TEXT MINING

LO2. PREPROCESSING

SUZAN VERBERNE 2022



TODAY'S LECTURE

Quiz about week 1

- Go from raw text to clean text
- Character encoding
- Edit distance (+ exercise)
- Regular expressions
- Tokenization and sentence splitting
- Lemmatization and stemming



- If we use words as features in a text classification task, the resulting vectors are
 - a. Low-dimensional and dense
 - b. Low-dimensional and sparse
 - c. High-dimensional and dense
 - d. High-dimensional and sparse



- If we use words as features in a text classification task, the resulting vectors are
 - a. Low-dimensional and dense
 - b. Low-dimensional and sparse
 - c. High-dimensional and dense
 - d. High-dimensional and sparse



- What does the long-tail distribution for text data refer to?
 - a. When we add a document to a collection, the number of unique terms will increase
 - b. In a given document, most of the terms will have a frequency of zero
 - c. In a given collection, there are many terms with a low frequency and few terms with a high frequency



- What does the long-tail distribution for text data refer to?
 - a. When we add a document to a collection, the number of unique terms will increase
 - b. In a given document, most of the terms will have a frequency of zero
 - c. In a given collection, there are many terms with a low frequency and few terms with a high frequency



- For which type of text processing task are capitalization and punctuation more useful?
 - a. For sequence labelling
 - b. For classification



- For which type of text processing task are capitalization and punctuation more useful?
 - a. For sequence labelling
 - b. For classification



- Which evaluation metric would you prioritize for the task of identifying terrorist threats on Twitter, and why?
 - a. Precision, because we want to be sure that we found all threats
 - b. Recall, because we want to be sure that we found all threats
 - c. Precision, because we don't want to accuse someone wrongly
 - d. Recall, because we don't want to accuse someone wrongly



- Which evaluation metric would you prioritize for the task of identifying terrorist threats on Twitter, and why?
 - a. Precision, because we want to be sure that we found all threats
 - b. Recall, because we want to be sure that we found all threats
 - c. Precision, because we don't want to accuse someone wrongly
 - d. Recall, because we don't want to accuse someone wrongly



WHO FOUND A TEAM MATE?

> (please raise your hands)



GO FROM RAW TEXT TO CLEAN TEXT



Suzan Verberne 2022

GATHERING RAW TEXT

Written text

- Digitized (scanned) documents (We need OCR = optical character recognition)
- Born-digital documents
 - text, html, pdf, MS Word documents

All text needs clean-up of some kind



DIGITAL INPUT NOT CLEAN

Scanned text and born-digital PDFs might contain:

- photos, tables, graphics
- layout or design information
- disclaimers, copyright statements
- headers, footers
- column and page breaks
- OCR errors
- character encoding errors

Semi-structured text: text with markup (HTML, XML, docx)



OPTICAL CHARACTER RECOGNITION (OCR)

TO THE TRVLY HO-NORABLE AND RIGHT WOR.

THY KNIGHT SIR THOMAS SMITH,

T R E A S V R E R for the Colonies and Companies of VIRGINIA: and Gouernour of Mulcouia, East-India, North-west Passage, and S 0 M M E R Ilands Companies.



ONORABLE SIR, the wifelt of Men, or rather the wifedome of God tells vs, that there Ecclef.3. t... is a time for all things: and that the great God, who at his owne will beganne Time it felfe, doth at his owne time beginne all things elfe: the foolifhneffe of men may aske and mufe why

was this fo foone, and that fo late ? but the wifedome of God knowes what is fit for every time : And furely amongst the fenfible fignes, and euident demonstrations of Gods all-gouerning prouidence, this is not the leaft, that he brings not forth his mightie works altogether, but makes every thing beautifull in his time. Eccl.2. M. And as in his creation he made not al at once, but produced them in their feuerall daies : fo in his gubernation, he reueileth not the knowledge of all things in one Age, but difcouers them in the feuerall ages of the World. And if man aske why God doth thus, holy David gives the answere; The Lord bath fo done his mar- Pfal. 111.4. vailous works, that they should be had in remembrance; for were they all in one age (fuch is our corruption) they would bee leffe observed and sooner forgotten, but being declared in their feuerall times, eucry Age finds matter to magnifie God; And therefore He whole glorious name is to be praised for ener, reueils fome meruailous thing in every generation, that fo his name may Pfal 72.19. be praised from Generation to Generation.

TO'THE TRVLY HO.

NORABLE AND RIGHT Wok THY KNIGHT SIR Tiiom^s SMITHY *YR E A S F R E R* for the Colonies and Conj panics of V i& G 1 N 1 A: aod Gouernsur of Mtif. couia, Eaft-lildiaNorth-weltPaffagc, and S 0 M M E 9 llands COM

N 0 R A B L 11 S 1 It, the W121 Of Men, Ot rather the wifedome of God tells vs, thatibere ~r a time for a# tkings: and that the great Gor4, who at his owne will beganne Time it fcl&, dothathisowne timebeginne all things elfe: the-foolithneffe ofinen may askc and tnufc why was this fo roone, and that fo late ? but the wifedome of God knowes what is fit for euery time : And furely among(I the fen fible fignes, and cuidcnt demonfirations of Gods all-gouerning prouidence, this is not the leaft, that he brings not fbrth his migh tie works altogether, but makgrvffY thing boc#tifs.11 in Ur time. rs.

And as in his creation he made not al at once, but produced them in their feuerall daies : fo in his gubernation, hcreuedethnotthc knowledge of all things in one Age, but difcouers them in the fcuerall ages of the World. And if man aske why God doth thus, holy Davidgiues the anfwerc; *The Lord both* **fo done** Ur marpfal.111.4.

HTML SOURCE

47 <title>'I am forever grateful': Prince Harry pays tribute to the Queen Prince Harry The Guardian</title>
48 <meta content="Duke of Sussex calls Oueen Elizabeth II his ‘quiding compass’:" name="description"/>
49 <meta charset="utf-8"/>
51 <meta content="width=device-width,minimum-scale=1,initial-scale=1" name="viewport"/>
52 <pre>set aname="theme-color" content="#052962" /></pre>
53 k rel="icon" href="https://static.guim.co.uk/images/favicon-32x32.ico">
55 <pre><link href="<u>https://assets.guim.co.uk/</u>" rel="preconnect"/></pre>
56 <link href="https://i.guim.co.uk" rel="preconnect"/>
57 <link href="https://j.ophan.co.uk" rel="preconnect"/>
58 <link href="https://ophan.theguardian.com" rel="preconnect"/>
59 <link href="https://sourcepoint.theguardian.com" rel="preconnect"/>
60 <link href="<u>https://assets.guim.co.uk</u>" rel="dns-prefetch"/>
61 <link href="https://i.guim.co.uk" rel="dns-prefetch"/>
62 <link href="https://j.ophan.co.uk" rel="dns-prefetch"/>
63 <link href="https://ophan.theguardian.com" rel="dns-prefetch"/>
64 <link href="https://api.nextgen.guardianapps.co.uk" rel="dns-prefetch"/>
65 <link href="https://hits-secure.theguardian.com" rel="dns-prefetch"/>
66 <link href="<u>https://interactive.guim.co.uk</u>" rel="dns-prefetch"/>
<pre>6/ <link href="https://phar.gu-web.net" rel="dns-pretetch"/></pre>
<pre>b8 <link href="https://static.theguardian.com" rel="uns-prefetch"/></pre>
by <link href="<u>https://support.theguardian.com</u>" rel="ans-pretetch"/>
70
/1 Script type application// types application// scheme aca/ "doid" "https://amp.thequardian.com/uk-news/2020/cen/12/orinea_harry_nave_tributa_ta_tba_nuean" "nuhlichar";//douganization" "acontext": "https://amp.thequardian.com/uk-news/2020/cen/12/orinea_harry_nave_tributa_ta_tba_nuean" "nuhlichar"; //douganization" "acontext": "https://amp.thequardian.com/uk-news/2020/cen/12/orinea_harry_nave_tributa_ta_tba_nuean" "nuhlichar"; //douganization" "acontext": "https://amp.thequardian.com/uk-news/2020/cen/12/orinea_harry_nave_tributa_ta_tba_nuean" "nuhlichar"; //douganization" "acontext": "https://amp.tka.nu/linka.tba_nuean" "nuhlichar"; //amp.tka.tba_nuean" "nuhlichar"; //amp.tka.tba_nuean"; "nuhlichar"; //amp.tka.tba_nuean"; "nuhlichar"; //amp.tka.tba_nuean"; "nuhlichar"; //amp.tka.tba_nuean; "nuhlichar"; "nuhlichar; "nuh
2 It grype - NewsArticle , guintext - Intips://schema.org , gui - Intips://amp.theguardian.com/uk-news/2022/sep/12/pinte-narty-pays-cribute-to-the-queen , publisher . guype - organization , guintext - Intips
75
<pre>76 <link href="https://amo.theguardian.com/uk-news/2022/sep/12/prince-harry-pays-tribute-to-the-gueen" rel="amohtml"/></pre>
77
78 <link as="font" crossorigin="" href="https://assets.guim.co.uk/static/frontend/fonts/guardian-headline_hoalts-not-hinted/GHGuardianHeadline_Medium.woff2?http3=true" rel="preload"/>
79 <link as="font" crossorigin="" href="https://assets.guim.co.uk/static/frontend/fonts/guardian-headline/noalts-not-hinted/GHGuardianHeadline-MediumItalic.woff2?http3=true" rel="preload"/>
80 <link as="font" crossorigin="" href="https://assets.guim.co.uk/static/frontend/fonts/guardian-headline/noalts-not-hinted/GHGuardianHeadline-Bold.woff2?http3=true" rel="preload"/>
81 <link as="font" crossorigin="" href="https://assets.guim.co.uk/static/frontend/fonts/guardian-textegyptian/noalts-not-hinted/GuardianTextEgyptian-Regular.woff2?http3=true" rel="preload"/>
82 <link as="font" crossorigin="" href="https://assets.guim.co.uk/static/frontend/fonts/guardian-textegyptian/noalts-not-hinted/GuardianTextEgyptian-Bold.woff2?http3=true" rel="preload"/>
83 <link as="font" crossorigin="" href="https://assets.guim.co.uk/static/frontend/fonts/guardian-textsans/noalts-not-hinted/GuardianTextSans-Regular.woff2?http3=true" rel="preload"/>
84 <link as="font" crossorigin="" href="https://assets.guim.co.uk/static/frontend/fonts/guardian-textsans/noalts-not-hinted/GuardianTextSans-Bold.woff2?http3=true" rel="preload"/>
<pre>86 <meta content="https://www.theguardian.com/uk-news/2022/sep/12/prince-harry-pays-tribute-to-the-queen" property="og:url"/></pre>
87 <meta content="PA Media" property="article:author"/>
88 <meta content="1200" property="og:image:width"/>
89 <meta al:ios:url"="" content="gnmguardian://uk-news/2022/sep/12/prince-harry-pays-tribute-to-the-queen?contenttype=Article&source=applinks" property="0g:image"/>
91 <meta content="https://www.facebook.com/theguardian" property="article:publisher"/>
<pre>%2 emeta property="0gitite" content="1 am forever graterul': Prince Harry pays tribute to the Queen"/></pre>
30 smleta property= 10:app_10 Content= 10044404020/ />
<pre>94 <meta _time"="" content="2022-09-12109:19:19:53.0002" property="artitle:modilate"/> 05 <meta _time"="" content="2022-09-12109:19:19:53.0002" property="artitle:modilate"/></pre>
35 smeta property- 09.image:neight content- /20 /2 De grade property- 09.image:neight content- /20 /2
so vmeta property- og description - content- buke of sussex carts queen crizabeth if his guiding compass //



https://www.theguardian.com/uk-news/2022/sep/12/prince-harry-pays-tribute-to-the-queen 16

CLEAN TEXT STORAGE

- Markup: meta-information in a text file that is clearly distinguishable from the textual content
 - In the case of XML and json, markup often provides useful information in text processing
 - We typically convert PDF to XML, using pdf-to-xml convertors
 - Also, benchmark data is often stored as XML

Character Encoding: the way that a computer displays text in a way that humans can understand.



CHARACTER ENCODING



Suzan Verberne 2022



> Character encoding: translates a string of 0s and 1s to a character

- > ASCII is a 7-bit encoding based on the English alphabet
 - 1100001 a
 - 1100010 b
 - 1100011 c

ASCII: American Standard Code for Information Interchange



HOW ABOUT THE REST OF THE WORLD?

- corpus linguistics (English)
- corpuslinguïstiek (Dutch)
- कोष भाषा विज्ञान (Hindi)
- (Arabic) الأسانيات الإحضا
- ➢ 語料庫語言學(Chinese)
- Hebrew) בלשנות קורפוס
- การ ศึกษา ภาษาศาสตร์ (Thai)
- Цорпус Лингуистицс (Serbian)

UNICODE

Universal standard for all writing systems (>100,000 characters)

Independent of platform, software, vendor

Interpretation of the character is done by the implementation (e.g. UTF-8) in the software (e.g. editor, printer or web browser) that determines the actual rendering (size, shape, font, style)

For maximum compatibility (forward and backward) we encode texts in UTF-8 when we read and write them



READ AND WRITE UTF-8 IN PYTHON 3

with open(filename,'r',encoding='utf-8') as raw: text = raw.read()

https://docs.python.org/3/library/functions.html#open



DATA CLEANING IN PRACTICE

- Digitalization, data conversion & cleaning are the first steps in the text mining process. These steps are:
 - necessary
 - time consuming
 - error prone
 - often complicated
- When building a text mining pipeline for new (raw) data, data collection and cleaning is a long and tedious step that is not to be underestimated!



EDIT DISTANCE

(SECTION 2.5 IN J&M)



Suzan Verberne 2022

WHY EDIT DISTANCE

For measuring string similarity:

- Spelling correction/normalization (think about normalizing the content of doctor's notes or user-generated content)
 - E.g. 'graffe' what word was meant?
 - \succ 'giraffe' differs by only one letter \rightarrow most likely
 - 'grail' or 'graf' differ in more letters
- Also needed for matching names or terms to databases that might contain spelling errors or typos
 - \succ Gleevic \rightarrow Gleevec (medication name)
 - ➢ Abbasiyantaeb → Abbasiantaeb

MINIMAL EDIT DISTANCE

- the minimum edit distance between two strings is defined as the minimum number of editing operations (insertion, deletion, substitution) needed to transform one string into another
- Levenshtein distance: insertion, deletion and substitution all have a cost of 1.
 - dog-do: 1
 - cat-cart: 1
 - cat-cut: 1
 - cat-act: ?
 - cat-act: 2

Iniversiteit

iden

COMPUTING MINIMAL EDIT DISTANCE

- "The space of all possible edits is enormous, so we can't search naively.
- However, lots of distinct edit paths will end up in the same state (string), so rather than recomputing all those paths, we could just remember the shortest path to a state each time we saw it.
- > We can do this by using dynamic programming.
- Dynamic programming is the name for a class of algorithms that apply a table-driven method to solve problems by combining solutions to sub-problems."

(J&M, section 2.5.1)





Operations: insert (cost 1), delete (cost 1), substitute (cost 1), copy (cost 0)



Suzan Verberne 2022

> Initialization D(i,0) = i D(0,j) = j > Recurrence Relation: For each i = 1...N D(i,j) = min - D(i-1,j) + 1 D(i,j-1) + 1 $D(i-1,j-1) + 1; if X(i) \neq Y(j)$ 0; if X(i) = Y(j)

> Termination: D(N,M) is distance

Universiteit

eiden



Operations: insert (cost 1), delete (cost 1), substitute (cost 1), copy (cost 0)



Suzan Verberne 2022



Operations: insert (cost 1), delete (cost 1), substitute (cost 1), copy (cost 0)



Suzan Verberne 2022

> Initialization D(i,0) = i D(0,j) = j > Recurrence Relation: For each i = 1...N D(i,j) = min - D(i-1,j) + 1 D(i,j-1) + 1 $D(i-1,j-1) + 1; if X(i) \neq Y(j)$ 0; if X(i) = Y(j)

> Termination: D(N,M) is distance

Universiteit

eiden

cost	operation	input	output	
1	substitute	С	f	
0	(copy)	а	а	
1	substitute	t	S	
1	substitute	S	t	
Total: 3				
cost	operation	input	output	
1	substitute	С	f	
0	(сору)	а	а	
1	deletion	t	*	
0	(сору)	S	S	
1	insertion	*	t	
Total: 3			33	

LEVENSHTEIN DISTANCE: EXERCISE

		a	n		a	С	t
	0	1	2	С	4	5	6
a	1						
	2						
C	3						
a	4						
t	5						



LEVENSHTEIN DISTANCE: EXERCISE

Suzan Verberne 2022

LEVENSHTEIN DISTANCE: EXERCISE

cost	operation	input	output
0	(сору)	а	а
1	insert	*	n
0	(сору)		
1	substitute	с	а
1	substitute	а	С
0	(сору)	t	t

REGULAR EXPRESSIONS

SECTION 2.1 IN J&M

Suzan Verberne 2022

WHY REGULAR EXPRESSIONS

- > For finding patterns, e.g.
 - > "[0-9][0-9][0-9][0-9] [A-Z][A-Z]"
 - "https?://\S+"
 - > "\S+@\S+"

WHY REGULAR EXPRESSIONS

Match/count patterns in a file/string

Or extract the matched pattern

Tool for making and debugging regular expressions: https://regex101.com/

Verbose regular expressions: <u>https://www.geeksforgeeks.org/verbose-in-python-regex/</u>

Suzan Verberne 2022

EXAMPLE (FROM THE BOOK)

- Suppose we want to write a regular expression to find all occurrences of the English article *the*
- A simple (but incorrect) pattern might be: /the/
- > Why is this incorrect?

EXAMPLE (FROM THE BOOK)

- Suppose we want to write a regular expression to find all occurrences of the English article *the*
- A simple (but incorrect) pattern might be: /the/
- > Why is this incorrect?
 - We want to find both the and The
 - But not: <u>the</u>ir, apo<u>the</u>cary, etc.
- Instead, if we would really want to cover all occurrences, and no non-relevant occurrences of the string, we need:

```
/(^|[^a-zA-Z])[tT]he([^a-zA-Z]|$)/
```


EXAMPLE (ALTERNATIVES)

from nltk.tokenize import word_tokenize
tokens = word_tokenize(raw_text)
print(tokens.count('the'))

import spacy import en_core_web_sm spacy_nlp = en_core_web_sm.load() spacy_doc = spacy_nlp(example_text) spacy_words = [token.text for token in spacy_doc] print(sum(token == 'the' for token in spacy_words))

TOKENIZATION AND SENTENCE SPLITTING

SECTION 2.2-2.4 & 2.4.5 IN J&M

Suzan Verberne 2022

DEFINITIONS

- Token An instance of a word or term occurring in a document
- Term A token when used as feature (or in an index), generally in normalized form (e.g. lowercased)

- Token count is the number of words (running) in a collection/document; this includes duplicates
- Vocabulary size is the number of unique terms; the feature size when we use words as features

TOKENIZATION

- Tokenization: Split text in tokens
 - (1) remove punctuation; (2) split on whitespaces characters
- Question: how would you want to tokenize these strings?
 - 1. Hewlett-Packard
 - 2. State-of-the-art
 - 3. aren't
 - **4**. C++
 - 5. cheap San Francisco-Los Angeles fares
 - 6. 20/03/1999
 - 7. 071 527 7043

TOKENIZATION WITH NLTK

from nltk.tokenize import word_tokenize tokens = word_tokenize(raw_text, language= 'english')

https://www.nltk.org/api/nltk.tokenize.html

TOKENIZATION WITH SPACY

<u>https://spacy.io/u</u> <u>sage/linguistic-</u> <u>features#how-</u> <u>tokenizer-works</u>

Customization for your task possible

STOP WORDS

- Stop words: extremely common words that don't carry any content
- Examples: a, an, and, are, as, at, be, by, for, from, has, he, in, is, it, its, of, on, that, the, to, was, were, will, with
- Stop word elimination used to be common for classification
- But... For which cases is this problematic?
 - "to be or not to be"
 - You might need stop words for multi-word terms, e.g. "King of Denmark", "Marks and Spencer"
- So term weighting is typically a better option (see lecture 4)

STOP WORDS

> Tasks for which we might want to remove stopwords

- topic modelling (see lecture 8)
- keyword extraction
- We never remove stopwords in:
 - sequence labelling tasks
 - classification tasks with small data (few/short documents)

SENTENCE SPLITTING

- We need sentence splitting for tasks that require sentence-level analysis, for example:
 - Finding the most similar sentences in a collection
 - Extracting the relations between two entities (e.g. a person and their birth year) within one sentence
 - Sentence-level sentiment analysis (e.g. analyzing reviews for positive and negative aspects)
 - Input for Transformer models with limited input length

SENTENCE SPLITTING

Many cases can be solved with a relatively simple approach:

Sentences end with . ? ! followed by whitespace

> Challenges:

- Abbreviations ('Mrs. Doe')
- Names/initials that include punctuation marks ('S. Verberne')
- Sentences without punctuation markers (headers/titles)
- Line endings inside sentences (PDF conversion)

SENTENCE SPLITTING WITH NLTK

from nltk.tokenize import sent_tokenize
sentences = sent_tokenize(document, language='english')

https://www.nltk.org/api/nltk.tokenize.html

SENTENCE SPLITTING WITH SPACY

https://spacy.io/usage/linguistic-features#sbd

Editable Code

spaCy v3.0 · Python 3 · via Binder

```
import spacy
```

```
nlp = spacy.load("en_core_web_sm")
doc = nlp("This is a sentence. This is another sentence.")
assert doc.has_annotation("SENT_START")
for sent in doc.sents:
    print(sent.text)
```


This is a sentence. This is another sentence.

LEMMATIZATION AND STEMMING

(SECTION 2.4.4 IN J&M)

Suzan Verberne 2022

BASIC WORD FORMS

> We might want to normalize specific word forms to the same term:

For example, it can be useful to have the term *bicycle* for each occurrence of either *bicycle* or *bicycles*

BASIC WORD FORMS

> We might want to normalize specific word forms to the same term:

- For example, it can be useful to have the term *bicycle* for each occurrence of either *bicycle* or *bicycles*
- Advantages:
 - reduces the number of features
 - generalizes better, especially for small datasets

BASIC WORD FORMS

- > We might want to normalize specific word forms to the same term:
 - For example, it can be useful to have the term *bicycle* for each occurrence of either *bicycle* or *bicycles*
- Advantages:
 - reduces the number of features
 - generalizes better, especially for small datasets
- Two types of basic word forms:
 - Lemma
 - Stem

LEMMA

Lemma: dictionary form of a word

> For example:

- Verbs: infinitive
 - 'think' for 'thinks', 'thinking', 'thought'
- Nouns: singular form
 - 'mouse' for 'mice'
 - 'computer' for 'computers'

Stem: the portion of a word that:

- is common to a set of (inflected) forms when all affixes are removed
- is not further analyzable into meaningful elements, being morphologically simple

- The stem is the part of the word that 'never' changes even when morphologically inflected. It is not necessarily an existing word:
 - 'comput' for forms 'computer', 'computing', 'computers', 'compute'

STEM VS LEMMA

- 'produced'
 - what is the lemma?
 - what is the stem?

STEM VS LEMMA

'produced'

- what is the lemma?
- what is the stem?

Answers:

- the lemma is 'produce'
- the stem is 'produc'

(think of 'producing' as one of the morphological forms of the verb)

We almost always prefer lemmas over stems. Stemming can be effective for very small collections.

EXAMPLE

- Sample text: Such an analysis can reveal features that are not easily visible from the variations in the individual genes and can lead to a picture of expression that is more biologically transparent and accessible to interpretation
- Lemmatizer: Such an analysis can reveal feature that be not easily visible from the variation in the individual gene and can lead to a picture of expression that be more biologically transparent and accessible to interpretation
- Stemmer: such an analysi can reveal featur that ar not easili visibl from the variat in the individu gene and can lead to a pictur of express that is more biolog transpar and access to interpret

SUZAN VERBERNE 2022

PYTHON PACKAGES AND TOOLS

Sklearn (<u>http://scikit-learn.org</u>) has built-in functionality for:

- Tokenization
- <u>https://scikit-</u> <u>learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVe</u> <u>ctorizer.html</u>
- Stop word removal (with option to supply your own list)
- NLTK (<u>http://www.nltk.org/</u>) and Spacy (<u>https://spacy.io/</u>) have functionality for:
 - Sentence splitting
 - Lemmatization and stemming
 - and additional pre-processing steps

HOMEWORK

- Read Jurafsky & Martin chapter 2: "Regular Expressions, Text Normalization, Edit Distance" (Brightspace)
- Complete this week's exercise: preprocessing with spacy. <u>https://course.spacy.io/en/chapter1</u> (note: only chapter 1, consisting of 12 parts)
- Contact address for the teaching assistants: <u>tmcourse@liacs.leidenuniv.nl</u>

AFTER THIS LECTURE...

- you know what issues to take into account when converting raw text to clean plain text
- you can explain the the difference between ASCII and Unicode; and why Unicode exists
- you can compute the minimal edit distance between two strings using the matrix algorithm
- you can use regular expressions in Python and you know what the possibilities and challenges of regular expressions are
- you can describe the challenges of tokenization and sentence splitting
 - you can explain the considerations for stop word removal
- you can define the difference between stemming and lemmatization

