

Counting Types for Massive JSON Datasets

BDA 2017, Nancy (présenté à DBPL 2017)

*Mohamed-Amine Baazizi, Dario Colazzo,
Giorgio Ghelli, Carlo Sartiani*

Counting types

- ▶ Can types count?
- ▶ Should they?

Type theory
perspective

- ▶ How to efficiently summarize the structure of large JSON datasets?
- ▶ How precise is the summary?

Database
perspective

The first problem

- ▶ Type inference for massive JSON datasets, BDA 2016/EDBT 2017

- ▶ We infer this type from a collection of JSON objects

```
{  
  title : Str ;  
  text : [ Str ] + Null ;  
  author : { address:T? ; affil:T? ; ... } ?  
  abstract : Str ?  
}
```

- ▶ How «optional» is the author?
- ▶ How frequently a text is Null?

Let us count

{ title : Str¹⁰⁰⁰ ;

text : ([Str⁸⁰⁰⁰]⁸⁰⁰ + Null²⁰⁰)¹⁰⁰⁰ ;

author : { add:T³⁰⁰ ; affil:T³⁰⁰ ; ... }⁸⁰⁰ ;

abstract : Str²⁰ ;

}¹⁰⁰⁰

The second problem

▶ How to capture correlation information?

▶ { **addr:T³⁰⁰**; **aff:T³⁰⁰**; r:T⁸⁰⁰ }⁸⁰⁰

Concision

▶ { **addr:T³⁰⁰**; **aff:T³⁰⁰**; r:T³⁰⁰ }³⁰⁰ + { r:T⁵⁰⁰ }⁵⁰⁰

▶ { **addr:T³⁰⁰**; r:T³⁰⁰ }³⁰⁰ + { **aff:T³⁰⁰**; r:T⁵⁰⁰ }⁵⁰⁰

▶ { **addr:T³⁰⁰**; r:T⁵⁰⁰ }⁵⁰⁰ + { **aff:T³⁰⁰**; r:T³⁰⁰ }³⁰⁰

▶ { **addr:T³⁰⁰**; r:T³⁰⁰ }³⁰⁰ + { **aff:T³⁰⁰**; r:T³⁰⁰ }³⁰⁰ + { r:T²⁰⁰ }²⁰⁰

Precision

The type system

- ▶ $B ::= \text{Null}^i \mid \text{Num}^i \mid \text{Str}^i \mid \text{Bool}^i$
- ▶ $R ::= \{ 1 : T, \dots, 1 : T \}^i$
- ▶ $A ::= [T]^i$
- ▶ $S ::= B \mid R \mid A$
- ▶ $T ::= S \mid 0 \mid T + T$
- ▶ **Examples**
 - ▶ Num^2 captures any multiset of two numbers
 - ▶ $[\text{Num}^4]^3$ a possible type for the multiset $\{ [1], [1], [1,2] \}^M$

The type inference algorithm

▶ Singleton

▶ $\vdash V : S$ $\vdash 3 : \text{Int}^1$

▶ Multiset

▶ $\vdash v_1, \dots, v_n :^M T$ $\vdash [1], [1], [1,2] :^M [\text{Num}^4]^3$

▶ Different abstraction levels

▶ $[1], [1], [1,2] :^M [\text{Num}^4]^3$

Concision

▶ $[1], [1], [1,2] :^M [\text{Num}^2]^2 + [\text{Num}^2]^1$

▶ $[1], [1], [1,2] :^M [\text{Num}^1]^1 + [\text{Num}^1]^1 + [\text{Num}^2]^1$

Precision

The type inference algorithm

▶ Singleton

Parametric inference

▶ $\overset{E}{\vdash} V : S$

$\vdash 3 : \text{Int}^1$

▶ Multiset

▶ $\overset{E}{\vdash} v_1, \dots, v_n :^M T \quad \vdash [1], [1], [1,2] :^M [\text{Num}^4]^3$

▶ Different abstraction levels

▶ $[1], [1], [1,2] :^M [\text{Num}^4]^3$

Concision

▶ $[1], [1], [1,2] :^M [\text{Num}^2]^2 + [\text{Num}^2]^1$

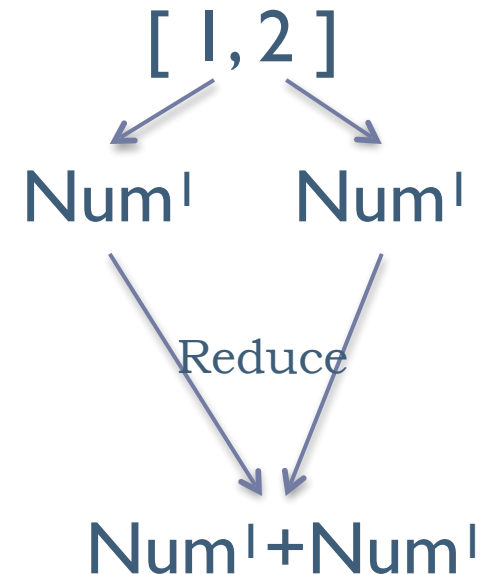
▶ $[1], [1], [1,2] :^M [\text{Num}^1]^1 + [\text{Num}^1]^1 + [\text{Num}^2]^1$

Precision

The type inference algorithm

$$\frac{\text{(TYPEARRAY)} \quad \vdash^E \{J_1, \dots, J_n\}^m :^m \mathbb{T}}{\vdash^E [J_1, \dots, J_n] : [\mathbb{T}]^1}$$

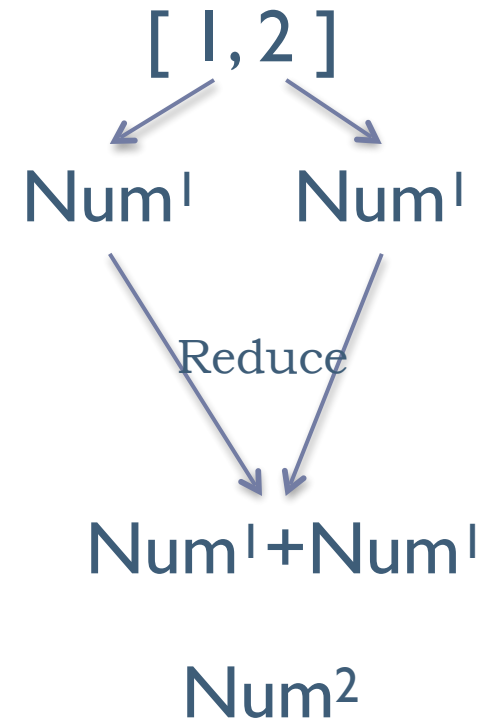
$$\frac{\text{(TYPEUNIONMULTISET)} \quad \vdash^E M_1 :^m \mathbb{T}_1 \quad \vdash^E M_2 :^m \mathbb{T}_2}{\vdash^E M_1 \cup^m M_2 :^m \text{Reduce}(\mathbb{T}_1, \mathbb{T}_2, E)}$$



The type inference algorithm

$$\frac{\text{(TYPEARRAY)} \quad \vdash^E \{J_1, \dots, J_n\}^m :^m \mathbb{T}}{\vdash^E [J_1, \dots, J_n] : [\mathbb{T}]^1}$$

$$\frac{\text{(TYPEUNIONMULTISET)} \quad \vdash^E M_1 :^m \mathbb{T}_1 \quad \vdash^E M_2 :^m \mathbb{T}_2}{\vdash^E M_1 \cup^m M_2 :^m \text{Reduce}(\mathbb{T}_1, \mathbb{T}_2, E)}$$



Parametric reduction

$$\begin{aligned}
 \text{Reduce}(\mathbb{T}_1, \mathbb{T}_2, E) = & \\
 & \oplus(\{ \text{Merge}(\mathbb{S}_1, \mathbb{S}_2, E) \mid \mathbb{S}_1 \in \circ\mathbb{T}_1, \mathbb{S}_2 \in \circ\mathbb{T}_2, E(\mathbb{S}_1, \mathbb{S}_2) \}^m \\
 & \cup^m \{ \mathbb{S}_1 \mid \mathbb{S}_1 \in \circ\mathbb{T}_1, \nexists \mathbb{S}_2 \in \circ\mathbb{T}_2. E(\mathbb{S}_1, \mathbb{S}_2) \}^m \\
 & \cup^m \{ \mathbb{S}_2 \mid \mathbb{S}_2 \in \circ\mathbb{T}_2, \nexists \mathbb{S}_1 \in \circ\mathbb{T}_1. E(\mathbb{S}_1, \mathbb{S}_2) \}^m)
 \end{aligned}$$

$$\text{Merge}(\mathbb{B}^m, \mathbb{B}^n, E) = \mathbb{B}^{m+n}$$

$$\begin{aligned}
 \text{Merge}(\mathbb{R}_1, \mathbb{R}_2, E) = & \\
 & \{ \{ (l, \text{Reduce}(\mathbb{T}_1, \mathbb{T}_2, E)) \mid (l, \mathbb{T}_1) \in \diamond(\mathbb{R}_1), (l, \mathbb{T}_2) \in \diamond(\mathbb{R}_2) \} \\
 & \cup \{ (l, \mathbb{T}_1) \mid (l, \mathbb{T}_1) \in \diamond(\mathbb{R}_1), \nexists \mathbb{T}_2. (l, \mathbb{T}_2) \in \diamond(\mathbb{R}_2) \} \\
 & \cup \{ (l, \mathbb{T}_2) \mid (l, \mathbb{T}_2) \in \diamond(\mathbb{R}_2), \nexists \mathbb{T}_1. (l, \mathbb{T}_1) \in \diamond(\mathbb{R}_1) \} \\
 & \}^{\#(\mathbb{R}_1)+\#(\mathbb{R}_2)}
 \end{aligned}$$

$$\text{Merge}([\mathbb{T}_1]^m, [\mathbb{T}_2]^n, E) = [\text{Reduce}(\mathbb{T}_1, \mathbb{T}_2, E)]^{m+n}$$

Equivalences of practical use

► Kind

$$\mathcal{K}(S_1, S_2) \Leftrightarrow \text{kind}(S_1) = \text{kind}(S_2)$$

$$\text{kind}(\text{Null}^i) = 0 \quad \text{kind}(\text{Bool}^i) = 1 \quad \text{kind}(\text{Num}^i) = 2$$

$$\text{kind}(\text{Str}^i) = 3 \quad \text{kind}(\mathbb{R}) = 4 \quad \text{kind}(\mathbb{A}) = 5$$

$$\{ \text{addr:T}^{300}; \text{aff:T}^{300}; r:\text{T}^{800} \}^{800}$$

► Label

$$\mathcal{L}([T_1], [T_2]) \quad \text{always}$$

$$\mathcal{L}([\mathbb{B}], [\mathbb{B}]) \quad \text{always}$$

$$\mathcal{L}(\mathbb{R}_1, \mathbb{R}_2) \Leftrightarrow \text{Keys}(\mathbb{R}_1) = \text{Keys}(\mathbb{R}_2)$$

$$\{ \text{addr:T}^{300}; r:\text{T}^{300} \}^{300} + \{ \text{aff:T}^{300}; r:\text{T}^{300} \}^{300} + \{ r:\text{T}^{200} \}^{200}$$

Kind reduction: twitter data

```
{ contributors: (Null9,599,980 + [Num20]20)9,600,000 ;  
  retweeted : Bool9,600,000 ;  
  retweeted_status {...} : {...}1,200,000 ;  
  deleted : {...}300,000 ;  
}9,900,000
```

Label reduction: twitter data

{ con: ...^{7,200,000}; ret: Bool^{7,200,000};... }^{7,200,000}

+{ con: ...^{1,200,000}; ret: Bool^{1,200,000};... }^{1,200,000}

+{ con: ...^{1,040,000}; ret: Bool^{1,040,000}; r_s: {}^{1,040,000};... }^{1,040,000}

+{ con: ...^{160,000}; ret: Bool^{160,000}; r_s: {}^{160,000};... }^{160,000}

+{ deleted: { }^{300,000};... }^{300,000} :

Label reduction: twitter data

{ con: ...^{7,200,000}; ret: Bool^{7,200,000};... }^{7,200,000}

+{ con: ...^{1,200,000}; ret: Bool^{1,200,000};... }^{1,200,000}

+{ con: ...^{1,040,000}; ret: Bool^{1,040,000}; r_s: {}^{1,040,000};... }^{1,040,000}

+{ con: ...^{160,000}; ret: Bool^{160,000}; r_s: {}^{160,000};... }^{160,000}

+{ deleted: { }^{300,000};... }^{300,000} :

```
{ contributors: (Null9,599,980 + [Num20]20)9,600,000 ;  
  retweeted : Bool9,600,000 ;  
  retweeted_status { ... } : { ... }1,200,000 ;  
  deleted : { ... }300,000 ;  
} 9,900,000
```

Kind reduction

Experiments

- ▶ Scala implementation
- ▶ Spark cluster of 5+1 nodes, 64GB, 100 cores
- ▶ 3 real life datasets :
 - ▶ Github (1m objects / 10GB / 14 sec)
 - ▶ Twitter (9.9m objects / 21GB / 53 sec)
 - ▶ Nytimes (1.2m objects / 21GB / 27 sec)
- ▶ Extraction of interesting features

Related work

- ▶ PL community: dependent types, probabilistic types
- ▶ Dataguide avec statistiques [Klettke et al. 2016]
- ▶ JavaScript library for MongoDB [Schmit 2017]
- ▶ Approaches w/o counting information
- ▶ No parametric approach, so far

[Klettke et al. 2016] Schema Extraction and Structural Outlier Detection for JSON-based NoSQL Data Stores, Technologie und Web (BTW)

[Schmidt. 2017]. mongodb-schema. (2017). <https://github.com/mongodb-js/mongodb-schema>.

To sum up

- ▶ An algorithm to summarize JSON data:
 - ▶ Well defined semantics
 - ▶ Parametric
 - ▶ Parallel
 - ▶ Yielding quantitative information
- ▶ What else may a counting type do?

Thank you!