

Text Data Mining – auf Grundlage von *Karine Megerdoomian (Handbook for Language Engineers)*

Seminar-Vortrag
im Rahmen der Computerlinguistik
Ausgearbeitet von
Daniel Gruber
(Juni 2004)

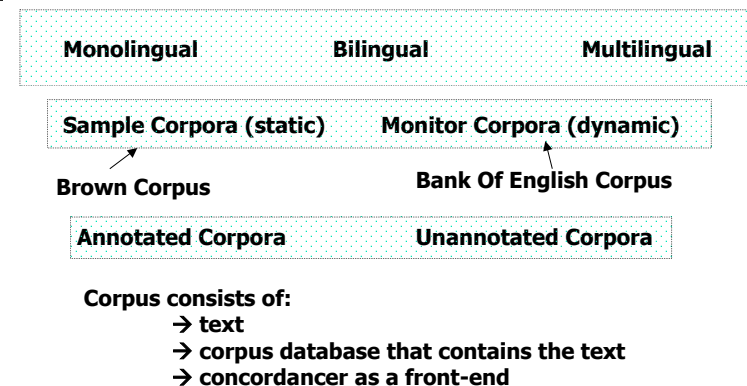
Corpus

- „A corpus is a *collection of text or speech* material that has been brought together according to be a certain set of predetermined criteria.“
 - Testing hypotheses about natural lang.
 - Extracting statistical and linguistic information

Corpus (Plural Corpora)

- A **collection** of written text or recorded speech
- Useful for *statistical* knowledge acquisition techniques
- Chomsky, 1957, Syntactic Structures:
„The corpus, if natural, will be so wildly skewed that the description [of language based on the corpus] would be no more than a mere list.“

Corpus – Types



Corpus – History

- Non-digital
 - Käding - 1897 – 100 million words [5000 analysts] → **shorthand/stenography**
 - Palmer – 1933 → language pedagogy
 - Eaton – 1940 → comparative linguistics
 - Fries, 1952 → corpus-based grammar
 - Quirk, 1961 → survey of English usage

Corpus – properties

- Machine readable form (digital)
- Representative of the domain under study
- Balanced sample
 - Text with specific parameters
 - (BNC, ANC)
- Finite (monitor corpus: non-finite [grows to reflect languages changes])

Corpora Resources I

- Brown Corpora (1 million tagged words)
 - Sample of written American English
- Lancaster-Oslo-Bergen corpus
- Penn Treebank
- British National Corpus (BNC)
- Project Gutenberg
 - <http://www.promo.net/pg/>
 - 6267 books free available
- Mannheimer Corpus Collection
 - <http://corpora.ids-mannheim.de/cosmas/>
 - (Demonstration einer statistischen Kookurenzanalyse)

Corpora Resources II

Brown Corpus	Lancaster-Oslo-Bergen	Penn Treebank	British National Corpus (BE)	Birmingham Collection of English Text	Reuters	Project Gutenberg	American National Corpus	Mannheimer Corpus Collection	UN Parallel Text Corpus
1967	1978		1994	1985	-	-			-
written AE	written BE		Written 90% spoken 10 %					German	English French Spanish
1.000.000 words		4.500.000 words	100.000.000 words	20.000.000 words	810.000 News	6267 Books		2 Bio. words	2.5 GB
		POS tagged	POS tagged						
		Wall Street Journal							UN electronic text archives (88-93)
									Multilingual

Corpus Analysis

- **Corpus linguistics** is the study of language through analysis of natural-occurring data. It involves computational methods and tools and develops theories of linguistics and language use.
- Annotation is mostly required for analyzing linguistic pattern
- Information Retrieval based on annotated corpora

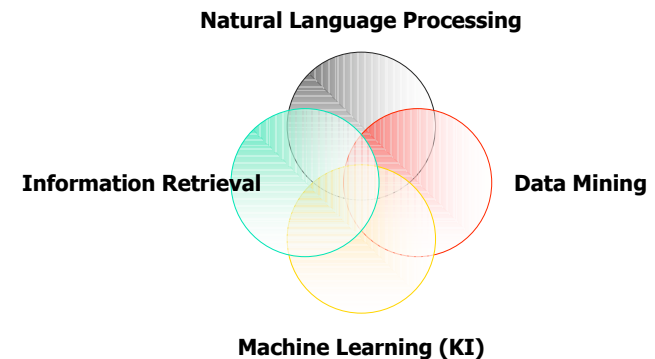
Text Mining or Text Data Mining

- **Computer – Internet:** Enormous growth in the volume of online text documents in multiple formats and languages
- *Goal: discover knowledge* from unstructured textual data
- **Text Mining vs. Data Mining**
 - Type of data under investigation
 - TM: **unstructured** natural-language documents
 - DM: **highly structured** data in data warehouses
- **Text Mining vs. Information Retrieval**
 - DM: derive **new** information from data
 - IR: extracts already **existing** information

Text Mining – Process

- (1) Information Extraction
 - TM systems include an IE module
 - Locates significant vocabulary items in NL documents (linguistic knowledge!)
- (2) Tokenization, stemming and tagging
- (3) Cluster a collection of documents (IR)
- (4) Categorize the clusters
- (5) Discover knowledge from the databases and visualization tools (DM)

Text Mining – Text Data Mining



Tokenization – pre-processing corpora I

- *Problem one:* Text can be enriched with markups (HTML, XML, ...), tables or other non-useful things...
- *Problem two:* We need a separated text to see **word and sentence boundaries**...

■ Sentence Boundaries

? ! . Are good choices, but how can we handle an abbreviation (like Dr.) or acronym (like I.B.M.)? What's with Numbers (112.211)?

Haplology: One Character has two simultaneous uses. [i.e. the period at the end could signal both a sentence break and an abbreviation]

90% of the periods in English are sentence boundary markers.

Tokenization – pre-processing corpora II

■ Sentence Boundaries

Solution: Regular Expressions with a list of abbreviations

Recent approaches:

Riley(1989): Statistical classification tree

Palmer and Hearst (1997): POS information used from a NN to predict sentence boundaries

Maximum entropy approach (probabilistic distribution of sentence boundaries in a text)

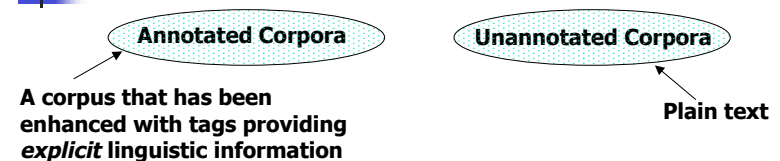
Tokenization – pre-processing corpora III

■ Word Segmentation

- > Occurrence of whitespaces or punctuation
- > May work in English

- Punctuation: Apostrophes can be a part of a word (rock'n roll) or marking possessive nouns (cat's)
- Contractions: Spanish: del → de el
- Hyphenation: Compound word (e-mail) or joining words (25-year-old)
- Whitespace: New York
- More: Numbers, Encoding (→ Unicode)

Corpus Annotation



■ How?

- GML (Generalized Markup Language)
- SGML (-> HTML)
- XML (subset of SGML)
- XCES (XML Corpus Encoding Standard [XML-scheme])
- RDF (Resource Description Framework)

Annotation Coverage

Annotation schemes

- Part of speech annotation
 - One of the first types of annotation
 - Most common annotation today
- Lemmatization (or stemming)
 - Based on morphological analysis
- Parsing
 - Parsed corpora: *treebanks* (→ *Penn Treebank*)
 - Annotate the syntactic phrases (sentence, verb phrase, noun phrase, prepositional phrase)
- Semantic annotation, discourse analysis (*greetings*, *apologies [sorry]*, *politeness [please]*), speech annotation

Tagset Design I

- Tagging is the foundation for further analysis
- Most common type is **part of speech (POS)** tagging
- Tagset: annotation tags used within the corpus
- Most widely used tagsets in:
 - Brown Corpus (<http://www.scs.leeds.ac.uk/ccalas/tagsets/brown.html>)
 - Penn Treebank (see next slide)
 - British National Corpus (BNC C5; distinguishes 61 categories; <http://www.natcorp.ox.ac.uk/what/c5spec.html>)

Tagset Design II

Penn Treebank Tagset

CC	Coordinating conjunction e.g. <i>and, but, or...</i>
CD	Cardinal Number
DT	Determiner
EX	Existential <i>there</i>
FW	Foreign Word
IN	Preposition or subordinating conjunction
JJ	Adjective
JJR	Adjective, comparative
JJS	Adjective, superlative
LS	List Item Marker
MD	Modal e.g. <i>can, could, might, may...</i>
NN	Noun, singular or mass
NNP	Proper Noun, singular

NNPS	Proper Noun, plural
NNS	Noun, plural
PDT	Predeterminer e.g. <i>all, both ... when they precede an article</i>
POS	Possessive Ending e.g. <i>Nouns ending in 's</i>
PRP	Personal Pronoun e.g. <i>I, me, you, he...</i>
PRPS	Possessive Pronoun e.g. <i>my, your, mine, yours...</i>
RB	Adverb Most words that end in <i>-ly</i> as well as degree words like <i>quite, too</i> and <i>very</i>
RBR	Adverb, comparative Adverbs with the comparative ending <i>-er</i> , with a strictly comparative meaning.
RBS	Adverb, superlative
RP	Particle
SYM	Symbol Should be used for mathematical, scientific or technical symbols
TO	<i>to</i>
UH	Interjection e.g. <i>uh, well, yes, my...</i>

Tagset Design III

Penn Treebank Tagset

VB	Verb, base form subsumes imperatives, infinitives and subjunctives	Punctuation Tags # \$ " () , . : ..
VBD	Verb, past tense includes the conditional form of the verb to be	
VBG	Verb, gerund or present participle	
VBN	Verb, past participle	
VBP	Verb, non-3rd person singular present	
VBZ	Verb, 3rd person singular present	
WDT	Wh-determiner e.g. <i>which, and that</i> when it is used as a relative pronoun	
WP	Wh-pronoun e.g. <i>what, who, whom...</i>	
WP\$	Possessive wh-pronoun e.g.	
WRB	Wh-adverb e.g. <i>how, where why</i>	

Tagging Methods

- Rule-based taggers (knowledge-based taggers)
 - Lexicon based
 - TAGGIT (Green and Rubin, 1971)
 - Unable to provide tags for construction that have not been recognized
- Probabilistic taggers
 - Need to be trained to build a probability matrix (word/grammatical class/probability)
 - Bigram analysis (probability: a word of a certain POS follow a word from another particular POS) (trigrams)
 - (hidden Markov model -> viterbi algorithm)
 - If an unknown word occurs the grammatical class can be found through the distributional information
- Hybrid taggers (CLAWS)

Text Mining, Corpus Building, and Testing

21

Tagging Methods: CLAWS

- BNC: 100.000.000 words
- Free WWW trail available <http://www.comp.lancs.ac.uk/ucrel/claws/trial.html>
- Tagsets: C5 (British National Corpus) and C7
- Output: horizontal or vertical (with columns)
- If you would like to use our trial service, please complete the form below. From an academic site, you can enter up to 10,000 words of English running text. From a non-academic site, you can enter up to 300 words of English running text. If you enter more, it will be cut off at the appropriate word limit. Input format guidelines are available. To tag the text you have entered click the button below the form.
- With C5 horizontal it looks like this:
- If_CJS you_PNP would_VM0 like_VVI to_TO0 use_VVI our_DPS trial_NN1 service_NN1 , please_AV0 complete_VVB the_AT0 form_NN1 below_AV0 , . From_PRP an_AT0 academic_AJ0 site_NN1 , you_PNP can_VM0 enter_VVI up_AV021 to_AV022 10,000_CRD words_NN2 of_PRF English_AJ0 running_AJ0 text_NN1 . . . From_PRP a_AT0 non-academic_AJ0 site_NN1 , you_PNP can_VM0 enter_VVI up_AV021 to_AV022 300_CRD words_NN2 of_PRF English_AJ0 running_AJ0 text_NN1 . . . If_CJS you_PNP enter_VVB more_AV0 , it_PNP will_VM0 be_VBI cut_VVN off_AVP at_PRP the_AT0 appropriate_AJ0 word_NN1 limit_NN1 . . . Input_NN1 format_NN1 guidelines_NN2 are_VBB available_AJ0 . . . To_TO0 tag_VVI the_AT0 text_NN1 you_PNP have_VHB entered_VVN click_VVB the_AT0 button_NN1 below_PRP the_AT0 form_NN1

Text Mining, Corpus Building, and Testing

22

Applications in Corpus Linguistics

Extracting information

Training purposes

Text Mining, Corpus Building, and Testing

23

Applications in Corpus Linguistics – Lexicon Acquisition

- Lexicon contains
 - Morphological
 - Syntactic
 - Semantic
 - Pragmatic information
- Lexicon used in
 - Information extraction
 - Document summarization
 - Machine translation
- Computational lexicon for taggers consists of the **lexemes** or **stem forms** of words
- If the tagger is used for a part of machine translation we need also the **translation**
- If incorrect information is in the lexicon the tagger cannot be as good as we want
- Machine readable dictionaries (MRD) containing lexical information
 - Static
 - Not available for many languages especially for translation purposes
- → Lexicon acquisition
 - Gather lexical information from text corpora
 - Corpora are now widely available
 - Can reflect dynamic and changing nature of language

Extracting information

Text Mining, Corpus Building, and Testing

24

Applications in Corpus Linguistics – Discourse Analysis

- Research topic: determining the psychological point of view (POV) of the author or a character in the text
- Current approaches:
 - Creating annotated corpora
 - Train a statistical discourse tagger module
 - Then the discourse tagger marks a corpus with POV expressions
 - → extremely difficult (no formal criteria)

Text Mining, Corpus Building, and
Testing

25

Applications in Corpus Linguistics – Word Sense Disambiguation

- Goal: select the *appropriate meaning* to a given word based on the linguistic context
- Importance
 - machine translation
 - information retrieval
 - parsing

Text Mining, Corpus Building, and
Testing

26

Applications in Corpus Linguistics – Machine Translation

- Several methods (full knowledge-based systems, interlingua methods, purely statistical approaches,..)
- MT includes
 - Multilingual lexicon
 - Tagger
 - Parser
 - Word sense disambiguation module
- All are based on corpora
- MT specific: **parallel corpora** (same text in several languages)
- **Text alignment**: create explicit link between the elements that are mutual translations (→ *aligned corpus*)
- Methods for aligning sentences
 - Comparing the lengths of textual units
 - Using Lexical content
 - Matching cognates (*verwandte*)

Text Mining, Corpus Building, and
Testing

27

Corpus Processing Tools

- UNIX/POSIX Tools
- Word Counts
- Concordances
- Collocations
- Testing and Evaluation

Text Mining, Corpus Building, and
Testing

28

Corpus Processing Tools – POSIX Tools

- POSIX: Portable Operating System Interface (...to be pronounced pahz-icks not poh-six)
- **POSIX Part 3: Shell and Utilities**
- **Standardization of tools so that they are widely available on several OS's.**
- A lot of utilities use *regular expressions* which are equivalent to *regular language* and equivalent to DFA (deterministic finite automata) and NFA (nondeterministic finite automata)
 - the word problem given from the expression can be solved efficiently
- Regular expressions are for instance:
 - A letter or a number
 - [...] := Class of characters – [^...] := complement class
 - Regular expression followed by +=: one or more
 - Regular expression followed by * := 0,1 or more (Kleene hull)
 - Regular expression | Regular expression := or (Regular expression)
- Helmut Herold, Linux-Unix Kurzreferenz
- Regular expressions in C: man regex

Text Mining, Corpus Building, and
Testing

29

Corpus Processing Tools – Word Counts

- Now we have some text on our POSIX-workstation
- We can use several utilities which are born to help us
- wc -w (for word count)
- tr (search and replace)
- cut (to handle text with columns)
- sort
- SED (stream editor) [sed 'script' file]
 - Script: [address1 [, address2]] function [args]
 - Addresses: number or /regex/
 - Functions: p (print), q (quit), s (replace)
 - sed '/Kommentar/!s/[^A-Za-z]/ /g' text.txt | wc -w
- AWK
 - Pattern { action}
 - Patterns: BEGIN, END, expression, regex, concatenated pattern, pattern1, pattern2 (all rows between pattern1 match and pattern2 match)
 - Builtin-variables: ARGV, ARGV, FILENAME, NF (number of fields)

Text Mining, Corpus Building, and
Testing

30

Corpus Processing Tools – Concordances

- Context of a word can be from interest
 - Especially for lexicographer
 - Intervening materials between verb and particle can be usefull for developing language grammars
- KWIC concordancing program (**Key Word in Context**)
 - Extracts all occurrences of the word of interest and displays it with the word in the center and the *surrounding context* on the two sides.

Text Mining, Corpus Building, and
Testing

31

IDS Copora „Kookkurrenzanalyse“ of „Curacao“

IDS Copora „Kookkurrenzanalyse“ of „Curacao“

Suchanfrage: Curacao (232 Treffs)

BelegNr.	TFR	Kollokationen	Häufigkeit	systematische Muster (lexematisch)
1+4:	500	Blue Wodka	4	75% Wodka ... und Blue Curacao
5+5:		Blue Sekt	5	60% Blue Curacao [...] Sekt
10+4:			41	87% Blue Curacao und
51+3:	472	Aruba Bonaire Eustatius	3	100% Aruba Bonaire [...] Curacao Saba St.
54+2:		Aruba Bonaire	12	93% die/dem Aruba [...] Bonaire (un
86+5:		Aruba Grenada	5	60% Barbados Grenada [...] Curacao [...]
71+24:		Aruba	24	75% die/dem Aruba (Bonaire)und Cur
35+3:	368	Bonaire Antillen Niederländische	3	86% Niederländischen Antillen [...]
98+17:		Bonaire	17	82% Antillen Aruba Bonaire (und) Curac
115+1:	177	Antillen	11	62% den Niederländischen Antillen (Arub
126+6:	141	Calister	6	100% Sängerin Izaline Calister aus Curac
132+6:	140	Isaline	6	100% Sängerin Isaline Calister aus Curac
138+9:	126	Orangensaft	9	55% Orangensaft und el Blue Curacao
147+5:	106	Eustatius	5	80% Bonaire Curacao Saba St Eustatius
152+7:	102	Grenada	7	85% Barbados Grenada (Aruba und) Curac
159+5:	90	Antilleninsel	5	100% der Antilleninsel Curacao
164+7:	83	Wodka	7	75% Wodka ... und Blue Curacao
171+4:	82	Willemstad	4	75% Willemstad (auf [...] Curacao
175+4:	79	Karibikinsel	6	100% der Karibikinsel Curacao
181+4:	79	Islands	4	50% Islands Curacao
185+7:	69	Sekt	7	71% Blue Curacao ... und Sekt
192+5:	69	Niederländischen	5	100% der Niederländischen Antillen Aruba
197+2:	67	deah	2	50% Curacao deah
198+4:	66	1878	4	100% 1878 die damals niederländische Ins
203+4:	64	Dominikanische	4	75% Curacao Dominica Dominikanische Rep
207+7:	62	niederländischen	7	57% der niederländischen Antilleninsel [...]
214+9:	60	Insel	9	88% die Insel [...] Curacao überfallen
223+6:	59	Karibik	6	50% Curacao [...] in der Karibik
229+4:	54	absetzte	4	100% nach Curacao absetzte
233+4:	54	Barbados	4	100% Barbados Grenada [...] Curacao Aruba
237+4:	53	Saba	4	100% Bonaire Curacao [...] Saba Bost.
241+6:	50	Inseln	6	65% Inseln Aruba und Curacao
247+4:	48	Zitronensaft	4	75% Curacao [...] Zitronensaft
251+3:	48	Dominica	3	100% Curacao Dominica Dominikanische Rep
254+3:	47	Bela	3	65% Curacao [...] Bela
257+3:	46	Cayman	3	100% Cayman Islands Curacao
260+3:	46	Sint	3	65% Curacao ... Sint
263+3:	44	Haarten	3	65% Curacao ... Haarten
266+3:	43	N V	3	100% N V Curacao
269+3:	41	Antigua	3	65% Curacao Antigua
272+4:	38	Venezuela	4	75% Curacao Venezuela
276+3:	36	Elwee	3	100% Elwee Curacao
279+5:	36	Amsterdam	5	60% Amsterdam Curacao

32

Collocations

- Concordance analysis is not very efficient for real quantitative analysis.
 - Not ordered by frequency
 - Number of hits can be very high
- „A *Collocation* is an **expression** that consists of a number of words within a short distance of each other.“
- Collocation analysis
- Compositionality

Testing and Evaluation

- Test corpora
 - Corpora created to be used for evaluating and testing statistical algorithms and the performance of NLP systems.
 - Typically annotated
- Truthfile = annotated test corpus ?
- Split *data set* at the beginning into a *training set* and a *test set*
- Measurements:
 - Information Retrieval: Precision and Recall
 - Accuracy and coverage
- Criteria for a test set or truthfile
 - Consistency (reflect linguistic phenomena)
 - Size of Tagset (larger tagset → larger corpora to train sm)

Symbolic and Statistical Paradigms in Computational Linguistics

- Historically: conflict between symbolic and statistical approaches (Chomsky: Syntactic Structures (1957))
- Reasons
 - Role of quantitative measures
 - Type of data to be investigated
- Computational Approaches (in the 60s):
 - Generative theory; building linguistic models based on formal grammars
- 90s (computers faster; disk space cheaper):
 - Renewed interest in statistical models
- 1994: Workshop: The Balancing Act (New Mexico State University); Goal: dialogue between researchers of both sides
- Since then hybrid approaches growing

Summary

- Introduction to corpus linguistics
- Chomsky „Syntactic Structures“ (1957)
- Corpus types and corpora resources
- Applications influences corpus collection and processing
- Tools and programs to analyse corpora



Discussion please...

Text Mining, Corpus Building, and
Testing