$\sigma_{\mathsf{C=e}} \Pi_{\mathsf{AC}}(\ \Pi_{\mathsf{AC}}(\mathsf{R}) \bowtie \Pi_{\mathsf{BC}}(\mathsf{R}) \cup \Pi_{\mathsf{AB}}(\mathsf{R}) \bowtie \Pi_{\mathsf{BC}}(\mathsf{R})) =$



Q: Suppose each tuple in R occurs 3 times (bag semantics). How many times occurs each tuple in the answer?

A. Set X=Y=Z=3

Sets of Contributing Tuples

$\sigma_{\mathsf{C=e}} \Pi_{\mathsf{AC}}(\ \Pi_{\mathsf{AC}}(\mathsf{R}) \bowtie \Pi_{\mathsf{BC}}(\mathsf{R}) \cup \Pi_{\mathsf{AB}}(\mathsf{R}) \bowtie \Pi_{\mathsf{BC}}(\mathsf{R})) =$



Trace only the set of input tuples that contributed to an output tuple

This is also a semi-ring! Which one?

CSE544 - Spring, 2013



Lineage [CuiWidomWiener 00 etc.]

Sets of contributing tuples



(Witness, Proof) why-provenance [BunemanKhannaTan 01] & [Buneman+ PODS08]

Sets of witnesses (w. =set of contributing tuples)

Semiring: (Why(*X*), ∪, ⊍, ∅, {∅})

Source: Tannen, EDBT 2010



Minimal witness why-provenance [BunemanKhannaTan 01]

Sets of minimal witnesses

Semiring: (PosBool(X), Λ , V, \neg , \bot)



Bags of sets of contributing tuples (of witnesses)

Semiring: (Trio(X), +, \cdot , 0, 1) (defined in [Green, ICDT 09])

Source: Tannen, EDBT 2010



Semiring: (B[X], +, ·, 0, 1)



Provenance polynomials [GKT, PODS 07] (N[X]-provenance) Bags of bags of contributing tuples **Semiring:** (N[X], +, ·, 0, 1)

Discretionary Access Control [LaPadula]

- Public = P
- Confidential = C
- Secret = S
- Top Secret = T
- No Such Thing... = 0



В	С	
b	С	X=C
b	е	Y=P
g	е	Z=T
	B b b	BCbcbege

Α	С	
а	С	$2 \cdot X^2 = ?$
а	е	X·Y = ?
d	е	$2 \cdot Y^2 + Y \cdot Z = ?$
f	е	$2 \cdot Z^2 + Y \cdot Z = ?$

Discretionary Access Control [LaPadula]

- Public = P
- Confidential = C
- Secret = S
- Top Secret = T
- No Such Thing... = 0



Α	С	
а	С	$2 \cdot X^2 = C$
а	е	$X \cdot Y = C$
d	е	$2 \cdot Y^2 + Y \cdot Z = C$
f	е	$2 \cdot Z^2 + Y \cdot Z = T$

(A, min, max, 0, P), where A = P < C < S < T < 0

But are there useful commutative semirings?

(B, ∧, ∨, ⊤, ⊥)	Set semantics
(ℕ, +, ·, 0, 1)	Bag semantics
(P(Ω), ∪, ∩, ∅, Ω)	Probabilistic events [FuhrRölleke 97]
(BoolExp(X), ∧, ∨, ⊤, ⊥)	Conditional tables (c-tables) [ImielinskiLipski 84]
(R ₊ ∞, min, +, 1, 0)	Tropical semiring (cost/distrust score/confidence need)
(A, min, max, 0, P) where A = P < C < S < T < 0	Access control levels [PODS8]

A provenance hierarchy



One semiring to rule them all... (apologies!)



A path downward from K_1 to K_2 indicates that there exists an **onto** (surjective) semiring homomorphism $h: K_1 \rightarrow K_2$

Using homomorphisms to relate models



Homomorphism?

h(x+y) = h(x)+h(y) h(xy)=h(x)h(y) h(0)=0 h(1)=1Moreover, for these homomorphisms h(x)=x