27/4/2024

# Summary of NLP (as of 27/4)

1.  Search Engine
    World's Application application

    Nutch → Crawler → web → extract

    BigData → Unstructured →

2.  Preprocessing ≀
    1. lower case
    2. punctuation
    3. stop words
    4. Numbers
    5. Tags (HTML, XML)
    6. Contraction <u>I've, shouldn't</u>
                        <u>I have   Shouldnot</u>
    7. lematatization ← → slow ─ acce. ─ large
    8. Stemming ← ─────────
                    └ fast ─ less accurate, large

    9. Tokenzation
    10. accents words  â, ä,

# Text to Numbers

| Count Vectorizer | Bag of Words | TF-IDF | Word embedding (word2vec, glove, Fasttext) |
|---|---|---|---|
| disadvantage | advantage | advantage | advantage |
| No Frequency within document | 1. Frequency within the document | 1. Frequency with corpus is also identify | 1. Context is preserved |
| NO Order | dis adwantage | disadvantage | disadvantage |
| NO Context | No Frequency with Corpus | 1. No context | 1. No order preserved |
| NO Frequnt within the Corpus | | 2. No order | 2. Same word can have different meaning based on Context |

1. Hi good morning,
2. How is your morning
3. Morning is beautiful

| | Hi | How | good | morning | your | beautiful |
|---|---|---|---|---|---|---|
| 1 | 1 | | 1 | 1 | | |
| 2 | | 1 | | 1 | 1 | |
| 3 | | | | 1 | | 1 |

I went to (HDFC) Bank.
I went to River (Bank)

I dont have any interest buying loan
The Interest rate is very high to get loan

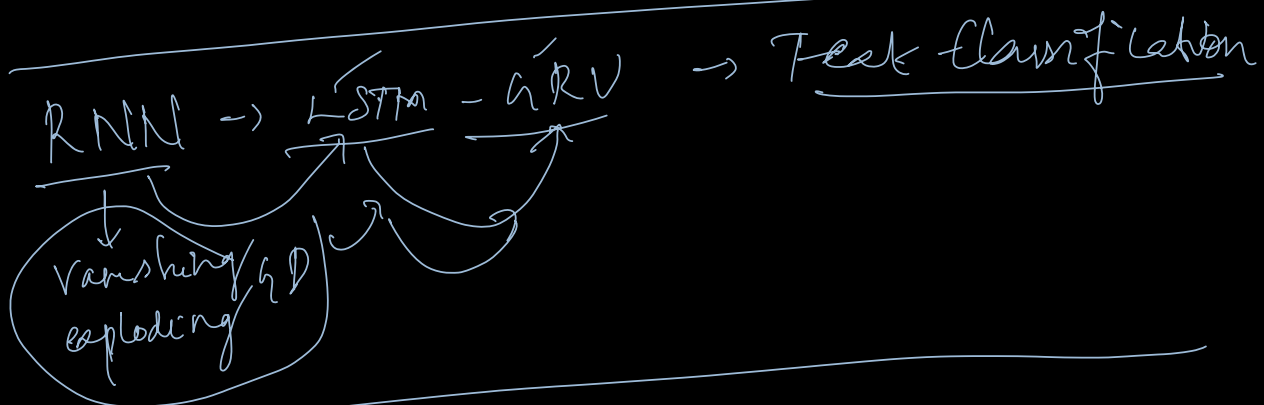Unsupervised → Text Clustering – k-means – bucket the similar document

Cosine Similarity : Recommendation engine
Content

Sentiment Analysis - lexicon-model -
Unsupervised textblob, Afinn, Vader

Text Classification
Bag of Words          logistic          Deep learning
embedding             machine
                      learning

RNN -> LSTM - GRU    -> Text Classification
      |
   Vanishing/GD
   exploding

Seq2Seq Model -> encoder - Decoder

              Machine Translation

Image Captioning