

Text Analysis and Retrieval

1. Introduction

Prof. Jan Šnajder

University of Zagreb
Faculty of Electrical Engineering and Computing (FER)

Academic Year 2022/2023



Creative Commons Attribution–NonCommercial–NoDerivs 3.0

v3.0

After this lecture, you'll...

- ① Know what this course is about and be glad that you've enrolled it
- ② Know what are the topics that we want to cover
- ③ Know what you need to do to earn your credits

Outline

1 Why this course?

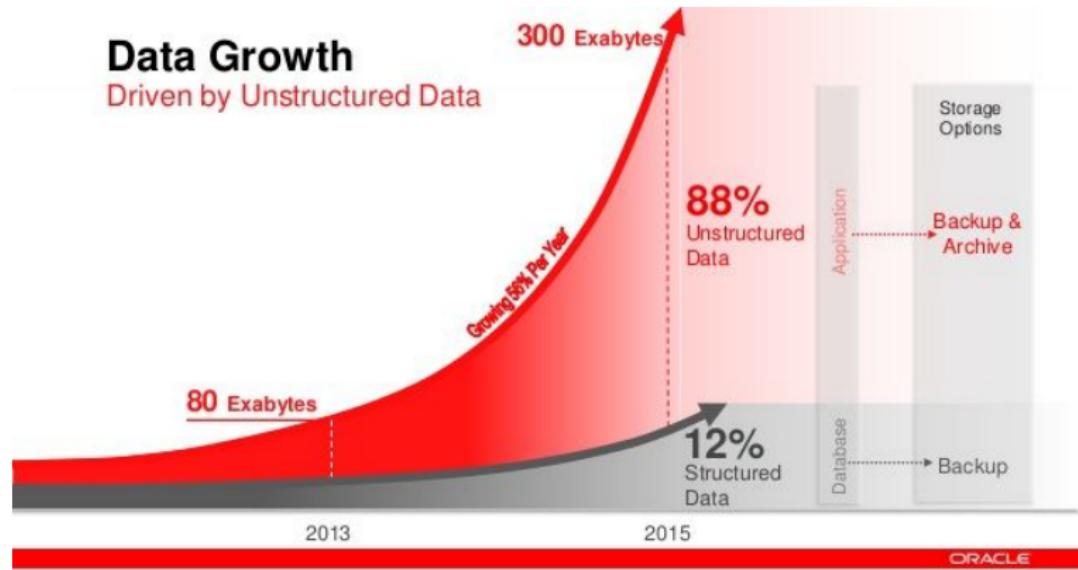
2 Topics

3 Organization

What is this course about?

- Text is everywhere (books, news, emails, tweets, reviews, . . .)
- Most **human knowledge** is stored in unstructured, textual form
- The amount of text data is vast and rapidly growing
- We need systems that address diverse **information needs** of the users and **extract information** from large volumes of unstructured data

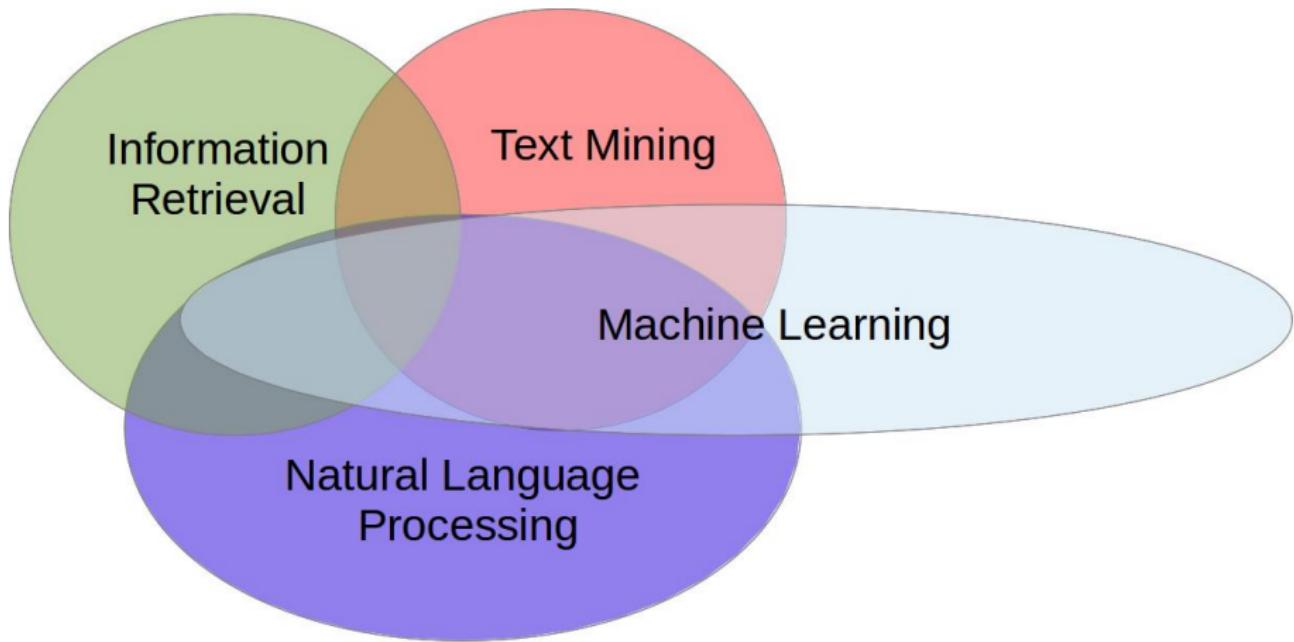
Unstructured data



Data Science



What's involved?

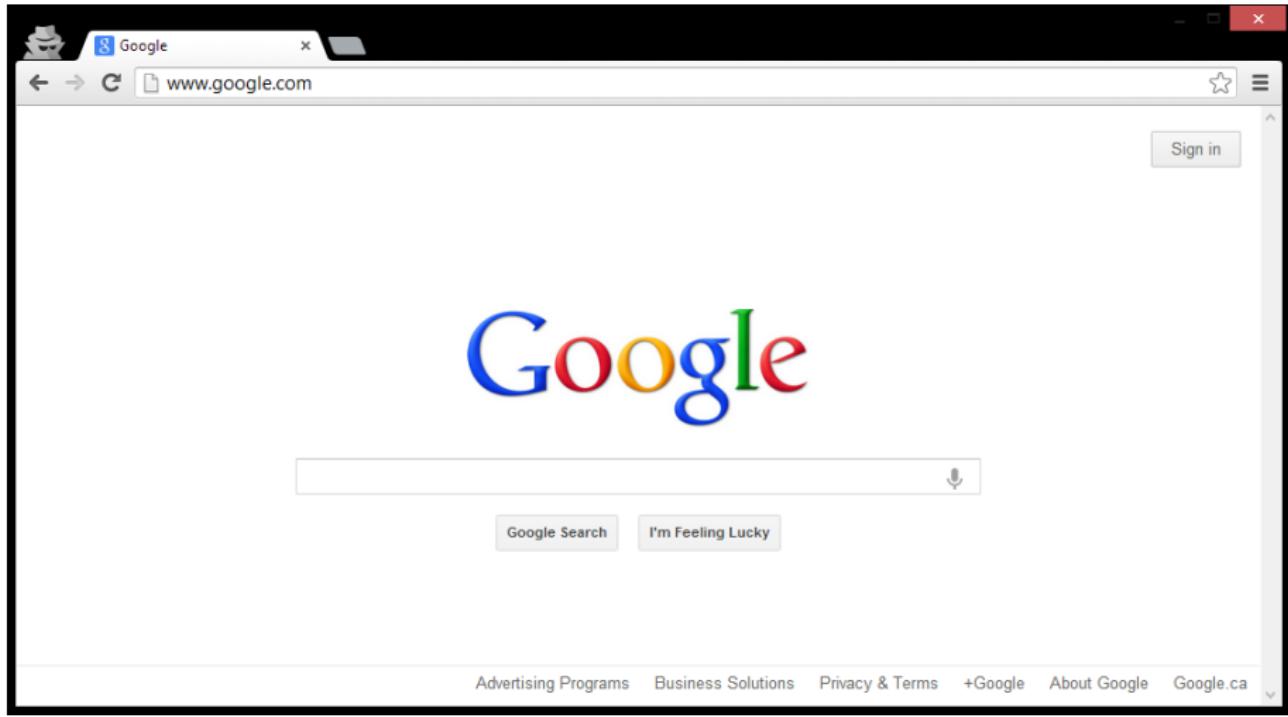


1. Information retrieval

Information Retrieval (Manning & Schütze 2008)

Information retrieval (IR) is finding material (usually documents) of an **unstructured** nature (usually text) that satisfies an **information need** from within large collections (usually stored on computers).

Web search



Domain-specific search

[Log in](#)[Help](#)[Post an ad \(free\)](#)[All Categories ▾](#) small parrot cage

in Earth ▾

[SEARCH](#)

[Bird Cage Suppliers Accepting Orders in Small Volume](#)

small-order.hktdc.com/ Reduce risk in sourcing! Place order in **small quantity** @ hktdc.com

[Used Bird and Parrot Cages, good condition](#)

PLS. EMAIL FOR APPT. TO VIEW ALL! Cleaning out my garage of used bird cages (\$5-40), iron parrot cages



[large parrot cage, 40"wide, 32"deep, 74" high, almost new](#)

green large parrot cage, just bought, too big for house



[parrot cage](#)

Parrot cage. Good condition Comes with two bowel feeders and wooden bench.

[Parrot Cage With toys and Dishes](#)

Parrot cage with toys and dishes. If interested please text me and I will send you a picture. Or just call me to



[Create your own ad](#) in Earth . It's easy and free!

[Parrot Cage](#)

Black Wrought iron Parrot cage, about 5 ft. tall, fits in corner. Good condition. 561-389-5347



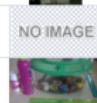
Follow, Like or add us to your Google + circles

[Large Parrot Cage For Sale](#)

LARGE WHITE 5 FOOT PARROT CAGE EXCELLENT CONDITION TOP PERCH WITH LOWER SHELF FOR STORAGE

[LARGE Parrot cage for sale](#)

Large Parrot Cage for sale with top perch and lower storage area comes with ladder for bird to climb back into



[wanted parrot cage 3 feet wide black willing to pay 150\\$](#)

Wanted Parrot cage willing to pay 150 Balck 3 feet wide



[Small Plastic cage with Hermit crab](#)

Purchased this summer by my oldest daughter for near \$35. (she is onto other activities.) Has large hermit

[PARROT CAGE, TOYS, BOOKS, ETC.](#)

ALL YOU NEED NOW IS THE PARROT: This Dometop cage provides your pet with a new home that is both

[Tweet](#) 0 [Follow](#)

[Like](#) 0

[+1](#) 0 [Follow](#) 1.2k

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 ... >

2. Text mining (aka text analytics)

Text mining (Hearst 2003)

Text mining (TM) is the discovery by computer of new, previously unknown information, by automatically extracting information from different written resources. A key element is the linking together of the extracted information to form new facts or new hypotheses to be explored further by more conventional means of experimentation.

Text mining

http://livingstories.googlelabs.com/https://afghanistan-PDFVRnRtBfalse, false, false, false, false;

The New York Times - Google | Feedback | Log out All stories | Subscribe to email updates | RSS feed

The War in Afghanistan

No updates since last visit

In a [major policy move](#), President Obama announced he would commit 30,000 more United States troops to Afghanistan but added he would begin drawing out American forces there starting in July 2011.

Mr. Obama's decision is likely to prove to a defining one for his Administration. It also poses several dilemmas. In staking out the reasons for escalating American commitment, he has also promised to end it. The White House policy also leaves an unanswered question: whether Afghan president Karzai will meet the challenge and how to respond if he does not. [Some lawmakers voiced skepticism](#).

 Video: The President's Afghanistan Address

[Read more...](#)

Nov 18, 2009 Nov 19, 2009 Nov 24, 2009 Nov 26, 2009 Nov 30, 2009 Dec 1, 2009

Clinton Visits Karzai on Eve of Inauguration Karzai Sworn In for Second Term as Afghan President Obama May Add 30,000 Troops in Afghanistan U.S. Seeks More Allied Troops for Afghanistan Obama's Speech on Afghanistan to Envision Exit Obama Adds Troops, but Maps Exit Plan

All coverage [Italy May Add 1,000 Extra Troops in Afghanistan](#) 1:07 PM **Timeline of important events**
The Global Response [Afghanistan Speech by Obama Wins Over Some Skeptics](#) 11:25 AM [Obama Adds Troops, but Maps Exit Plan](#)
By James Dao [Dec 1, 2009](#)
Casualties [President Obama intended his speech at West Point to rally Americans behind his plan to send 30,000 more troops to Afghanistan and to set an 18-month timetable for starting a withdrawal. And interviews suggest that, while opinions on the war remained wildly diverse, Mr...](#) 11:25 AM [Dec 1, 2009](#)
Opinion
The Afghan Elections
U.S. Policy [Related](#)
[Afghans See Sharp Shift in U.S. Tone - Feature](#)
The Troop Debate
Events
Articles
People

Graphic



Changing Opinions on Afghanistan

Nov 18, 2009 Nov 19, 2009 Nov 24, 2009 Nov 26, 2009 Nov 30, 2009 Dec 1, 2009

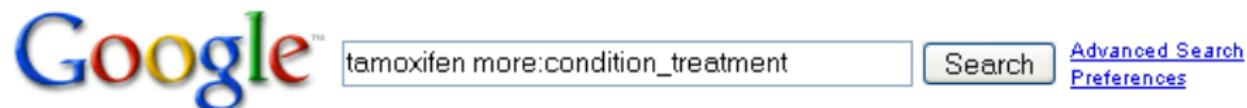
Decrease Increase Keep same

All news [Read more...](#) Jump to: [Multimedia](#)
Events
Articles
People

[Afghan Plan Faces Sharp Questioning From Senators](#) Dec 2, 2009 11:45 AM
By Brian Knowlton, Elizabeth Bumiller

Šnajder (UNIZG FER) TAR – Intro Academic Year 2022/2023 12 / 49

Information retrieval + Text mining



Refine results for **tamoxifen**:

Overview	Drug uses	Research overview	From medical authorities
Symptoms	Side effects	Practice guidelines	Alternative medicine
Tests/diagnosis	Interactions	Patient handouts	For health professionals
Treatment	Warnings/recalls	Continuing education	For patients
Causes/risk factors		Clinical trials	Support groups

Tamoxifen: Q & A - National Cancer Institute

The benefits of **tamoxifen** as a **treatment** for breast cancer are firmly ... How long should a patient take **tamoxifen** for the **treatment** of breast cancer? ...

www.cancer.gov/cancertopics/factsheet-Therapy/tamoxifen - 50k - [Cached](#) - [Similar pages](#)

Tamoxifen

3. Natural language processing

Natural language processing

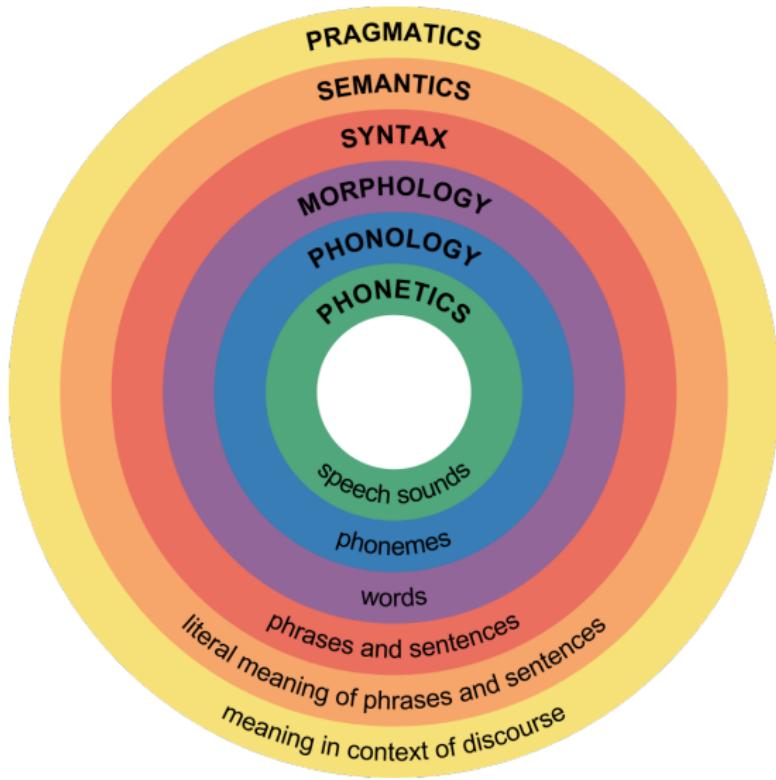
(Wikipedia)

Natural language processing (NLP) is a field of computer science, artificial intelligence, and linguistics concerned with the **interactions between computers and human (natural) languages**. As such, NLP is related to the area of human–computer interaction.

Why is TAR a challenge?

- Text is written in **natural language**. Natural language is a tough nut to crack. It is **complex, ambiguous, vague**, and relies on **commonsense knowledge**.
- Full understanding of natural language is an **AI-complete problem**.
- On top of this: dealing with large amounts of data poses serious **technical challenges**.

Language complexity



Language vagueness

Rezultati pretrage za *mali stan*

1 2 3 4 5 6 7-12 SLJEDEĆA »

6091 oglasi

Sortiraj Relevantnosti ▾

Vau Vau Njuškalo oglasi



Srdoči, odličan **mali stan** u novogradnji



Stan u stambenoj zgradi, Prizemlje

Stambena površina: **24.44 m²**

Objavljen: 03.03.2014.

46.500 € ~ 355.685 kn



KVATERNIKOV TRG - UREĐEN **mali stan**



Stan u stambenoj zgradi, 5 kat

Stambena površina: **40.00 m²**

Objavljen: 06.03.2014.

70.000 € ~ 535.441 kn



MALI STAN U FAŽANI !!



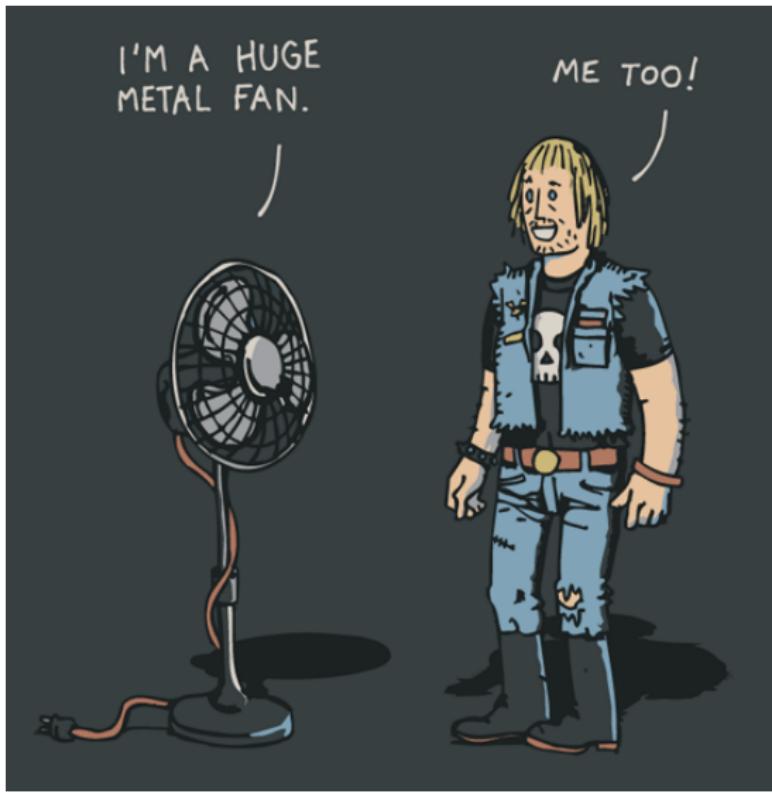
Stan u stambenoj zgradi, ---

Stambena površina: **35.00 m²**

Objavljen: 05.03.2014.

48.500 € ~ 370.984 kn

Language ambiguity



Word sense ambiguity

jaguar photos - Google Search

https://www.google.ie/search?q=jaguar+photos&oq=jaguar+photos&aqs=chrome..69i57j0l5.1948j0j7&sourceid=chrome&espv=210&es_sm=91&ie=UTF-8

Apps Postman Other Bookmarks

jaguar photos in Classic Cars

Web Images More Search tools

About 143,000,000 results (0.46 seconds)

Images for jaguar - Report images



YOUR COMMUNITY RESULTS
Powered by bryansk

Jaguar Photos - Jaguar Car Pictures - Motor Trend Magazine
jaguar, photos
www.motortrend.com/new_cars/07/jaguar/photos/index.html Communities: [Classic Cars](#), [My Searches](#)

Jaguar Photos & 2010 2009 Jaguar Pictures at Automobile Magazine
jaguar, photos
www.automobilemag.com/new_car_photos/01/jaguar/index.html Communities: [Classic Cars](#), [My Searches](#)

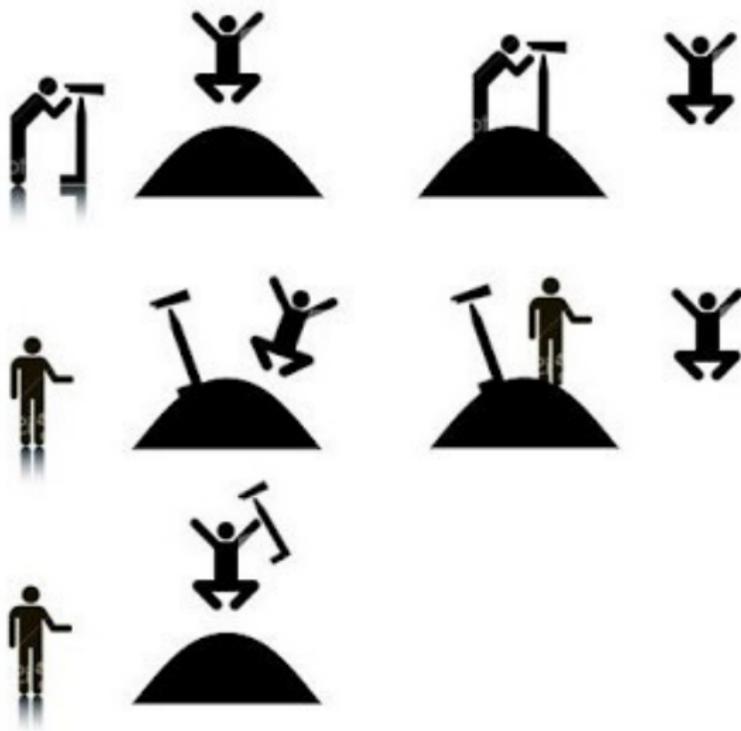
Jaguar photos
We have compiled some stunning Jaguar photos below for your pleasure. Feel free to send us your comments, feedback and if you own nice Jaguar photos, ...
www.jagplanet.com/photos Communities: [Classic Cars](#), [My Searches](#)

Jaguar photos - Panthera onca - ARKive
www.arkive.org/species/mammals/jaguar/
View all of ARKive's Jaguar photos - *Panthera onca*. ... Female jaguar with cub in birth den ... Female jaguar and ten week old cub crossing shallow creek.

Jaguars, Jaguar Pictures, Jaguar Facts - National Geographic
animals.nationalgeographic.com/animals/mammals/jaguar/

Structural ambiguity

"I saw the man on the hill with a telescope"



Outline

1 Why this course?

2 Topics

3 Organization

Course aims and target audience

Aims

Provide a **systematic overview** of both traditional and advanced methods for text analysis and retrieval. The course is divided into **two parts**. The first part covers the **fundamentals**: NLP, IR, and ML. The second part focuses on **concrete applications** in information extraction and text analytics, with an emphasis on methods based on statistical natural language processing and machine learning.

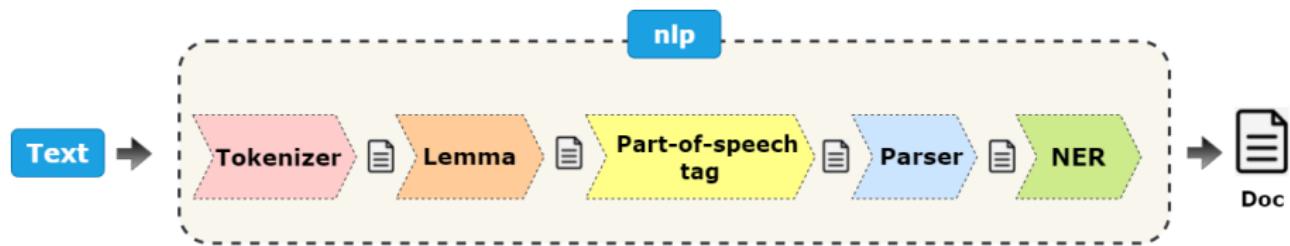
Audience

Students interested in gaining an understanding of and practical experience with basic information retrieval and text analysis techniques. Knowledge of at least basic machine learning is required. No previous knowledge of natural language processing is required.

Topics

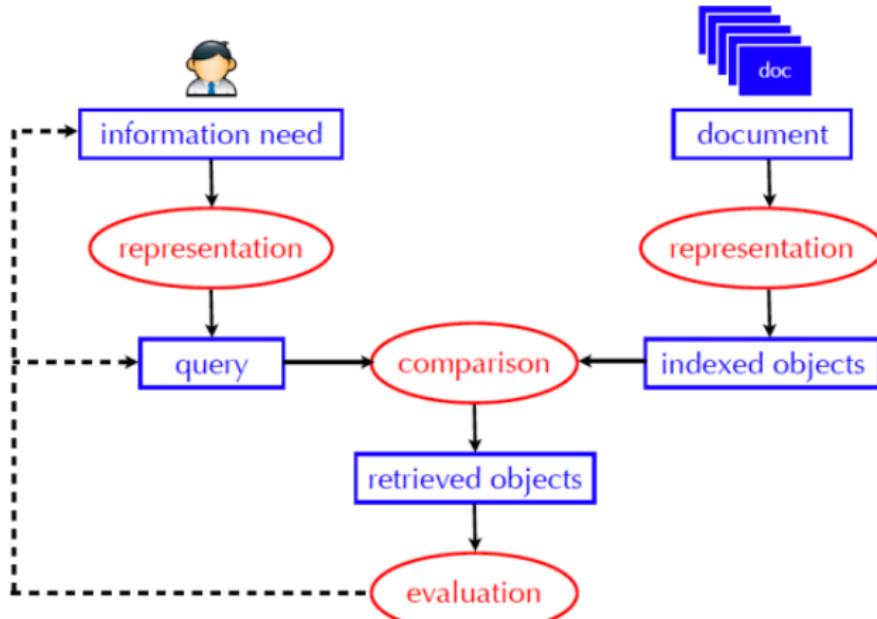
- ① Introduction
- ② Basics of NLP
- ③ Basics IR
- ④ Machine learning for NLP
- ⑤ Semantics
- ⑥ Neural NLP: Recurrent models
- ⑦ Neural NLP: Transformers
(two weeks break)
- ⑧ Information extraction
- ⑨ Question answering
- ⑩ Summarization
- ⑪ Natural language inference
- ⑫ Sentiment analysis
- ⑬ Author profiling; Wrap-up

Basic of NLP

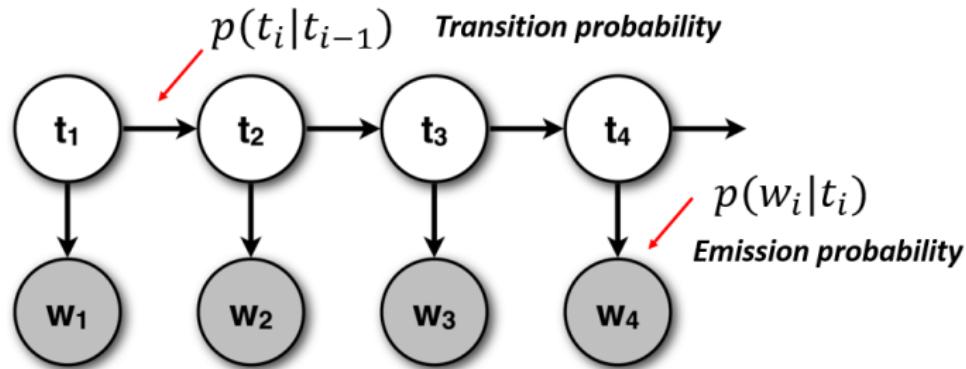


Information retrieval

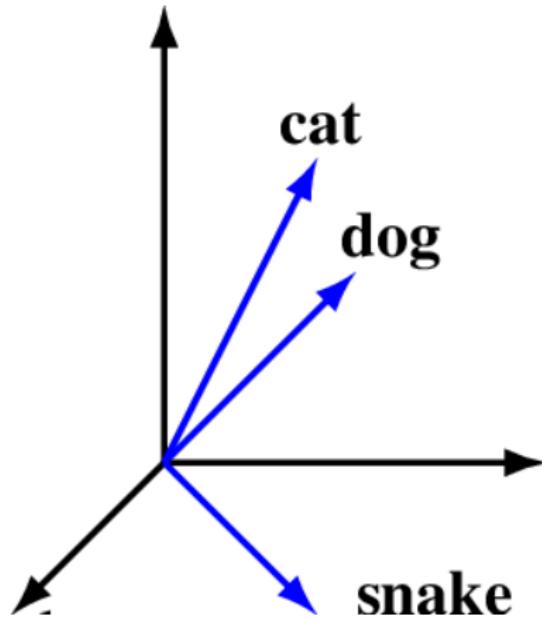
Basic Information Retrieval Process



Machine learning for NLP

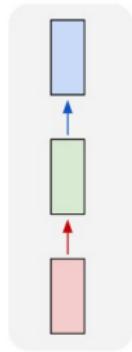


Semantics

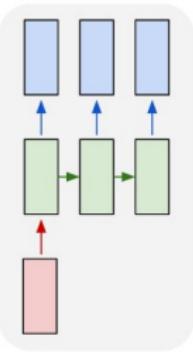


Neural NLP: Recurrent models

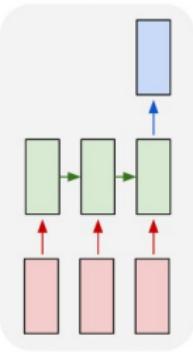
one to one



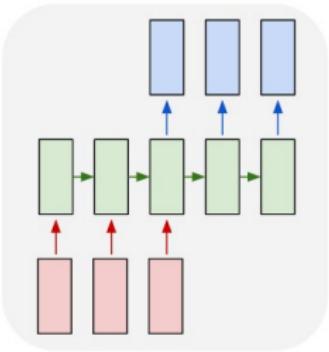
one to many



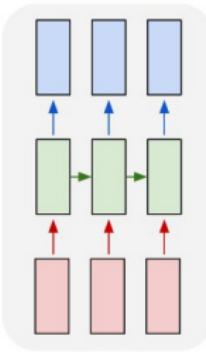
many to one



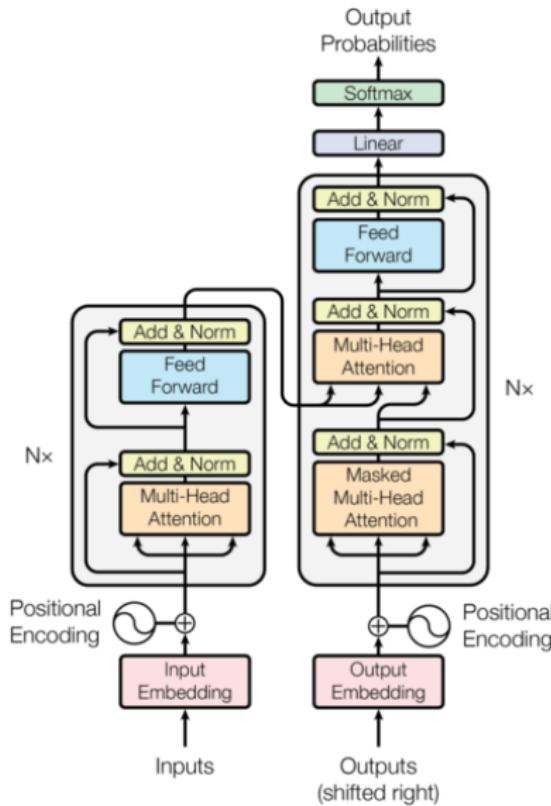
many to many



many to many



Neural NLP: Transformers



Outline

1 Why this course?

2 Topics

3 Organization

Teaching staff

Lecturers:

- Prof. Jan Šnajder

TAs:

- Josip Jukić

Student TAs:

- Tin Ferković, Ivan Martinović, Luka Pavlović, Marko Rajnović, Matea Vasilj, Janko Vidaković

Course web site:

<http://www.fer.hr/predmet/apt>

Powered by...



TakeLab

Text Analysis and Knowledge Engineering Lab

<http://takelab.fer.hr>



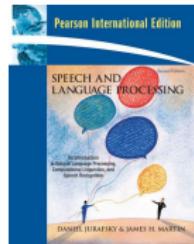
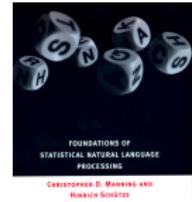
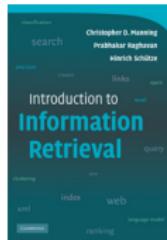
Learning outcomes

On successful completion of the course, you will be able to:

- ① Summarize the application areas, trends, and challenges in text analysis and retrieval
- ② Describe the fundamental techniques of text analysis and retrieval
- ③ Use linguistic preprocessing tools
- ④ Design and implement a text analysis/retrieval system
- ⑤ Apply machine learning algorithms to text analysis tasks
- ⑥ Evaluate a text analysis/retrieval system
- ⑦ Organize and formulate a system description paper
- ⑧ Describe, review, analyze, and criticize the main text analysis methods present in scientific papers

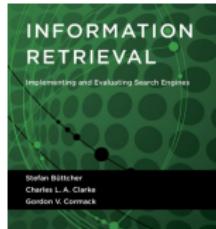
Textbooks

- C. D. Manning, P. Raghavan, H. Schütze.
Introduction to Information Retrieval, CUP, 2008
- C. D. Manning, H. Schütze: *Foundations of Statistical Natural Language Processing*,
MIT Press, 1999
- D. Jurafsky, J. H. Martin. *Speech and Language Processing*, 3rd ed., 2023.

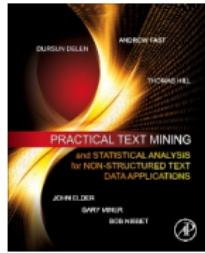


Textbooks

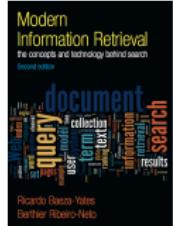
- S. Buettcher, C. L. A. Clarke, G. V. Cormack. *Information Retrieval: Implementing and Evaluating Search Engines*, MIT Press, 2010



- G. Miner, J. Elder IV, T. Hill, R. Nisbet, D. Denlen, A. Fast. *Practical Text Mining and Statistical Analysis for Non-structured Text Data Applications*, Academic Press, 2012

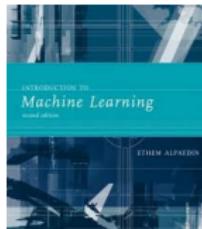


- R. Baeza-Yates, B. Ribeiro-Neto. *Modern Information Retrieval: The Concepts and Technology Behind Search*, ACM Press, 2011

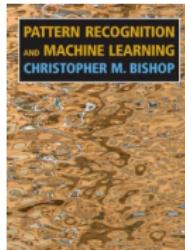


Textbooks – Machine learning

- Ethem Alpaydin: *Introduction to Machine Learning*, MIT Press, 2009



- Christopher Bishop: *Pattern Recognition and Machine Learning*, Springer, 2007



- Kevin P. Murphy: *Machine Learning: A Probabilistic Perspective*, MIT Press, 2012



Exams

- No exams

Theory: study assignments and in-class quizzes

- Together, in class:
 - 45 minute introduction to a topic
 - study assignment handed out
- You, at home:
 - do the study assignment
 - prepare for the in-class quiz and self-check the learning outcomes
- Together, in the next class:
 - 15 minute QA session to resolve unclarities
 - 15 minute in-class quiz (**6 points**)
 - 15 minute discussion session to dive deeper (**1 bonus point**)

Applications: paper reading and discussions

- Together, in class:
 - 15 minute introduction to a topic
 - pointer to the paper
- You, at home:
 - read the paper
 - submit the paper summary and questions (1 point)
- Together, in the next class:
 - 15 minute paper comprehension quiz (5 points)
 - 30 minutes going through the paper together
 - 15 minutes paper scoring
 - 15 discussion session (1 bonus point)

Lab assignments

- Hands-on experience of using the latest NLP tools in Python
- Complements the lectures and prepares you for the project
- You need to solve and submit your assignments – **no demonstrations**
- 3 assignments for **3 points each**

Lab schedule

- Assignment 1: Preprocessing (week 3)
- Assignment 2: Classification and sequence labeling (week 5)
- Assignment 3: Neural NLP (week 7)

Project assignments

- Groups of 3 students work together on a focused TAR topic:
 - ① Study related work
 - ② Prepare the data (collect, clean, annotate)
 - ③ Implement the system
 - ④ Evaluate the system
 - ⑤ Write up a project report
 - ⑥ Give a 10-minute presentation + demonstration
- Expected workload: 44 hours per person over 11 weeks
- Score based on the results, project report, and presentation
 - Group score, to be distributed among and by the group members
 - Maximum differential between members' points is 2× (e.g., if one member has 50, the others must have at least 25)
- Project assignments ⇒ soon available at the course website

Project report

- In a form of a scientific paper
- Min. 3 pages and max. 4 pages + 1 page for references
- Structure:
 - ① Intro: explain the problem, the motivation, and what you did
 - ② Related work: short task overview, similar algorithms/systems
 - ③ Description of your system
 - ④ Evaluation of your system
 - ⑤ Conclusion
 - ⑥ References
- Preferably in English (but we won't insist ☺)
- To be typeset in \LaTeX (we provide a template)
- Reviewed by teaching staff + TAs

Project report

- Published on-line in the now renowned “TAR Course Project Report” series (with your consent)
 - Take a look at [TAR 2021 PR](#), [TAR 2020 PR](#), [TAR 2019 PR](#),
[TAR 2018 PR](#), [TAR 2017 PR](#), [TAR 2016 PR](#), and [TAR 2014 PR](#)
- We highly encourage you to write your report in English and let us publish it on-line
 - ⇒ gives great EXPOSURE to your work
 - ⇒ future employers will FIGHT over you! ☺

Grade breakdown

	Continuous		Exam	
	Threshold	% grade	Threshold	% grade
Theory	50%	18%		
Labs	50%	9%		
Applications	50%	18%		
Project		55%		

If you fail to pass theory or reading threshold, you may retake the quiz (but you may retake at most 2 quizzes).

Grading



Excellent (5)	89
Very good (4)	76
Good (3)	63
Acceptable (2)	50

"I would have had higher grades, but I donated them to charity."

Project work schedule

- March 4: Topics announced
- March 4–13: Bidding for topics (in teams)
- March 14: Topics assigned
- Apr: Project progress checkpoint “alpha”
- May: Project progress checkpoint “beta”
- June: Project submission deadline (no extension!)
- June: Project presentations

ECTS breakdown

- This course is worth **5 ECTS credits = 150 hours of work**
- 44 hrs for the project
- 26 hrs for paper reading
- 26 hrs for the lectures
- 36 hrs for theory self-study
- 10 hrs for lab sessions
- 8 hrs for project presentations

You now...

- ① Know what this course is about and be glad that you've enrolled it
- ② Know what are the topics that we want to cover
- ③ Know what you need to do to earn your credits