

COURSE NAVIGATION

- NI-DDW - Web Data Mining
- Classification
- Home Work
  - 1. Data Acquisition - Web Crawler/Scraper
  - 2. Text Mining
  - 3. Social Network Analysis
  - 4. Web Analytics/Web Usage Mining
  - 5. Indexing + Document Retrieval
  - 6. Recommender systems
- Lectures
- Seminar projects (optional)
- Tutorials

2. Text Mining

Task

- Find any suitable textual data for processing which will contain at least 500 sentences.
  - you can manually collect texts from BBC/CNN/New York Times, or
  - use the crawler from the first homework/tutorial and extend it to crawl particular website and collect content for this task, or
  - use any other suitable texts (e.g. OpenData speech datasets)
- Perform following NLP tasks
  - POS tagging
  - NER with entity classification (using nltk.ne\_chunk)
  - NER with custom patterns
    - e.g. every match of: adjective (optional) and proper noun (singular/plural) is matched as the entity
    - see slides 31 or 38 from lecture 4 for some NLTK examples using RegexpParser or custom NER
  - NER + classification using existing language model (e.g. NER models from Hugging Face)
- Implement your custom entity classification
  - For each detected entity (using both nltk.ne\_chunk and custom patterns)
    - Try to find a page in the Wikipedia
    - Extract the first sentence from the summary
    - Detect category from the sentence as a noun phrase
      - Example:
        - for „Wikipedia“ entity the first sentence is „Wikipedia (/wɪkɨˈpiːdiə/ or /ˌwɪkiˈpiːdiə/ WIK-i-PEE-dee-ə) is a free online encyclopedia that aims to allow anyone to edit articles.“
        - you can detect pattern „... is/VBZ a/DT free/JJ online/NN encyclopedia/NN ...“
        - the output can be „Wikipedia“: „free online encyklopedia“
  - For unknown entities assign default category e.g. „Thing“

Wikipedia package in Python:

```
import wikipedia
results = wikipedia.search("Wikipedia")
print(results)
page = wikipedia.page("Wikipedia")
print(page.summary)
```

Instructions for submitting

In your repository provide the following information:

- Description of the data you used for processing
- For each processing step (POS, NER based on ne\_chunk, custom NER, entity classification, language model) list the main results (e.g. top entities)
- Compare results of entity classification approaches
  - nltk-based classification
  - wikipedia-based classification using nltk entities as the input
  - wikipedia-based classification using custom patterns as the input
  - language model classifications
- Provide your implementation
- Comment on
  - issues during the design/implementation
  - ideas for extensions/improvements/future work