<u>Lecture 1:</u>          <u>Introduction to NLP</u>

# Natural Language Processing or NLP is a ~~mai~~ multidisciplinary field combining math linguistics and computer science, the goal is to get computer to do useful things with natural language data.

Common NLP tasks and applications:
_____

→ Translation

→ Summarization

→ Question answering

→ speech recognition

→ classification

→ assisted writing.

# some reasons why NLP is hard:

1. Content
2. Inflections
3. Content
4. misspelled or misused words.
5. physical gesture.

# Fundamentals of NLP

→ Preprocessing:

    → tokenization
    → simplification technique
    → tagging and
    → simple rules based approaches.

→ Basic vectorization:
  → turning text to numbers.
  → measuring similarity between docs.

→ Modelling overview:
  → types of machine learning algorithms
  → vs. models
  → evaluation

→ first steps into classification:
  → classifying text using Navie Bayes.
  → evaluation with precision and recall

→ Topic modeling:
  → Automatically finding topics in docs using latent dirichlet allocation.

# Deep learning for NLP:

→ word vectors: capturing word meaning, the concept of embedding.

→ Recurrent Neural Networks: Capturing sequence information and generating learn language.

→ Neural Networks: what they are, how they work and details around training

→ Sequence2 sequence seq2seq and attention: training a neural network to transform one sequence to another.

→ Transformation: The dominant mainstream architecture today. Pretraining and transfer learning

## Processing :

1. tokenization : The process of segmenting our documents into tokens.

### Code:

```
// insurting library libraries

() !pip install -u spacy == 3. *

() ! python -m spacy info

() import spacy
() ! python -m spacy download en_Cone_web-sm

() nlp = spacy.load ('en_Cone-web-sm')

() type (nlp)

//sample sentense testing

S = "I eat nice"
doc = nlp (s)
print ([t.text for t in doc]   //output !
                                 ["I", "eat", "nice"
```

# Basic programming:

case folding, stop word removal, stemming,
lemmatization.

## Case folding

sentense: "Mr. Cook went into the
kitchen to Cook dinner"

### without cf:

{ Cook, dinner, into,
kitchen, mr, the,
to, went, cook}

### with cf

{cook, dinner,
into, kitchen, mr,
the, to, went}

### Code:

```
Print([t.lower_ for t in doc])

// to skip for first word

Print([t.lower_ if not t.is_sent_start else
       t for t in doc])
```

## Stop word removal:

stop words → {the, a, of, an, this, that}

### Code:

```
Print (nlp. Defaults. stop_words)
Print (len (nlp. Defaults. stop_words))

Print ([t for t in doc if not t.is_stop])
```

## stemming

removing word suffixes or prefixes.

Banking  }
           }  Bank
Banks     }

# Lematization

Reduce a word down to its lemma or dictionary form.

Did
Done
Doing
} Do

Code: [(t.text, t.lemma_) for t in doc]

# Week-2 || Advance processing

Part of speech (POS) Tagging:
  {noun, verb, adjective, .....}

Example:
  "John watched an old movie at the cinema."

  Prop. N / Verb / Det / Noun / ADP / DET / ADJ / Pun Noun / Punc

## || Code

→ [(t.text, t.pos_) for t in doc]

→ [(t.text, t.tag_) for t in doc]

## Named Identity Recognition (NER)

  ..{a person, a location, an organization .....]

Named entity: anything that can be referred by a proper ~~noun~~ name. They often have a proper Noun (PROPN) pos tag.

Example:

Person → PER
Location → LOC
Geopolitical Entity ⟷ GPE
Organization → ORG

Usefulness:

→ Organizing / Categorizing Corpus

→ Question answering

→ Critical in information extraction.

Challange:

An entity can speak multiple tokens.

Hamilton

us president?    City    f₁ driven?

## Code:

```
doc = nlp(s)          // s is a sentence or string

[(t.text, t.ent_type_) for t in doc)

// entity non zero check

[(t.text, t.ent_type_) for t in doc if
              t.ent_type != 0]

// iterating through ents.

[(ent.text, ent.label) for ent in doc.ents)
```

## Parsing

Determining the syntactic structure of
a sentence.

type:
→ Constituency parcing
→ Dependency parcing

# Constituency parsing using CFG:

NP — Noun phrase
VP — verb phrase
PP → Prepositional phrase
NN → Noun
PPP → Personal pronoun
NNP → Proper Noun
VB → verb (Base form
DT → Determiner
IN → Preposition

| Production Rules | Lexicon |
|---|---|
| S → NP VP | DT → the \| a \| this \| that |
| NP → PRP \| NNP \| DT NN | PRP → I \| she \| he |
| VP → VB \| VB NP \| VP PP | IN → in \| at |
| PP → IN NP | NN → book \| hotel \| room. |

" she enrolled in th course at university,"

| She |
|-----|
| NP |

| enrolled | in the | course at univers |
|----------|--------|-------------------|
| | VP | |

" she enrolled in the Course at university. "

| she |
|-----|
| NP |

| enrolled in the Course at university |
|--------------------------------------|
| VP |

| · |
|---|
| · |

| She |
|-----|
| PRP |

| enrolled |
|----------|
| VBD |

| in the Course |
|---------------|
| PP |

| at the university |
|-------------------|
| PP |

| in |
|----|
| IN |

| the course |
|------------|
| NP |

| at |
|----|
| IN |

| the univer |
|------------|
| NP |

| the |
|-----|
| DT |

| course |
|--------|
| NN |

| the |
|-----|
| DT |

| university |
|------------|
| NN |

# Rule

S → NP VP

NP → DET N | DET ADJ N

VP → V NP

## Lexicon

DET → a | the

ADJ → beautiful | Perching

N → bird | birds | grain | grains

V → peck | pecks | pecking

→ "The bird pecks the grain."

# Parsing

→ Determine the cymetric structure of sentense.

" She enrolled

" The quick brown fox jump over the lazy dog "