

# Comments and answers - Assignment 3

Project Text Analysis

Leon Wetzel, Teun Buijse and Roman Terpstra

---

## Exercise 1

### Part 3

Path Similarity

```
('car-automobile', 1.0)
('coast-shore', 0.5)
('monk-slave', 0.2)
('moon-string', 0.1111111111111111)
('food-fruit', 0.09090909090909091)
('journey-car', 0.05)
```

Leacock-Chodorow Similarity

```
('car-automobile', 3.6375861597263857)
('coast-shore', 2.9444389791664407)
('monk-slave', 2.0281482472922856)
('moon-string', 1.4403615823901665)
('food-fruit', 1.2396908869280152)
('journey-car', 0.6418538861723948)
```

Process finished with exit code 0

As you can see the ranking is quite different, mostly when it comes to medium-level similarity pairs. The similarity measures in the NLTK corpus reader do not consider these to be very similar, while the subjects did, in our eyes, successfully recognize that these are more similar than monk-slave for example. This could be explained by the fact that car-automobile, for example, can be used interchangeably. While food-fruit have very different meaning but belong to the same category. ## Exercise 2

### Part 1

To start the server, we use the following bash command:

```
shell script $ java -mx4g -cp "*" edu.stanford.nlp.pipeline.StanfordCoreNLPServer -port 9000 -
timeout 15000 -serverProperties server.properties
```

We count 2 ORGANIZATION's, 1 LOCATION and 20 PERSON's. Not all the shown named entities are correct. For example, *Augusta Ada King* and *Ada* are named as ORGANIZATION, which is not correct in this particular context.

### Part 2

There appear to be 6 other models for named entities. We take a closer look at the 4 and 7 classes type of models.

- english.conll.4class.distsim.prop
- regexner.patterns
- english.muc.7class.distsim.crf.ser.gz

- english.conll.4class.distsim.crf.ser.gz
- english.all.3class.distsim.prop
- english.muc.7class.distsim.prop

Let's start off with the 4 classes. We alter `server.properties` by replacing the old value of `ner.model` with `edu/stanford/nlp/models/ner/english.conll.4class.distsim.crf.ser.gz`. We then feed the text file to the application, which returns more ORGANIZATION labels, although one can conclude that more often than not a ORGANIZATION label is not properly used in this context. In addition, the MISC label is introduced, being a label often present at indications of nationality.

Let's use the 7 classes model. The `ner.model` in `server.properties` will be changed to `edu/stanford/nlp/models/ner/english.muc.7class.distsim.crf.ser.gz`. We can see that next to our familiar named entities, we also see a new named entity DATE. The ORGANIZATION entity is also more present than the original model.

### ### Part 3

For this part of the exercise, we use the `edu/stanford/nlp/models/ner/english.conll.4class.distsim.crf.ser.gz` model.

```
C:\Users\leonw\Documents\PTA\venv\Scripts\python.exe
C:/Users/leonw/Documents/PTA/week3/exercise2.py
countess 0
independence 0
Lord 0
Countess ORGANIZATION
Ada ORGANIZATION
Byron ORGANIZATION
question 0
attempt 0
disease 0
England LOCATION
augusta 0
Ada PERSON
Byron PERSON
bent 0
byron 0
club 0
instrument 0
Lord LOCATION
others 0
December 0
earth 0
baron 0
king 0
working 0
Isabella PERSON
programmer 0
lord 0
godhead 0
Engine ORGANIZATION
engineer 0
algorithm 0
history 0
mother 0
War ORGANIZATION
friendship 0
overlord 0
Greek ORGANIZATION
employment 0
relationship 0
end 0
mathematics 0
```

approach 0  
Luigi PERSON  
concern 0  
charles 0  
year 0  
class 0  
November 0  
oeuvre 0  
greek 0  
poet 0  
adenosine 0  
forefather 0  
endowment 0  
engine 0  
science 0  
King ORGANIZATION  
skill 0  
work 0  
Babbage PERSON  
england 0  
world 0  
anne 0  
pastime 0  
talent 0  
interest 0  
relate 0  
Independence ORGANIZATION  
cock 0  
article 0  
sight 0  
access 0  
person 0  
feat 0  
logic 0  
study 0  
request 0  
set 0  
note 0  
church 0  
machine 0  
society 0  
marriage 0  
kinship 0  
imagination 0  
individual 0  
insanity 0  
november 0  
overture 0  
charlemagne 0  
father 0  
Anne PERSON  
Between ORGANIZATION  
engineering 0  
analyst 0  
\_ 0  
wedlock 0  
computer 0  
child 0  
Lovelace ORGANIZATION  
Biography 0  
month 0  
war 0  
doubt 0  
vision 0  
deaminase 0

wife 0  
stage 0  
creature 0  
Lovelace PERSON  
programâ 0  
lovelace 0  
company 0  
eminence 0  
locomotive 0  
car 0  
december 0  
technology 0  
calculator 0  
mathematician 0  
workplace 0  
Analyst 0  
motion 0  
writer 0  
sake 0  
mentality 0  
mind-set 0  
calendar 0  
bill 0  
chiefly 0  
tool 0  
Metaphysician MISC  
death 0  
Charles PERSON  
capability 0  
Menabrea PERSON  
universe 0  
path 0  
isabella 0  
don 0  
Ada LOCATION  
Byron LOCATION  
Analytical ORGANIZATION  
adult 0  
founder 0  
Notes 0  
€ 0  
populace 0  
' ' 0  
hardening 0  
== 0  
campaign 0  
worldly 0  
effort 0  
Augusta ORGANIZATION

Process finished with exit code 0