

Towards linguistics with syntax-analysed data and command-line tools

Veronika Laippala
veronika.laippala@utu.fi



TURKUNLP
.ORG



**UNIVERSITY
OF TURKU**

Topic of the day

- Using syntax-analysed data as the basis of linguistic analysis
 - A lot of potential for linguistics
- Command line
 - Much more flexible than corpus software
 - Allows for more advanced and efficient processing
 - Standard solutions eventually make their way to software, but new and more complex methods are not (necessarily) available there
- Apply these to some typical linguistic research settings

Steps for today

1. “Theory”
2. “Practice”

Note:

- See the references for more information on the topics!
- Feel free to select the parts you do from the notebooks (we cannot do all of them during this session)
- The practice section done on Google colab
- You’ll need Google credentials

Steps for today

1. “Theory”

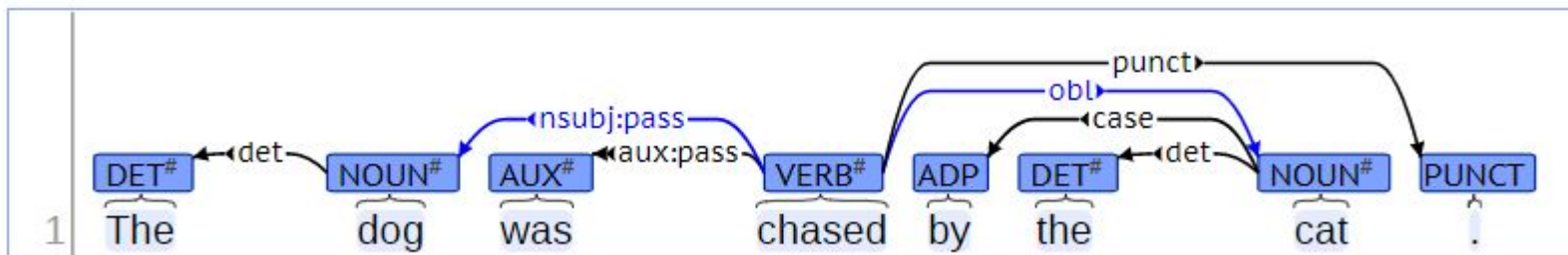
- a. Universal dependencies
- b. Type A and Type B studies (Biber and Jones 2009)
- c. Keyness (and text classification)

2. “Practice”

- a. Processing data on command line in the syntax-analysed Conllu format
- b. Keyness
 - i. “Standard”
 - ii. Text dispersion
 - iii. Linear classifier

Universal dependencies

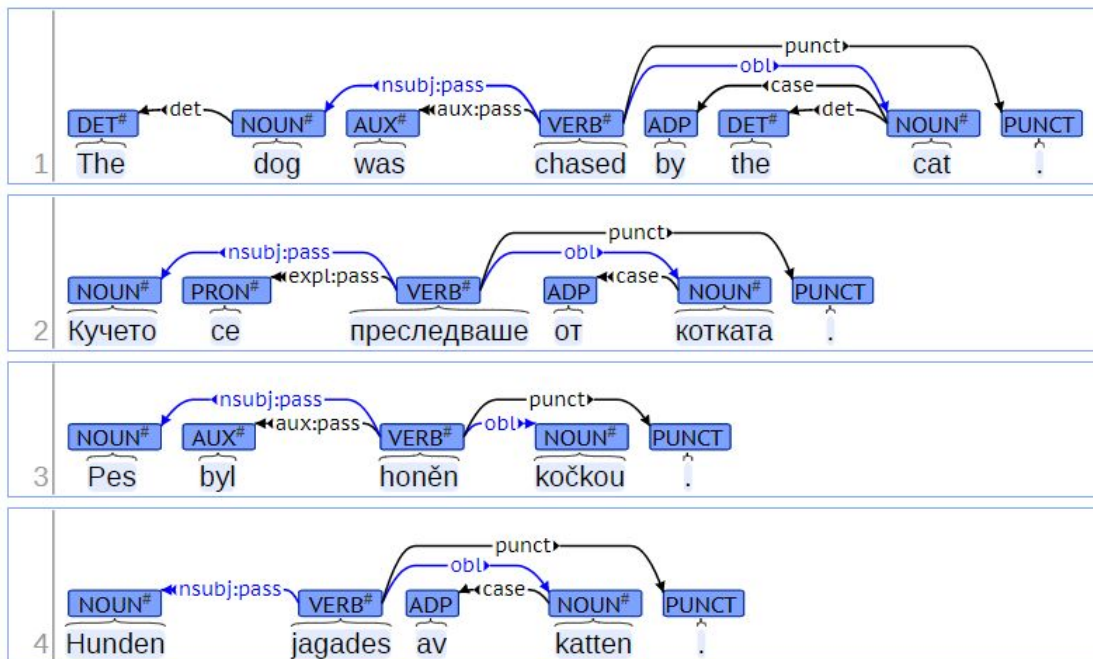
- Scheme for presenting dependency syntax



- Typed dependency relations between words
- Basic dependency representation forms a tree
- Exactly one word is the head of the sentence
- All other words are dependent on another word in the sentence


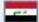




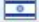























Universal dependencies









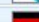














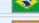









- Cross-linguistically valid






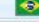


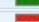









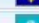

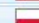























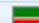


















Current UD Languages

Information about language families (and)

▶		Afrikaans
▶		Akkadian
▶		Akuntsu
▶		Albanian
▶		Amharic
▶		Ancient Greek
▶		Ancient Hebrew
▶		Apurina
▶		Arabic
▶		Armenian
▶		Assyrian
▶		Bambara
▶		Basque
▶		Beja
▶		Belarusian
▶		Bengali
▶		Bhojpuri
▶		Breton
▶		Bulgarian
▶		Buryat
▶		Cantonese
▶		Catalan
▶		Cebuano
▶		Chinese
▶		Chukchi
▶		Classical Chinese
▶		Coptic
▶		Croatian
▶		Czech
▶		Danish

▶		Dutch
▶		English
▶		Erzya
▶		Estonian
▶		Faroese
▶		Finnish
▶		French
▶		Frisian Dutch
▶		Galician
▶		German
▶		Gothic
▶		Greek
▶		Guajajara
▶		Guarani
▶		Hebrew
▶		Hindi
▶		Hindi English
▶		Hittite
▶		Hungarian
▶		Icelandic
▶		Indonesian
▶		Irish
▶		Italian
▶		Japanese
▶		Javanese
▶		Kaapor
▶		Kangri
▶		Karelian
▶		Karo
▶		Kazakh
▶		Khunsari
▶		Kiche
▶		Komi Permyak
▶		Komi Zyrian
▶		Korean
▶		Kurmanji
▶		Latin

▶		Ligurian
▶		Lithuanian
▶		Livvi
▶		Low Saxon
▶		Madi
▶		Makurap
▶		Maltese
▶		Manx
▶		Marathi
▶		Mbya Guaraní
▶		Moksha
▶		Mundurucu
▶		Naija
▶		Nayini
▶		Neapolitan
▶		North Sami
▶		Norwegian
▶		Old Church Slavonic
▶		Old East Slavic
▶		Old French
▶		Old Turkish
▶		Persian
▶		Polish
▶		Pomak
▶		Portuguese
▶		Romanian
▶		Russian
▶		Sanskrit
▶		Scottish Gaelic
▶		Serbian
▶		Skolt Sami
▶		Slovak
▶		Slovenian
▶		Soi
▶		South Le
▶		Spanish

▶		Soi
▶		South Levantine Arabic
▶		Spanish
▶		Swedish
▶		Swedish Sign Language
▶		Swiss German
▶		Tagalog
▶		Tamil
▶		Tatar
▶		Teco
▶		Telugu
▶		Thai
▶		Tupinamba
▶		Turkish
▶		Turkish German
▶		Ukrainian
▶		Umbrian
▶		Upper Sorbian
▶		Urdu
▶		Uyghur
▶		Vietnamese
▶		Warlpiri
▶		Welsh
▶		Western Armenian
▶		Wolof
▶		Xibe
▶		Yakut
▶		Yoruba

<https://universaldependencies.org/>

Design principles of the treebanks

- Needs to have a solid linguistic foundation
- Be transparent and accessible to non-specialists
- Support well downstream language understanding tasks

Disclaimer

- Despite of the common guidelines, adaptation can vary between different treebanks
 - Language-specific additions
 - Language-specific omissions
 - Treebank level annotation practices

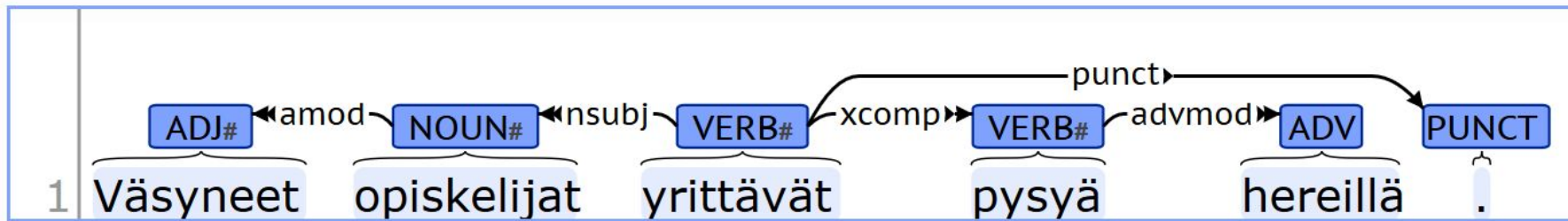
Treebank + machine learning = parser

1. Segmentation
2. Tokenization and sentence splitting
3. Morphological annotation + part-of-speech
4. Lemmatization
5. Dependency syntax

Turku-neural-parser-pipeline

<http://turkunlp.org/Turku-neural-parser-pipeline/>

A neural parsing pipeline for segmentation, morphological tagging, dependency parsing and lemmatization with pre-trained models for more than 50 languages. Top ranker in the CoNLL-18 Shared Task.



"Väsyneet"

ADJ

Case: Nom, Degree: Pos, Number: Plur

Lemma: väsynyt

Postag: A

LAS 91%

Token segmentation ~99%

POS ~98%

Lemmas ~95%

Trankit is a light-weight Transformer-based Toolkit for multilingual Natural Language Processing (NLP). It provides a pipeline for fundamental NLP tasks over 100 languages, and 90 pretrained pipelines for 56 languages.

Trankit can be easily installed via pip: `pip install trankit`

For more information, please check out our [github repo](#), [documentation](#), and [technical paper](#).

Usage

```
1 from trankit import Pipeline
2 # initialize a pipeline on English
3 p = Pipeline(lang='english', gpu=True, cache_dir='./cache')
4
5 doc = '''Michael helped shoot the majority of my firm's website
6 and we could not have been happier.'''
7
8
9 # perform all tasks on the input
10 all = p(doc)
11
12 sents = p.sspllit(doc) # sentence segmentation
13 tokens = p.tokenize(doc) # tokenization
```

	Nominals	Clauses	Modifier words	Function Words
Core arguments	<u>nsubj</u> <u>obj</u> <u>iobj</u>	<u>csubj</u> <u>ccomp</u> <u>xcomp</u>	https://universaldependencies.org/u/dep/all.html	
Non-core dependents	<u>obl</u> <u>vocative</u> <u>expl</u> <u>dislocated</u>	<u>advcl</u>	<u>advmod</u> * <u>discourse</u>	<u>aux</u> <u>cop</u> <u>mark</u>
Nominal dependents	<u>nmod</u> <u>appos</u> <u>nummod</u>	<u>acl</u>	<u>amod</u>	<u>det</u> <u>clf</u> <u>case</u>
Coordination	MWE	Loose	Special	Other
<u>conj</u> <u>cc</u>	<u>fixed</u> <u>flat</u> <u>compound</u>	<u>list</u> <u>parataxis</u>	<u>orphan</u> <u>goeswith</u> <u>reparandum</u>	<u>punct</u> <u>root</u> <u>dep</u>

id wordform lemma upos pos morpho head deptype

1	They	they	PRON	PRP	Case=Nom Number=Plur	2	nsubj
2	buy	buy	VERB	VBP	Number=Plur Person=3 Tense=Pres	0	root
3	and	and	CONJ	CC	_	4	cc
4	sell	sell	VERB	VBP	Number=Plur Person=3 Tense=Pres	2	conj
5	books	book	NOUN	NNS	Number=Plur	2	obj
6	.	.	PUNCT	.	_	2	punct

<https://universaldependencies.org/format.html>

What to do with conllu data then?

Two research designs in corpus linguistics

(see Biber & Jones 2009, Biber, Conrad & Reppen 1998)

Table 1.1 *Association patterns in language use*

-
-
- | | |
|------|---|
| A. | Investigating the use of a linguistic feature (lexical or grammatical) |
| (i) | Linguistic associations of the feature |
| – | lexical associations (associations with particular words) |
| – | grammatical associations (associations with particular grammatical constructions) |
| (ii) | Non-linguistic associations of the feature |
| – | distribution across registers |
| – | distribution across dialects |
| – | distribution across time periods |
| B. | Investigating varieties or texts (e.g., registers, dialects, historical periods) |
| (i) | Linguistic association patterns |
| – | individual linguistic features or classes of features |
| – | co-occurrence patterns of linguistic features |
-
-

What to do with conllu data then?

- Type A
 - *Zero or pronominal subject? (Helasvuoto & Kyröläinen 2016)*
 - *Near-synonyms (Biber, Conrad and Reppen 1998: 93)*
- Type B
 - *Filtering -- filter away function words etc.*
 - *Normalisation with lemmatisation*
 - *Lexico-grammatical characteristics of texts (Biber 1988)*

Today focus on type B

- Characteristics of entire texts
- Focus on *keyness*
 - A loose theoretical framework for analyzing important characteristics of a set of texts (Scott 1997, Scott & Tribble 2006)
 - Keywords ~ a group of words that function as a key to the text
 - Tell about the style and aboutness of the texts

Transportation/ lodging	Tourism	Narrative/ description	Physical features	Places/attractions	Food/drink
airport	adventure	afternoon	beach	city	beer
biking	arrived	amazing	beaches	gardens	delicious
boat	attractions	around	cliffs	museum	dinner
booked	destination	beautiful	hills	park	lunch
bus	explore	day	island	places	restaurant
ferry	exploring	enjoyed	islands	shops	restaurants
flight	guide	famous	mountain	town	
flights	holiday	hour	mountains	village	
headed	locals	located	river	villages	
hike	photo	lovely	rocks		
hiking	photos	nearby	sea		
hostel	sights	night	trees		
hostels	tour	north	water		
hotel	tourist	scenery			
hotels	tourists	scenic			
journey	tours	south			
ride	travel	spectacular			
road	travellers	steep			
streets	travelling	stunning			
trail	trip	sun			
trails	visit	sunny			
walk	visited	sunset			
walked	visiting	swimming			
walking	visitors	weather			

Biber & Egbert 2018

In practice

- Two corpora: target + reference
- Comparison of the words in the target corpus and in the reference corpus
- Keywords the ones that are over- or underrepresented in the target corpus
- Result a list of keywords
- Filtering + lemmatisation useful in the preprocessing

Methods used to extract keywords

“The standard”:

Frequency list of
target corpus
words

Frequency list of
reference corpus
words

Challenges with keyword analysis

- Text should be the unit of observation!
- (Can you think of why?)
- How to evaluate the keywords?

Text dispersion (Egbert & Biber 2019)

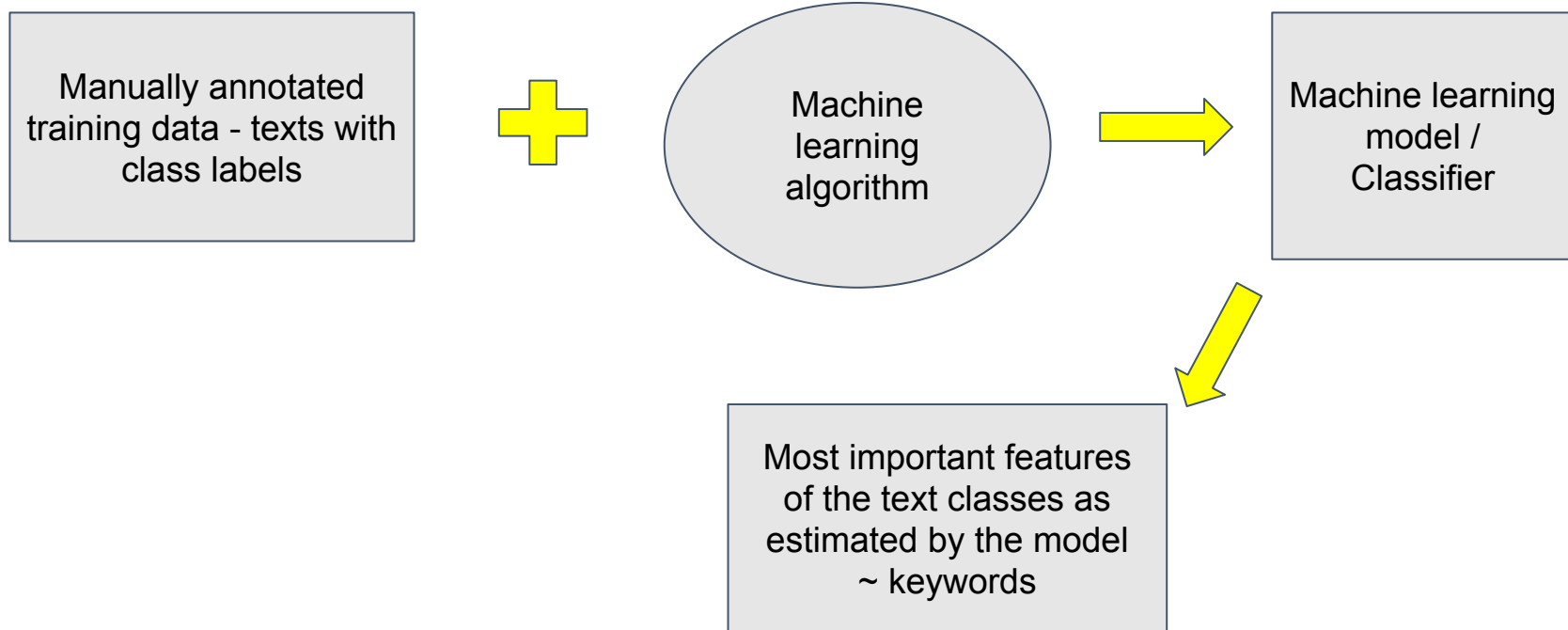
**Document
frequencies** of the
target corpus
words

**Document
frequencies** of the
reference corpus
words

Challenges with keyword analysis

- How to evaluate globally how well the keywords reflect **the target and reference corpora**?
- How to evaluate locally how well the keywords reflect **individual documents**?
- One possible solution to these challenges is to set up the task as text classification between the target and reference corpora

Text classification

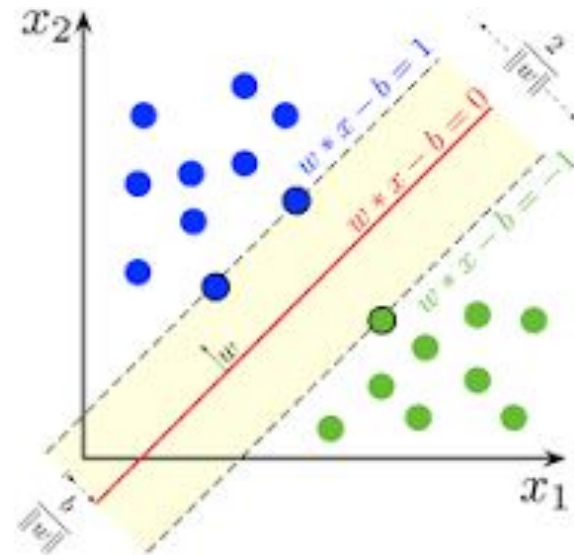


Text classification in keyword analysis

- Answers many of the challenges listed previously
 - Text as a unit of observation
 - Global evaluation: model performance
 - Local evaluation: how confident is the classifier about the text?
- Machine learning not new in linguistics
 - Random forests are often applied to study e.g. the geographical preference of -ing or to complements (e.g. Dehors & Gries 2016)

Support vector machines

- Good with sparse data (such as texts)!
- Feature importance often used to explain predictions in NLP (e.g. Sharoff et al. 2010)



Text classification

- Split data to train and test sets
 - ML focuses on *predicting* new instances
 - → Evaluation done on the test set that has not been seen during the training
- Typical evaluation measures precision, recall, f1-score (balanced and harmonized mean of precision+recall)
- Featurization - how are the data represented to the classifier?
 - Words, but also lemmas, syntactic information, etc.
 - → Comparison of the importance of words / other features for the classes in the data (see e.g. Laippala et al. 2021)

References

Biber, D., Conrad, S. and Reppen, R (1998). *Corpus Linguistics: Investigating Language Structure and Use*. Cambridge.

Biber, D., & Jones, J. K. (2009). *Quantitative methods in corpus linguistics*. In *Corpus Linguistics: An International Handbook* (Vol. 2, pp. 1286-1304). De Gruyter Mouton.

Biber, D., Egbert, J. (2018) *Register Variation Online*. Cambridge.

Deshors, Sandra & Gries, Stefan. (2016). 2016. Deshors, Sandra C. & Stefan Th. Gries. *Profiling verb complementation constructions across New Englishes: A two-step random forests analysis of -ing vs. to- complements*. *International Journal of Corpus Linguistics* 21(2): 192-218.. *International Journal of Corpus Linguistics*. 21. 10.1075/ijcl.21.2.03des.

Egbert, Jesse & Biber, Doug. (2019). *Incorporating text dispersion into keyword analyses*. *Corpora*. 14. 77-104. 10.3366/cor.2019.0162.

Helasvuo, Marja-Liisa and Kyröläinen, Aki-Juhani. "Choosing between zero and pronominal subject: modeling subject expression in the 1st person singular in Finnish conversation " *Corpus Linguistics and Linguistic Theory*, vol. 12, no. 2, 2016, pp. 263-299.

Laippala, Veronika, Jesse Egbert, Douglas Biber, Aki-Juhani Kyröläinen (2021) *Exploring the role of lexis and grammar for the stable identification of register in an unrestricted corpus of web documents. Language resources and evaluation*.

Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Yoav Goldberg, Jan Hajič, Christopher D. Manning, Ryan McDonald, Slav Petrov, Sampo Pyysalo, Natalia Silveira, Reut Tsarfaty, Daniel Zeman. *Universal Dependencies v1: A Multilingual Treebank Collection*. In *Proceedings of LREC*. 2016.

Marie-Catherine de Marneffe, Christopher D. Manning, Joakim Nivre, Daniel Zeman. *Universal Dependencies*. *Computational Linguistics*. 2021.

Mike Scott, Chris Tribble (2006): *Textual Patterns: Key Words and Corpus Analysis in Language Education*

Sharoff, S., Z. Wo, and K. Markert. 2010. *The web library of babel: evaluating genre collections*. *Proceedings of the Seventh Conference on International Language Resources and Evaluation*, 3063–3070.