

NLP - Natural Language processing

Agenda:

- ① Introduction to NLP
- ② Application of NLP
- ③ Various NLP Tasks
- ④ What is Language
- ⑤ Approaches to NLP
- ⑥ Challenges in NLP
- ⑦ NLP pipeline

① My name is Barry and
I am a data scientist

②

There was a tiger. He was very
cruel in his youth. He was a man-
eater. In his old age, he became
very weak. He could no longer
hunt down any animal. So he
suffered much for want of food.
ad a gold bangle in his
possession. He sat in an open
place in the forest with that
bangle.

Tomorrow's Agenda: (practical)

- ① Text preprocessing Technique
- ② Text representation
- ③ Word embeddings

NLP pipeline

① Data Acquisition:

- Available Data (CSV, TXT, PDF, XLS)
- Others Data (DB, Internet, API, Scraping)
- No data (create your own data)

Note: If you have less data then use Data Augmentation

→ Replace with synonyms

ex: I am a Data Scientist ✓

I am a AI Engineer ✓

→ biagram flip

exp: ① I am Barry ✓
② Barry is my name ✓

→ back translate

→ add Additional noise

exp: I am ~ Data scientist,
I love this job → extra

② Text preparation:

- ① cleanup: HTML tags, emoji, spelling correction
- ② Basic preprocessing
- ③ Advance preprocessing

Basic preprocessing:

* Tokenization { sentence
word

I am Barry,
I am a PS
list = ['I', 'am', 'Barry']
list = ['I am Barry',
I am a PS]

* optional steps:

- ① stop word removal
- ② stemming - less use
- ③ lematization - good
- ④ Removing punctuation (-!?.)
- ⑤ lower case
- ⑥ language Detection

'I' 'am' 'Barry'
[0 1 2]

[playing, played, plays]
↓
[play]

Barry, barry
barry



Text summar
in Eng
Bengli →

Advanced pre-processing:

- ① parts of speech tagging
- ② punting
- ③ coreference resolution

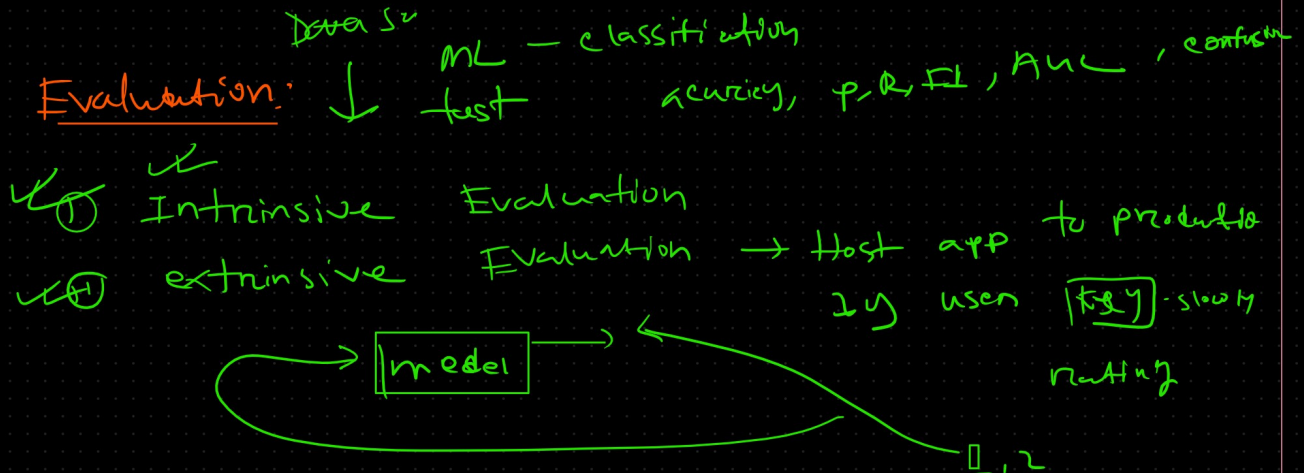
③ Feature Engineering

- Text Vectorization
- TFIDF, Bag of word, BoW, one Hot, word to vec
- Encoding

④ Modelling

- ① heuristic
- ② ML
- ③ DL
- ④ cloud API (AWS, GCP, Azure) → paid one

⑤ Evaluation:



⑥ Deployment:

monitoring, retraining

common terms:

- ① corpus (entire text)
- ② vocabulary (unique word)
- ③ Document (one row)
- ④ word (single word)

Harrison was in fact a scion of the Virginia planter aristocracy. He was born at Berkeley in 1773. He studied classics and history at Hampden-Sydney College, then began the study of medicine in Richmond.