

# Trimming a consistent DL knowledge base, relying on linguistic evidence

supervisors : Laure Vieu et Nathalie Aussenac-Gilles

Julien Corman

IRIT

## Problem

Intuition

Similarity

Plausibility

Trimming

Experiments

Extensions

Appendices

MF

Trimming :  
illustration

Trimming :  
assumptions

# 1 Problem

## 2 Intuition

## 3 Similarity

## 4 Plausibility

## 5 Trimming

## 6 Experiments

## 7 Extensions

## 8 Appendices

### ■ MF

### ■ Trimming : illustration

### ■ Trimming : assumptions

# DL and OWL

## Description Logics (DL):

- Decidable fragments of FOL
- $\mathcal{ALC}$  = “modal fragment” of FOL : unary and binary predicates only (called *atomic concepts* and *roles*), no identity, no function, restrictions on quantification (see appendix).
- Extensions : nominals, cardinality restriction, role subsumption, role composition, inverse roles, ...
- Algorithms and libraries for different tasks/problems : consistency, entailment, modularity, minimal conflicts, ...

Problem

Intuition

Similarity

Plausibility

Trimming

Experiments

Extensions

Appendices

MF

Trimming :  
illustration

Trimming :  
assumptions

# DL and OWL

## Description Logics (DL):

- Decidable fragments of FOL
- $\mathcal{ALC}$  = “modal fragment” of FOL : unary and binary predicates only (called *atomic concepts* and *roles*), no identity, no function, restrictions on quantification (see appendix).
- Extensions : nominals, cardinality restriction, role subsumption, role composition, inverse roles, ...
- Algorithms and libraries for different tasks/problems : consistency, entailment, modularity, minimal conflicts, ...

## OWL 2

- Knowledge representation language, W3C recommendation.
- Equivalent to the DL  $\mathcal{SROIQ}^{(D)}$
- Several syntaxes, among which a (hardly readable) RDF serialization.

Problem

Intuition

Similarity

Plausibility

Trimming

Experiments

Extensions

Appendices

MF

Trimming :  
illustration

Trimming :  
assumptions

# Example 1

## DBpedia

```
hasKeyPerson(Virgin Holidays, CEO).  
hasKeyPerson(Caixa Bank, CEO).  
hasOccupation(Peter Munk, CEO).
```

Problem

Intuition

Similarity

Plausibility

Trimming

Experiments

Extensions

Appendices

MF

Trimming :  
illustration

Trimming :  
assumptions

# Example 1

## DBPedia

`hasKeyPerson`(*Virgin Holidays*, *CEO*).  
`hasKeyPerson`(*Caixa Bank*, *CEO*).  
`hasOccupation`(*Peter Munk*, *CEO*).  
`hasKeyPerson`(*BrookField Office Properties*, *Peter Munk*).  
 $\top \sqsubseteq \forall \text{hasKeyPerson}.\text{Person}.$

Problem

Intuition

Similarity

Plausibility

Trimming

Experiments

Extensions

Appendices

MF

Trimming :  
illustration

Trimming :  
assumptions

# Example 1

## DBPedia

```
hasKeyPerson(Virgin Holidays, CEO).  
hasKeyPerson(Caixa Bank, CEO).  
hasOccupation(Peter Munk, CEO).  
hasKeyPerson(BrookField Office Properties, Peter Munk).  
 $\top \sqsubseteq \forall \text{hasKeyPerson}.\text{Person}.$ 
```

- Intuitively absurd : violates for instance “No individual (*CEO* here) can be both a person and the occupation of a person”.
- More pragmatically, may lead to erroneous inferences : e.g. *Virgin Holidays* and *Caixa Bank* have the same Person as a keyPerson.
- But logically consistent and coherent.

## Example 2

### DBPedia

`doctoralAdvisor`(*Thaddeus S.C. Lowe*, *Smithsonian Institution*).  
`doctoralAdvisor`(*Nick Katz*, *Bernard Dwork*).  
`owningCompany`(*Smithsonian Networks*, *Smithsonian Institution*).  
 $T \sqsubseteq \forall \text{doctoralAdvisor}.\text{Person}.$

- Still logically consistent and coherent.



## Example 2

### DBPedia

`doctoralAdvisor`(*Thaddeus S.C. Lowe*, *Smithsonian Institution*).  
`doctoralAdvisor`(*Nick Katz*, *Bernard Dwork*).  
`owningCompany`(*Smithsonian Networks*, *Smithsonian Institution*).  
 $T \subseteq \forall \text{doctoralAdvisor}.\text{Person}.$

- Still logically consistent and coherent.
- These are not just “factual” errors, like `director`(*Citizen Kane*, *Woody Allen*).
- Source of the problem :
  - genuine typos
  - incompatible understandings/uses of a same DL individual/concept/role.



# OWL data : consistent/coherent by default

- One of the following is necessary for an OWL 2 dataset to be inconsistent/incoherent :

- `owl:complementOf` or `owl:disjointWith`
- `owl:negativeObjectPropertyAssertion`
- `owl:disjointObjectProperties`, `owl:AsymmetricProperty` or `owl:irreflexiveObjectProperty`.
- `owl:oneOf`
- `owl:Nothing`
- `owl:objectMaxCardinality`
- etc...

## Problem

Intuition

Similarity

Plausibility

Trimming

Experiments

Extensions

Appendices

MF

Trimming :  
illustration

Trimming :  
assumptions

# OWL data : consistent/coherent by default

- One of the following is necessary for an OWL 2 dataset to be inconsistent/incoherent :
  - `owl:complementOf` or `owl:disjointWith`
  - `owl:negativeObjectPropertyAssertion`
  - `owl:disjointObjectProperties`, `owl:AsymmetricProperty` or `owl:irreflexiveObjectProperty`.
  - `owl:oneOf`
  - `owl:Nothing`
  - `owl:objectMaxCardinality`
  - etc...

Rarely used (source : LODStats (LODCloud sample))

- `owl:subClassOf` : > 89 000 occ.
- `owl:complementOf` : 2 occ.
- `owl:disjointWith` : 33 occ.

# Proposal

## Problem

### Intuition

### Similarity

### Plausibility

### Trimming

### Experiments

### Extensions

### Appendices

#### MF

#### Trimming : illustration

#### Trimming : assumptions

- Automatically gathered linguistic evidence in order to **detect** and **repair** such violations of common sense.
  - **Detect** : identify consequences of a set  $\Gamma$  of axioms which are unlikely to hold if the rest of  $Cn(\Gamma)$  does.
  - **Repair** : suggest axioms to be preferably discarded or amended
- Linguistic input : web pages

# Proposal

## Problem

### Intuition

### Similarity

### Plausibility

### Trimming

### Experiments

### Extensions

### Appendices

#### MF

#### Trimming : illustration

#### Trimming : assumptions

- Automatically gathered linguistic evidence in order to **detect** and **repair** such violations of common sense.
  - **Detect** : identify consequences of a set  $\Gamma$  of axioms which are unlikely to hold if the rest of  $Cn(\Gamma)$  does.
  - **Repair** : suggest axioms to be preferably discarded or amended
- Linguistic input : web pages
- Main hypothesis (distributional evidence) : individuals which share linguistic contexts tend to instantiate the same concepts.  
Inspiration : ontology population/named entity classification (Tanev and Magnini, ...)

- 1 Problem
- 2 Intuition**
- 3 Similarity
- 4 Plausibility
- 5 Trimming
- 6 Experiments
- 7 Extensions
- 8 Appendices
  - MF
  - Trimming : illustration
  - Trimming : assumptions

# Linguistic evidence : intuition

## Example

$\Gamma = \{$   
   $\text{doctoralAdvisor}(\textit{Thaddeus S.C. Lowe}, \textit{Smithsonian Institution}),$   
   $\text{doctoralAdvisor}(\textit{Nick Katz}, \textit{Bernard Dwork}),$   
 $\top \subseteq \forall \text{doctoralAdvisor. Person}, \dots \}$

- $\Gamma \models \text{Person}(\textit{Smithsonian Institution})$
- $\Gamma \models \text{Person}(\textit{Bernard Dwork})$



# Linguistic evidence : intuition

## Example

$\Gamma = \{$   
   $\text{doctoralAdvisor}(\textit{Thaddeus S.C. Lowe}, \textit{Smithsonian Institution}),$   
   $\text{doctoralAdvisor}(\textit{Nick Katz}, \textit{Bernard Dwork}),$   
 $\top \subseteq \forall \text{doctoralAdvisor. Person}, \dots \}$

- $\Gamma \models \text{Person}(\textit{Smithsonian Institution})$   
   $\Gamma \models \text{Person}(\textit{Bernard Dwork})$
- Assume also that :  
   $\Gamma \models \text{Person}(\textit{Margaret Atwood})$   
   $\Gamma \models \text{Person}(\textit{Peter Munk})$   
   $\Gamma \models \text{Person}(\textit{Thaddeus S.C. Lowe}) \dots$
- Does “the Smithsonian institution” behave like terms denoting other instances of `Person` according to  $\Gamma$  ?  
  Does “Bernard Dwork” behave like terms denoting other instances of `Person` according to  $\Gamma$  ?

# Linguistic evidence : intuition

## Example

$\Gamma = \{$   
  `doctoralAdvisor`(*Thaddeus S.C. Lowe*, *Smithsonian Institution*),  
  `doctoralAdvisor`(*Nick Katz*, *Bernard Dwork*),  
 $\top \sqsubseteq \forall \text{doctoralAdvisor. Person, ...} \}$

- $\Gamma \models \text{Person}(\textit{Smithsonian Institution})$   
   $\Gamma \models \text{Person}(\textit{Bernard Dwork})$
- “#the Smithsonian Institution was born”  
  “Bernard Dwork was born”

# Linguistic evidence : intuition

## Example

$\Gamma = \{$   
   $\text{doctoralAdvisor}(\textit{Thaddeus S.C. Lowe}, \textit{Smithsonian Institution}),$   
   $\text{doctoralAdvisor}(\textit{Nick Katz}, \textit{Bernard Dwork}),$   
   $\top \sqsubseteq \forall \text{doctoralAdvisor. Person}$   
   $\text{Organization} \sqsubseteq \neg \text{Person}, \dots \}$

- $\Gamma \models \text{Person}(\textit{Smithsonian Institution})$
- $\Gamma \models \text{Person}(\textit{Bernard Dwork})$
- $\Gamma \models \neg \text{Organization}(\textit{Smithsonian Institution})$
- $\Gamma \models \neg \text{Organization}(\textit{Bernard Dwork})$

# Linguistic evidence : intuition

## Example

$\Gamma = \{$   
  `doctoralAdvisor`(*Thaddeus S.C. Lowe*, *Smithsonian Institution*),  
  `doctoralAdvisor`(*Nick Katz*, *Bernard Dwork*),  
   $\top \sqsubseteq \forall \text{doctoralAdvisor. Person}$   
   $\text{Organization} \sqsubseteq \neg \text{Person}, \dots \}$

- $\Gamma \models \text{Person}(\textit{Smithsonian Institution})$   
   $\Gamma \models \text{Person}(\textit{Bernard Dwork})$   
   $\Gamma \models \neg \text{Organization}(\textit{Smithsonian Institution})$   
   $\Gamma \models \neg \text{Organization}(\textit{Bernard Dwork})$
- “the Smithsonian Institution was established”, “the Smithsonian Institution’s workforce”  
  “#Bernard Dwork was established”, “#Bernard Dwork’s workforce”

# Linguistic evidence : intuition

## Example

$\Gamma = \{$   
  `doctoralAdvisor`(*Thaddeus S.C. Lowe*, *Smithsonian Institution*),  
  `doctoralAdvisor`(*Nick Katz*, *Bernard Dwork*),  
   $\top \sqsubseteq \forall \text{doctoralAdvisor. Person}$   
   $\text{Organization} \sqsubseteq \neg \text{Person}, \dots \}$

- Linguistic contexts may help identify :
  - plausible consequences of  $\Gamma$  :  $\text{Person}(\textit{Bernard Dwork})$ ,  
   $\neg \text{organization}(\textit{Bernard Dwork})$
  - implausible consequences of  $\Gamma$  :  $\text{Person}(\textit{Smithsonian Institution})$ ,  $\neg \text{organization}(\textit{Smithsonian Institution})$

# Choices

- Focus on  $\Psi_{\Gamma}$  : consequences of  $\Gamma$  the form  $A(e)$  or  $\neg A(e)$ , with  $A$  an atomic concept and  $e$  an individual.
- Linguistic terms labeling concepts and roles are never used (only terms labeling individuals).

## Individual labels rather than concept labels ?

- Concept labels tend to be more polysemous : e.g. “Group”, “Function”, “Element”, ...
- Lack of linguistic occurrences for :
  - Ad hoc concepts labels : ex (eGov ontologies) : “Triple path”, “Structuring event type” (0 google occ.)
  - Abstract concepts : e.g. “perdurant”

## Unary rather than binary predicates ?

- labels already known  $\Rightarrow$  lack of linguistic cooccurrences.

- 1 Problem
- 2 Intuition
- 3 Similarity**
- 4 Plausibility
- 5 Trimming
- 6 Experiments
- 7 Extensions
- 8 Appendices
  - MF
  - Trimming : illustration
  - Trimming : assumptions

# Similarity

- Distributional hypothesis : represent a term  $t$  by its linguistic contexts
- A context  $c$  :
  - sequence of words preceding/surrounding/following an occurrence of the term, possibly lemmatized
  - syntactic dependency, ...
  - ignoring punctuation, determiners, ...
- A terms  $t$  is represented as a vector  $\mathbf{v}^t$  of frequencies with each observed context.



# Similarity

- Distributional hypothesis : represent a term  $t$  by its linguistic contexts
- A context  $c$  :
  - sequence of words preceding/surrounding/following an occurrence of the term, possibly lemmatized
  - syntactic dependency, ...
  - ignoring punctuation, determiners, ...
- A terms  $t$  is represented as a vector  $\mathbf{v}^t$  of frequencies with each observed context.
- Weighting observed frequencies :
  - $PMI(c, e) = -\log \frac{p(c, e)}{p(e) \cdot p(c)}$
  - self-information (Giulano and Gliozzo) :  $\text{self}(c) = -\log(p(c))$ , with  $p(c)$  obtained from an external language model...

# Similarity

Problem

Intuition

Similarity

Plausibility

Trimming

Experiments

Extensions

Appendices

MF

Trimming :  
illustration

Trimming :  
assumptions

- Distributional hypothesis : represent a term  $t$  by its linguistic contexts
- A context  $c$  :
  - sequence of words preceding/surrounding/following an occurrence of the term, possibly lemmatized
  - syntactic dependency, ...
  - ignoring punctuation, determiners, ...
- A term  $t$  is represented as a vector  $\mathbf{v}^t$  of frequencies with each observed context.
- Weighting observed frequencies :
  - $PMI(c, e) = -\log \frac{p(c, e)}{p(e) \cdot p(c)}$
  - self-information (Giulano and Gliozzo) :  $\text{self}(c) = -\log(p(c))$ , with  $p(c)$  obtained from an external language model...
- Reducing vector dimensions : latent semantic analysis (SVD), latent Dirichlet allocation, skip-gram model, ...
- Similarity  $\text{sim}(t_1, t_2)$  given by some distance (cosine, ...) between  $\mathbf{v}^{t_1}$  and  $\mathbf{v}^{t_2}$ .

- 1 Problem
- 2 Intuition
- 3 Similarity
- 4 Plausibility**
- 5 Trimming
- 6 Experiments
- 7 Extensions
- 8 Appendices
  - MF
  - Trimming : illustration
  - Trimming : assumptions

# Plausibility of $A(e) \in \text{Cn}(\Gamma)$

Notation :

- $\text{sim}(e, e')$  : similarity between distributional representations of terms denoting  $e$  and  $e'$ .
- $\text{inst}_\Gamma(A) = \{e' \mid \Gamma \models A(e')\}$
- $S = \text{inst}_\Gamma(A) \setminus \{e\}$  : *support set* for  $A(e)$ .
- $\text{sim}(e, S) \doteq \sum_{e' \in S} \frac{\text{sim}(e, e')}{|S|}$
- $X_{e, |S|}^\Gamma$  (random variable) : expected average similarity between  $e$  and  $|S|$  random individuals of  $\text{inst}_\Gamma(\top) \setminus \{e\}$ .

Plausibility score  $\text{sc}_\Gamma(A(e))$

- $\text{sc}_\Gamma(A(e)) = p(X_{e, |S|}^\Gamma \leq \text{sim}(e, S))$
- Measures how surprisingly high the similarity between  $e$  and individuals of  $S$  is.
- Based on the similarity between  $e$  and all individuals.

Problem

Intuition

Similarity

Plausibility

Trimming

Experiments

Extensions

Appendices

MF

Trimming :  
illustration

Trimming :  
assumptions

# Support set $S$

- $S = \text{inst}_\Gamma(A) \setminus \{e\}$  : *support set* for  $A(e)$ .
- What about  $\text{inst}_\Gamma(\neg A)$  ?
- Linguistically unrealistic : no reason to think that two instances of  $\neg A$  should behave similarly.

## Example

- $\Gamma \models \neg \text{Person}(WW2)$
- $\Gamma \models \text{Person}(\textit{Thelonious Monk})$
- $\text{sim}(\textit{Smithsonian Institution}, WW2) > \text{sim}(\textit{Smithsonian Institution}, \textit{Thelonious Monk}) ???$

# Support set $S$

- $S = \text{inst}_\Gamma(A) \setminus \{e\}$  : *support set* for  $A(e)$ .
- What about  $\text{inst}_\Gamma(\neg A)$  ?
- Linguistically unrealistic : no reason to think that two instances of  $\neg A$  should behave similarly.

## Example

- $\Gamma \models \neg \text{Person}(WW2)$
- $\Gamma \models \text{Person}(\textit{Thelonious Monk})$
- $\text{sim}(\textit{Smithsonian Institution}, WW2) > \text{sim}(\textit{Smithsonian Institution}, \textit{Thelonious Monk}) ???$

- Support set for  $\neg A(e)$  :  $S = \text{inst}_\Gamma(A)$

# Plausibility of $\neg A(e) \in \text{Cn}(\Gamma)$

Notation :

- $\text{sim}(e, e')$  : similarity between distributional representations of  $e$  and  $e'$ .
- $\text{inst}_\Gamma(A) = \{e' \mid \Gamma \models A(e')\}$
- $S = \text{inst}_\Gamma(A)$  : *support set for  $\neg A(e)$* .
- $X_{e,|S|}^\Gamma$  (random variable) : expected average similarity between  $e$  and  $|S|$  random individuals of  $\text{inst}_\Gamma(\top) \setminus \{e\}$ .

Plausibility score  $\text{sc}_\Gamma(\neg A(e))$

- $\text{sc}_\Gamma(A(e)) = p(X_{e,|S|}^\Gamma \geq \sum_{e' \in S} \frac{\text{sim}(e, e')}{|S|})$
- Measures how surprisingly *low* the similarity between  $e$  and individuals of  $S$  is.
- Based on the similarity between  $e$  and all individuals.

# Expected similarity

- $X_{e,|S|}^\Gamma$  : expected average similarity between  $e$  and  $|S|$  random individuals of  $\text{inst}_\Gamma(\mathcal{T}) \setminus \{e\}$ .
- Intuition : *ceteris paribus*, the lower  $|S|$ , the less informative  $\text{sim}(e, S)$  should be.
- The lower  $|S|$ , the more uniform de distribution of  $X_{e,|S|}^\Gamma$  should be.

## Distribution of $X_{e,|S|}^\Gamma$

- $m \doteq \text{sim}(e, \text{inst}_\Gamma(\mathcal{T}) \setminus \{e\})$
- $X_{e,|S|}^\Gamma \sim \text{Beta}(m|S| + 1, (1 - m)|S| + 1)$

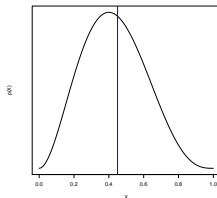


# Expected plausibility : example

- $\psi_1 = A(e), \psi_2 = B(e)$
- $m = \text{sim}(e, \text{inst}_\Gamma(\mathbb{T}) \setminus \{e\}) = 0.4$

$\psi_1$

- $S = \text{inst}_\Gamma(A) \setminus \{e\}$
- $|S| = 5$
- $\text{sim}(e, S) = 0.45$
- $\text{sc}_\Gamma(\psi_1) = 0.558$

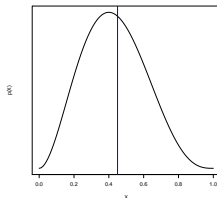


# Expected plausibility : example

- $\psi_1 = A(e), \psi_2 = B(e)$
- $m = \text{sim}(e, \text{inst}_\Gamma(T) \setminus \{e\}) = 0.4$

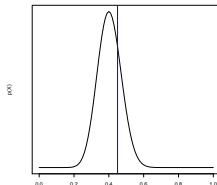
$\psi_1$

- $S = \text{inst}_\Gamma(A) \setminus \{e\}$
- $|S| = 5$
- $\text{sim}(e, S) = 0.45$
- $\text{sc}_\Gamma(\psi_1) = 0.558$



$\psi_2$

- $S = \text{inst}_\Gamma(B) \setminus \{e\}$
- $|S| = 50$
- $\text{sim}(e, S) = 0.45$
- $\text{sc}_\Gamma(\psi_2) = 0.754$



- 1 Problem
- 2 Intuition
- 3 Similarity
- 4 Plausibility
- 5 Trimming**
- 6 Experiments
- 7 Extensions
- 8 Appendices
  - MF
  - Trimming : illustration
  - Trimming : assumptions

# Trimming

- An input KB  $K$ .
- Objective : use plausibility scores to decide which axioms should be preferably discarded or amended within  $K$ .
- Equivalently, select the optimal  $\Gamma_1, \dots, \Gamma_n \in 2^K$ .

Linguistic compliance comp :  $2^K \mapsto \mathbb{R}$

$$\text{comp}(\Gamma) = \sum_{\psi \in \Psi_\Gamma} \frac{\text{scr}(\psi)}{|\Psi_\Gamma|}$$

# Trimming

- An input KB  $K$ .
- Objective : use plausibility scores to decide which axioms should be preferably discarded or amended within  $K$ .
- Equivalently, select the optimal  $\Gamma_1, \dots, \Gamma_n \in 2^K$ .

Linguistic compliance comp :  $2^K \mapsto \mathbb{R}$

$$\text{comp}(\Gamma) = \sum_{\psi \in \Psi_\Gamma} \frac{\text{scr}(\psi)}{|\Psi_\Gamma|}$$

- $\prec$  : strict partial order over  $2^K$  :  $\Gamma_1 \prec \Gamma_2$  iff either  $\text{comp}(\Gamma_1) < \text{comp}(\Gamma_2)$ , or  $(\text{comp}(\Gamma_1) = \text{comp}(\Gamma_2) \text{ and } \Gamma_1 \subset \Gamma_2)$ .
- Assumption : focus on syntax (see appendix).
- Output  $\mathcal{O}$  : intersection, or possibly disjunction of the subbases which are maximal wrt  $\prec$ .

# Trimming : practical limits

- Maximizing comp is not trivial :
  - $\text{comp}(\Gamma)$  is not directly a function of  $\Gamma$ , but of  $\Psi_\Gamma$  : so there may be an optimal  $\Psi' \subseteq \Psi_K$ , and no  $\Gamma$  such that  $\Psi_\Gamma = \Psi'$ .
  - For  $\psi \in \Psi_{\Gamma_1} \cap \Psi_{\Gamma_2}$ ,  $\text{scr}_{\Gamma_1}(\psi) \neq \text{scr}_{\Gamma_2}(\psi)$  in general, because the support sets for  $\psi$  differ in  $\Gamma_1$  and  $\Gamma_2$ .
- The output  $\mathcal{O}$  can be very weak, e.g. if  $|\mathcal{O}| < 0.5 * |K|$

# Trimming : practical limits

- Maximizing comp is not trivial :
  - $\text{comp}(\Gamma)$  is not directly a function of  $\Gamma$ , but of  $\Psi_\Gamma$  : so there may be an optimal  $\Psi' \subseteq \Psi_K$ , and no  $\Gamma$  such that  $\Psi_\Gamma = \Psi'$ .
  - For  $\psi \in \Psi_{\Gamma_1} \cap \Psi_{\Gamma_2}$ ,  $\text{scr}_{\Gamma_1}(\psi) \neq \text{scr}_{\Gamma_2}(\psi)$  in general, because the support sets for  $\psi$  differ in  $\Gamma_1$  and  $\Gamma_2$ .
- The output  $\mathcal{O}$  can be very weak, e.g. if  $|\mathcal{O}| < 0.5 * |K|$

## More plausible scenarios

- Search space previously circumscribed : e.g. discard at most  $n$  axioms.
- (Iteratively) discard the worst axiom (see evaluation).

# Alternatives to comp

Linguistic compliance  $\text{comp}_K : 2^K \mapsto \mathbb{R}$

$$\text{comp}_K(\Gamma) = \sum_{\psi \in \Psi_\Gamma} \frac{\text{sc}_K(\psi)}{|\Psi_\Gamma|}$$

- More amenable to optimizations.
- Ex (trivial) : a subbase  $\Gamma_1$  with  $\max_{\psi \in \Psi_{\Gamma_1}} \text{sc}_K(\psi) < \text{comp}_K(\Gamma_2)$  for some already evaluated subbase  $\Gamma_2$ .  
 $\Rightarrow$  No subbase of  $\Gamma_1$  can be optimal wrt  $\prec$ .
- Drawback : potentially higher number of optimal subbases.



# Alternatives to comp

## Lexicographic ordering $\preceq_{lex} \subseteq 2^K \times 2^K$

- Instead of plausibilities mean, penalize subbases whose consequences have a low plausibility (see appendix)
- Then  $\prec$  is defined by  $\Gamma_1 \prec \Gamma_2$  iff either  $\Gamma_1 \prec_{lex} \Gamma_2$ , or ( $\Gamma_1 =_{lex} \Gamma_2$  and  $\Gamma_1 \subset \Gamma_2$ ).

## Lexicographic ordering $\preceq_{lex_K} \subseteq 2^K \times 2^K$

- Identical to  $\preceq_{lex}$ , but using  $sc_K$  instead of  $sc_\Gamma$  for plausibility.
- Closer to traditional KB debugging / belief base revision : identify undesired consequences within  $K$  before trimming.

- 1 Problem
- 2 Intuition
- 3 Similarity
- 4 Plausibility
- 5 Trimming
- 6 Experiments**
- 7 Extensions
- 8 Appendices
  - MF
  - Trimming : illustration
  - Trimming : assumptions

# Input 1 : real data

## Real input KB

- Source : LOD
  - Evaluation procedure : manually verify if consequences with lowest plausibility and discarded axioms are actually erroneous.
  - Advantage : plausible data
  - Drawback : subjective evaluation (low inter annotator agreement)
- 
- Dataset  $K_{DBP}$  : 5721 (logical) axioms automatically extracted from DBPedia (see appendix).
  - 1095 individuals
  - ABox + TBox
  - expressivity :  $\mathcal{AL}^{(D)}$

## Input 2 : artificially degraded data

### Artificially degraded KB

- Source : higher quality KB
- Degrading procedure : randomly select an axiom  $\phi$  of  $K$ , and generate  $\phi'$  by replacing  $\text{sign}(\phi)$  with random elements of  $\text{sign}(K)$ . The syntactic structure remains unchanged.
- Requirements : the resulting base  $K' = K \cup \{\phi'\}$  must be consistent, and  $|\Psi_K| < |\Psi_{K'}|$ .
- Assumption : random axioms are very likely be absurd, and so random consequences to be outliers within  $\Psi_{K'}$ .

## Input 2 : artificially degraded data

### Artificially degraded KB

- Source : higher quality KB
- Degrading procedure : randomly select an axiom  $\phi$  of  $K$ , and generate  $\phi'$  by replacing  $sign(\phi)$  with random elements of  $sign(K)$ . The syntactic structure remains unchanged.
- Requirements : the resulting base  $K' = K \cup \{\phi'\}$  must be consistent, and  $|\Psi_K| < |\Psi_{K'}|$ .
- Assumption : random axioms are very likely be absurd, and so random consequences to be outliers within  $\Psi_{K'}$ .
- Evaluation : automatically retrieve the generated axioms and consequences within  $K'$  and  $\Psi_{K'}$  respectively.
- Drawback : artificial data
- Advantage : objective evaluation

## Input 2 : artificially degraded KB

Problem

Intuition

Similarity

Plausibility

Trimming

**Experiments**

Extensions

Appendices

MF

Trimming :  
illustration

Trimming :  
assumptions

- Dataset  $K_F$  : 1028 axioms automatically extracted from the NEON fisheries ontology (see appendix).
- 71 individuals
- ABox + TBox
- expressivity :  $\mathcal{SI}$

# Linguistic input

## Corpora

- Web pages retrieved with a search engine, using individuals' labels as queries.
- $K_{DBP} : \approx 57\,000$  pages,  $K_F : \approx 6\,300$  pages

Problem

Intuition

Similarity

Plausibility

Trimming

**Experiments**

Extensions

Appendices

MF

Trimming :  
illustration

Trimming :  
assumptions

# Linguistic input

## Corpora

- Web pages retrieved with a search engine, using individuals' labels as queries.
- $K_{DBP} : \approx 57\,000$  pages,  $K_F : \approx 6\,300$  pages

## Linguistic contexts

- **LP** : (customized) sequences of surrounding lemma-POS (shifting window), frequencies weighted with PMI  
Limit : “more results about  $X$ ”, “more about  $X$  on Twitter”, ...
- **NP** : Ngrams preceding or following the term, frequencies weighted with PMI
- **NS** : Ngrams, frequencies weighted with self-information (querying the Microsoft Web N-gram corpus).
- **NPS** : Ngrams + PMI + self-information.

- Similarities : cosine distances



# Evaluation : plausibility

- Input :  $K_F$
- Generation of 100 random axioms  $\phi_1, \dots, \phi_{100}$  out of  $K_F$ .
- $K_1, \dots, K_{100}$  : 100 input KBs, such that  $K_i = K_F \cup \{\phi_i\}$ .
- For each  $K_i$ , order  $\Psi_{K_i}$  by plausibility.
- $\Psi_{K_i}^{rand} = \Psi_{K_i} - \Psi_{K_F}$ .

Problem

Intuition

Similarity

Plausibility

Trimming

**Experiments**

Extensions

Appendices

MF

Trimming :  
illustration

Trimming :  
assumptions

# Evaluation : plausibility

- Input :  $K_F$
- Generation of 100 random axioms  $\phi_1, \dots, \phi_{100}$  out of  $K_F$ .
- $K_1, \dots, K_{100}$  : 100 input KBs, such that  $K_i = K_F \cup \{\phi_i\}$ .
- For each  $K_i$ , order  $\Psi_{K_i}$  by plausibility.
- $\Psi_{K_i}^{rand} = \Psi_{K_i} - \Psi_{K_F}$ .

|     | rank         | p-val  |
|-----|--------------|--------|
| LP  | 4.15 / 216.1 | <0.001 |
| NP  | 9.73 / 216.1 | <0.001 |
| NS  | 7.33 / 216.1 | <0.001 |
| NPS | 5.59 / 216.1 | <0.001 |

Average ranking among  $\Psi_{K_i}$  of the lowest-ranked formula of  $\Psi_{K_i}^{rand}$ ,  
and p-value for the rankings of all formulas of all  $\Psi_{K_i}^{rand}$

- For most  $K_i$  (75/100),  $|\Psi_{K_i}^{rand}| = 1$ . In most of these cases (57/75), the only formula in  $\Psi_{K_i}^{rand}$  was also the one with lowest plausibility in  $\Psi_{K_i}$ .

# Evaluation : trimming

- For each  $K_i$ , the set  $\Delta_i = \Gamma_{i,1}, \dots, \Gamma_{i,1029}$  of all immediate subbase of  $K_i$  was computed.
- Within  $\Delta_i$ , all  $\Gamma_{i,j}$  such that  $\Psi_{\Gamma_{i,j}} \neq \Psi_{K_i}$  were ordered according to  $\prec$ .
- Weighting : LP (lemmaPos + PMI)

|                             | rank         | p-val     |
|-----------------------------|--------------|-----------|
| $\text{comp}(\Gamma)$       | 7.86 / 80.03 | $< 0.001$ |
| $\text{comp}_{K_i}(\Gamma)$ | 8.05 / 80.03 | $< 0.001$ |
| $\preceq_{lex}$             | 6.51 / 80.03 | $< 0.001$ |
| $\preceq_{lex_{K_i}}$       | 2.47 / 80.03 | $< 0.001$ |

Average ranking of the randomly generated statement  $\phi_i$  for each  $K_i$ ,  
and p-value for the rankings of all  $\phi_i$

# Evaluation : iterated trimming, $K_F$

- $K' = K_F$  extended with 20 random axioms
- $|K'| = 1028 + 20 = 1048$

|     |                   | val. | prec. & rec. | p-val (prop. test) |
|-----|-------------------|------|--------------|--------------------|
| NPS | comp              | 9    | 0.45         | $< 0.001$          |
|     | $\text{comp}_K$   | 9    | 0.45         | $< 0.001$          |
|     | $\preceq_{lex}$   | 3    | 0.15         | $< 0.002$          |
|     | $\preceq_{lex_K}$ | 9    | 0.45         | $< 0.001$          |
| LP  | comp              | 10   | 0.5          | $< 0.001$          |
|     | $\text{comp}_K$   | 10   | 0.5          | $< 0.001$          |
|     | $\preceq_{lex}$   | 5    | 0.25         | $< 0.001$          |
|     | $\preceq_{lex_K}$ | 10   | 0.5          | $< 0.001$          |

Table: Randomly generated axioms among the first 20 discarded ones

# Evaluation : iterated trimming, $K_{DBP}$

■  $|K_{DBP}| = 5721$

|     |                 | val. | prec. |
|-----|-----------------|------|-------|
| NPS | comp            | 7    | 0.35  |
|     | $\preceq_{lex}$ | 3    | 0.15  |
| LP  | comp            | 11   | 0.55  |
|     | $\preceq_{lex}$ | 5    | 0.25  |

Table: Actually erroneous axioms among the 20 first discarded ones

1 Problem

2 Intuition

3 Similarity

4 Plausibility

5 Trimming

6 Experiments

**7 Extensions**

8 Appendices

- MF

- Trimming : illustration

- Trimming : assumptions

Problem

Intuition

Similarity

Plausibility

Trimming

Experiments

**Extensions**

Appendices

MF

Trimming :  
illustration

Trimming :  
assumptions

# Extensions

## Complex concepts

- Most DLs allow the construction of arbitrary complex DL concepts, e.g.  $\exists \text{doctoralAdvisor}.\top$
- They could (in theory) be used instead of  $A$ .
- If  $\Psi_{\Gamma}^*$  is the set of all resulting consequences, no finite subset  $\Psi'$  of  $\Psi_{\Gamma}^*$  is such that  $\Psi_{\Gamma}^* \subseteq \text{Cn}(\Psi')$ .  
 $\Rightarrow$  Need to choose among these concepts.
- Some complex concepts are not relevant linguistically, e.g.  $(\text{Moldavian} \sqcup \text{Muslim}) \sqcap \text{Lawyer} \sqcap \exists \text{hasFather}.\forall \text{livesIn}.\text{Apartment}$

Problem  
Intuition  
Similarity  
Plausibility  
Trimming  
Experiments  
**Extensions**  
Appendices  
MF  
Trimming :  
illustration  
Trimming :  
assumptions

# Extensions

Problem  
Intuition  
Similarity  
Plausibility  
Trimming  
Experiments  
**Extensions**  
Appendices  
MF  
Trimming :  
illustration  
Trimming :  
assumptions

## Complex concepts

- Most DLs allow the construction of arbitrary complex DL concepts, e.g.  $\exists \text{doctoralAdvisor.T}$
- They could (in theory) be used instead of  $A$ .
- If  $\Psi_\Gamma^*$  is the set of all resulting consequences, no finite subset  $\Psi'$  of  $\Psi_\Gamma^*$  is such that  $\Psi_\Gamma^* \subseteq \text{Cn}(\Psi')$ .  
 $\Rightarrow$  Need to choose among these concepts.
- Some complex concepts are not relevant linguistically, e.g.  $(\text{Moldavian} \sqcup \text{Muslim}) \sqcap \text{Lawyer} \sqcap \exists \text{hasFather}.\forall \text{livesIn.Appartment}$

$e \neq e'$

- Set  $\Psi_\Gamma = \{\psi = e \neq e' \mid \Gamma \models \psi\}$
- Penalize  $\text{comp}(\Gamma)$  if  $\sim(e, e')$  is high.



- 1 Problem
- 2 Intuition
- 3 Similarity
- 4 Plausibility
- 5 Trimming
- 6 Experiments
- 7 Extensions
- 8 **Appendices**
  - MF
  - Trimming : illustration
  - Trimming : assumptions

# Modal fragment (MF) of FOL ( $= \mathcal{ALC}$ )

- Problem
- Intuition
- Similarity
- Plausibility
- Trimming
- Experiments
- Extensions
- Appendices
  - MF**
  - Trimming : illustration
  - Trimming : assumptions

- If  $A$  is a unary predicate, then  $A(x) \in \text{MF}$ .
- MF is closed under boolean operators.
- If  $\phi \in \text{MF}$ ,  $y$  does not appear in  $\phi$ , and  $R$  is a binary predicate, then :
  - $\exists y(R(x, y) \wedge \phi[x/y]) \in \text{MF}$
  - $\forall y(R(x, y) \rightarrow \phi[x/y]) \in \text{MF}$

# Trimming with $\preceq_{lex_K}$

## Example

$\Omega = \{$   
(1) `doctoralAdvisor`(*Thaddeus S.C. Lowe*, *Smithsonian Institution*),  
(2) `doctoralAdvisor`(*Nick Katz*, *Bernard Dwork*),  
(3)  $\top \sqsubseteq \forall \text{doctoralAdvisor. Person}$   
(4)  $\text{Organization} \sqsubseteq \neg \text{Person}$   
 $\}$

- Assume `doctoralAdvisor`, *Bernard Dwork* and *Smithsonian Institution* do not appear in  $\Gamma \setminus \Omega$ .
- Trimming :
  - discarding axioms in order to give up implausible consequences, but retain plausible ones.
  - no axiom should be unnecessarily discarded
- Only one solution here : discarding (1).

# Trimming : assumptions

- $\prec$  : strict partial order over  $2^K$  :  $\Gamma_1 \prec \Gamma_2$  iff either  $\text{comp}(\Gamma_1) < \text{comp}(\Gamma_2)$ , or  $(\text{comp}(\Gamma_1) = \text{comp}(\Gamma_2) \text{ and } \Gamma_1 \subset \Gamma_2)$ .
- Minimize **syntactic** information loss whenever possible, i.e.  $\Gamma_1$  and  $\Gamma_2$  viewed as bases, not theories.

In particular :

- If  $\text{Cn}(\Gamma_1) = \text{Cn}(\Gamma_2)$ , but  $\Gamma_1 \not\subseteq \Gamma_2$  and  $\Gamma_2 \not\subseteq \Gamma_1$ , then  $\Gamma_1$  and  $\Gamma_2$  are not comparable wrt  $\prec$ .
- Redundancies should be preserved when possible : if  $\text{Cn}(\Gamma_1) = \text{Cn}(\Gamma_2)$  and  $\Gamma_1 \subset \Gamma_2$ , then  $\Gamma_1 \prec \Gamma_2$  still holds.

# Lexicographic ordering $\preceq_{lex}$

- $\omega_\Gamma \doteq \omega_\Gamma^1, \dots, \omega_\Gamma^{|\Psi_\Gamma|}$  : formulas of  $\Psi_\Gamma$  order by increasing score  $sc_\Gamma$
- $sc_\Gamma(\omega_\Gamma) = sc_\Gamma(\omega_\Gamma^1), \dots, sc_\Gamma(\omega_\Gamma^{|\Psi_\Gamma|})$
- $\preceq_{lex}$  defined by  $\Gamma_1 \preceq_{lex} \Gamma_2$  iff either :
  - $sc_{\Gamma_1}(\omega_{\Gamma_1}) = sc_{\Gamma_2}(\omega_{\Gamma_2})$ , or
  - there is a  $1 \leq i \leq |\Psi_{\Gamma_2}|$  such that  $sc_{\Gamma_1}(\omega_{\Gamma_1}^i) = sc_{\Gamma_2}(\omega_{\Gamma_2}^i)$  for all  $1 \leq j < i$ , and either  $sc_{\Gamma_1}(\omega_{\Gamma_1}^i) < sc_{\Gamma_2}(\omega_{\Gamma_2}^i)$  or  $|\Psi_{\Gamma_1}| = i - 1$