

# Natural Language processing applications

Omar Alrassam

# Outlines

- What is NLP?
- The importance of NLP
- Applications of NLP
- NPL in Kurdish language
- Future work
- Conclusion
- References

# What is NLP?

- **Natural Language Processing**, or (NLP) for short is a tract of Artificial Intelligence and Linguistics, devoted to make computers understand the statements or words written in human languages.
- Natural language processing is the area of study dedicated to the automatic manipulation of speech and text by software.
- Natural Language Processing, is broadly defined as the automatic manipulation of natural language, like speech and text, by software.
- The study of natural language processing has been around for more than 50 years and grew out of the field of linguistics

# The importance of NLP

- NLP is important for scientific, economic, social, and cultural reasons. NLP is experiencing rapid growth as its theories and methods are deployed in a variety of new language technologies. For this reason it is important for a wide range of people to have a working knowledge of NLP. Within industry, this includes people in human-computer interaction, business information analysis, and web software development. Within academia, it includes people in areas from humanities computing and corpus linguistics through to computer science and artificial intelligence. (To many people in academia, NLP is known by the name of “Computational Linguistics.”)

# Some Applications of NLP

- Spelling checker
- Machine translation
- Text summarization
- Automatic Question Answering
- Speech Recognition
- Text-to-Speech Synthesis
- Information Retrieval

# Spelling Checker

Spelling Correction is a process of detecting and sometimes providing suggestions for incorrectly spelled words in a text. Spell Checker is an application program that flags words in a document that may not be spelled correctly. Spell Checker may be stand-alone capable of operating on a block a text

# Spelling Checker cont

## \* Error Detection Techniques

1- Dictionary Lookup Technique

2- N-gram Analysis Technique  $P(w_n | w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})}$

3- Neural Net Techniques

# Some important terms in NLP

- Corpus

1- One of the first things required for natural language processing (NLP) tasks is a corpus.

2- A large collection of written or spoken language, that is used for studying the language.

- Stop Words: A stop word is a commonly used word (such as “the”, “a”, “an”, “in”) that a search engine has been programmed to ignore, both when indexing entries for searching and when retrieving them as the result of a search query.



# Some important terms in NLP

- Example of stop words

Sample text with Stop Words	Without Stop Words
GeeksforGeeks – A Computer Science Portal for Geeks	GeeksforGeeks , Computer Science, Portal ,Geeks
Can listening be exhausting?	Listening, Exhausting
I like reading, so I read	Like, Reading, read

# Some important terms in NLP

- Stemming

- Example

- Cars, car's, cars' □ car

- Am, is, are □ be

- the boy's cars are different colors

- the boy car be differ color

# Challenges in Kurdish Text Processing

**According to the study of (Sheykh Esmaili, 2014)**

- Despite having a large number of speakers, the Kurdish language is among the less-resourced languages.
- Spoken by more than 30 million people in western Asia.
- Has many dialects, the two major ones
  - Northern Kurdish (Kurmanji)
  - Central Kurdish (Sorani)

# Challenges in Kurdish Text Processing

## 1- Dialect Diversity

- Grammar (F,M)
- Writing system

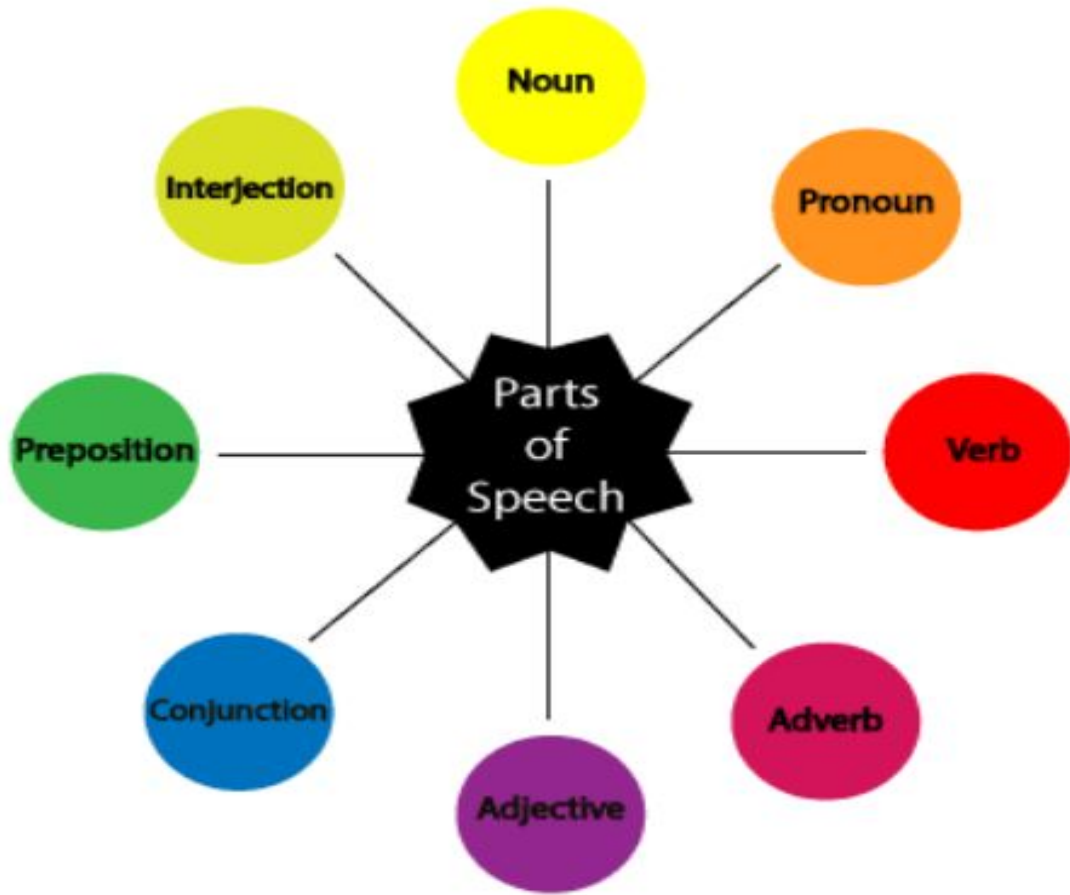
**Kurmanji using Latin-based letters**

**Sorani using Arabic-based letters**

**As a result Difficulties in Segmentation and  
Tokenization process**

# Challenges in Kurdish Text Processing

## 2- part-of-speech detection.



# Challenges in Kurdish Text Processing

3- Arabic alphabet does not have capitalization

and therefore it is more difficult to recognize

sentence boundaries as well as recognizing

Named Entities.

# Challenges in Kurdish Text Processing

4- Lack of Language Resources Kurdish is a resource-scarce language for which the only linguistic resource available on the Web is raw text. More concretely, in spite the few attempts in building corpus and lexicon, Kurdish still does not have any large-scale and reliable general/domainspecific corpus.

Furthermore, no test collection {which is essential in evaluation of Information Retrieval systems} or stemming algorithm has been developed for Kurdish so far.

# Challenges in Kurdish Text Processing

## Updating words

- 1\_ کارئزما \_\_\_\_\_ کهسی پیروژ\_1
- 2\_ ریفورم \_\_\_\_\_ چاکسازی
- 3\_ ناسیونالیزم \_\_\_\_\_ نهتهوهگه‌ری
- 4\_ ستاندار \_\_\_\_\_ پیوانه‌یی
- 5\_ پرۆتوکول \_\_\_\_\_ ریکه‌وتن
- 6\_ پرۆلیتار \_\_\_\_\_ چینی کرێکاران
- 7\_ کاپیتالیزم \_\_\_\_\_ سه‌رمایه‌داری
- 8\_ فیودالیزم \_\_\_\_\_ ده‌ر به‌گایه‌تی
- 9\_ ئەنارشیزم \_\_\_\_\_ ئاژاوه‌چیتی
- 10\_ ئایدیالیزم \_\_\_\_\_ فه‌لسه‌فه‌ی میسالی (خه‌یالی)
- 11\_ ریالیزم \_\_\_\_\_ فه‌لسه‌فه‌ی واقه‌گه‌ری (واقعی)
- 12\_ یوتوپیا \_\_\_\_\_ ئەم زار اوایی زیاتر بۆ خه‌یالی فه‌لسه‌فی فه‌یله‌سوفان به‌کادیت و تایبه‌ته به فه‌یله‌سوفان
- 13\_ ئەلته‌ر ناتیف - به‌دیل (جیگره وه )
- 14\_ ئەکادیمی - زانستی
- 15\_ ئەنتی - دژ



# Challenges in Kurdish Text Processing

- 16\_ بایلو جی - زیندهوهر ناسی
- 17\_ کومیدیا - بهز مه سات
- 18\_ جینوساید - قرکردن، قه لاجوکردن
- 19\_ دیموگرافی - دانیشتوانی
- 20\_ ئەکتیف - چالاک، کارا
- 22\_ ستاندار - پیوانهیی
- 23\_ ئەنتوگرافیا \_\_\_ وهسفی گهلان، گهلناسی
- 24\_ ئوتوکراتی \_\_\_ حوکمی تاکه کەس
- 25\_ ئورگان \_\_\_ ئەندام، زمانحال
- 26\_ ئۆپزیسیون \_\_\_ بهر هه لستکار
- 27\_ ئینسکلۆپیدیا \_\_\_ موسوعه (زانباری نامه )
- 28\_ ئیکۆلوژی \_\_\_ ژینگه ناسی
- 29\_ ئیکۆنومیزم \_\_\_ ئابوریناسی
- 30\_ هیومانیزم \_\_\_ مرفه گه رایی

# Challenges in Kurdish Text Processing

31\_ کەپیتالیزم \_\_ سەرمايهدارى

32\_ کاندید \_\_ پالیۆراو

33\_ کامپ \_\_ نیشینگه

34\_ کۆنکریتی \_\_ بەلگه‌دار

35\_ دوکیۆمێنت \_\_ بەلگه‌نامه

36\_ کۆمسیۆن \_\_ لیژنه

37\_ دیسپلین \_\_ ریکخەر

38\_ دیموکراسی \_\_ حوکمی گهل

39\_ دۆسییه \_\_ فایل

40\_ دیمۆگرافی \_\_ دانیشتووی

41\_ دیالیکت \_\_ شیوه‌زار

42\_ دیفاکتۆ \_\_ ئەمری واقیع

43\_ یراکتیک \_\_ جێبه‌جێکردن

# Challenges in Kurdish Text Processing

پروژه — کرده\_46

47\_ ناسیونالیزم — نهتهوهگهری

48\_ کۆلونیالیزم — داگیرکەر

49\_ کۆلونیالیکرا — داگیرکراو

50\_ ماتریالیزم — ماددییهت

# Challenges in Kurdish Text Processing



# Future work

## Deep learning Algorithm

- **Deep learning** (also known as **deep structured learning** or **hierarchical learning**) is part of a broader family of [machine learning](#) methods based on artificial neural networks.

# Conclusion

- In this seminar the importance of NLP is discussed and the challenging of Kurdish text was explained.
- Also, focusing on the Text processing in Kurdish language was the main aim of this seminar.