

Harmonizing Tagsets for Multilingual Corpora via Concept Lattice

Alexandr Rosen

Institute of Theoretical and Computational Linguistics
Faculty of Arts, Charles University in Prague, Czech Republic

Dialogue 2010, Bekasovo
29 May 2010

The problem

- In a **multilingual parallel corpus**, each language comes with a potentially incompatible **morphosyntactic tagset**.
- The **variety of tagsets** is a problem for **human users** and for **applications**, too.

en	in IN	the DT	remotest JJS	exurbs NNS
de	in APPR	den ART	abgelegensten ADJA	Außenbezirken NN
nl	in 600	dit 370	schitterende 103	appartement 000
fr	dans PRP	les DET:ART	plus lointaines ADV ADJ	banlieues NOM
sp	en PREP	las ART	zonas NC	más remotas ADV ADJ
it	da PRE	queste PRO:demo	lingue NOM	babeliche ADJ
ru	v Sp-l	samych P—pl	otdaljonnych Afp-plf	rajonach Ncmpln
cs	v RR-6	těch PDXP6	nejodlehlejších AAFP6—3A	zástavbách NNFP6—A
bg	na R	tova Pde-os-n	prijatelstvo Ansi	dviženie Ncnsi
pl	w prep:loc:nwok	tym adj:sg:loc:m3:pos	wspaniałym adj:sg:loc:m3:pos	apartamencie subst:sg:loc:m3
hu	a ART	szép ADJ	katalán ADJ	lányba NOUN(CAS(ILL))

Reasons for the variety of tagsets

- Different **labels** for the same concept
- Typological differences among **languages**
- Different choice of a **theory** or its **aspect**

A solution

- Map each tag in every language to a **common abstract tagset**, a kind of tagset interlingua.
- Use *Formal Concept Analysis* to build a **hierarchy** of categories, linked to language-specific tags.

Linguistically motivated choices can be reflected in the common tagset

- In some tagsets, **adjectives**, **ordinal numerals** and **possessive pronouns** are all in a single class.
- In other tagsets, they are different **lexical** classes.
- Anyway they share the property of being a **syntactic** or **inflectional** adjective.

→ We propose a 3-way **cross-classification** of word classes.

Cross-classification

	lexical	inflectional	syntactic
ordinal numeral <i>pátý</i>	num	adj	adj
cardinal numeral <i>pět</i>	num	noun	noun
personal prn <i>ty</i>	prn	noun	noun
possessive prn <i>tvůj</i>	prn	adj	adj
relative prn <i>který</i>	prn	adj	noun
interrogative prn <i>který</i>	prn	adj	noun/adj
adverbial pple <i>volající</i>	verb	pple	advb

Further support for the cross-classification: sets of morphological categories appropriate to cross-classified POS

- Czech possessive pronouns show **pronoun–antecedent** agreement as lexical pronouns
- and **attribute–noun** agreement as syntactic adjectives

(1) Jana přišla, ale **jejího** syna jsem neviděl.
Jana_{FEM,3,NOM} came but her_{FEM,3} son_{MASC,ACC} I haven't seen
'Jana has come, but I haven't seen her son.'

Multiple inheritance hierarchy

- Sets of **categories** as partially ordered **types** in multiple inheritance **ISA hierarchy**.
- Each type **denotes** a set of language-specific **tags**.
- ‘Translational’ **ambiguities** due to mismatching tags are **not resolved but properly represented**.

Formal Concept Analysis (FCA)

- A **logical formalism** and **tool** to build the **hierarchy** and to do **reasoning** about categories.
- FCA classifies **objects (tags)** according to their **properties (categories)**.

Formal context for adjectives and cardinal/ordinal numerals

- A table to identify **objects** and their **attributes**.

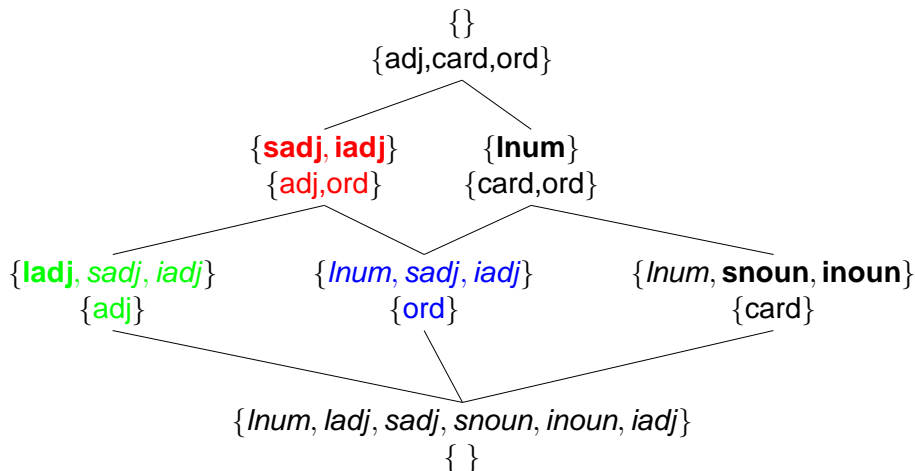
	<i>ladj</i>	<i>lnum</i>	<i>iadj</i>	<i>inoun</i>	<i>sadj</i>	<i>snoun</i>
adj	✓		✓		✓	
ord		✓	✓		✓	
card		✓		✓		✓

Formal concepts derived from the formal context table

1	$\langle\{\text{adj,ord,card}\}, \{\}\rangle$
2	$\langle\{\text{ord,card}\}, \{l\text{num}\}\rangle$
2	$\langle\{\text{adj,ord}\}, \{i\text{adj},s\text{adj}\}\rangle$
3	$\langle\{\text{adj}\}, \{l\text{adj},i\text{adj},s\text{adj}\}\rangle$
3	$\langle\{\text{ord}\}, \{l\text{num},i\text{adj},s\text{adj}\}\rangle$
3	$\langle\{\text{card}\}, \{l\text{num},i\text{noun},s\text{noun}\}\rangle$
4	$\langle\{\}\rangle$

- **Concept:** a set of **objects** and a set of **attributes**.
- Objects of a concept are also **part of its superconcept**.
- Concepts are **partially ordered** by specificity (roughly: the more attributes, the more specific).

Concept lattice for adjectives and ordinal numerals



Reasoning about attributes

- Possible **implications**: *ladj* \Rightarrow *sadj* or *snoun* \Rightarrow *Inum*
- ... useful for **general queries** (“show me all adjectives”)
- ... or to match **incompatible language-specific tags**

Problems

- Attributes in FCA formal contexts interpreted in **conjunction**, for **disjunctive attributes** we have to introduce a more **general attribute**, perhaps in a pre-processing step.
- The main task: manually specifying **formal context for all tags**.