

ProvSQL Tutorial

Introduction to Semiring Provenance

Pierre Senellart



PSL 



inria informatics mathematics



institut
universitaire
de France

DesCartes school, 10 October 2022

Provenance management

- Data management **all about query evaluation**

Provenance management

- Data management **all about query evaluation**
- What if we want **something more** than the query result?
 - Where does the result come from?
 - Why was this result obtained?
 - How was the result produced?
 - What is the probability of the result?
 - How many times was the result obtained?
 - How would the result change if part of the input data was missing?
 - What is the minimal security clearance I need to see the result?
 - What is the most economical way of obtaining the result?
 - How can a result be explained in layman terms?

Provenance management

- Data management **all about query evaluation**
- What if we want **something more** than the query result?
 - Where does the result come from?
 - Why was this result obtained?
 - How was the result produced?
 - What is the probability of the result?
 - How many times was the result obtained?
 - How would the result change if part of the input data was missing?
 - What is the minimal security clearance I need to see the result?
 - What is the most economical way of obtaining the result?
 - How can a result be explained in layman terms?
- **Provenance management**: along with query evaluation, record **additional bookkeeping information** allowing to answer the questions above

Data model

- **Relational data model:** data decomposed into relations, with labeled attributes. . .

Data model

- **Relational data model:** data decomposed into relations, with labeled attributes...

name	position	city	classification
John	Director	New York	unclassified
Paul	Janitor	New York	restricted
Dave	Analyst	Paris	confidential
Ellen	Field agent	Berlin	secret
Magdalen	Double agent	Paris	top secret
Nancy	HR director	Paris	restricted
Susan	Analyst	Berlin	secret

Data model

- **Relational data model**: data decomposed into relations, with labeled attributes...
- ... with an extra **provenance annotation** for each tuple (think of it first as a tuple id)

name	position	city	classification	prov
John	Director	New York	unclassified	<i>t₁</i>
Paul	Janitor	New York	restricted	<i>t₂</i>
Dave	Analyst	Paris	confidential	<i>t₃</i>
Ellen	Field agent	Berlin	secret	<i>t₄</i>
Magdalen	Double agent	Paris	top secret	<i>t₅</i>
Nancy	HR director	Paris	restricted	<i>t₆</i>
Susan	Analyst	Berlin	secret	<i>t₇</i>

Commutative semiring $(K, 0, 1, \oplus, \otimes)$

- Set K with distinguished elements $0, 1$
- \oplus **associative, commutative** operator, with identity 0_K :
 - $a \oplus (b \oplus c) = (a \oplus b) \oplus c$
 - $a \oplus b = b \oplus a$
 - $a \oplus 0 = 0 \oplus a = a$
- \otimes **associative, commutative** operator, with identity 1_K :
 - $a \otimes (b \otimes c) = (a \otimes b) \otimes c$
 - $a \otimes b = b \otimes a$
 - $a \otimes 1 = 1 \otimes a = a$
- \otimes **distributes** over \oplus :

$$a \otimes (b \oplus c) = (a \otimes b) \oplus (a \otimes c)$$

- 0 is **annihilating** for \otimes :

$$a \otimes 0 = 0 \otimes a = 0$$

Example semirings

- $(\mathbb{N}, 0, 1, +, \times)$: **counting** semiring
- $(\{\perp, \top\}, \perp, \top, \vee, \wedge)$: **Boolean** semiring
- $(\{unclassified, restricted, confidential, secret, top\ secret\}, top\ secret, unclassified, \min, \max)$: **security** semiring
- $(\mathbb{N} \cup \{\infty\}, \infty, 0, \min, +)$: **tropical** semiring
- $(\{\text{Boolean functions over } \mathcal{X}\}, \perp, \top, \vee, \wedge)$: semiring of **Boolean functions** over \mathcal{X}
- $(\mathbb{N}[\mathcal{X}], 0, 1, +, \times)$: semiring of integer-valued **polynomials** with variables in \mathcal{X} (also called **How**-semiring or **universal** semiring)
- $(\mathcal{P}(\mathcal{P}(\mathcal{X})), \emptyset, \{\emptyset\}, \cup, \uplus)$: **Why**-semiring over \mathcal{X}
($A \uplus B := \{a \cup b \mid a \in A, b \in B\}$)

Semiring provenance [Green et al., 2007]

- We **fix** a semiring $(K, 0, 1, \oplus, \otimes)$
- We assume provenance annotations are **in K**
- We consider a query q from the **positive relational algebra** (selection, projection, renaming, cross product, union; joins can be simulated with renaming, cross product, selection, projection)
- We define a semantics for the provenance of a tuple $t \in q(D)$ **inductively** on the structure of q

What can we do with it?

counting semiring: count the number of times a tuple can be derived, multiset semantics

Boolean semiring: determines if a tuple exists when a subdatabase is selected

security semiring: determines the minimum clearance level required to get a tuple as a result

tropical semiring: minimum-weight way of deriving a tuple (think shortest path in a graph)

Boolean functions: Boolean provenance, with applications to **probabilistic databases**

integer polynomials: universal provenance, see further

Why-semiring: Why-provenance [Buneman et al., 2001], set of combinations of tuples needed for a tuple to exist

ProvSQL: Provenance within PostgreSQL

[Senellart et al., 2018]

- **Lightweight** extension/plugin for PostgreSQL ≥ 9.5
- Provenance annotations stored as **UUIDs**, in an extra attribute of each provenance-aware relation
- A provenance circuit **relating UUIDs** of elementary provenance annotations and arithmetic gates stored as table
- All computations done in the **universal semiring** (more precisely, extensions of it to support more operations)
- **Probability computation** from the provenance circuits, via various methods

Bibliography I

- Peter Buneman, Sanjeev Khanna, and Wang Chiew Tan. Why and where: A characterization of data provenance. In *Database Theory - ICDT 2001, 8th International Conference, London, UK, January 4-6, 2001, Proceedings.*, 2001.
- Todd J Green, Grigoris Karvounarakis, and Val Tannen. Provenance semirings. In *PODS*, 2007.
- Pierre Senellart, Louis Jachiet, Silviu Maniu, and Yann Ramusat. ProvSQL: provenance and probability management in postgresql. In *VLDB*, 2018. Demonstration.