

Manual annotation complexity grid

Karën Fort

karen.fort@univ-lorraine.fr / https://members.loria.fr/KFort





Dealing with the complexity of manual annotation

Analysing the complexity of an annotation campaign

What to annotate?

How to annotate?

Weight of the context

Synthesis: using the right tool, at the right place

WYMR: What You Must Remember

Manual annotation

Why is it easy/difficult? What part should be helped/automatized?

Manual annotation

Why is it easy/difficult? What part should be helped/automatized?

Gene renaming [Fort et al., 2012] :

The yppB gene complemented the defect of the recG40 strain. yppB and ypbC and their respective null alleles were termed recU and "recU1" (recU:cat) and recS and "recS1" (recS:cat), respectively. The recU and recS mutations were introduced into rec-deficient strains representative of the alpha (recF), beta (addA5 addB72), gamma (recH342), and epsilon (recG40) epistatic groups.

Manual annotation

Why is it easy/difficult? What part should be helped/automatized?

Structured named entities [Grouin et al., 2011] :

```
pers.ind pers.ind
name.first name.last
Lionel et Sylviane Jospin
```

What is complex?



Dealing with the complexity of manual annotation

Analysing the complexity of an annotation campaign

What to annotate?

How to annotate?

Weight of the context

Synthesis: using the right tool, at the right place

WYMR: What You Must Remember

Complexity dimensions

- ► 5 independent dimensions:
 - ▶ 2 related to the localisation of annotations
 - ▶ 3 related to the characterisation of annotations
- ▶ 1 not independent: the context



- ► Scale from 0 (null complexity) to 1 (maximal complexity) to allow for the comparison between campaigns
- ▶ Independent from the volume to annotate and the number of annotators

Example: gene renaming

- Identification of gene names in the source signal:
 The yppB gene complemented the defect of the recG40 strain. yppB and ypbC
 and their respective null alleles were termed "recU" and "recU1" (recU:cat)
 and "recS" and "recS1" (recS:cat), respectively.
- Identification of gene couples expressing a renaming relation:
 The yppB gene complemented the defect of the recG40 strain. yppB and ypbC
 and their respective null alleles were termed "recU" and "recU1" (recU:cat)
 and "recS" and "recS1" (recS:cat), respectively.

Dealing with the complexity of manual annotation

Analysing the complexity of an annotation campaign

What to annotate?

How to annotate?

Weight of the context

Synthesis: using the right tool, at the right place

WYMR: What You Must Remember

Discrimination

Parts-of-speech [Marcus et al., 1993], pre-annotated : I/PRP do/VBP n't/RB feel/VB very/RB ferocious/JJ ./.

Gene renaming[Fort et al., 2012], no pre-annotation:

The yppB:cat and ypbC:cat null alleles rendered cells sensitive to DNA-damaging agents, impaired plasmid transformation (25- and 100-fold), and moderately affected chromosomal transformation when present in an otherwise Rec+ B. subtilis strain. The yppB gene complemented the defect of the recG40 strain. yppB and ypbC and their respective null alleles were termed recU and "recU1" (recU:cat) and recS and "recS1" (recS:cat), respectively. The recU and recS mutations were introduced into rec-deficient strains representative of the alpha (recF), beta (addA5 addB72), gamma (recH342), and epsilon (recG40) epistatic groups.

Discrimination

Parts-of-speech [Marcus et al., 1993], pre-annotated : I/PRP do/VBP n't/RB feel/VB very/RB ferocious/JJ ./.

Gene renaming[Fort et al., 2012], no pre-annotation:

The yppB:cat and ypbC:cat null alleles rendered cells sensitive to DNA-damaging agents, impaired plasmid transformation (25- and 100-fold), and moderately affected chromosomal transformation when present in an otherwise Rec+ B. subtilis strain. The yppB gene complemented the defect of the recG40 strain. yppB and ypbC and their respective null alleles were termed recU and "recU1" (recU:cat) and recS and "recS1" (recS:cat), respectively. The recU and recS mutations were introduced into rec-deficient strains representative of the alpha (recF), beta (addA5 addB72), gamma (recH342), and epsilon (recG40) epistatic groups.

 \Rightarrow more difficult if the units to annotate are scattered, in particular if the segmentation is not obvious.

Discrimination

The discrimination weight is all the more high as the proportion of what *should* be annotated as compared to what *could* be annotated is low.

Definition

$$\textit{Discrimination}(\textit{Flow}) = 1 - \frac{|\textit{Annotations}(\textit{Flow})|}{\sum_{i=1}^{\textit{LevelSeg}} |\textit{UnitsObtainedBySeg}_i(\textit{Flow})|}$$

⇒ Need for a reference segmentation

Parts-of-speech[Marcus et al., 1993] : I/PRP do/VBP n't/RB feel/VB very/RB ferocious/JJ ./.

 $Discrimination_{PTB_{POS}} = 0$

Gene renaming[Fort et al., 2012] :

The yppB:cat and ypbC:cat null alleles rendered cells sensitive to DNA-damaging agents, impaired plasmid transformation (25- and 100-fold), and moderately affected chromosomal transformation when present in an otherwise Rec+ B. subtilis strain. The yppB gene complemented the defect of the recG40 strain. yppB and ypbC and their respective null alleles were termed recU and "recU1" (recU:cat) and recS and "recS1" (recS:cat), respectively. The recU and recS mutations were introduced into rec-deficient strains representative of the alpha (recF), beta (addA5 addB72), gamma (recH342), and epsilon (recG40) epistatic groups.

 $Discrimination_{Identification} = 0,9$ $Discrimination_{Renaming} = 0,95$

► extending or shrinking the discriminated unit:
Madame Chirac → Monsieur et Madame Chirac

- ► extending or shrinking the discriminated unit:
 Madame Chirac → Monsieur et Madame Chirac
- ▶ decompose a discriminated unit into several elements: le préfet Érignac → le préfet Érignac

- ► extending or shrinking the discriminated unit:
 Madame Chirac → Monsieur et Madame Chirac
- ▶ decompose a discriminated unit into several elements: le préfet Érignac → le préfet Érignac
- or group together several discriminated units into one unique annotation: Sa Majesté le roi Mohamed VI → Sa Majesté le roi Mohamed VI

Definition

$$Delimitation(Flow) = min\left(rac{Substitutions + Additions + Deletions}{|Annotations(Flow)|}, 1
ight)$$

 $Delimitation_{Identification} = 0$ $Delimitation_{Renaming} = 0$

 $Delimitation_{PTB_{POS}} = 0$

 $\begin{array}{l} \textit{D\'elimitation}_{\textit{EN}_{\textit{TypesSubtypes}}} = 1 \\ \textit{D\'elimitation}_{\textit{EN}_{\textit{Components}}} = 0, 3 \end{array}$

Dealing with the complexity of manual annotation

Analysing the complexity of an annotation campaign

What to annotate?

How to annotate?

Weight of the context

Synthesis: using the right tool, at the right place

WYMR: What You Must Remember

Expressiveness of the annotation language

Definition

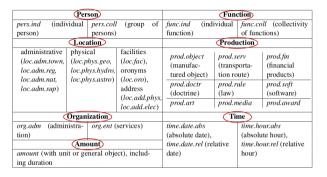
The degrees of expressiveness of the annotation language are the following:

- ▶ 0.25: type languages
- ▶ 0.5: relational languages of arity 2
- ▶ 0.75: relational languages of arity higher than 2
- ► 1: higher-order languages

 $Expressiveness_{Identification} = 0.25$ $Expressiveness_{Renaming} = 0.25$

Person					Function				
1	lividual	1		of	J	ividual	func.col	,	
person) persons)					function) of functions)				
Location					Production				
administrative (loc.adm.town, loc.adm.reg, loc.adm.nat, loc.adm.sup)		facilities (loc.fac), oronyms (loc.oro), address (loc.add.ploc.add.e	ohys,	tured object) tion prod.doctr (doctrine) (law		oorta- oute) rule	prod.fin (financial products) prod.soft (software) prod.award		
Organization					Time				
org.adm (administra- tion) org.ent (sec			services)	time.date.abs (absolute date), time.date.rel (re			time.hour.abs (absolute hour), tive time.hour.rel (relat		
amount (with unit or general object), including duration					date)		hour)		

Types and sub-types used for structured NE annotation [Grouin et al., 2011]



Level 1: pers, func, loc, prod, org, time, amount \rightarrow 7 possibilities (degree of freedom = 6).

Person					Function				
pers.ind (ind	ividual pers.coll		(group	of	func.ind (ind	ividual	func.coll	(collectivity	
person)	person) persons)		,		function) of functions)			ons)	
Location					Production				
administrative (loc.adm.town, loc.adm.nat, loc.adm.sup) physical (loc.phys.as loc.phys.as		hys.geo, iys.hydro,	facilities (loc.fac), oronyms (loc.oro), address (loc.add.phys), loc.add.elec)		prod.object (manufac- tured object) prod.doctr (doctrine) prod.art	(transportation route) prod.rule (law) (prod.fin (financial products) prod.soft (software) prod.award	
Organization					Time				
org.adm (administra- org.ent (tion)		(services)		time.date.abs (absolute date),		time.hour.abs (absolute hour),			
Amount					time.date.rel (relative time.		time.hou	hour:rel (relative	
amount (with unit or general object), including duration					date)		hour)		

Level 1: pers, func, loc, prod, org, time, amount \rightarrow 7 possibilities (degree of freedom = 6).

Level 2: prod.object, prod.serv, prod.fin, prod.soft, prod.doctr, prod.rule, prod.art, prod.media, prod.award o 9 possibilities (degree of freedom = 8).

Person					Function				
pers.ind (ind	ividual	pers.coll	(group	of	func.ind (individual func.coll (collect		(collectivity		
person)		persons)			function)	of functions)		ons)	
Location					Production				
administrative (loc.adm.rown, loc.adm.reg, loc.adm.nat, loc.adm.sup)		facilities (loc.fac), oronyms (loc.oro), address (loc.add.ploc.add.e	ohys,	prod.object (manufac- tured object) prod.doctr (doctrine) prod.art	prod.s (transp tion re prod.r (law) prod.r	oorta- oute)	prod.fin (financial products) prod.soft (software) prod.award		
Organization					Time				
org.adm (administration) org.ent ((services)		time.date.abs (absolute date), time.date.rel (relative		time.hour.abs (absolute hour), time.hour.rel (relative		
amount (with unit or general object), including duration				date)		hour)			

Level 1: pers, func, loc, prod, org, time, amount \rightarrow 7 possibilities (degree of freedom = 6).

Level 2: prod.object, prod.serv, prod.fin, prod.soft, prod.doctr, prod.rule, prod.art, prod.media, prod.award o 9 possibilities (degree of freedom = 8).

Level 3: loc.adm.town, loc.adm.reg, loc.adm.nat, $loc.adm.sup \rightarrow$ 4 possibilities (degree of freedom = 3).

24 / 35

Degree of freedom

$$\nu = \nu_1 + \nu_2 + \ldots + \nu_m$$

where ν_i is the maximal degree of freedom the annotator has when choosing the i^{th} sub-type $(\nu_i = n_i - 1)$.

Dimension (size) of the tagset

$$Dimension(Flow) = min(rac{
u}{ au}, 1)$$

where τ is the threshold from which we consider the tagset to be very large (experimentally determined).

$$Dimension_{Identification} = 0$$

 $Dimension_{Renaming} = 0.04$
 $Dimension_{NE_{TypesSubtypes}} = 0.34$

Degree of ambiguity: residual ambiguity

Using the traces left by the annotators:



```
[...] <EukVirus>3CDproM</EukVirus> can process both structural and nonstructural precursors of the <EukVirus uncertainty-type = "toogeneric"><taxon>poliovirus</taxon> polyprotein</EukVirus> [...].
```

Définition

$$AmbiguityRes(Flow) = \frac{|Annotations_{amb}|}{|Annotations|}$$

$$AmbiguityRes_{Identification} = 0.04$$

 $AmbiguityRes_{Renaming} = 0.02$

Degree of ambiguity: theoretical ambiguity

Proportion of the units to annotate that corresponds to ambiguous vocables.

Definition

$$AmbiguityTh(Flow) = \frac{\sum_{voc_i=1}^{|Voc(Flow)|} (Ambig(voc_i) * freq(voc_i, Flow))}{|Units(Flow)|}$$

with

$$Ambig(voc_i) = \left\{ egin{array}{ll} 1 & ext{if} & |Tags(voc_i)| > 1 \ 0 & ext{else} \end{array}
ight.$$

$$AmbiguityTh_{Identification} = 0.01$$

ightarrow Does not apply to renaming relations

Dealing with the complexity of manual annotation

Analysing the complexity of an annotation campaign

What to annotate?

How to annotate?

Weight of the context

Synthesis: using the right tool, at the right place

WYMR: What You Must Remember

Context to take into account

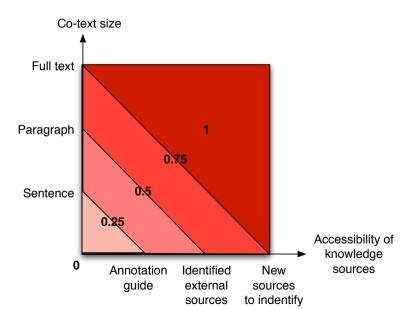
- **size of the window** to take into account in the source signal:
 - ► The sentence:

 I/PRP do/VBP n't/RB feel/VB very/RB ferocious/JJ ./.



- ▶ number of knowledge elements to be rallied or degree of accessibility of the knowledge sources that are consulted:
 - annotation guidelines
 - nomenclatures (Swiss-Prot)
 - new sources to be found (Wikipedia, etc.)

Weight of the context



Dealing with the complexity of manual annotation

Analysing the complexity of an annotation campaign

What to annotate?

How to annotate?

Weight of the context

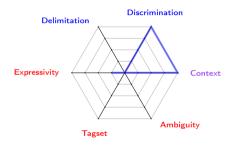
Synthesis: using the right tool, at the right place

WYMR: What You Must Remember

Synthesis of the complexity dimensions



Classification of it pronouns as anaphoric or impersonal



Gene names identification



Manual annotation and NLP:

- usage
- ► cost

Manual annotation:

- definition
- organization
- complexity grid

Practice: Compute the complexity dimensions for your annotations

Assignment: Reading a research paper (graded)

Read carefully and email your teacher the synopsis (2p max, pdf) by next week

```
[Dandapat et al., 2009] :
Cheap and Fast – But is it Good?
Evaluating Non-Expert Annotations for Natural Language Tasks
https://aclanthology.org/D08-1027.pdf
```

Dandapat, S., Biswas, P., Choudhury, M., and Bali, K. (2009).

Complex linguistic annotation - no easy way out! a case from bangla and hindi POS labeling tasks.

In Proceedings of the third ACL Linguistic Annotation Workshop, Singapour.

Fort, K., François, C., Galibert, O., and Ghribi, M. (2012).

Analyzing the impact of prevalence on the evaluation of a manual annotation campaign.

In Proceedings of the International Conference on Language Resources and Evaluation (LREC), Istanbul, Turquie.

7 pages.

Grouin, C., Rosset, S., Zweigenbaum, P., Fort, K., Galibert, O., and Quintard, L. (2011).

Proposal for an extension of traditional named entities: From guidelines to evaluation, an overview.

In <u>Proceedings of the 5th Linguistic Annotation Workshop</u>, pages 92–100, Portland, Oregon, USA.

Poster.

Marcus, M., Santorini, B., and Marcinkiewicz, M. A. (1993).

Building a large annotated corpus of English: The Penn Treebank.

Computational Linguistics, 19(2):313-330.