# Natural Language Processing
## *Session 1*

## Nick Kadochnikov

University of Chicago – MS Applied Data Science

# About the Instructor (Nick!!!)

**Nick Kadochnikov**; kadochnikov@uchicago.edu
Associate Clinical Professor at UChicago PSD
20+ years of hands-on experience in data analytics to solve a variety of business problems
  - Social media and digital marketing analytics
  - Population health, medical prognosis, healthcare interoperability
  - Customer targeting and segmentation
  - Recommender systems, propensity to buy, cross-sell and up-sell modeling
  - Product development analytics
  - Fraud prevention
- Multiple data mining packages:
  - Python, Spark, SAS, SPSS, ILOG CPLEX, Netezza, Hadoop, Hive, Pig
- Worked with large volumes of structured and unstructured data (billions of records), including: transactional, financial, firmographic / demographic, organizational and macroeconomic data
- Worked in multiple countries across the globe
- Education:
  - MS in Global Marketing Management, Virginia Commonwealth University
  - MS and BS in Economics, St. Petersburg State University

Personal:
  - Love everything high speed: auto racing, go-karts, skiing and rollerblading
  - Passion for Renaissance and Baroque architecture

## WORK EXPERIENCE

Recent work data includes company industry & sector classification, document classification, and information extraction

## TA EXPERIENCE

NLP, Conversation AI Capstone, Reinforcement Learning

## EDUCATION

M.S. in Analytics | University of Chicago | 2022

CFA charter | CFA Institute | 2019

B.S. in Physics | University of Chicago | 2014

## IGNAS GRABAUSKAS

Data Scientist @ STBLaw

# About Me

Avid learner and a data science enthusiast with close to 5 years of experience working at the intersection of data and technology across various domains including Supply Chain, Investment Banking and Legal Consulting. Passionate about NLP and truly appreciate the importance of regular expressions.

# SWATHI GANESAN

AI Engineer

# Education

**2015 - 2019**
SASTRA UNIVERSITY, India
B.Tech in Computer Sci. & Engg.

**2022 - 2023**
University of Chicago
M.S. in Applied Data Science

# Experience

**2019 - 2022**
Decision Scientist - Mu Sigma, India

**Jun '23 - Sep '23**
Data Science Intern - William Blair, Chicago

**Currently**
AI Engineer - Harbor Global, Chicago

# Personal

Love exploring new cultures and experiences through food, travel and books.
I enjoy cooking and grocery shopping.
Trader Joe's is my happy place.
<Book recommendation: Becoming Trader Joe>

3

## PROFESSIONAL BACKGROUND

Management Consultant @PwC

Data Analyst @One Mount Group

AI Engineer @Harbor Global

## EDUCATION

BSc in Information System @Boston University

MSc in Applied Data Science @University of Chicago

## PERSONAL

My Spotify Wrapped 2023

**LINH LE**

linhcle@uchicago.edu | Chicago, IL

# Class rules

- I am not a "professor"
- My role is not to "teach", but to share expertise and facilitate learning
- Please turn your Zoom video-on
- Please join the class on-time

THE UNIVERSITY OF CHICAGO

# Syllabus highlights

- Use of Generative AI
  - In this course, students are allowed to use AI tools (such as ChatGPT) on all assignments.
- Academic Integrity
  - Sharing course content is strictly prohibited. Under no circumstances are students permitted to:
  - Download lecture videos;
  - Share sample code on open-source websites;
  - Plagiarize individual assignments;
  - Discuss course content, whether verbally, or electronically, with other cohorts of MS in Applied Data Science students.



The gift of feedback

# Class materials

- Book
  - Speech and Language Processing, 2nd Edition by Daniel Jurafsky and James Martin
    - https://www.pearson.com/us/higher-education/program/Jurafsky-Speech-and-Language-Processing-2nd-Edition/PGM181706.html
  - You can also use free online PDF chapters from the draft version of the 3rd edition:
    - https://web.stanford.edu/~jurafsky/slp3/
- Software:
  - Python with Jupyter Notebooks
  - Word processing and presentation software
- Assignments
  - Hands-on experience with software and techniques
  - Your TA will be able to assists you each week as needed

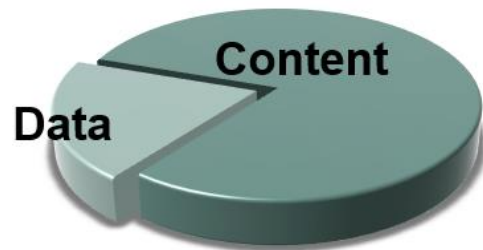THE UNIVERSITY OF CHICAGO

# Session 1 Agenda

- NLP applications
  - Text Mining, Survey mining, Social Media Analytics
- Basic text processing
- Regular expressions

# Introduction to Text Analytics

# What is Text Analytics?

- Over 80% of information being stored is unstructured
- Text analytics unlocks the power of that information for a variety of functions and applications



PC 143 (Hunter)
15 June 2006 23:47
Suspect identified himself as John Setsuko. Matched description given by night club doorman (IC1, Male, Ag 22-24 yrs, blue Everton shirt). Stopped whilst driving White Ford Mondeo, W563 WDL. Address given as 22 East Dene Ridge, Copdock, Ipswich. Searched at scene and found in possession of 1oz Cannabis Resin and lockable pocket knife.
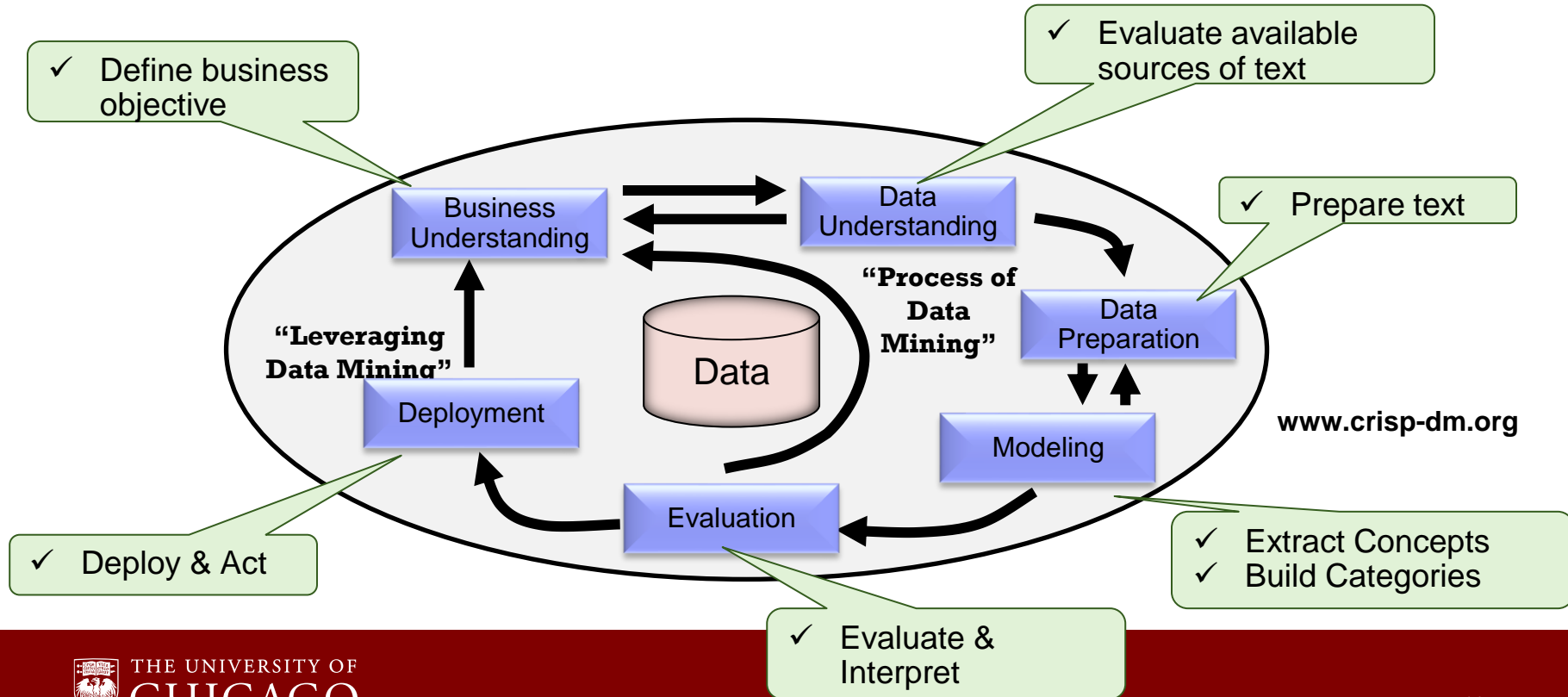
| Arresting_Officer | PC 143 |
| Arrest_Date_Time | 15/06/2006 : 23:47 |
| Suspect_Forename | John |
| Suspect_Surname | Setsuko |
| Suspect_VRN | W563WDL |
| Suspect_Vehicle_Color | White |
| Suspect_Vehicle_Make | Ford Mondeo |
| Suspect_Addr_Street | 22 East Dene Ridge |
| Suspect_Addr_Town | Ipswich |
| Evidence_1_Description | 1 oz Cannabis Resin |
| Classification | Drug possession |

THE UNIVERSITY OF CHICAGO

Identify email address in police report

Identify car make / model in police report

# CRISP-DM Methodology applies to text mining



Define business objective

Evaluate available sources of text

Prepare text

Business Understanding

Data Understanding

"Process of Data Mining"

Data Preparation

"Leveraging Data Mining"

Data

Deployment

Modeling

www.crisp-dm.org

Deploy & Act

Evaluation

Extract Concepts
Build Categories

Evaluate & Interpret

THE UNIVERSITY OF CHICAGO

# Data Mining vs. Text Mining

- In traditional data mining application you can either train the model on target variable (supervised modeling), or let the model find natural patterns in the data (i.e. unsupervised clustering)

- Same concepts apply to NLP problems!
  - You can build NLP classifiers
    - Target variable can be topic, sentiment, etc.
  - You can build NLP clusters
    - i.e. topic modeling, document / sentence clustering, etc.

# Text Mining and Data Preparation

## Data Mining

- Data cleaning
  - Selecting relevant data
  - Data quality – errors
  - Interpreting and handling missing data
- Data transformation
  - Get the data into the right form to ask the relevant questions…
- An iterative process

## Text Mining

- Data cleaning
  - Synonyms, abbreviations, specialized vocabulary, common typos
- Data transformation
  - Simplify text:
    - Stemming, lemmatizations, etc.
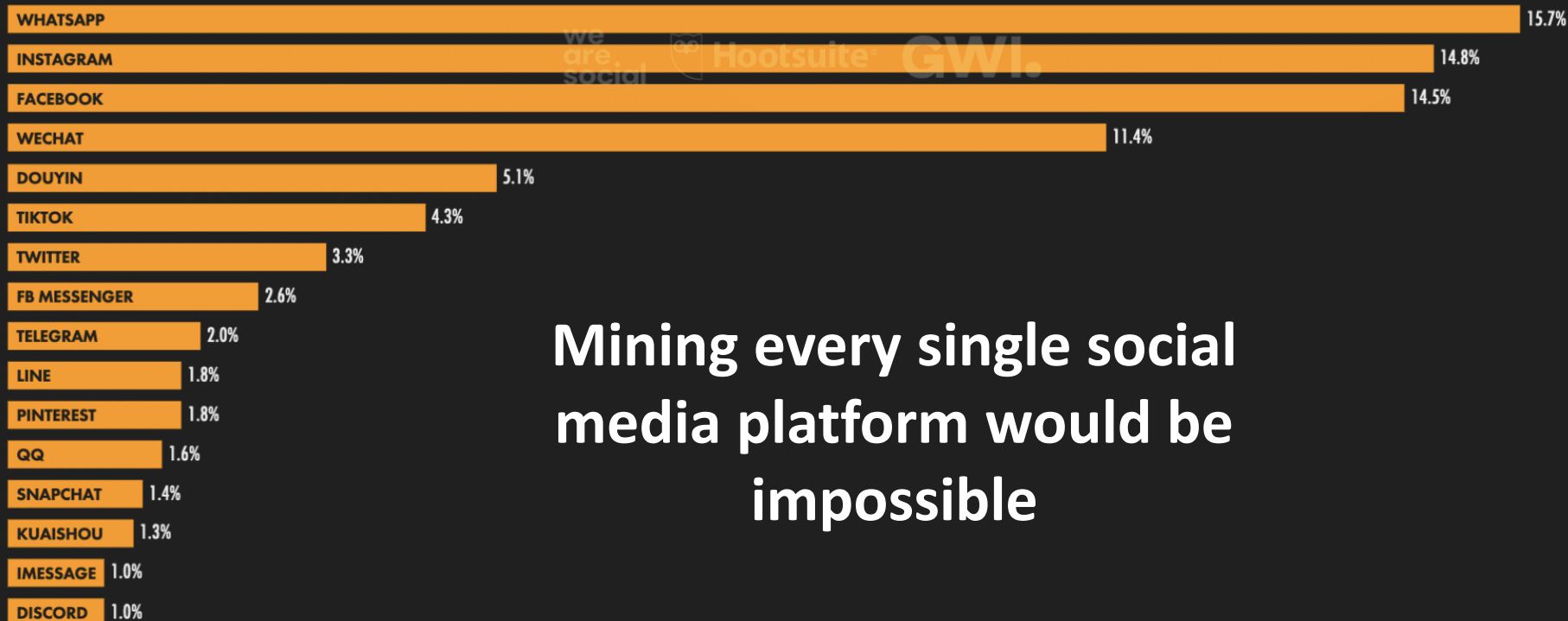
THE UNIVERSITY OF CHICAGO

# Social Media Analytics

What is the difference between Social Media Analytics
and
traditional text mining / natural language processing (NLP)?

What about survey mining

# The solution comes in form of data aggregation services

| | Twitter | Reviews/ message boards | Blogs | Forums | Online news articles | Video com-ments |
|---|---|---|---|---|---|---|
| **Data Aggregators** | | | | | | |
| **Twitter (Gnip)** | • | • | • | • | • | • |
| **Topsy** | • | | | | | |
| **DataSift** | • | • | | • | • | • |
| **Boardreader** | | • | • | • | • | • |
| **Trendiction** | • | • | • | • | • | • |

THE UNIVERSITY OF CHICAGO

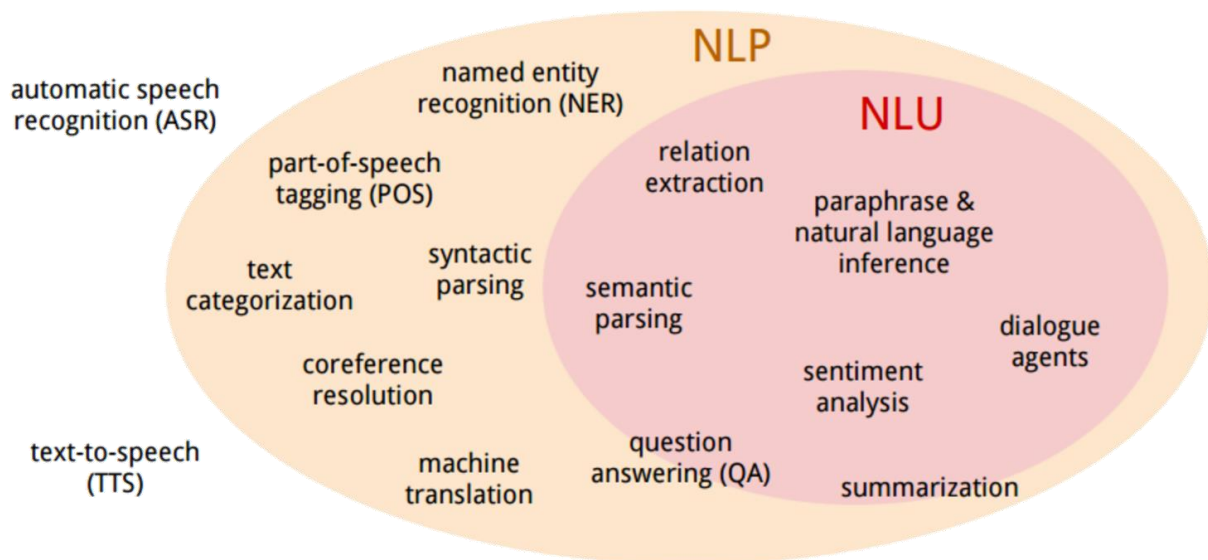# Supply Chain Risks: Distributions of mentions by topics
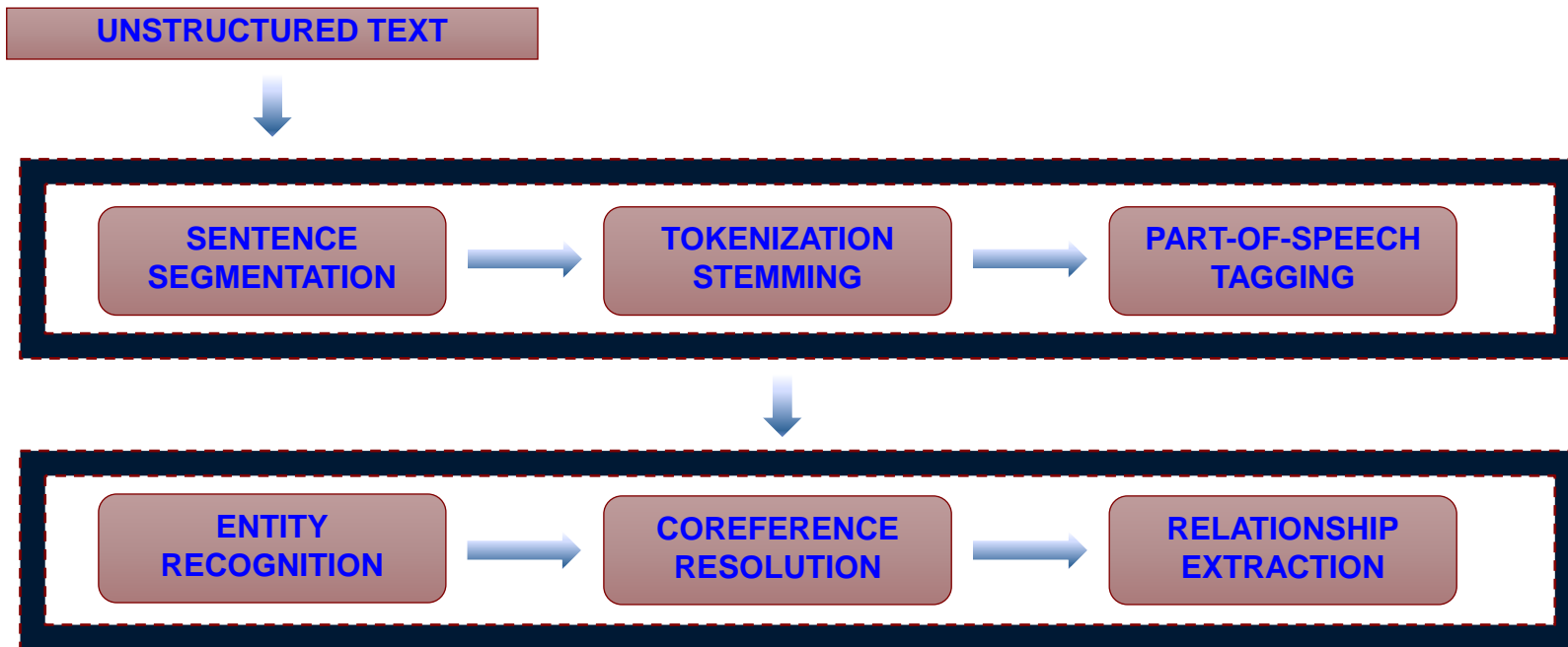
# NLP Tasks and Process

# NLP: Big Picture

Natural Language Processing (NLP) is a field of computer science and artificial intelligence, concerned with making computers process natural (human) language

Computational Linguistics (CL) is the field of using computers to understand language

# Basic Text Processing & Information Extraction (IE)



UNSTRUCTURED TEXT

SENTENCE SEGMENTATION → TOKENIZATION STEMMING → PART-OF-SPEECH TAGGING

ENTITY RECOGNITION → COREFERENCE RESOLUTION → RELATIONSHIP EXTRACTION

THE UNIVERSITY OF CHICAGO

24

Identify top three most frequent beverages
consumed in the book

# Basic Text Processing

## Regular Expressions

Dan Jurafsky

# **Regular expressions**

- A formal language for specifying text strings

- How can we search for any of these?

  - woodchuck

  - woodchucks

  - Woodchuck

  - Woodchucks

# Regular Expressions: Disjunctions

- Letters inside square brackets []

| Pattern | Matches |
|---------|---------|
| `[wW]oodchuck` | Woodchuck, woodchuck |
| `[1234567890]` | Any digit |

- Ranges `[A-Z]`

| Pattern | Matches | |
|---------|---------|---|
| `[A-Z]` | An upper case letter | Drenched Blossoms |
| `[a-z]` | A lower case letter | my beans were impatient |
| `[0-9]` | A single digit | Chapter 1: Down the Rabbit Hole |

# Regular Expressions: Negation in Disjunction

- Negations `[^Ss]`
  - Caret means negation only when first in []

| Pattern | Matches | |
|---------|---------|---|
| `[^A-Z]` | Not an upper case letter | `Oyfn pripetchik` |
| `[^Ss]` | Neither 'S' nor 's' | `I have no exquisite reason"` |
| `[^e^]` | Neither e nor ^ | `Look here` |
| `a^b` | The pattern a caret b | `Look up a^b now` |

# Regular Expressions: More Disjunction

- Woodchucks is another name for groundhog!
- The pipe | for disjunction

| Pattern | Matches |
|---|---|
| groundhog\|woodchuck | |
| yours\|mine | yours<br>mine |
| a\|b\|c | = [abc] |
| [gG]roundhog\|[Ww]oodchuck | |

Photo D. Fletcher

# Regular Expressions: ?   *   +   .

| Pattern | Matches | | | | |
|---|---|---|---|---|---|
| colou?r | Optional previous char | color | colour | | |
| oo*h! | 0 or more of previous char | oh! ooh! | oooh! | ooooh! | |
| o+h! | 1 or more of previous char | oh! ooh! | oooh! | ooooh! | |
| baa+ | | baa | baaa | baaaa | baaaaa |
| beg.n | | begin | begun | begun | beg3n |

Stephen C Kleene

Kleene *,   Kleene +

Dan Jurafsky

# Regular Expressions: Anchors ^ $

| Pattern | Matches |
|---------|---------|
| ^[A-Z] | <u>P</u>alo Alto |
| ^[^A-Za-z] | <u>1</u>    "<u>Hello</u>" |
| \.$ | The end<u>.</u> |
| .$ | The end<u>?</u>  The end<u>!</u> |

# **Example**

- Find me all instances of the word "the" in a text.

    `the`                  **Misses capitalized examples**

    `[tT]he`            **Incorrectly returns `other` or `theology`**

    `[^a-zA-Z][tT]he[^a-zA-Z]`

# Regular Expressions in Python (re)

# **Errors**

- The process we just went through was based on fixing two kinds of errors
  - Matching strings that we should not have matched (there, then, other)
    - False positives (Type I)
  - Not matching things that we should have matched (The)
    - False negatives (Type II)

# **Errors cont.**

- In NLP we are always dealing with these kinds of errors.

- Reducing the error rate for an application often involves two antagonistic efforts:

  - Increasing accuracy or precision (minimizing false positives)
  - Increasing coverage or recall (minimizing false negatives).

# **Summary**

- Regular expressions play a surprisingly large role
  - Sophisticated sequences of regular expressions are often the first model for any text processing text
- For many hard tasks, we use machine learning classifiers
  - But regular expressions are used as features in the classifiers
  - Can be very useful in capturing generalizations

37

# Thank You!