

	<h2>Advanced Language Technologies</h2>
	<p>Information and Communication Technologies Research Area "Knowledge Technologies" <u>Jožef Stefan International Postgraduate School</u> Winter 2009 / Spring 2010</p> <p>Lecture II. Computer Corpora</p> <p><u>Tomaž Erjavec</u></p>

	<h2>Overview of the lecture</h2>
	<ol style="list-style-type: none"> 1. Background 2. Corpus compilation and markup 3. Morphosyntactic tagging

	<h2>Background</h2>
	<ul style="list-style-type: none"> • What is a corpus? • Using corpora • Characteristics of a corpus • Typology of corpora • History • Slovene language corpora

	<h2>A corpus is:</h2>
	<ul style="list-style-type: none"> ■ a large collection of texts ■ in digital format ■ language "as it is" ■ a sample of the language it is meant to represent ■ used for describing language (descriptive/empirical linguistics)

	<h2>A more precise definition</h2>
	<ul style="list-style-type: none"> ■ Corpus (plural corpora) is Latin for <i>body</i> ■ Guidelines of the Expert Advisory Group on Language Engineering Standards, EAGLES: <ul style="list-style-type: none"> - Corpus : A collection of pieces of language that are selected and ordered according to explicit linguistic criteria in order to be used as a sample of the language. - Computer corpus : a corpus which is encoded in a standardised and homogeneous way for open-ended retrieval tasks. Its constituent pieces of language are documented as to their origins and provenance. ■ For computer scientists: a dataset

	<h2>Using corpora</h2>
	<ul style="list-style-type: none"> ■ Applied linguistics: <ul style="list-style-type: none"> - <i>Lexicography</i>: making dictionaries (first users of corpora) - <i>Translation studies</i>: translation equivalents with contexts translation memories, machine aided translations - <i>Language learning</i>: real-life examples, curriculum development ■ Corpus linguistics: <ul style="list-style-type: none"> - linguistics based not on introspection, but on observation of real data ■ <i>Language technology</i>: <ul style="list-style-type: none"> - testing set for developed methods; - <i>training set</i> for inductive learning (<u>statistical Natural Language Processing</u>)

Characteristics of a (good) corpus

- *Quantity*:
the bigger, the better
- *Quality*:
the texts are authentic; the mark-up is validated
- *Simplicity*:
the computer representation is understandable, with the markup easily separated from the text
- *Documented*:
the corpus contains bibliographic and other meta-data

Typology of corpora I.

- Medium:
 - *written language*
 - *spoken language* (spoken, but in writing / transcription)
 - *speech corpora* (actual speech signal)
- Content:
 - *reference corpora* (representative), e.g. BNC
 - *sub-language corpora* (specialised), e.g. COLT
- Structure:
 - corpora with *integral* texts
 - corpora or of text *samples* (historical and legal reasons)
e.g. Brown

Typology of corpora II

- Time:
 - *static corpora*
 - *monitor corpora* (language change)
- Languages:
 - *monolingual corpora*
 - multilingual *parallel corpora* (e.g. Hansard, Europarl, JRC Acquis)
 - multilingual *comparable corpora*
- Annotation:
 - *plain text corpora*
 - *annotated corpora*

Reference corpora

- **Characteristics:**
 - a sample of the “complete” language
 - large, expensive, detailed and explicit design criteria
 - typically of contemporary language
 - documented and annotated
 - legally clean, available
- **Criteria for including texts:**
 - representativeness:
corpus includes “all” text types
 - balance:
the sizes of text type samples are in proportion to their “importance” for the speakers of the language
- **methodology v.s. practical constraints**

History

- History of Computational linguistics:
- 1950 -- 1960: empiricism
weak computers:
frequency lists
 - 1970 -- 1980: cognitive modeling (generative approaches, artificial intelligence)
deep analysis / “basic science”: computational linguistics
 - 1990 -- ...: empiricist revival, also combined approaches
quantity / usefulness:
language technologies
 - 2000 -- ...: The Web
- History of computer corpora:
- First milestones: Brown (1 million words) 1964; LOB (also 1M) 1974
 - The spread of reference corpora: Cobuild Bank of English (monitor, 100..200..M) 1980; BNC (100M) 1995; Czech CNC (100M) 1998; Croatian HNK (100M) 1999...
 - Slovene language reference corpora: FIDA (100M), Nova Beseda (100M...) 1998; FIDA+ (600M) 2006.
 - EU corpus oriented projects in the '90: NERC, MULTEXT-East,...
 - Language resources brokers: LDC 1992, ELRA 1995
 - Web as Corpus (2000)
 - more, larger, for more languages, with diverse annotations

Slovene language corpora

- Monolingual reference corpora:
- ZRC SAZU: Beseda, 1998; Nova beseda, 2000-
 - DZS, Amebis, FF, IJS: FIDA, 1998, FidaPlus, 2006
 - IJS, FF: JOS corpora
- Parallel corpora:
- IJS: MULTEXT-East 1998-, SVEZ-IJS, 2004, JRC-ACQUIS, 2006
 - SVEZ: EuroKorpus
 - FF: TRANS, 2002
 - UP: Turist Corpus, 2008
- Speech corpora:
- Laboratory for Digital Signal Processing, University of Maribor: SpeechDat, ONOMASTICA...
 - Laboratory of Artificial Perception, Systems and Cybernetics, University of Ljubljana: SQL, GOPOLIS,...

	<h2>II. Compilation and markup of corpora</h2>
	<ul style="list-style-type: none"> • Steps in the preparation of a corpus • What annotation can be added to the text • Computer coding of corpora • Markup Methods

	<h2>Before making your own corpus</h2>
	<p>check if an appropriate corpus is already available</p> <ul style="list-style-type: none"> ■ google ■ corpora@lists.uib.no ■ LDC, ELRA

	<h2>Steps in the preparation of a corpus</h2>
	<ol style="list-style-type: none"> 1. Choosing the component texts and acquiring digital originals 2. Up-translation to standard format 3. Linguistic annotation 4. Documentation 5. Use and Dissemination

Getting the text

1. Choosing the component texts:
linguistic and non-linguistic criteria;
availability; simplicity; size
2. Copyright
sensitivity of source (financial and
privacy considerations); agreement
with providers; usage, publication
3. Acquiring digital originals
OCR; digital originals; Web
 - BootCat

Processing

1. Conversion to common format
consistency; character set encodings;
structure
 - Web as Corpus: Wacky tools
2. Documentation
e.g. TEI header; Open Archives etc.
3. Linguistic annotation
language dependent methods; errors

Use and dissemination

- Using the corpus:
 - concordancer (linguists)
e.g. [FidaPLUS](#), [SKE](#), [iKorpus](#)
 - statistics extraction
 - development of new methods for analysis
- Dissemination:
 - legalities (source copyright, corpus use
agreement)
 - mode: concordancer or dataset

Computer coding of corpora

- Encoding must ensure
 - durability
 - interchange between computer platforms
 - interchange between applications
- Basic standard: *Extended Markup Language, XML*
 - a number of companion standards and technologies: XSLT, XML Schema, ISO Relax NG, XPath, XQuery, ...
- The vocabulary of annotations for corpora and other language resources are defined by the *Text Encoding Initiative, TEI*
- XML/TEI used much wider than just for corpora:
 - annotation of dictionaries: English-Slovene, Japanese-Slovene (from jaSlo)
 - for annotating text-critical editions

Corpus annotation

- Annotation = interpretation
- Documentation about the corpus (example)
 - Document structure (example)
 - Basic linguistic markup: sentences, words (example), punctuation, abbreviations (example)
 - Lemmas and morphosyntactic descriptions (example)
 - Syntax (example)
 - Alignment (example)
 - Terms, semantics, anaphora, pragmatics, intonation,...

Example: TEI header

```
<teiHeader id="ecmr.H" type="text" lang="sl-en" creator="ET" status="update"
date.created="1999-04-13" date.updated="1999-06-22" >
<fileDesc>
<titleStmt>
<title lang="sl">Ekonomsko ogledalo; 13 &scaron;tevilik 98/99</title>
<title lang="en">Slovenian Economic Mirror; 13 issues, 98/99</title>
<respstmt>
<name>Andrej Skubic, FF</name>
<resp lang="sl">Zagotovitev digitalnega originala, poravnava</resp>
<resp lang="en">Provision of digital original, alignment</resp>
<name>Tomaž Erjavec, IJS</name>
<resp lang="sl">Tokenizacija, pretvorba v TEI</resp>
<resp lang="en">Tokenisation, conversion to TEI</resp>
</respstmt>
</titleStmt> ...
```

Example: text structure

```
<quote id="Osl.1.8.18" rend="center;it">
<lg id="Osl.1.8.18.1">
  <l id="Osl.1.8.18.1.1">Tam pod kostanjevim drevesom</l>
  <l id="Osl.1.8.18.1.2">izdala si me,</l>
  <l id="Osl.1.8.18.1.3">izdal sem te,</l>
  <l id="Osl.1.8.18.1.4">ne da bi trenila z očesom.</l>
</lg>
</quote>
<p id="Osl.1.8.19">
  <s id="Osl.1.8.19.1">Trije možje se niso niti ganili.</s>
  <s id="Osl.1.8.19.2">Toda ko je <name>Winston</name>
  znova pogledal v Rutherfordov propadli obraz, je opazil, da so
  njegove oči polne solz.</s> ...
```

Example: morphosyntactic tagging

```
<s id="Osl.1.2.2.1">
<w lemma="biti" ana="Vcps-sma">Bil</w>
<w lemma="biti" ana="Vcip3s-n">je</w>
<w lemma="jasen" ana="Afpmsnn">jasen</w><c>,</c>
<w lemma="mrzel" ana="Afpmsnn">mrzel</w>
<w lemma="aprilski" ana="Aopmsn">aprilski</w>
<w lemma="dan" ana="Ncmsn">dan</w>
<w lemma="in" ana="Ccs">in</w>
<w lemma="ura" ana="Ncfpn">ure</w>
<w lemma="biti" ana="Vcip3p-n">so</w>
<w lemma="biti" ana="Vmpps-pfa">bile</w>
<w lemma="trinajst" ana="Mcnpl">trinajst</w><c>.</c>
</s>
```

Example: alignment

```
<linkGrp id="Osl.en.1" type="body" targtype="s"
domains="Oen Osl">
<link xtargets="Osl.1.2.2.1 ; Oen.1.1.1.1">
<link xtargets="Osl.1.2.2.2 ; Oen.1.1.1.2">
<link xtargets="Osl.1.2.3.1 ; Oen.1.1.2.1">
<link xtargets="Osl.1.2.3.2 ; Oen.1.1.2.2">
...
<link xtargets="Osl.1.2.6.5 ; Oen.1.1.5.5">
<link xtargets="Osl.1.2.6.6 ; Oen.1.1.5.6 Oen.1.1.5.7">
<link xtargets="Osl.1.2.6.7 ; Oen.1.1.5.8">
...
```

Methods for linguistic markup

- *hand annotation*: documentation, first steps
generic (XML, spreadsheet) editors or specialised editors
- *semi-automatic*: morphosyntactic and other linguistic annotation
cyclic approach: machine, hand, validate, correct, machine, ...
- *machine, with hand-written rules*: tokenisation
regular expression
- *machine, with inductively built models from annotated data*:
"supervised learning"; HMMs, decision trees, inductive logic programming,...
- *machine, with inductively built models from un-annotated data*:
"unsupervised learning"; clustering techniques
- **overview of the field**

III. Morphosyntactic tagging

- Better known as part-of-speech (PoS) tagging
- Tagging is the task of labeling each word in a sequence of words with its appropriate part-of-speech
- Words are often ambiguous with respect to their POS:
 - *saw* → singular noun
 - *saw* → past tense of verb *see*
- Purposes and applications (examples):
 - pre-processing step for further analyses:
 - lemmatisation
 - syntactic structure, etc.
 - text indexing, e.g. nouns are more useful than verbs
 - pronunciation in speech processing

Steps in tagging

- for each word token in text the tagger needs to know all its possible tags (ambiguity class)
→ a morphological lexicon
- given the context in which the word appears in, the tagger must decide in the correct tag:
 - he saw/V a man carrying a saw/N
- so, tagging performs limited syntactic disambiguation

Example: Penn Treebank

Under/IN the/DT proposal/NN ,/, Delmed/NNP
would/MD issue/VB about/IN 123.5/CD
million/CD additional/JJ Delmed/NNP
common/JJ shares/NNS to/TO
Fresenius/NNP at/IN an/DT average/JJ
price/NN of/IN about/IN 65/CD cents/NNS
a/DT share/NN ,/, though/IN under/IN
no/DT circumstances/NNS more/JJR than/IN
75/CD cents/NNS a/DT share/NN ./.

PoS taggers

- Most taggers induce the language model from a hand-annotated corpus
- Typically, two resources are induced:
 - lexicon, giving the ambiguity class of a word and their frequencies in the training corpus
 - the tag of a word in text depends on its local context

Tagging with Markov Models

- Sequence of tags in a text is regarded a Markov chain
- Limited horizon: A word's tag only depends on the previous tag: $p(x_{i+1} = \tilde{t} \mid x_1, \dots, x_i) = p(x_{i+1} = \tilde{t} \mid x_i)$
- Time invariant: This dependency does not change over time: $p(x_{i+1} = \tilde{t} \mid x_i) = p(x_2 = \tilde{t} \mid x_1)$
- Task: Find the most probable tag sequence for a sequence of words
- Maximum likelihood estimate of tag t^* following \tilde{t} :
 $p(t^* \mid \tilde{t}) = f(\tilde{t}, t^*) / f(\tilde{t})$
- Optimal tags for a sentence:
 $t'_{1,n} = \arg \max p(t_{1,n} \mid w_{1,n}) = \prod p(w_i \mid t_i) p(t_i \mid t_{i-1})$

Most popular Markov model tagger

- TnT (Trigrams 'n Tags)
- induces lexicon and tag trigrams from the training corpus
- has heuristics to tag unknown words
- has no problem with large tagsets
- fast in training and tagging
- freely available for non-commercial use
- but only as a Linux executable
- OS alternative: hunpos

TreeTagger

- uses decision trees
- relatively fast
- comes with lots of models for various languages
- executables freely available
<http://www.ims.uni-stuttgart.de/projekte/corplex/TreeTagger/>

Transformation-based Tagging (TbT)

- Basic idea: transform an imperfect tagging into one with fewer errors by changing wrong tags
- Features that trigger changes can be conditioned on words and on more context and are user specified
- Components:
 - specification of transformations
 - learning algorithm: constructs a ranked list of transformations
- A transformation consists of two parts:
 - triggering environment + rewrite rule
- Examples:
 - if previous tag is TO and current tag is NN then change it to VB
 - if one of previous two words is *n?* and current tag is VBP then change it to VB
 - if next tag is JJ and current tag is JJR then change it to RBR
 - if one of previous three tags is MD and current tag is VBP then change it to VB

Yet another Tagger

For a while, trying out new approaches to tagging was in fashion

- Maximum Entropy taggers
- Support Vector Machine taggers
- Memory based taggers
- ...

Tagsets

- A tagset is a set of part-of-speech tags
- Classical 8 classes (Thrax, 100 BC): noun, verb, article, participle, pronoun, preposition, adverb, conjunction
- But all tagset use more tags than that!
- Criteria:
 - specificity: degree to which humans use the tagset uniformly on the same text
 - accuracy: evaluation of output on tagged text
 - suitability for intended application

Tagsets for English

- For English, there exist several tagsets: Brown, CLAWS, Penn, ...
- English tagsets include PoS + some other morphological (inflectional) properties: 30-80 tags
- Penn Treebank Tagset for English: 37 tags, e.g.
 - JJ adjective, positive
 - JJR adjective, comparative
 - JJS adjective, superlative
 - NN non-plural common noun
 - NNS plural common noun
 - NNP non-plural proper name
 - NNPS plural proper name
 - IN preposition
 - ...

Morphosyntactic tagsets

- For inflectionally rich languages (such as Slavic languages), tagsets contain much more information than just PoS
- Slovene, Czech, etc. > 1000 different morphosyntactic tags
 - gender, number, case, animacy, definiteness, ...
- Efforts to standardise tagsets across languages:
 - Eagles
 - MULTEXT
 - MULTEXT-East

MULTEXT-East

- EU project in '90s: development of language resources for Central and East-European languages
- also development of morphosyntactic specifications, lexica and annotated corpus
- Parallel annotated corpus:
Orwell's 1984
- Several later releases, V3 in 2004, V4 in 2010
- Web site: <http://nl.ijs.si/ME/>

MULTEXT-East morphosyntactic specifications

- Specify
 - what morphosyntactic features particular languages distinguish,
 - what their names and values are,
 - how they can be mapped to tags (morphosyntactic descriptions, MSDs)
- e.g. that *Ncms* is:
 - a valid for Slovene
 - is equivalent to *PoS:Noun, Type:common, Gender:masculine, Number:singular*
- <http://nl.ijs.si/ME/V3/msd/html/>

	<h2>JOS morphosyntactic specifications</h2>
	<ul style="list-style-type: none"> ■ only for Slovene ■ based on MULTTEXT-East but changed some features and lexical assignments ■ also moved to 100% XML/TEI encoding ■ bi-lingual (Slovene and English) ■ also made annotated corpora: <ul style="list-style-type: none"> – jos100k (hand validated) – jos1M (partially hand validated) ■ http://nl.ijs.si/jos/

	<h2>Conclusions</h2>
	<ul style="list-style-type: none"> ■ What is a corpus ■ How to make it ■ How to annotate it
