

# First-Year Project 4: Natural Language Processing

Christian Hardmeier

27 April 2021

# Natural Language Processing

# Machine Translation

Translate text Translate .docx & .pptx files

Translate from **Latvian** (detected) ▾

Mašintulkšana, saīsināti MT, ir datorlingvistikas apakšnozare, kas nodarbojas ar automatizētu vienā valodā rakstīta teksta tulkošanu citā valodā.

MT programmas atšķirībā no citiem ar datora izmantošanu saistītiem paņēmieniem – datorizētas tulkošanas (computer-assisted translation, computer-aided translation — CAT) un interaktīvas tulkošanas (tulkošanas ar datora palīdzību) – veic tulkošanu ar minimālu cilvēka līdzdalību vai bez tās.

Translate into **Danish** ▾

Maskinoversættelse, forkortet MT, er et delområde inden for computerlingvistik, der beskæftiger sig med automatisk oversættelse af tekst skrevet på ét sprog til et andet sprog.

I modsætning til andre computerstøttede oversættelser (CAT) og computerstøttede oversættelsesteknikker (CAT) udfører MT-programmer oversættelser med ringe eller ingen menneskelig indblanding.

▶ 🔍

🔗 ⚡

# Writing Aids

## Demo document

### The basics

Mispellings and grammatical errors can effect your credibility. The same goes for misused commas, and other types of punctuation. Not only will Grammarly underline these issues in red, it will also showed you how to correctly write the sentence.

### But wait...there's more?

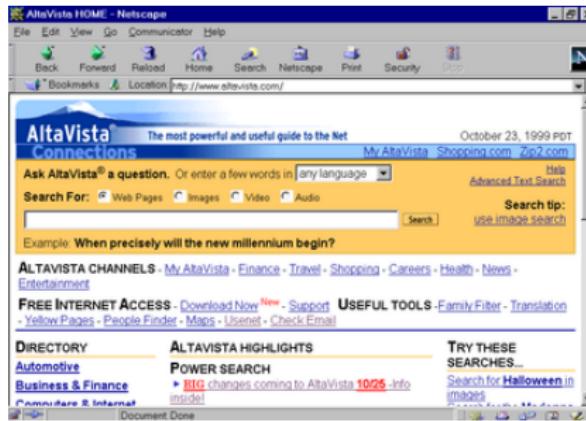
Blue underlines mean a clarity issue has been spotted by Grammarly. You'll find suggestions that can possibly help you revise a wordy sentence in an effortless manner.

Grammarly will also inspect your vocabulary carefully and suggest the best word with green underlines to make sure you don't have to analyze your writing too much.

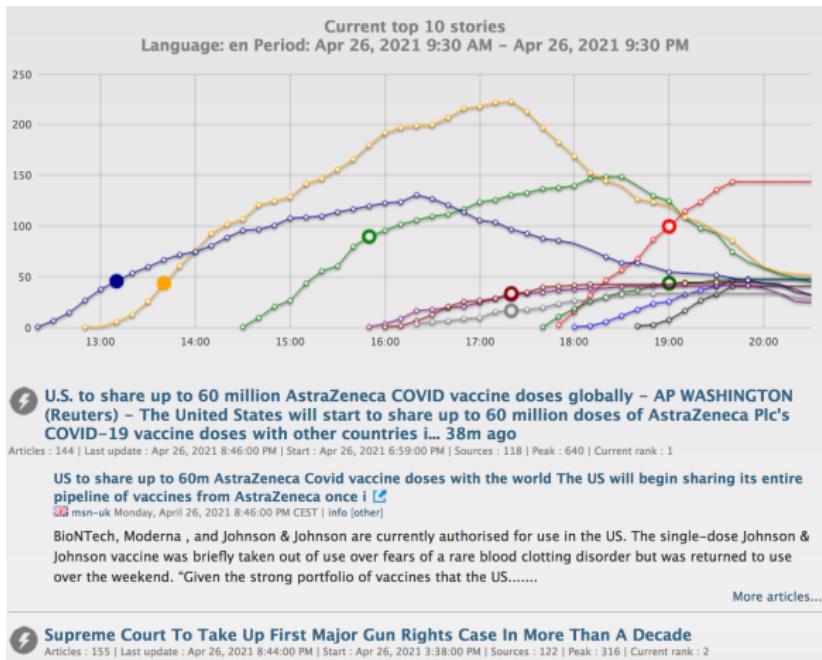
The screenshot shows the Grammarly interface with the following details:

- All alerts**: 16 alerts
- Hide Assistant**
- Overall score**: 61 (See performance)
- Goals**: 3 of 5 set
- All alerts** (highlighted):
  - Correctness**: 6 alerts (red bar)
  - Clarity**: A bit unclear (blue bar)
  - Engagement**: A bit bland (green bar)
  - Delivery**: Serious issues (purple bar)
- SPELLING**: Mispellings → **Misspellings**
  - The word **Misspellings** is not in our dictionary. If you're sure this spelling is correct, you can add it to your personal dictionary to prevent future alerts.
- effect**: Correct your spelling
- commas**: Remove the comma
- punctuation**: Remove a space
- ,**: Add the word(s)

# Information retrieval



# Information aggregation/extraction



<https://emm.newsbrief.eu/>

# Natural Language

## **Where Young College Graduates Are Choosing to Live**

*New York Times, 20 Oct 2014*

When young college graduates decide where to move, they are not just looking at the usual suspects, like New York, Washington and San Francisco. Other cities are increasing their share of these valuable residents at an even higher rate and have reached a high overall percentage, led by Denver, San Diego, Nashville, Salt Lake City and Portland, Ore., according to a report published Monday by City Observatory, a new think tank.

And as young people continue to spurn the suburbs for urban living, more of them are moving to the very heart of cities — even in economically troubled places like Buffalo and Cleveland. The number of college-educated people age 25 to 34 living within three miles of city centers has surged, up 37 percent since 2000, even as the total population of these neighborhoods has slightly shrunk.

Some cities are attracting young talent while their overall population falls, like Pittsburgh and New Orleans. And in a reversal, others that used to be magnets, like Atlanta and Charlotte, are struggling to attract them at the same rate.

# **Moxifloxacin Loaded Nanoemulsions having Tocopheryl Succinate as Integral Component Improves Pharmacokinetics and Enhances Survival in E Coli Induced Complicated Intra Abdominal Infection.**

Shukla P, Verma AK, Dwivedi P, Yadav A, Gupta PK, Rath SK, Mishra PR.

## **Abstract**

In the present work a novel nanoemulsion laden with moxifloxacin has been developed for effective management of complicated Intra-abdominal infections. Moxifloxacin nanoemulsion fabricated using high pressure homogenization was evaluated for various pharmaceutical parameters, pharmacokinetics and pharmacodynamics in rats with *E coli* induced sepsis. The developed nanoemulsion MONe6 (size  $168\pm28$  d.nm and ZP  $-24.78\pm0.45$ mV respectively) was effective for intracellular delivery and sustaining the release of MOX. MONe6 demonstrated improved plasma ( $AUC_{MONe6}/MOX = 2.38$  fold) and tissue pharmacokinetics of MOX ( $AUC_{MONe6}/MOX = 2.63$  and 1.47 times in lung and **liver** respectively). Calculated PK/PD index correlated well with reduction in bacterial burden in plasma as well as tissues. Enhanced survival on treatment with MONe6 (65.44 %) and as compared to control group (8.22 %) was result of reduction in lipid peroxidation, neutrophil migration and cytokine levels (TNF- $\alpha$  and IL1b) as compared to untreated groups in rat model of *E coli* induced sepsis. Parenteral nanoemulsions of MOX hold a promising advantage in the therapy of *E. coli* induced complicated intra-abdominal infections and is helpful in the prevention of further complication like septic shock and death.

Athough people enjoyment differs from stage to stage, youth enjoys more as compared with the older people is definitely reasonable to agree with the statement. I observed in my life at part, that young people enjoys more.

In modern world, due to the development of technology, they are more fascinated to new kind of objects, such as mobiles, computer, internet, etc., and by using them that were available enjoying life.

Youngers are more affectionated to the day to day development of technology, by using those things they enjoying more. Older people are not that much aware of these technologies. Because technology is developing rapidly in these day to day life. Younger ones are more energetic than older people, they can go and do whatever they want. They make their friends as a part to enjoy in most of the public places, they usually attends party and meetings, and have capability to think and improve by adapting new technologies they enjoys more.



**Barbara Plank** @barbara\_plank · Apr 21

...

Mirella Lapata is introducing the first [#EACL2021NLP](#) keynote talk by Marco Baroni

He has contributed in so many ways to [#NLProc](#) - looking forward to his talk!



1

6



**Omnia Zayed** @OmniaHZayed · Apr 20

...

Come and join our amazing team in the lovely city of Galway! 3 PhD positions are available @unlp\_nuig of @DSlatNUIG @nuigalway on the intersection between NLP and Multimodal Data Analysis. The positions are funded by @insight\_centre. Check the advert below for more info.  
[#nlproc](#)

# What distinguishes natural language data?

# Tabular data

D	E	F	G	H	I	J	K	L
Longitude	Latitude	Police_Force	Accident_Severity	Number_of_Vehicles	Number_of_Casualties	Date	Day_of_Week	Time
-0.153842	51.508057	1	3	2	3	18/02/2019	2	17:50
-0.127949	51.436208	1	3	2	1	15/01/2019	3	21:45
-0.124193	51.526795	1	3	2	1	01/01/2019	3	01:50
-0.191044	51.546387	1	2	1	1	01/01/2019	3	01:20
-0.200064	51.541121	1	3	2	2	01/01/2019	3	00:40
0.020461	51.548879	1	3	2	3	01/01/2019	3	02:45
-0.099071	51.367605	1	3	1	1	01/01/2019	3	01:35
-0.088978	51.489509	1	3	3	5	01/01/2019	3	02:10
0.141957	51.572326	1	3	2	1	01/01/2019	3	01:15
-0.243769	51.399529	1	3	3	1	01/01/2019	3	04:30
0.070738	51.556734	1	3	1	1	01/01/2019	3	01:15
-0.021065	51.533238	1	1	1	1	01/01/2019	3	03:00
-0.087182	51.549218	1	3	1	1	01/01/2019	3	02:45
0.070277	51.557075	1	2	2	1	01/01/2019	3	04:10
-0.170889	51.49621	1	3	3	1	01/01/2019	3	00:20
-0.007064	51.506832	1	2	1	1	01/01/2019	3	05:55
-0.376691	51.509481	1	3	2	1	01/01/2019	3	00:50
-0.393167	51.533115	1	3	2	1	01/01/2019	3	02:45
-0.111369	51.387849	1	2	2	1	01/01/2019	3	05:54
0.048712	51.547165	1	3	2	1	01/01/2019	3	07:30
-0.336704	51.482519	1	3	2	1	01/01/2019	3	09:41
-0.065175	51.592686	1	3	2	1	01/01/2019	3	14:18
-0.1046	51.507028	1	3	1	1	01/01/2019	3	15:20

## Structured representation

- ▶ Set of **data points** with associated **attributes** (or features).
- ▶ Known **metadata** describing the format and meaning of each attribute:

<i>Longitude</i>	Fractional degrees	E of Greenwich
<i>Latitude</i>	Fractional degrees	N of Equator
<i>Police_Force</i>	Integer	Index into list
<i>Number_of_Vehicles</i>	Integer	Vehicle count

- ▶ Observations are **interpreted** in a particular way.

# Image data



## Unstructured 2D signal

- ▶ 2-dimensional spatial organisation.  
Every pixel has 8 neighbours.
- ▶ Each individual pixel carries almost no information.  
Adding noise to any or all pixels doesn't change the picture.
- ▶ Structure and meaning arise from pixels in context.
- ▶ Strong correlations between nearby elements.

Det mörkeras mer och mer: December =  
mörker i augusti. De svarta rosabladerna  
ha redan släcktsit sig. Men kortsäpporna  
nå minn bord lyda i gälla skräckand förg  
i allt det grå som för att närmast om att  
de en gång blifvo uppförna för att skrygg  
en sjuk och galen frestes svarmod. Men

på svenska. Det är på det hela taget ett förbannat språk vi ha. Orden trampa varandra på tårna och knuffa varandra i ränstenen. Och allt blir så påtagligt och rått. Inga halvtoner, inga lättanta sydningar och mjuka övergångar. Ett språk som tycks vara skapat till bruk för den outrotliga pöbelvanan att plumpa ut med sanningen i alla väder.

Det mörknar mer och mer: decembermörker i augusti. De svarta rosenbladen ha redan skrynklat sig. Men kortlapparna på mitt bord lysa i gälla skrattande färger i allt det grå som för att påminna om att de en gång blevo uppfunna för att skingra en sjuk och galen furstes svårmod. Men jag fasar vid blotta tanken på arbetet att samla ihop dem och vända de aviga rätt och blanda dem till en ny patiens, jag kan bara sitta och se på dem och lyssna till hur "hjärterknekten och spaderdamen viska dystert om sin begravda kärlek", som det står i samma sonett.

*Le beau valet de cœur et la dame de pique  
causent sinistrement de leurs amours défunts.*

Jag kunde ha lust att gå upp i det smutsiga gamla rucklet där snett över och dricka porter med flickorna. Röka en sur pipa och dra en spader med värdinnan och ge henne goda råd för hennes reumatism. Hon var här i förra veckan och klagade sin nöd, fet och präktig. Hon hade en tjock guldbrosch under isterhakan och betalade en femma kontant. Hon skulle bli smickrad av en kontravisit.

Det ringde på tamburdörren. Nu öppnar Kristin... Vad kan det vara? Jag har ju sagt till att jag inte tar emot i dag... En detektiv?... Som låtsas vara sjuk, uppträder som patient... Kom in du, min gubbe, jag skall nog sköta om dig...

Kristin gläntade på dörren och slängde ett brev med svarta kanter på mitt bord. Inbjudning att bevista jordfästningen...

\*

— Min handling, ja... "Vil Monsieur have den Historie paa heroiske Vers, saa koster det 8 Skilling..."

25 augusti.

Jag såg i drömmen gestalter från min ungdom. Jag såg henne som jag kysste en midsommarnatt för länge sedan, då jag var ung och icke hade dödat någon. Jag såg också andra unga flickor av dem som hörde till vår krets den tiden; en som gick och läste det året jag blev student och som alltid ville tala med mig om religionen; en annan, som var äldre än jag och som gärna stod och viskade med mig i skymningen bakom en jasminhäck i vår trädgård. Och en annan, som alltid gjorde narr av mig, men som blev så ond och häftig och föll i krampgråt en gång då jag gjorde narr av henne... De gingo bleka i en blek skymning, deras ögon stodo vidöppna och förskrämnda, och de gjorde tecken åt

## **Where Young College Graduates Are Choosing to Live**

*New York Times, 20 Oct 2014*

When young college graduates decide where to move, they are not just looking at the usual suspects, like New York, Washington and San Francisco. Other cities are increasing their share of these valuable residents at an even higher rate and have reached a high overall percentage, led by Denver, San Diego, Nashville, Salt Lake City and Portland, Ore., according to a report published Monday by City Observatory, a new think tank.

And as young people continue to spurn the suburbs for urban living, more of them are moving to the very heart of cities — even in economically troubled places like Buffalo and Cleveland. The number of college-educated people age 25 to 34 living within three miles of city centers has surged, up 37 percent since 2000, even as the total population of these neighborhoods has slightly shrunk.

Some cities are attracting young talent while their overall population falls, like Pittsburgh and New Orleans. And in a reversal, others that used to be magnets, like Atlanta and Charlotte, are struggling to attract them at the same rate.

## Written text: Unstructured 1D signal with discrete elements

- ▶ One-dimensional spatial (linear) organisation.  
Every token has 2 neighbours.
- ▶ Discrete elements:  
Changing a word is *categorical*, not gradual.
- ▶ Hierarchical structure at multiple levels.  
Structural details poorly understood  
(competing linguistic theories).
- ▶ Strong correlations between nearby elements.
- ▶ Individual points (words/tokens) inherently meaningful,  
but ambiguous.

# What makes language processing difficult?

# Ambiguity

Lexical ambiguity: **Bank**



[https://commons.wikimedia.org/wiki/File:European\\_Central\\_Bank\\_041107.jpg](https://commons.wikimedia.org/wiki/File:European_Central_Bank_041107.jpg)

[https://commons.wikimedia.org/wiki/File:River\\_Erosion\\_-\\_geograph.org.uk\\_-\\_358650.jpg](https://commons.wikimedia.org/wiki/File:River_Erosion_-_geograph.org.uk_-_358650.jpg)

# Ambiguity

## Structural ambiguity

I saw the man in the park with the telescope.

- ▶ Who has the telescope?

## Referential ambiguity

Anna has a little sister. She loves her very much.

- ▶ Who loves whom?

# Vagueness

- ▶ Language is often *vague* or *underspecified*,  
and it would be unnatural to be totally precise in every situation.
- ▶ What is bigger?
  - ▶ A large dog,
  - ▶ or a small elephant?
- ▶ At what time does the afternoon end and the evening start?

# Variation

- ▶ Languages
  - ▶ (even in one text: *code switching*)
- ▶ Register
  - ▶ Formal vs. informal
  - ▶ Written vs. spoken
- ▶ Domain
  - ▶ What *is* the text about?

## World knowledge

- ▶ Humans use *world knowledge* to interpret language.
- ▶ Cues from *context* of language use:
  - ▶ Textual context
  - ▶ Audiovisual context
  - ▶ Situational context
  - ▶ Cultural context

## World knowledge

- Well, what? He's not happy?
  - He can't be, can he, if he's, you know, messing around.
  - You gonna see him again?
  - Do you think I should?
  - Want me to be honest?
  - No. No.
  - When you gonna take **that thing** off?
  - **It's** too tight. I've got to get **it** cut off.
  - Mm.
- What is *that thing*?

Lantana (2001)

# Project Overview

# Overview

- ▶ Social media data from Twitter
- ▶ One tweet per line. User handles anonymised as `@user`.
- ▶ Predict speaker's *intentions* or *state of mind* (pragmatics).
- ▶ Supervised classification:
  - ▶ Given an input, predict a label.
  - ▶ Trained on data with manually annotated labels.
- ▶ 7 seven different tasks – pick 2 of them!
- ▶ *Binary* classification: 2 classes  
*Multiclass* classification: more than 2 classes

# Binary classification tasks

- ▶ Irony Detection
  - ▶ *Off to bed can't wait to feel this hangover.*
  - ▶ *Jeez it's a lovely morning out!!*  #Ireland #December
- ▶ Offensive Language Detection
  - ▶ *@user @user @user She is a walking talking lie.. that's why.*
  - ▶ *@user Eric holder is Obama's straw man of corruption*
- ▶ Hate Speech Detection
  - ▶ Lots of really disgusting hate speech
  - ▶ Personal attacks, often sexualised

# Emotion Recognition

- ▶ ANGER

*Ppl like that irritate my soul*

- ▶ JOY

*Happy Birthday @user #cheer #cheerchick #jeep #jeepgirl  
#IDriveAJeep #jeepjeep #Cheer*

- ▶ OPTIMISM

*The point of living, and being an optimist, is to be foolish enough to believe the best is yet to come' - Peter Ustinov  
#optimism #quote*

- ▶ SADNESS

*@user wow I'm just really sadden by that. Terrible*

# Emoji Prediction

20 labels – 😍, 😂, ..., 🎄, 📸, 😜



*Man these are the funniest kids ever!! That face!*  
*#HappyBirthdayBubb @ FLIPnOUT Xtreme*



*Sundays are all about the cute babies and dogs! #Ballard*  
*#sundaymarket #littlestmodel...*



*Christmas is up!! (@ The Dog House in Seattle, WA)*



*A spicy Volcano Roll just erupted in my mouth! Delectable!*

# Sentiment Analysis

- ▶ NEGATIVE

*Thanks manager for putting me on the schedule for Sunday*

- ▶ NEUTRAL

*Who wants to be my date to the White Sox vs Red Sox game  
Tuesday*

- ▶ POSITIVE

*Happy Birthday Nick J May you live long and Happy :)*

# Stance Detection

3 labels – *favour, neutral, against*  
in relation to 5 target topics:

- ▶ Abortion
- ▶ Atheism
- ▶ Climate change
- ▶ Feminism
- ▶ Hillary Clinton

## Stance Detection: Atheism

- ▶ NONE

*@user Old age has not made you any wiser or more mature.  
For shame!*

- ▶ FAVOUR

*If you regularly base your thoughts on superstitions, you might  
not be able to think well. #freethinker*

- ▶ AGAINST

*Daily time in God's Word yields lasting freedom. #assurance*

# Where do those labels come from?

- ▶ Manually labelled (annotated) data: 6 out of 7 tasks
  - ▶ Very popular method to encode human knowledge in small to medium datasets
  - ▶ Requires clear task definition and guidelines
  - ▶ Quality control with multiple annotators
- ▶ *Fortuitous* data: Emoji classification

*@user it's all good* 😊

{ TEXT: *@user it's all good*  
LABEL: 😊

# What will we learn?

- ▶ Basic properties of natural language
- ▶ Preprocessing: How to prepare texts for automatic processing?
- ▶ Annotation:
  - ▶ How to label up data for supervised classification?
  - ▶ How to evaluate the quality of manual annotations?
- ▶ Basic NLP feature extraction and classification
- ▶ Evaluation of classifiers in NLP

# Preprocessing

## Data acquisition

From a customer/  
employer

From a data provider (LDC, ELRA,  
ELRC)

Web scraping

Human informants

### Data formats

Presentation-oriented  
formats (HTML,  
DOC/DOCX, PDF)

Encodings (“code  
pages”; Unicode,  
Latin-1, ...)

Metadata conserva-  
tion:

*Where does the  
data come from?*

*Who produced it?*

*Where can I find  
more context?*

### Cleaning

Language  
identification

Page head-  
ers/footers

Spelling errors

### Segmentation

Documents

Chapters/sections

Sentences

Words

på svenska. Det är på det hela taget ett förbannat språk vi ha. Orden trampa varandra på tårna och knuffa varandra i ränstenen. Och allt blir så påtagligt och rått. Inga halvtöner, inga lättä antydningar och mjuka övergångar. Ett språk som tycks vara skapat till bruk för den outrotliga pöbelvanan att plumpa ut med sanningen i alla väder.

Det mörknar mer och mer: decembermörker i augusti. De svarta rosenbladen ha redan skrynklat sig. Men kortlapparna på mitt bord lysa i gälla skrattande färger i allt det grå som för att påminna om att de en gång blevo uppfunkna för att skingra en sjuk och galen furstes svårmod. Men jag fasar vid blotta tanken på arbetet att samla ihop dem och vända de aviga rätt och blanda dem till en ny patiens, jag kan bara sitta och se på dem och lyssna till hur "hjärterknekt och spaderdam viska dystert om sin begravda kärlek", som det står i samma sonett.

*Le beau valet de cœur et la dame de pique  
causent sinistrement de leurs amours défunts.*

Jag kunde ha lust att gå upp i det smutsiga gamla rucklet där snett över och dricka porter med flickorna. Röka en sur pipa och dra en spader med värdinnan och ge henne goda råd för hennes reumatism. Hon var här i förra veckan och klagade sin nöd, fet och präktig. Hon hade en tjock guldbrosch under isterhakan och betalade en femma kontant. Hon skulle bli smickrad av en kontravisit.

Det ringde på tamburdörren. Nu öppnar Kristin ... Vad kan det vara? Jag har ju sagt till att jag inte tar emot i dag ... En detektiv? ... Som låtsas vara sjuk, uppträder som patient ... Kom in du, min gubbe, jag skall nog sköta om dig ...

Kristin gläntade på dörren och slängde ett brev med svarta kanter på mitt bord. Inbjudning att bevisa jordfästningen ...

\*

— Min handling, ja ... "Vil Monsieur have den Historie paa heroske Vers, saa koster det 8 Skilling ..." .

25 augusti.

Jag såg i drömmen gestalter från min ungdom. Jag såg henne som jag kysste en midsommarnatt för länge sedan, då jag var ung och icke hade dödat någon. Jag såg också andra unga flickor av dem som hörde till vår krets den tiden; en som gick och läste det året jag blev student och som alltid ville tala med mig om religionen; en annan, som var äldre än jag och som gärna stod och viskade med mig i skymningen bakom en jasminhäck i vår trädgård. Och en annan, som alltid gjorde narr av mig, men som blev så ond och häftig och föll i krampgråt en gång då jag gjorde narr av henne ... De gingo bleka i en blek skymning, deras ögon stodo vidöppna och förskrämda, och de gjorde tecken åt

på svenska. Det är på det hela taget ett förbannat språk vi ha. Orden trampa varandra på tårna och knuffa varandra i ränstenen. Och allt blir så påtagligt och rått. Inga halvtöner, inga lättä antydningar och mjuka övergångar. Ett språk som tycks vara skapat till bruk för den outrotliga pöbelvanan att plumpa ut med sanningen i alla väder.

Det mörknar mer och mer: decembermörker i augusti. De svarta rosenbladen ha redan skrynklat sig. Men kortlapparna på mitt bord lysa i gälla skrattande färger i allt det grå som för att påminna om att de en gång blevo uppfunkna för att skingra en sjuk och galen furstes svårmad. Men jag fasar vid blotta tanken på arbetet att samla ihop dem och vända de aviga rätt och blanda dem till en ny patiens, jag kan bara sitta och se på dem och lyssna till hur "hjärterknekt och spaderdam viska dystert om sin begravda kärlek", som det står i samma sonett.

Le beau valet de cœur et la dame de pique  
causent sinistrement de leurs amours défunts.

Jag kunde ha lust att gå upp i det smutsiga gamla rucklet där snett över och dricka porter med flickorna. Röka en sur pipa och dra en spader med värdinnan och ge henne goda råd för hennes reumatism. Hon var här i förra veckan och klagade sin nöd, fet och präktig. Hon hade en tjock guldbrosch under isterhakan och betecklade en femme kontant. Hon skulle bli smickrad.

Det ringde på tamburdörren. Nu öppnar Kristin ... Vad kan det vara? Jag har ju sagt till att jag inte tar emot i dag ... En detektiv? ... Som låtsas vara sjuk, uppträder som patient ... Kom in du, min gubbe, jag skall nog sköta om dig ...

Kristin gläntade på dörren och slängde ett brev med svarta kanter på mitt bord. Inbjudning att bevista jordfästningen ...

\*

---

— Min handling, ja ... "Vil Monsieur have den Historie paa heroske Vers, saa koster det 8 Skilling ..."

---

25 augusti.

Jag såg i drömmen gestalter från min ungdom. Jag såg henne som jag kysste en midsommarnatt för länge sedan, då jag var ung och icke hade dödat någon. Jag såg också andra unga flickor av dem som hörde till vår krets den tiden; en som gick och läste det året jag blev student och som alltid ville tala med mig om religionen; en annan, som var äldre än jag och som gärna stod och viskade med mig i skymningen bakom en jasminhäck i vår trädgård. Och en annan, som alltid gjorde narr av mig, men som blev så ond och häftig och föll i krampgråt en gång då jag gjorde narr av henne ... De gingo till landet i blek skymning, deras ögon stodofämnda, och de gjorde tecken åt

## Chapters/sections

på svenska. Det är på det hela taget ett förbannat språk vi ha. Orden trampa varandra på tårna och knuffa varandra i ränstenen. Och allt blir så påtagligt och rått. Inga halvtöner, inga lättä antydningar och mjuka övergångar. Ett språk som tycks vara skapat till bruk för den outrotliga pöbelvanan att plumpa ut med sanningen i alla väder.

Det mörknar mer och mer: decembermörker i augusti. De svarta rosenbladen ha redan skrynklat sig. Men kortlapparna på mitt bord lysa i gälla skrattande färger i allt det grå som för att påminna om att de en gång blevo uppfunkna för att skingra en sjuk och galen furstes svårmod. Men jag fasar vid blotta tanken på arbetet att samla ihop dem och vända de aviga rätt och blanda dem till en ny patiens, jag kan bara sitta och se på dem och lyssna till hur "hjärterknekt och spaderdam viska dystert om sin begravda kärlek", som det står i samma sonett.

Le beau valet de cœur et la dame de pique  
causent sinistrement de leurs amours défunts.

Jag kunde ha lust att gå upp i det smutsiga gamla rucklet där snett över och dricka porter med flickorna. Röka en sur pipa och dra en spader med värdinnan och ge henne goda råd för hennes reumatism. Hon var här i förra veckan och klagade sin nöd, fet och präktig. Hon hade en tjock guldbrosch under isterhakan och betalade en femma kontant. Hon skulle bli smickrad a

Det ringde på tamburdörren. Nu öppnar Kristin ... Vad kan det vara? Jag har ju sagt till att jag inte tar emot i dag ... En detektiv? ... Som låtsas vara sjuk, uppträder som patient ... Kom in du, min gubbe, jag skall nog sköta om dig ...

Kristin gläntade på dörren och slängde ett brev med svarta kanter på mitt bord. Inbjudning att bevisa jordfästningen ...

\*

— Min handling, ja ... "Vil Monsieur have den Historie paa heroske Vers, saa koster det 8 Skilling ..."

25 augusti.

Jag såg i drömmen gestalter från min ungdom. Jag såg henne som jag kysste en midsommarnatt för länge sedan, då jag var ung och icke hade dödat någon. Jag såg också andra unga flickor av dem som hörde till vår krets den tiden; en som gick och läste det året jag blev student och som alltid ville tala med mig om religionen; en annan, som var äldre än jag och som gärna stod och viskade med mig i skymningen bakom en jasminhäck i vår trädgård. Och en annan, som alltid gjorde narr av mig, men som blev så ond och häftig och föll i krampgråt en gång då jag gjorde narr av henne ... De gingo bleka i en blek skymning, deras ögon stodo krämnda, och de gjorde tecken åt

## Metainformation

på svenska. Det är på det hela taget ett förbannat språk vi ha. Orden trampa varandra på tårna och knuffa varandra i ränstenen. Och allt blir så påtagligt och rått. Inga halvtöner, inga lättä antydningar och mjuka övergångar. Ett språk som tycks vara skapat till bruk för den outrotliga pöbelvanan att plumpa ut med sanningen i alla väder.

Det mörknar mer och mer: decembermörker i augusti. De svarta rosenbladen ha redan skrynklat sig. Men kortlapparna på mitt bord lysa i gälla skrattande färger i allt det grå som för att påminna om att de en gång blevo uppfunkna för att skingra en sjuk och galen furstes svårmad. Men jag fasar vid blotta tanken på arbetet att samla ihop dem och vända de aviga rätt och blanda dem till en ny patiens, jag kan bara sitta och se på dem och lyssna till hur "hjärterknekt och spaderdam viska dystert om sin begravda kärlek", som det står i samma sonett.

Le beau valet de cœur et la dame de pique  
causent sinistrement de leurs amours défunts.

Jag kunde ha lust att gå upp i det smutsiga gamla rucklet där snett över och dricka porter med flickorna. Röka en sur pipa och dra en spader med värdinnan och ge henne goda råd för hennes reumatism. Hon var här i förra veckan och klagade sin nöd, fet och präktig. Hon hade en tjock guldbrosch under isterhakan och betalade en kontant. Hon skulle bli smickrad av en

Det ringde på tamburdörren. Nu öppnar Kristin ... Vad kan det vara? Jag har ju sagt till att jag inte tar emot i dag ... En detektiv? ... Som låtsas vara sjuk, uppträder som patient ... Kom in du, min gubbe, jag skall nog sköta om dig ...

Kristin gläntade på dörren och slängde ett brev med svarta kanter på mitt bord. Inbjudning att bevista jordfästningen ...

\*

— Min handling, ja ... "Vil Monsieur have den Historie paa heroske Vers, saa koster det 8 Skilling ..."

25 augusti.

Jag såg i drömmen gestalter från min ungdom. Jag såg henne som jag kysste en midsommarnatt för länge sedan, då jag var ung och icke hade dödat någon. Jag såg också andra unga flickor av dem som hörde till vår krets den tiden; en som gick och läste det året jag blev student och som alltid ville tala med mig om religionen; en annan, som var äldre än jag och som gärna stod och viskade med mig i skymningen bakom en jasminhäck i vår trädgård. Och en annan, som alltid gjorde narr av mig, men som blev så ond och häftig och föll i krampgråt en gång då jag gjorde narr av henne ... De gingo till ka i en blek skymning, deras ögon stodof förskrämda, och de gjorde tecken åt

## Paragraphs

på svenska. Det är på det hela taget ett förbannat språk vi ha. Orden trampa varandra på tårna och knuffa varandra i ränstenen. Och allt blir så påtagligt och rått. Inga halvtöner, inga lättä antydningar och mjuka övergångar. Ett språk som tycks vara skapat till bruk för den outrotliga pöbelvanan att plumpa ut med sanningen i alla väder.

Det mörknar mer och mer: decembermörker i augusti. De svarta rosenbladen ha redan skrynklat sig. Men kortlapparna på mitt bord lysa i gälla skrattande färger i allt det grå som för att påminna om att de en gång blevo uppfunkna för att skingra en sjuk och galen furstes svårmad. Men jag fasar vid blotta tanken på arbetet att samla ihop dem och vända de aviga rätt och blanda dem till en ny patiens, jag kan bara sitta och se på dem och lyssna till hur "hjärterknekt och spaderdam viska dystert om sin begravda kärlek", som det står i samma sonett.

Le beau valet de cœur et la dame de pique  
causent sinistrement de leurs amours défunts

Jag kunde ha lust att gå upp i det smutsiga gamla rucklet där snett över och dricka porter med flickorna. Röka en sur pipa och dra en spader med värdinnan och ge henne goda råd för hennes reumatism. Hon var här i förra veckan och klagade sin nöd, fet och präktig. Hon hade en tjock guldbrosch under isterhakan och batelede en sommar kontant. Hon skulle bli

Det ringde på tamburdörren. Nu öppnar Kristin ... Vad kan det vara? Jag har ju sagt till att jag inte tar emot i dag ... En detektiv? ... Som låtsas vara sjuk, uppträder som patient ... Kom in du, min gubbe, jag skall nog sköta om dig ...

Kristin gläntade på dörren och slängde ett brev med svarta kanter på mitt bord. Inbjudning att bevisa jordfästningen ...

\*

— Min handling, ja ... "Vil Monsieur have den Historie paa heroske Vers, saa koster det 8 Skilling ..."

25 augusti.

Jag såg i drömmen gestalter från min ungdom. Jag såg henne som jag kysste en midsommarnatt för länge sedan, då jag var ung och icke hade dödat någon. Jag såg också andra unga flickor av dem som hörde till vår krets den tiden; en som gick och läste det året jag blev student och som alltid ville tala med mig om religionen; en annan, som var äldre än jag och som gärna stod och viskade med mig i skymningen bakom en jasminhäck i vår trädgård. Och en annan, som alltid gjorde narr av mig, men som blev så ond och häftig och föll i krampgråt en gång då jag gjorde narr av henne ... De sista åren i mitt liv, deras ögon stod och de gjorde tecken åt

Foreign-language material

## Fine-grained segmentation

på svenska. Det är på det hela taget ett förbannat språk vi ha. Orden trampa varandra på tårna och knuffa varandra i rännstenen. Och allt blir så påtagligt och rått. Inga halvtoner, inga lätta antydningar och mjuka övergångar. Ett språk som tycks vara skapat till bruk för den outrotliga pöbelvanan att plumpa ut med sanningen i alla väder.

## Fine-grained segmentation

på svenska. | Det är på det hela taget ett förbannat språk vi ha. | Orden trampa varandra på tårna och knuffa varandra i rännstenen. | Och allt blir så påtagligt och rått. | Inga halvtoner, inga lätta antydningar och mjuka övergångar. | Ett språk som tycks vara skapat till bruk för den outrotliga pöbelvanan att plumpa ut med sanningen i alla väder.

Sentence boundaries

## Fine-grained segmentation

på svenska. Det är på det hela taget ett förbannat språk vi ha. Orden trampa varandra på tårna och knuffa varandra i rännstenen. Och allt blir så påtagligt och rått. Inga halvtoner, inga lätta antydningar och mjuka övergångar. Ett språk som tycks vara skapat till bruk för den outrotliga pöbelvanan att plumpa ut med sanningen i alla väder.

Word/token boundaries

## Fine-grained segmentation

på svenska. Det är på det hela taget ett förbannat språk vi ha. Orden trampa varandra på tårna och knuffa varandra i rännstenen. Och allt blir så påtagligt och rått. Inga halvtoner, inga lätta antydningar och mjuka övergångar. Ett språk som tycks vara skapat till bruk för den outrotliga pöbelvanan att plumpa ut med sanningen i alla väder.

Word/token boundaries

But... does that make sense?

svenska.

påtag-

ligt

halvtoner,

övergångar.

# Tokenisation

- ▶ *Tokenisation or word segmentation* is the task of splitting a text into minimal processing units
  - ▶ for further linguistic analysis
  - ▶ as input to machine learning solutions
- ▶ The units should be *meaningful* for the model or code that processes them.
- ▶ Models learn relations between tokens, so the segmentation should produce units that have meaningful relations.
- ▶ Often, “words” are a good level of segmentation to aim at.

# What is a word?

“the smallest meaningful sequence of sounds that can be uttered in isolation”

- ▶ Definition from spoken language.
- ▶ In written language, we often assume that words are delimited by spaces or punctuation.
- ▶ The devil is in the detail:
  - ▶ *don't* –
  - ▶ *I'm* –
  - ▶ *dishwasher* –
  - ▶ *washing machine* –

## Reasons to prefer a particular segmentation

- ▶ Matching existing linguistic theory/annotation
- ▶ Matching existing models
- ▶ Trade-off between
  - ▶ Meaningfulness
  - ▶ Frequency
  - ▶ Number of parameters in models to train

So how to do  
word segmentation?

# Regular expressions

- ▶ “Language” for pattern matching in texts.
- ▶ Efficient and expressive.
- ▶ Available in many programming languages and tools.
- ▶ Part of the Python standard library:

```
import re
```

- ▶ Basic principle:  
A regex *pattern matches a string (or a part of it).*

## Basic patterns

A cat sat on the mat. Next to it sat a rat.

- ▶ A letter (or any symbol without special meaning) matches itself.

a t A at \_

- ▶ Square brackets: Match one of the characters.

[cs] [.\_]

- ▶ Negated square brackets: Match anything but those characters

[^a-z] (*beware: äöå etc.!*)

# Basic patterns

A cat sat on the mat. Next to it sat a rat.

- ▶ A letter (or any symbol without special meaning) matches itself.

a t A at \_

- ▶ Square brackets: Match one of the characters.

[cs] [.\_]

- ▶ Negated square brackets: Match anything but those characters

[^a-z] (*beware: äöå etc.!*)

## Basic patterns

A cat sat on the mat. Next to it sat a rat.

- ▶ A letter (or any symbol without special meaning) matches itself.

a t A at \_

- ▶ Square brackets: Match one of the characters.

[cs] [.\_]

- ▶ Negated square brackets: Match anything but those characters

[^a-z] (*beware: äöå etc.!*)

# Basic patterns

A cat sat on the mat. Next to it sat a rat.

- ▶ A letter (or any symbol without special meaning) matches itself.

a t A at \_

- ▶ Square brackets: Match one of the characters.

[cs] [.\_]

- ▶ Negated square brackets: Match anything but those characters

[^a-z] (*beware: äöå etc.!*)

# Basic patterns

A\_cat\_sat\_on\_the\_mat. Next\_to\_it\_sat\_a\_rat.

- ▶ A letter (or any symbol without special meaning) matches itself.

a t A at \_

- ▶ Square brackets: Match one of the characters.

[cs] [.]

- ▶ Negated square brackets: Match anything but those characters

[^a-z] (*beware: äöå etc.!*)

# Basic patterns

A **cat** **sat** on the mat. Next to it **sat** a rat.

- ▶ A letter (or any symbol without special meaning) matches itself.

a t A at ↴

- ▶ Square brackets: Match one of the characters.

[cs] [.]

- ▶ Negated square brackets: Match anything but those characters

[^a-z] (*beware: äöå etc.!*)

# Basic patterns

A\_cat\_sat\_on\_the\_mat. Next\_to\_it\_sat\_a\_rat.

- ▶ A letter (or any symbol without special meaning) matches itself.

a t A at \_

- ▶ Square brackets: Match one of the characters.

[cs] [.]

- ▶ Negated square brackets: Match anything but those characters

[^a-z] (*beware: äöå etc.!*)

# Basic patterns

A\_cat\_sat\_on\_the\_mat.\_Next\_to\_it\_sat\_a\_rat.

- ▶ A letter (or any symbol without special meaning) matches itself.

a t A at \_

- ▶ Square brackets: Match one of the characters.

[cs] [.]

- ▶ Negated square brackets: Match anything but those characters

[^a-z] (*beware: äöå etc.!*)

## Special elements

- ▶ ^, \$ – beginning and end of line
- ▶ \s – whitespace
- ▶ \S – not whitespace
- ▶ \w – word-internal character
- ▶ \W – not word-internal character
- ▶ and more: see library documentation

**Note:** \w use a specific definition of *word*,  
which may or may not be what you want!

## Combining elements

- ▶  $a^*$  – zero or more occurrences of the previous item  $a$  (Kleene star)
- ▶  $a^+$  – one or more occurrences of the previous item  $a$
- ▶  $a|b$  – alternatives
- ▶  $(abc)$  – grouping, so you can write  $(abc)^*$  or  $(ab|c)^+$

# Regex for string splitting

```
>>> import re  
>>> line = 'A cat sat on the mat. His name was Måns.'
```

- ▶ If you can easily describe the *delimiter*: `re.split`

```
>>> re.split(' ', line)  
['A', 'cat', 'sat', 'on', 'the', 'mat.', 'His', 'name', 'was', 'Måns.']}
```

- ▶ If you can easily describe the *tokens*: `re.findall`

```
>>> re.findall(r'\w+', line)  
['A', 'cat', 'sat', 'on', 'the', 'mat', 'His', 'name', 'was', 'Måns']
```

- ▶ Find a match *at the beginning of the string*: `re.match`
- ▶ Find a match *somewhere in the string*: `re.search`

# Separate patterns for tokens and delimiters

*code example*

## Tips for designing a tokeniser

- ▶ Start with a smallish subset of your data and tokenise it manually.
- ▶ Design regular expressions to match the desired tokenisation.
- ▶ Run it over a larger set, identify problems and refine.
- ▶ Looking through vocabulary lists can help you find problems (especially tokens that only occur once or twice).

## Intrinsic evaluation

Evaluating tokeniser “as tokeniser”.

- ▶ Requires a correct *gold standard segmentation*.
- ▶ *Word boundary recall*:

$$R = \frac{|C \cap G|}{|G|}$$

$R$ : word boundary recall

$C$ : word boundaries found by tokeniser

$G$ : word boundaries in gold standard

## Extrinsic evaluation

Measure how the performance of a downstream task depends on tokenisation.

- ▶ In the last part of the project, you will build a classifier.
- ▶ This allows you to do an extrinsic evaluation of the tokeniser.
- ▶ Wait until the end...

## Comparing tokenisations

- ▶ You can *compare* your tokenisation to that of another implementation and assess it qualitatively.
- ▶ Use the right tools to compare large outputs efficiently:
  - ▶ `difflib.SequenceMatcher` in the Python standard library
  - ▶ `diff` utility on Unix command line
- ▶ Your TAs can help you get going!

## Exercise

- ▶ Work on creating the tokenisation of your datasets.  
(Section 1 in the project description)
- ▶ Tokenisation is a prerequisite for everything that comes after!
- ▶ Your TAs are here to help you! Make good use of that!