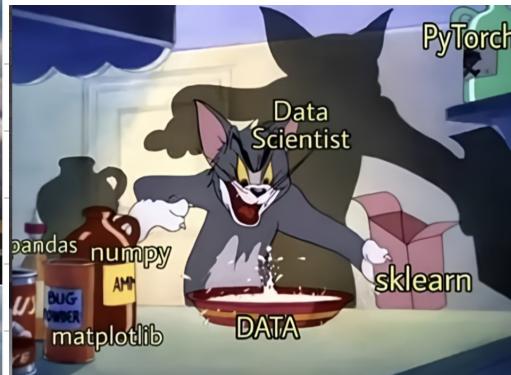


Session -1

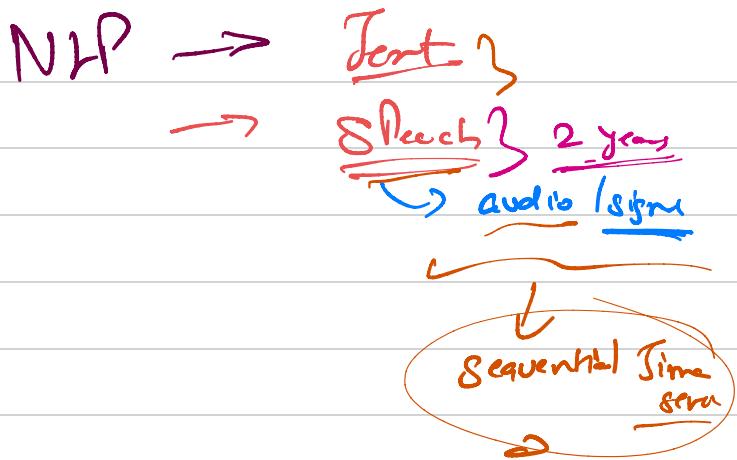
INTRO To NLP

April 12, 2024



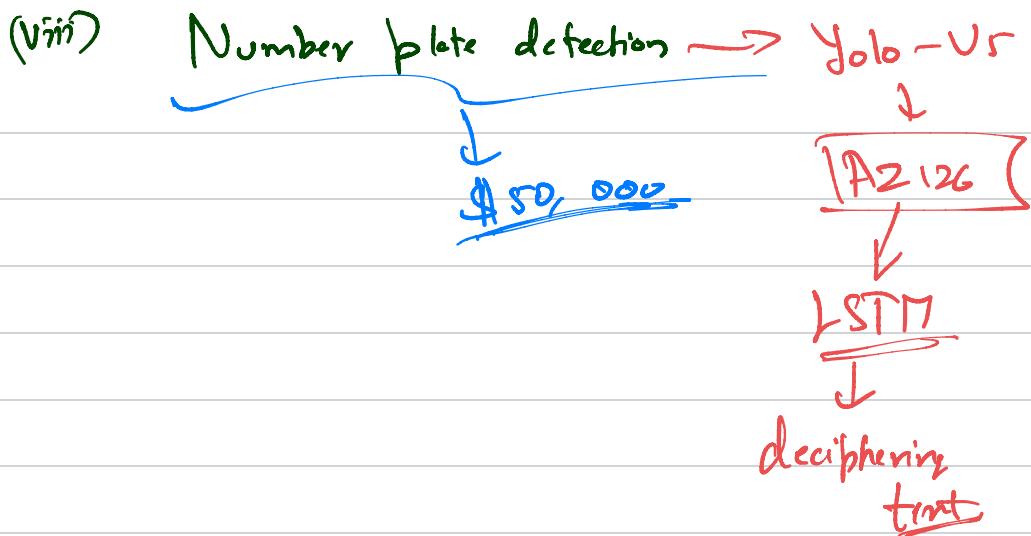
AGENDA

- ① Overview of NLP module + applications.
- ② Basic terms & preprocessing
- ③ Task: Sentiment Classifications
- ④ Encode text as Vector
- ⑤ Build a Simple model.



→ Projects → End of module

- (i) Sent Classification (Sentiment)
- (ii) L2T (CPT)
- (iii) NMT (neural machine translation)
- (iv) Chatbots
- (v) Sent → Coming From L2T
 - Plagiarism → ?
- (vi) Sent Summarization → Pointer generator
 - Abstraction
 - Text summarization
- (vii) CC - Closed captions (Youtube)
- (viii) Spelling checker



CV \rightarrow Spatial Kind \rightarrow forming Spatial ds

RNN \rightarrow " " " " Sequential
Structure

~~are~~, can, can't, shall
will

NLP (couple of decade)

```

graph LR
    NLP[NLP] --> RB[Rule-based]
    NLP --> H[Heuristic]
    RB --> B[Bow]
    RB --> T[TF-IDF]
  
```

DL \rightarrow Word2Vec \rightarrow Word \rightarrow Vector

\hookrightarrow 2012-2013 40-50 / 100

→ Nurse → Woman
↓
Female

Doctor → male

Remove
bias

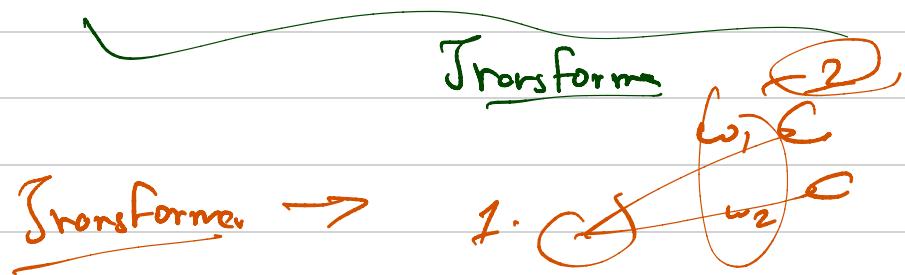
→ Probabilistic Techinique → LDA, CRF
✓
15-20 years

RNN → LSTM, GRU → last 10 years

→ Transformer → Attention
→ Key idea

→ BERT - GPT

RTX-4090 → 24 GB
→ 2-3 tokens INR
→ VRAM
→ 80 GB
→ 165 million \$
→ 25,000 AI models
→ OpenAI
↓
2 months



GPT-4.5 \rightarrow 1.8 Trillion

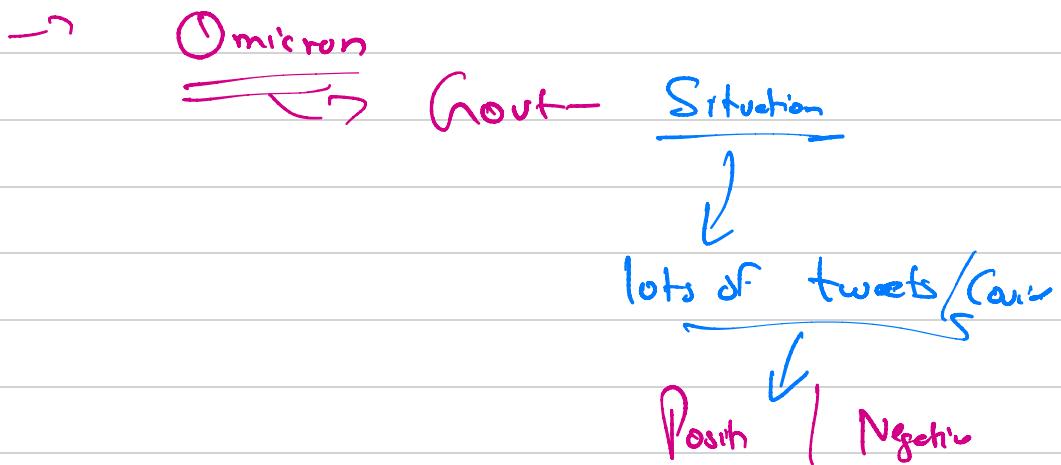
H. Brain \rightarrow 1000 trillion

10⁴¹ ~~x~~ True alpha

min amt of switch

word

Sent - Classification



→ Some pre-processing: !!

"I love WFH,
thank you!!"

①

Soldenization

②

Stop words

break down sentence
words

③

Stemming

Lemmatisation

Crude algos

running → run

U-fast

Carrying → car

gone → go
Wont → go

P-Statement:

banned keywords:

Product 1

→ allowed (or not)

Cannabis → indica
marijuana → con't allow

g - cr - → Product
marijuana → 1/0

①

Spelling mistakes → garlic

②

Hindi + English + Tamil + Telugu

③

Synonym nam

Granulation Porut stone

80%

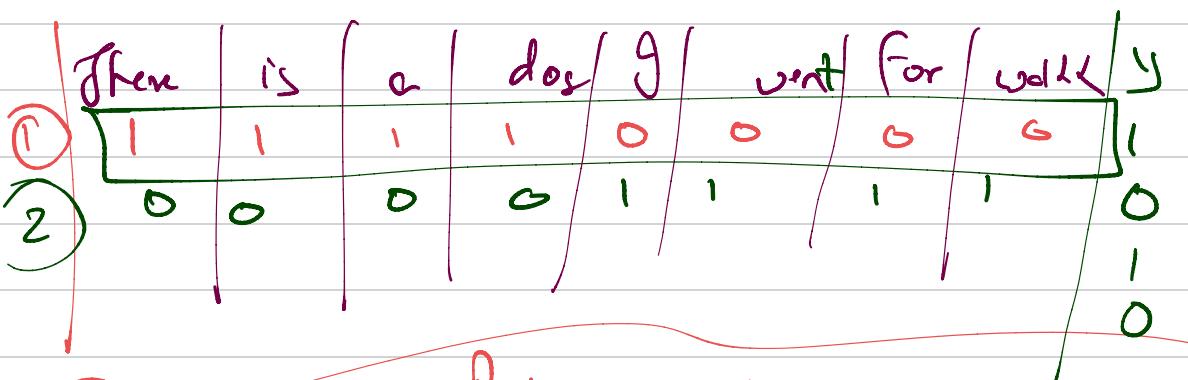


Character level - LSTM network
naive bayes } - 99.6%

Bow →

There is a dog

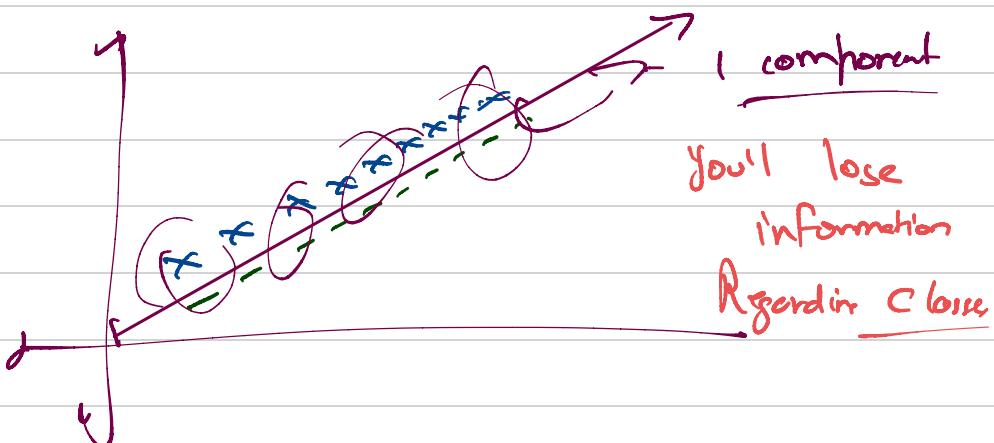
g went for walk.



Problem with representation

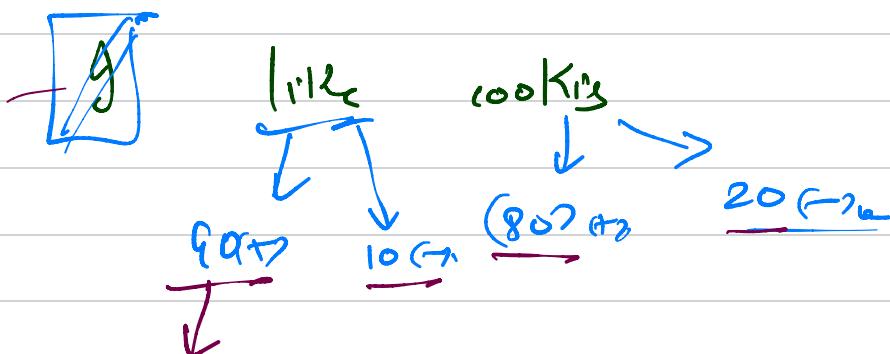
U^T U sparse

→ Need to reduce dimensionality.



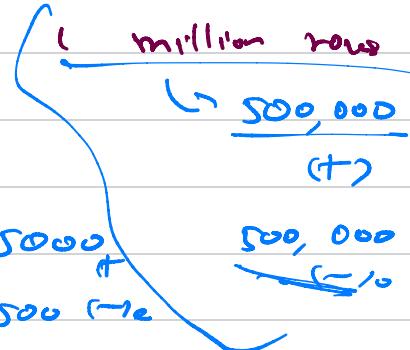
Orange → across all the tweets
 → AD → tweets → Freq of
 → → tweets → Freq of
Orange

Orange
 , AD freq, word freq
 ↓
 + the words in a given sent



80% of 100K's, 30 → d dimension of
 +
 2 dimensions
10K
10K, 2

Trainin. →



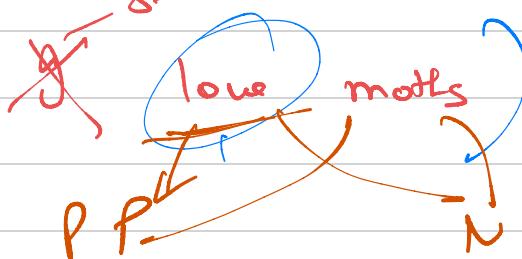
~~This is my dot~~

→ 5000

→ 500

500,000
(-)

8:00



test dot

→ Model limitation

You rely → on v. big dataset

