

Text Mining

- Part of Speech Tagging (POS tagging)
- Named Entity Recognition (NER)
- Relation Extraction



Zhaopeng Xing
Doctoral Student at Carolina
Health Informatics Program in
Biomedical Health Informatics



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



The Big Picture – Information extraction

Noun,

National

From Wikipedia, the

"NIH" redirect

The National Institutes of Health (NIH) is the primary agency for biomedical research within the U.S. Department of Health and Human Services. It conducts its own research and provides research funding to other institutions.

As of 2013, the NIH has over 17,000 employees, including 11,000 scientists, and as of 2003, the NIH budget was over \$26 billion.^[6]

The NIH comprises 27 scientific departments, including the National Cancer Institute, the National Institute of Mental Health, and the National Institute of Diabetes and Digestive and Kidney Diseases.

In 2019, the NIH contributed to the

Verb, A noun

Agency overview	
Formed	1887
Preceding agency	Hygienic Laboratory
Headquarters	Bethesda, Maryland, U.S.
Employees	20,262 (2012), ^[1] including 6,000 research scientists (2019). ^[2]
Annual budget	▲ US\$39 billion (2019) ^[2] ▲ US\$37 billion ^[3] (2018) ^[4]
Agency executives	Francis Collins, Director Lawrence Tabak, Principal Deputy Director
Parent agency	Department of Health & Human Services
Child agencies	National Cancer Institute National Institute of Allergy and Infectious Diseases National Institute of Diabetes and Digestive and Kidney Diseases National Heart, Lung, and Blood Institute National Library of Medicine
Website	NIH.gov

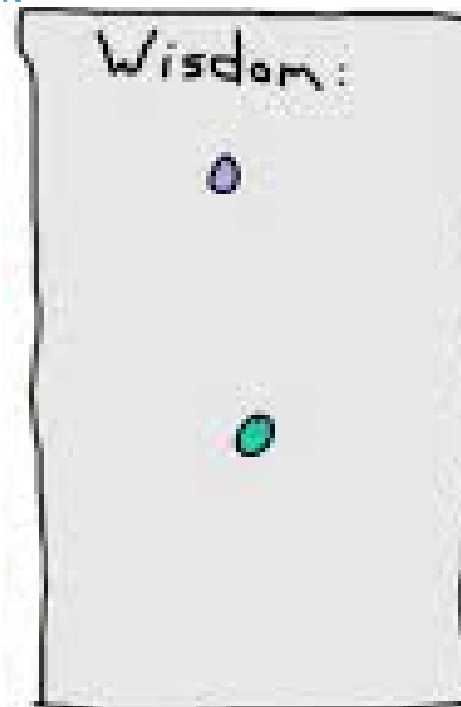
Number Noun, Organization.
A part of

United States government responsible for the health and well-being of the United States. It is located in Bethesda, Maryland. The NIH conducts research and provides major biomedical research funding.

It has over 17,000 employees and more than 4,000 postdoctoral research fellows. It is the largest biomedical research institution in the world,^[5] while, in 2013, it had a budget of approximately \$37 billion annually in the U.S., or about US\$26.4 billion.

It is responsible for many of the major medical advances of the 20th century, including the use of lithium to manage bipolar disorder, the discovery of HIV, and human papillomavirus (HPV).^[7]

It also publishes the Index Medicus, which measured the largest number of citations in the medical literature.





Outline

Goal

- What is Part of Speech Tagging (POS tagging) and how it is used in text mining?
- What is Named Entity Recognition (NER) and how it is used in text mining?
- What is the dependency structure and how to extract predicate-argument relation?

Part of Speech Tagging (POS tagging)

Named Entity Recognition (NER)

Relation Extraction

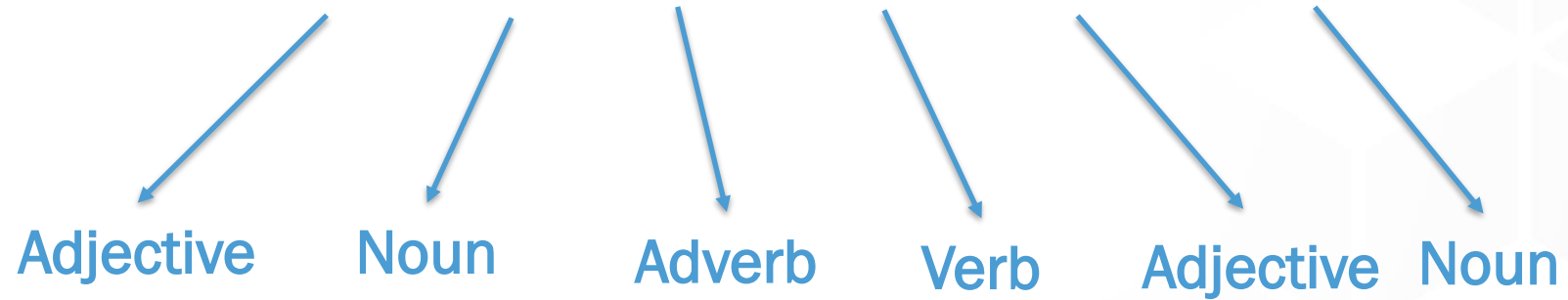


THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



POS tagging

“My dog also likes spicy sausages.”





Non-trivial function: Disambiguation

Like

*Please let me know how you would **like** to proceed.* Verb

***Like** most people, I'd prefer to have enough money not to work.*

Preposition

Having the same characteristics as...

*She looks **like** she is about to cry.*
...as if...

Conjunction



Penn treebank and Universal tag set

Penn treebank

Number	Tag	Description	19.	PRPS	Possessive pronoun
1.	CC	Coordinating conjunction	20.	RB	Adverb
2.	CD	Cardinal number	21.	RBR	Adverb, comparative
3.	DT	Determiner	22.	RBS	Adverb, superlative
4.	EX	Existential <i>there</i>	23.	RP	Particle
5.	FW	Foreign word	24.	SYM	Symbol
6.	IN	Preposition or subordinating conjunction	25.	TO	<i>to</i>
7.	JJ	Adjective	26.	UH	Interjection
8.	JJR	Adjective, comparative	27.	VB	Verb, base form
9.	JJS	Adjective, superlative	28.	VBD	Verb, past tense
10.	LS	List item marker	29.	VBG	Verb, gerund or present participle
11.	MD	Modal	30.	VBN	Verb, past participle
12.	NN	Noun, singular or mass	31.	VBP	Verb, non-3rd person singular present
13.	NNS	Noun, plural	32.	VBZ	Verb, 3rd person singular present
14.	NNP	Proper noun, singular	33.	WDT	Wh-determiner
15.	NNPS	Proper noun, plural	34.	WP	Wh-pronoun
16.	PDT	Predeterminer	35.	WP\$	Possessive wh-pronoun
17.	POS	Possessive ending	36.	WRB	Wh-adverb
18.	PRP	Personal pronoun			

Universal tags

Open class words	Closed class words	Other
ADJ	ADP	PUNCT
ADV	AUX	SYM
INTJ	CCONJ	X
NOUN	DET	
PROPN	NUM	
VERB	PART	
	PRON	
	SCONJ	



POS tags – Where we can use it?

Syntactic chunking

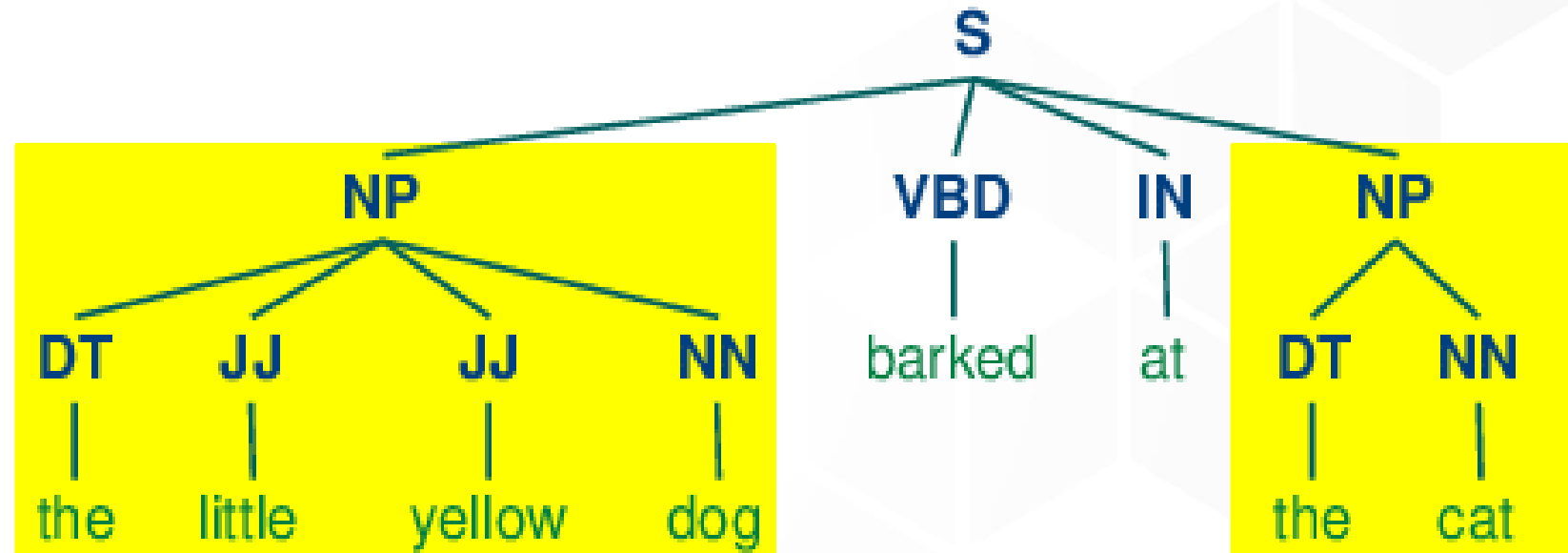
Noun phrase

Verb phrase

Preposition phrase

Adjective phrase

Adverb phrase





POS tags – Where we can use it?

Sentiment analysis

- Adjective and adverb as sentiment indicators

Positive: Good selection of food , fairly reasonable prices. Typical American Chinese food. Waitress was attentive

Negative: Clearly one of the worst restaurants in the country, not because of their food or service but rather their misleading/fraudulent business practices



More to think about

What are some other use cases of POS tagging in information extraction? (Maybe google “POS tagging”)

Task 1: POS tags and noun phrase extraction

Notice: Your work in Jupyter notebook will not be saved. Please download it.

Part of Speech Tagging (POS tagging)

Named Entity Recognition (NER)

Relation Extraction



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



Named Entity Recognition

Find and classify information units from text

[Who] did **[What]** at **[When]** in **[Where]**



NER – Example

“The **NIH** was founded in **1887** and is now part of the **United States Department of Health and Human Services**. The **NIH** is located in **Maryland, U.S.** and has **nearly 1,000** scientists and support staff. The **NIH** obtained **US\$39** billion from **Congress** in **2019**”

ORG	NIH, the United States Department of Health and Human Services, Congress
DATE	1887, 2019
MONEY	US\$39 billion
CARDINAL	Nearly 1,000
GPE	Maryland, U.S.





- Patient history, e.g., [Problem], [Test], [Treatment] in biomedical NER

She has had **<reatment> physical therapy </reatment>** and recovered completely from that .
<test> Initial examination </test> showed **<problem> bruising </problem>** around the left eye , normal lung examination , normal heart examination , normal neurologic function with a baseline decreased mobility of **<problem> her left arm </problem>** .

Problem	Bruising, her left arm
Test	Initial examination
Treatment	Physical therapy



More to think about

In the domain you are interested in, what types of entities are important to extract?

Task 2: Extract named entities

Notice: Your work in Jupyter notebook will not be saved. Please download it.



Part of Speech Tagging (POS tagging)
Named Entity Recognition (NER)

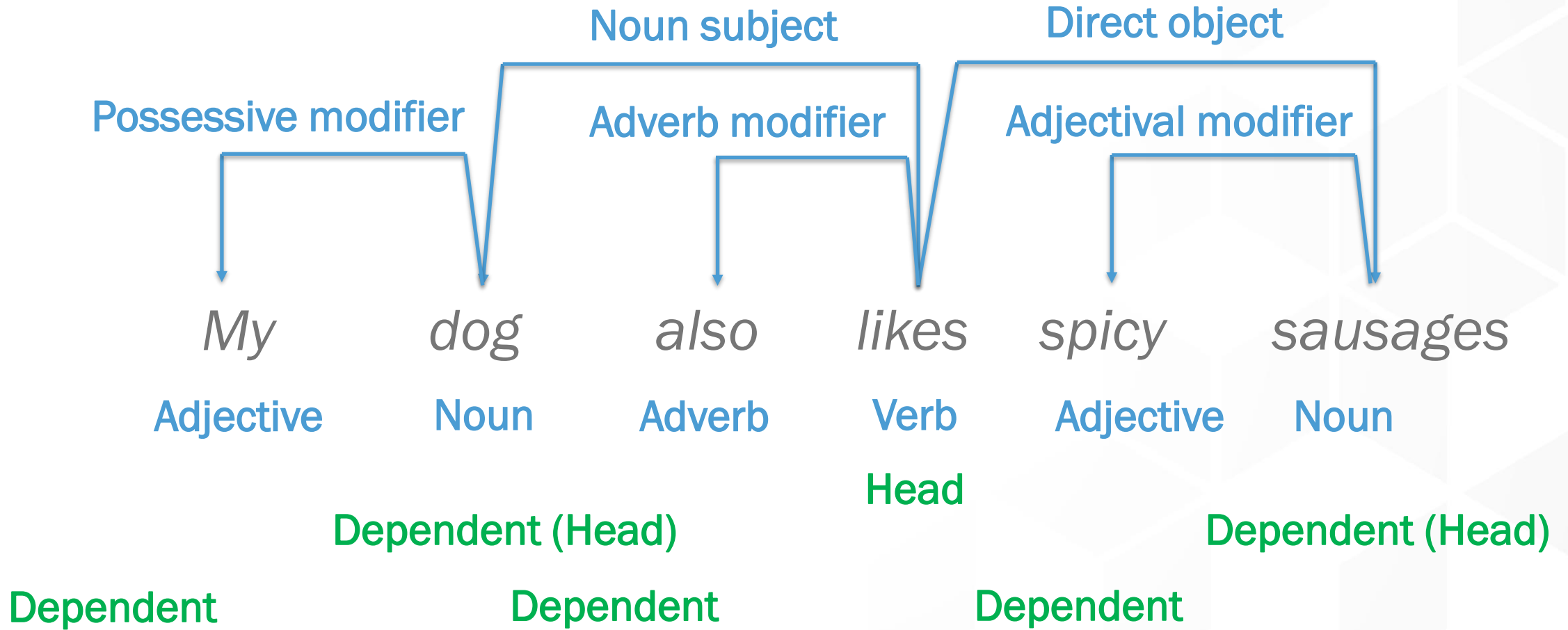
Relation Extraction



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



Syntactic Dependency Structure





More to think about

Task 3: How about the syntactic dependency structure of the following complex sentence?

I remember that you have given Tom a gift

Bell makes and distributes computer products.

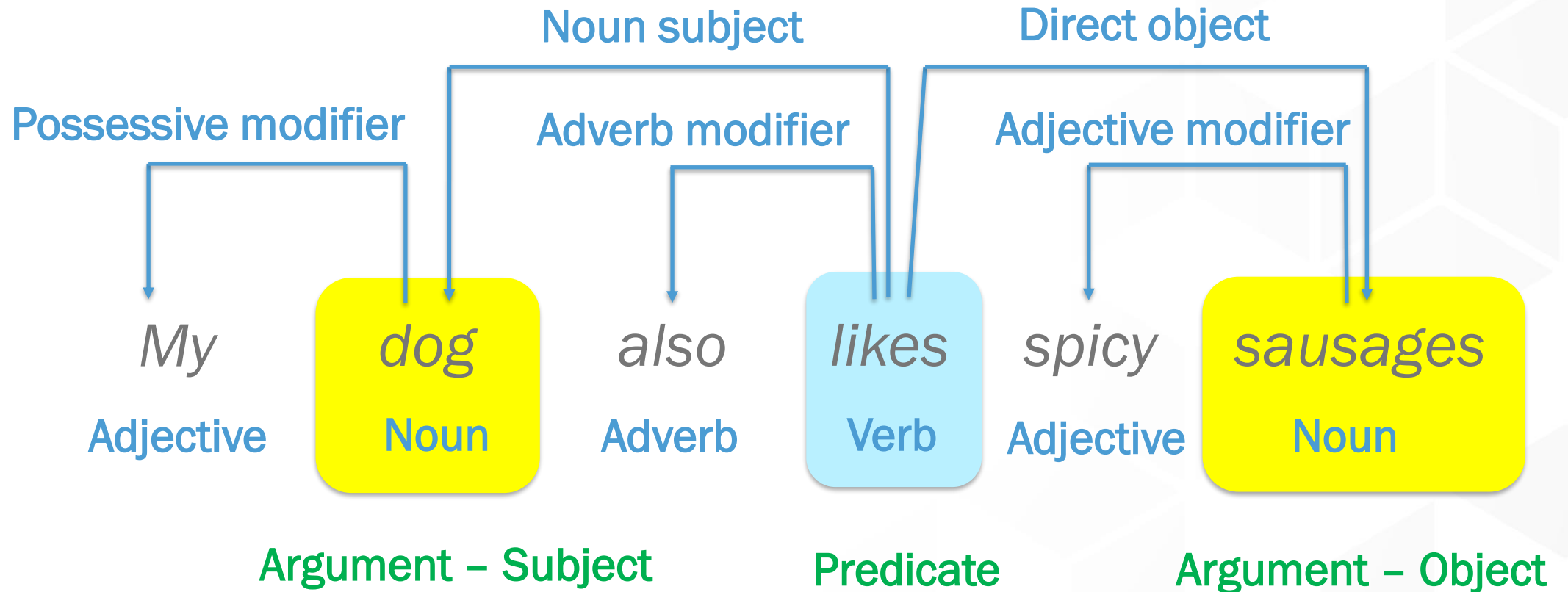
The NIH is located in Maryland, U.S. and it contains nearly 1,000 scientists and support staff.

Notice: Your work in Jupyter notebook will not be saved. Please download it.





Predicate–argument structure





Biomedical relation extraction

- Evidence-based practice, e.g., Drug efficacy, interaction, side effects

“We used **hemofiltration** to **treat** a **patient with digoxin overdose** that was **complicated** by **refractory hyperkalemia**”

Subjects	Predicate	Objects
Hemofiltration	TREATS	Patient, digoxin overdose
Hyperkalemia	COMPLICATES	Digoxin overdose
Digoxin overdose	PROCESS OF	Patient



More to think about

Task 4 (Optional): extract subject-predicate-object relation by analyzing the syntactic dependency

Notice: Your work in Jupyter notebook will not be saved. Please download it.

- You can play with the biomedical [relation extraction tool](#) with the text you are interested in or the sample text below.

This study demonstrates that netilmicin is a safe and effective antibiotic that can be used as a first choice treatment of acute bacterial conjunctivitis.





Wrap-up and take-aways

- Part of Speech Tagging
 - Named Entity Recognition
 - Relation Extraction
-
- Unstructured data → structured data
 - Information objects and relation



CAROLINA HEALTH INFORMATICS PROGRAM

We are grateful to the following organizations for their support:

UNITED HEALTH FOUNDATION®



Program for Precision Medicine in Health Care

Interested in an MPS or a Ph.D. in digital health?

Drop us a line: chip@unc.edu.



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



Reference

Jurafsky, D., & Martin, J. H. (2016). *Speech and Language Processing*. (3rd ed.)

- Overview:

Task of information extraction,

<https://people.cs.umass.edu/~mccallum/courses/inlp2007/lect20-ie.ppt.pdf>

- POS part:

Fan, J., Prasad, R., Yabut, R. M., Loomis, R. M., Zisook, D. S., Mattison, J. E., & Huang, Y. (2011). Part-of-speech tagging for clinical text: wall or bridge between institutions? *AMIA Annual Symposium Proceedings, 2011*, 382–391.

- NER:

Ramshaw, L. A., & Marcus, M. P. (1999). Text Chunking Using Transformation-Based Learning. In S. Armstrong, K. Church, P. Isabelle, S. Manzi, E. Tzoukermann, & D. Yarowsky (Eds.), *Natural language processing using very large corpora* (Vol. 11, pp. 157–176). Springer Netherlands. https://doi.org/10.1007/978-94-017-2390-9_10



Reference (cont.)

Rindflesch, T.C. and Fiszman, M. (2003). The interaction of domain knowledge and linguistic structure in natural language processing: interpreting hypernymic propositions in biomedical text. *Journal of Biomedical Informatics*, 36(6):462-477.

Kundeti, S. R., Vijayananda, J., Mujjiga, S., & Kalyan, M. (2016). Clinical named entity recognition: Challenges and opportunities. *2016 IEEE International Conference on Big Data (Big Data)*, 1937–1945. <https://doi.org/10.1109/BigData.2016.7840814>

Wu, Y., Jiang, M., Xu, J., Zhi, D., & Xu, H. (2017). Clinical named entity recognition using deep learning models. *AMIA Annual Symposium Proceedings*, 2017, 1812–1819.

Wang, Y., Wang, L., Rastegar-Mojarad, M., Moon, S., Shen, F., Afzal, N., Liu, S., Zeng, Y., Mehrabi, S., Sohn, S., & Liu, H. (2018). Clinical information extraction applications: A literature review. *Journal of Biomedical Informatics*, 77, 34–49. <https://doi.org/10.1016/j.jbi.2017.11.011>

- Relation extraction:

Nivre, J. (2005). Dependency grammar and dependency parsing. *MSI Report*, 5133, 1–32.

Wang, Y., Wang, L., Rastegar-Mojarad, M., Moon, S., Shen, F., Afzal, N., ... Liu, H. (2018). Clinical information extraction applications: A literature review. *Journal of Biomedical Informatics*, 77, 34–49.



Figure Citation

Text information extraction

<https://www.gapingvoid.com/>

POS Tagging

<http://acl.ldc.upenn.edu/J/J93/J93->

https://www.ling.upenn.edu/courses/Fall_2003/ling001/penn_treebank_pos.

<https://universaldependencies.org/u/pos/index.html>

NER

<http://text-machine.cs.uml.edu/cliner/samples/cliner-sample-output.pdf>

NLP applications in clinical information extraction

Name	Description	Website
cTAKES	Open-source NLP system based on UIMA framework for extraction of information from electronic health records unstructured clinical text	http://ctakes.apache.org/
MetaMap	National Institutes of Health (NIH)-developed NLP tool that maps biomedical text to UMLS concepts	https://metamap.nlm.nih.gov/
MedLEE	NLP system that extracts, structures, and encodes clinical information from narrative clinical notes	http://zellig.cpmc.columbia.edu/medlee/
KnowledgeMap		https://medschool.vanderbilt.edu/cpm/center-precision-medicine-blog/kmci-knowledgemap-concept-indexer
Concept Indexer (KMCI)	NLP system that identifies biomedical concepts and maps them to UMLS concepts	
HITEx	Open-source NLP tool built on top of the GATE framework for various tasks such as principal diagnoses extraction and smoking status extraction	https://www.i2b2.org/software/projects/hitex/hitex_manual.html
MedEx	NLP tool used to recognize drug names, dose, route, and frequency from free-text clinical records	https://medschool.vanderbilt.edu/cpm/center-precision-medicine-blog/medex-tool-finding-medication-information
MedTagger	Open-source NLP pipeline based on UIMA framework for indexing based on dictionaries, information extraction, and machine learning–based named entity recognition from clinical text	http://ohnlp.org/index.php/MedTagger
ARC	Automated retrieval console (ARC) is an open-source NLP pipeline that converts unstructured text to structured data such as Systematized Nomenclature of Medicine – Clinical Terms (SNOMED CT) or UMLS codes	http://blulab.chpc.utah.edu/content/arc-automated-retrieval-console
Medtex	Clinical NLP software that extracts meaningful information from narrative text to facilitate clinical staff in decision-making process	https://aehrc.com/research/projects/medical-free-text-retrieval-and-analytics/#medtex
CLAMP	NLP software system based on UIMA framework for clinical language annotation, modeling, processing and machine learning	https://sbmi.uth.edu/ccb/resources/clamp.htm
MedXN	A tool to extract comprehensive medication information from clinical narratives and normalize it to RxNorm	http://ohnlp.org/index.php/MedXN
MedTime	A tool to extract temporal information from clinical narratives and normalize it to the TIMEX3 standard	http://ohnlp.org/index.php/MedTime
PredMED	NLP application developed by IBM to extract full prescriptions from narrative clinical notes	

(Wang et al., 2018)

Text Mining

- Part of Speech Tagging (POS tagging)
- Named Entity Recognition (NER)
- Relation Extraction



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



Zhaopeng Xing
Doctoral Student at Carolina
Health Informatics Program in
Biomedical Health Informatics