

Corpus Annotation

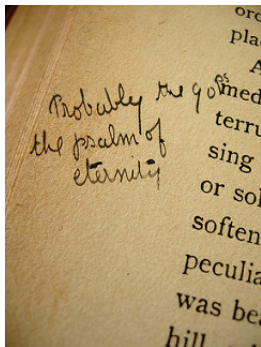
Béatrice Daille

Faculté des Sciences et Techniques de Nantes
Département Informatique

2018-2019

- Corpus : exemples, définition, caractérisation
- Métadonnées : Dublin Core, TEI
- Web comme corpus
- **Annotation des corpus**
- Exploration des corpus : phrasèmes, collocations

Annotation : définition



L'annotation recouvre à la fois le processus consistant à apposer (ad-) une note sur un support, l'ensemble des notes ou chaque note particulière qui en résulte et ce, sans préjuger a priori de la nature du support considéré (texte, vidéo, images, etc.), du contenu sémantique de la note (note chiffrée, valeur choisie dans un référentiel fermé ou texte libre), de son positionnement global ou local, ni de son objectif (visée évaluative ou caractérisante, simple commentaire discursif).

Objectifs

Utilité des annotations

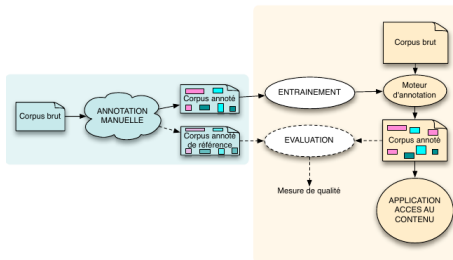
TALN

- ① Acquisition : classifieur construit par apprentissage supervisé
- ② Évaluation : annotation de référence pour évaluer un algorithme

Linguistique

- ① Exploration d'un corpus
- ② Étude d'un phénomène linguistique à grande échelle

Annotation pour le TAL



Architecture schéma d'annotation

Identique à l'architecture d'un SGBD

- ① Niveau externe : visualisé et manipulé par les utilisateurs
- ② Niveau intermédiaire : schéma logique (spécification des annotations)
- ③ Niveau physique : codage et stockage des annotations (annotations déportées ou insérées, XML et formats linéaires)

Des ressources diverses

Diversité

- médias : texte, parole, musique, vidéo
- champs d'applications : linguistique, domaines spécialisés

Complexité

- 1960-1990 : morpho-syntaxe
- 1990-2000 : morpho-syntaxe et syntaxe (corpus arborés), entités nommées simples, sémantique (WordNet, FrameNet)
- depuis 2000 : annotations sémantiques variées (opinions, émotions, etc.), entités nommées structurées, etc.

Taille

- Brown Corpus, 1 million de mots
- BNC, 100 millions de mots

Types d'annotation

Extrême diversité des annotations

- niveau document/texte ou signal
- expressivité (booléen, catégories, relations)
- complexité : phénomènes linguistiques complexes, ontologies

Annotations au niveau du document

Intégration dans les méta-données

Profilage du texte : genre, discours, positionnement, niveau de langue, domaine

Annotations au niveau du texte

Enrichissements linguistiques :

- Segmentation : phrasèmes, entités nommées
- Catégories grammaticales
- Lemmatisation
- Analyse syntaxique

Catégories grammaticales universelles

Jeu de catégories grammaticales validées sur 22 langues
Tables de correspondance pour 27 jeux d'étiquettes

VERB verbes (tous les temps et modes)

NOUN noms (commun et propre)

PRON pronoms

ADJ adjectifs et numéraux ordinaux

ADV adverbes

ADP adpositions (prépositions and postpositions)

CONJ conjonctions

DET déterminants (regroupent les articles et adjectifs possessifs, démonstratifs, interrogatifs, exclamatifs ou numéraux cardinaux)

NUM chiffres

PRT particules ou autres mots fonctionnels

X autres : mots étrangers, typos, abréviations, interjections

. ponctuation

<http://code.google.com/p/universal-pos-tags/>

"A Universal Part-of-Speech Tagset" by Slav Petrov, Dipanjan Das and Ryan McDonald LREC 2012

Exemple annotation catégories grammaticales

Text				
Cyndi savoredthe soup.				
^0...^5...^10...^15...^20				
Annotations				
Id	Type	SpanStart	Span End	Features
1	token	0	5	pos=NP
2	token	6	13	pos=VBD
3	token	14	17	pos=DT
4	token	18	22	pos=NN
5	token	22	23	
6	name	0	5	name_type=person
7	sentence	0	23	constituents=[1],[2],[3],[4],[5]

Corpus arborés

- Grammaires de constituants
- simple parenthésage des constituants
[Ce guide][[leur] permet [de [se familiariser [. . .]
- étiquetage des constituants
[N Ce_DEDEMMS guide_NCOMS N][
[P leur_PPCA6MP P] permet_V_VINIP3
[P de_PREPD [Vi se_PPPE6MP familiariser_VPRN [. . .]

Corpus arborés

- Grammaires de dépendances
- indication des relations de dépendances
Ce @DN>
guide @NV2>
leur @PV>
permet
[...]
- Corpus Suzanne

Corpus arborés

- Grammaires de constituants et grammaires de dépendances
- indication des relations fonctionnelles
 [N <Sujet> Ce_DEDEMMS guide_NCOMS N]
 [V [P <ObjetIndirect> leur_PPCA6MP P] permet V_VINIP3 [...]]
- relations logiques ou profondes
 [N Ce_DEDEMMS guide_NCOMS N]
 [V [P [P <8> leur_PPCA6MP P] permet V_VINIP3
 [P de_PREPD [N <8> N] [Vi se_PPPE6MP familiariser_VPRN [...]]
- Pen Treebank

Corpus arborés : Penn Treebank

- 4,8 M de tokens en anglais américain (Brown corpus ré-annoté)
- annotation en morpho-syntaxe : 48 catégories (12 ponctuations et symboles)
- annotation en syntaxe : 15 catégories

Dépendances universelles

37 dépendances élémentaires

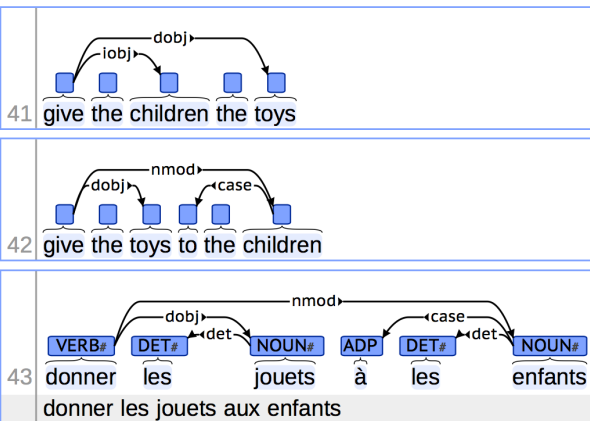
Clausal Argument Relations	Description
NSUBJ	Nominal subject
DOBJ	Direct object
IOBJ	Indirect object
CCOMP	Clausal complement
XCOMP	Open clausal complement
Nominal Modifier Relations	Description
NMOD	Nominal modifier
AMOD	Adjectival modifier
NUMMOD	Numeric modifier
APPOS	Appositional modifier
DET	Determiner
CASE	Prepositions, postpositions and other case markers
Other Notable Relations	Description
CONJ	Conjunct
CC	Coordinating conjunction

Figure 14.2 Selected dependency relations from the Universal Dependency set. (de Marneffe et al., 2014)

<http://universaldependencies.org/u/dep/index.html>

"Universal Stanford Dependencies : A cross-linguistic typology" by Marie-Catherine de Marneffe, Timothy Dozat, Natalia Silveira, Katri Haverinen, Filip Ginter, Joakim Nivre, Christopher D. Manning. LREC 2014

Exemples



Annotations au niveau du texte

Phénomènes linguistiques complexes :

- Lexèmes, termes
- Anaphores
- Dialogue
- Expressions d'évaluations

ISO-TimeML

Expressions temporelles

les événements dans le temps présents dans les articles de presse, les biographies, Wikipédia, etc.

TIMEX3 attribut TYPE

DATE le 15 janvier 2019, hier, l'été 2006

TIME 11h30, treize heures quinze, la nuit dernière

DURATION 2 mois, trois jours,

SET deux fois par mois, tous les 2 jours

Tout l'empan de l'expression référentielle est marqué.

ISO-TimeML

Les événements

la notion d'événement recouvre tous les types de situations (états, activités, achèvements, etc.) EVENT attribut CLASS

OCCURRENCE vendre, construire, tuer

STATE amoureux, fatigué

REPORTING dire, annoncer

I-Action essayer, promettre

I-State croire, vouloir

Aspectual commencer, terminer, arrêter

Perception voir, entendre, sentir

autres attributs : TYPE, TENSE, ASPECT, POLARITY, MODALITY
balise sur la tête du groupe événementiel, en excluant les auxiliaires, les modifieurs, les adverbes de négation, les clitiques, etc.

ISO-TimeML

Les relations

TLINK relation temporelle entre un événement et une date.

ALINK relation aspectuelle. Par exemple, entre un verbe aspectuel (commencer, cesser, continuer) et son complément événementiel.

SLINK relation de subordination. Par exemple, la relation qui existe entre un verbe modal (falloir, devoir) ou de perception (voir, entendre) et son complément événementiel.

MLINK relation pour mesurer la durée d'un événement

attribut **RELTYPE** pour spécifier la relation :

TLINK : BEFORE, IS_INCLUDED

ISO-TimeML

Les marqueurs de temporalité

SIGNAL

prépositions temporelles sur, dans, avant

conjonction de subordination temporelles avant

caractères spéciaux - et /

le marqueur est identifié lors d'une relation TLINK

Exemple

Jean est <EVENT id="e1" class="OCCURRENCE" pos="VERB"
 tense="PAST" vform="PASTPART" > *né* </EVENT>
 <SIGNAL id="s1" > *avant* < /SIGNAL>
I' <EVENT id="e2" class="OCCURRENCE" pos="NOUN" >
introduction < /EVENT> *de l'euro.*
 <TLINK id="l1" eventID="e1" relatedToEvent="e2" signalID="s1"
 relType="BEFORE" / >

Annotations au niveau du texte

Correction

- Normalisation : détection des erreurs, formes non standards
- Anonymisation

Coût annotation

Penn Tree Bank (PTB)

PTB1 (Marcus et al. 1993)

- correction de l'étiquetage morpho-syntaxique : 3 000 mots à l'heure, 3 H par jour
- correction de l'étiquetage syntaxique : 750 mots à l'heure, 3 H par jour
- + courbe d'apprentissage de 1 mois (étiquettes morpho-syntaxiques) à 2 mois (syntaxe)

Coût annotation

Prague Dependency Treebank

- 1996-2004 (Böhmová et al. 2001)
- construit à partir du CNC (Czech National Corpus)
- 3 niveaux de structures
 - ① morphologique : 1,8 millions de mots
 - ② analytique : syntaxe en dépendances
 - ③ sémantique

Coût annotation

Prague Dependency Treebank

Version 1.0

- annotation manuelle des niveaux morphologique et analytique
- temps : 5 ans
- nombre de personnes : 22 dont 17 simultanément pendant certaines périodes
- coût : 600 000 \$

Coût annotation

GENIA

- (Kim et al. 2008)
- 400 000 mots annotés en microbiologie
- 5 annotateurs à mi-temps, 1 coordinateur senior 1 coordinateur junior pendant 1,5 an

Coût annotation

ESTER

- 100 heures de parole transcrite (campagne d'évaluation ESTER, systèmes de transcription, 2008)
- 1 h de parole = entre 20 à 60 H de transcription

Consensus au cœur de l'annotation

- guide d'annotation
- réunion avec les annotateurs, plusieurs phases : mise au point, vérification de l'applicabilité du guide d'annotation
- mesurer le consensus, garantie d'une cohérence

Dimensions de complexité de l'annotation

(Fort et al. 2012)

- ① discrimination des unités à annoter
- ② délimitation des unités à annoter
- ③ expressivité du langage d'annotation
- ④ jeu d'étiquettes
- ⑤ ambiguïté
- ⑥ contexte textuel à consulter, connaissances à mobiliser

Les guides d'annotation

Des ressources précieuses

Limiter le coût de l'annotation

- outils d'annotation : WebAnno, Glozz, GATE, Knowtator, Callisto, Prat, etc. → pas toujours adaptés
- propagation d'étiquettes, pré-annotation, apprentissage actif
- myriadisation : Amazon Mechanical Turk et jeux ayant un but

Mesures d'accord entre annotateurs

- Mesures d'accord entre deux annotateurs
 - accord observé A_0
 - accord attendu A_e
 - S de Scott A_e^S
 - π de Scott A_e^π
 - κ de Cohen A_e^κ
- Interprétation de la valeur des accords
- Accord inter-annotateurs personnalisé
- Analyse de l'accord inter-annotateur global

Accord observé

Pourcentage d'accord ou accord observé A_o

$$A_o = \frac{1}{i} \sum_{i \in I} \text{agr}_i \quad (1)$$

avec i : nombre total d'éléments annotables
et agr : accord sur l'annotation

Limitation :

- Pas de prise en compte du hasard : un petit nombre de catégories donne un meilleur accord observé
- Pas de compensation en fonction de la distribution des éléments dans les catégories : une catégorie prépondérante influence largement l'accord observé

Exemple accord observé

Annotation du registre de langue

Matrice de confusion : 2 annotateurs : Paul et Marie, 3 catégories : Familier, Courant, Soutenu, 3 fichiers

		Paul	Paul	Paul	
		F	C	S	Somme
Marie	F	0	0	0	0
Marie	C	1	2	0	3
Marie	S	0	0	0	0
Somme		1	2	0	3

$$A_o = \frac{1}{3} \cdot 2 = 0,66$$

Accord attendu

Prise en compte que les annotateurs classent un élément quelconque dans une catégorie au hasard A_e

$$S_{\pi\kappa} = \frac{A_o - A_e}{1 - A_e} \quad (2)$$

S : annotations réalisées au hasard suivent une distribution uniforme dans les catégories (Bennett et al. 1954)

$$A_o^S = \frac{1}{k} \quad (3)$$

k : nombre de catégories k

Calcul sur exemple précédent : $A_o^S = \frac{1}{3} = 0,33$

$$S = \frac{0,66 - 0,33}{1 - 0,33} = 0,49$$

- Problème : un petit nombre de catégories donne un meilleur accord observé

π de Scott

Accord attendu en supposant que la répartition des éléments dans les catégories n'est pas homogène A_e^π (Scott, 1955)

Estimation de la répartition moyenne réalisée par les annotateurs

$$A_e^\pi = \sum_{k \in K} \left(\frac{n_k}{2i} \right)^2 \quad (4)$$

$$\pi = \frac{A_o - A_e^\pi}{1 - A_e^\pi} \quad (5)$$

avec n_k : nombre d'affectations à k par les deux annotateurs

i : nombre d'éléments annotables i

k : une catégorie

Exemple calcul Scott

		Paul	Paul	Paul	
		F	C	S	
Marie	F	a	b	c	C_{Marie}^F
Marie	C	d	e	f	C_{Marie}^C
Marie	S	g	h	j	C_{Marie}^S
		C_{Paul}^F	C_{Paul}^C	C_{Paul}^S	

		P	P	P	
		F	C	S	
M	F	0	0	0	0
M	C	1	2	0	3
M	S	0	0	0	0
		1	2	0	

$$C_{Marie}^F = a + b + c$$

$$\mathbf{n}_F = C_{Marie}^F + C_{Paul}^F$$

$$A_e^\pi = \frac{(\frac{n_F}{2})^2 + (\frac{n_C}{2})^2 + (\frac{n_S}{2})^2}{i^2} = \frac{(\frac{1+0}{2})^2 + (\frac{3+2}{2})^2 + 0}{9} = 0,72$$

$$\pi = \frac{0,66 - 0,72}{1 - 0,72} = -0,23$$

κ de Cohen

Modélisation du hasard A_o^κ (Cohen 1960)

Répartition des éléments entre catégories peut être différent pour chaque annotateur

$$A_e^\kappa = \sum_{k \in K} \frac{n_{c1k}}{i} \cdot \frac{n_{c2k}}{i} \quad (6)$$

$$\kappa = \frac{A_o - A_e^\kappa}{1 - A_e^\kappa} \quad (7)$$

avec n_{c1k} : nombre d'affectations à k par l'annotateur $c1$

i : nombre d'éléments annotables i

Exemple calcul Cohen

		Paul	Paul	Paul
		F	C	S
Marie	F	a	b	c
Marie	C	d	e	f
Marie	S	g	h	j

		Paul	Paul	Paul
		F	C	S
Marie	F	0	0	0
Marie	C	1	2	0
Marie	S	0	0	0

$$A_o = \frac{a+e+i}{a+b+\dots+j} = \frac{2}{3} = 0,66$$

$$\mathbf{n}_{c1F} = \frac{C_{Marie}^F}{i} = \frac{a+b+c}{a+b+\dots+j} = \frac{0}{3} = 0$$

$$\mathbf{n}_{c1C} = \frac{C_{Marie}^C}{i} = \frac{d+e+f}{a+b+\dots+j} = \frac{3}{3} = 1$$

$$A_e^\kappa = \mathbf{n}_{c1F} \cdot \mathbf{n}_{c2F} + \mathbf{n}_{c1C} \cdot \mathbf{n}_{c2C} + \mathbf{n}_{c1S} \cdot \mathbf{n}_{c2S} = 0 + 0,66 + 0 = 0,66$$

$$\kappa = \frac{0,66-0,66}{1-0,66} = 0$$

$$\mathbf{n}_{c2F} = \frac{C_{Paul}^F}{i} = \frac{a+d+g}{a+b+\dots+j} = \frac{1}{3} = 0,33$$

$$\mathbf{n}_{c2C} = \frac{C_{Paul}^C}{i} = \frac{b+e+h}{a+b+\dots+j} = \frac{2}{3} = 0,66$$

Interprétation

Valeur prise dans $[-\frac{A_e}{1-A_e}, 1]$ avec 1 accord total

- Ordre $S \geq \pi$ et $\pi \leq \kappa$
- si $\pi < \kappa$ et π et κ proches : peu de biais entre les annotateurs

$0 \leq 0,2$	faible	Landis et Koch (1977)
$0,2 \leq 0,4$	passable	
$0,4 \leq 0,6$	moyen	
$0,6 \leq 0,8$	substantiel	
$0,8 \leq 0,1$	parfait	
$0 \leq 0,67$	à ignorer	Krippendorff (1980)
$0,67 \leq 0,8$	indécis	
$0,8 \leq 1$	bon	
$0 \leq 0,4$	faible	Green (1997)
$0,4 \leq 0,75$	correct/bon	
$0,75 \leq 1$	élevé	

Interprétation

S : ne dépend que du nombre de catégories

→ n'est pas sensible à la répartition des éléments dans les catégories

π et κ : lorsque les catégories sont fortement disproportionnées, et même si fort accord sur la catégorie prédominante, très sensibles aux désaccords sur les catégories minoritaires

multi π ou Fleiss- κ

généralisation de π de Scott à plusieurs annotateurs

le nombre de paires d'annotateurs en accord par rapport à toutes les paires de jugements possibles pour l'élément considéré

$$agr_i = \frac{1}{c(c-1)} \sum_{k \in K} \mathbf{n}_{ik}(\mathbf{n}_{ik} - 1) \quad (8)$$

$$A_o^{\pi-m} = \frac{1}{i} \sum_{i \in I} agr_i \quad (9)$$

$$A_e^{\pi-m} = \frac{1}{(ic)^2} \sum_{k \in K} (\mathbf{n}_k)^2 \quad (10)$$

$$\pi - m = \frac{A_o^{\pi-m} - A_e^{\pi-m}}{1 - A_e^{\pi-m}} \quad (11)$$

c : nombre d'annotateurs c

Exemple de calcul de Fleiss- κ

4 annotateurs, 3 catégories, 3 fichiers

n_{ik}	F	C	S	P_c
Fic1	0	4	0	$P_{Fic1} = \frac{1}{4,3}(0 + 4 * 3 + 0)$
Fic2	1	3	0	$P_{Fic2} = \frac{1}{4,3}(1 * 0 + 3 * 2 + 0)$
Fic3	0	3	1	$P_{Fic3} = \frac{1}{4,3}(0 + 3 * 2 + 0)$
Somme	1	10	1	12
P_k	$P_F = \frac{1}{12}$ 0,08	$P_C = \frac{10}{12}$ 0,83	$P_S = \frac{1}{12}$ 0,08	

$$\sum_{i \in I} P_i = 1 + 0,5 + 0,5 = 2$$

$$A_e^{\pi-m} = \frac{1^2 + 10^2 + 1^2}{(3 \cdot 4)^2} = \frac{112}{144} = 0,77$$

$$A_o^{\pi-m} = \frac{1}{3} \cdot 2 = 0,67$$

$$\pi - m = \frac{0,67 - 0,77}{1 - 0,77} = -0,30$$

Hypothèses

- Les annotateurs annotent tous les éléments → pas toujours le cas en situation réelle
- La distance entre les catégories est identique → des classes peuvent être plus proches que d'autres
possibilité de pondérer les catégories : indice α (Krippendorff, 2004)
définition d'une matrice de pondération adaptée aux catégories
par exemple : pour l'annotation du registre, un désaccord 0 entre familier et soutenu ; 0,5 entre courant et soutenu ou courant et familier

$$A_e^{\pi-pm} = \frac{1}{(ic)^2} \sum_{k \in K} \sum_{k' \in K} \text{pond}(k, k') \mathbf{n}_k \mathbf{n}'_k$$

- Tous les annotateurs n'ont pas la même expertise → comparaison des accords entre paires d'annotateurs et éventuellement sur certaines catégories

Références

La bible

Handbook of Linguistic Annotation. Ide and Pustejovsky (eds). Springer 2017.

Mesures d'accords

Ron Artstein et Massimo Poesio : Inter-coder agreement for computational linguistics. Computational Linguistics , 34(4) :555–596, 2008. ISSN 0891-2017.