

Chap - 22 (NLP)

NLP is a branch of AI that helps to communicate with intelligent system using human language instead of programming language.



Purpose of NLP

- i) Language Understanding: Help computers understand human language. Ex: Google search, Siri
- ii) Language Generation: make computers generate text
Ex: Chat GPT writes essay, paragraph
- iii) Machine Translation: Translate one language to another. Ex: Google translator
- iv) Sentiment Analysis: detect emotion or opinion
Ex: social media reviews

Subject:.....

Date:.....

v) Speech recognition: Convert spoken language to text. Ex: Voice assistant (Siri)

vi) Text summarization: Large text \rightarrow Short text (meaningful)

vii) Question Answering: Answer question asked in natural language.

viii) Chatbots: conversational system respond like

Ex: ChatGPT

N-gram character model

N-gram character model is a sequence of N characters that appear together.

⇒ 1-gram (Unigram) → one character ("a")

⇒ 2-gram (Bigram) → two characters ("th")

⇒ 3-gram (Trigram) → Three characters ("the")

A N-gram model learns how often these character groups appear in a language, helping to predict the next character.

* A N-gram model is defined as a **Markov** chain of

order $N-1$.

⇒ next character depends on the previous $N-1$ characters.

Bigram → previous one Trigram → previous 2

$$P(c_1:N) = \prod_{i=1}^N P(c_i | c_{i-1:i-1}) = \prod_{i=1}^N P(c_i | c_{i-2:i-1})$$

Uses

- i) Language identification (\rightarrow detect language)
- ii) Spelling Correction
- iii) Genre Classification (\rightarrow identify the type of text (news, legal, scientific))
- iv) Named Entity Recognition (NER) (\rightarrow find Names of people, drugs, place etc.)

Smoothing n-gram models

Problem: When we make an N-gram model, we count how often words or letters come together.

\Rightarrow But sometimes, in our training text, some word combinations never appear.

Ex: "th" \rightarrow this, that (common)

But "ht" \rightarrow maybe never show

model makes its probability = 0

But "http" word is available in reality.

$$\frac{(n-i+1)(n)}{\prod_{i=1}^n} = \frac{(n-i+1)}{\prod_{i=1}^n} = \frac{n}{(n-i)}$$

Solution: \rightarrow smoothing

\hookrightarrow we don't let anything have probability 0

Types of Smoothing

i) Laplace Smoothing: Just add 1 to every count

Ex: 'ht' appear 0 times, make it 1 time.

\rightarrow simple but not very accurate.

ii) Backoff model: If big model doesn't know something then backs off to smaller model.

Try 3-gram \rightarrow "The cat" if missing

" 2-gram \rightarrow "cat"

" 1-gram \rightarrow "the"

iii) Linear Interpolation: Combining trigram, bigram, unigram together.

$$P(c_i | c_{i-2:i-1}) = \lambda_3 P(c_i | c_{i-2:i-1}) + \lambda_2 P(c_i | c_{i-1}) + \lambda_1 P(c_i)$$

where, $\lambda_3 + \lambda_2 + \lambda_1 = 1$

$[\lambda \rightarrow \text{importance scores}]$

Model evaluation

There are many ways to build n-gram model (unigram, bigram, trigram). We need to know which one works best.

↳ solution model evaluation.

Step-1: split the corpus

⇒ Training set — Teach the model

⇒ Validation set — Test how well it learned

Step-2: how to measure "goodness"

↳ option-1: Task-based evaluation —

Ex: language identification

→ Accuracy predicts (%)

↳ option-2: probability-based evaluation —

→ give probability for validation test

→ Higher probability → better model

Problem → for long texts probability became tiny

→ solution: perplexity +
[
entropy & A]

Subject:.....

Date:.....

Step-3 : Perplexity

$$\text{perplexity}(c_1:N) = \frac{1}{P(c_1:N)^{1/N}}$$

Model predicts well \rightarrow probability \uparrow perplexity \downarrow

Model confused \rightarrow probability \downarrow perplexity \uparrow

N-gram Word model

"Which word is likely come after another word?"

\Rightarrow Unigram \rightarrow single words ("The", "cat")

\Rightarrow Bigram \rightarrow 2 words ("The cat", "runs fast")

\Rightarrow Trigram \rightarrow 3 words ("The cat runs", "Cat runs fast")

problem: Word model can see new words, that's not in the training model.

Solution: Use unknown token <UNK> which is replace with unknown word.

Then calculate probability.

Ex:SentenceModel

I love pizza

I love pizza

I love xyzburgers

I love <UNK>

→ gives perfect result in Trigram.

→ gives perfect result in Trigram.

Perplexity : Trigram > bigram > unigramText classification

Text classification is the process of teaching a computer to "automatically" sort a piece of text into a category like spam, positive/negative or news.

Example : Spam Detection

⇒ split emails into two folders

→ Spam folder → all spam example

→ Inbox (Ham) → all real emails

Subject:.....

Date:.....

⇒ Clues for Spam

<u>Spam</u>	<u>Ham</u>
"Buy Now", "Cheap", "Viagra", "Meeting at 3 pm"	
All uppercase (ALL CAPS)	Normally written
Weirdly written (you'd-serve)	Normal words

⇒ Computer build two model

$$P(\text{Message} | \text{spam}), P(\text{Message} | \text{ham})$$

use bayes rule,

$$\underset{c \in \{\text{spam}, \text{ham}\}}{\operatorname{argmax}} P(c | \text{message}) = \underset{c \in \{\text{spam}, \text{ham}\}}{\operatorname{argmax}} P(\text{message} | c) P(c),$$

⇒ Counting total ham and spam message

⇒ result

Information Retrieval (IR)

Finding documents that are relevant to a user's information need.

Ex: In google — [AI Book]

Components

i) Corpus of documents: The collection of documents the system search.

Ex: a paragraph, a page, text, file

ii) Query in a query language: specifies what the user wants.

Ex: [AI book] — keyword

iii) Result set: System thinks which documents is suitable for search.

iv) Presentation: How results shown. list or visualization.

IR Scoring Functions

- a methods used by search engines to rank documents based on how relevant they are to a user's query.
- BM25 is one of the best.

BM25: a formula that gives each search page a score.

High score → more relevant

BM25 three factors

i) TF → How many times a word appears

Ex: If a query is "farming"

Page A: 10 times ✓

Page B: 1 time ✗

Inverse
document
frequency

ii) IDF → How rare a word is!

Ex:

"in" → appears everywhere → low importance

"farmers" → rare → high importance

iii) length documents →

long Page may mention all words by chance → low score

short, focused Pages mention all query word → high score.

Subject:.....

Date:.....

$$BM25(d_j, q_1:N) = \sum_{i=1}^N IDF(q_i) \cdot \frac{TF(q_i, d_j) \cdot (k+1)}{TF(q_i, d_j) + k \cdot (1 - b + b \cdot \frac{|d_j|}{L})}$$

where,

$|d_j|$ = length of document d_j

L = avg document length in the corpus

$$L = \sum_i \frac{|d_i|}{N}$$

$$\begin{aligned} k &= 2 \\ b &= 0.75 \end{aligned} \quad \text{parameters}$$

$$IDF(q_i) = \log \frac{N - DF(q_i) + 0.5}{DF(q_i) + 0.5}$$

Subject:.....

Date:.....

IR System evolution

→ 1970 - 1980 : 2nd gen.

→ 1980 - 1990 : 3rd gen.

→ 1990 - 2000 : 4th gen.

→ 2000 - 2010 : 5th gen.

→ 2010 - 2020 : 6th gen.

→ 2020 - 2030 : 7th gen.

→ 2030 - 2040 : 8th gen.

→ 2040 - 2050 : 9th gen.

b. ARCHITECTURE

a. HUB

→ 1970 - 1980 : HUB

→ 1980 - 1990 : HUB

→ 1990 - 2000 : HUB

→ 2000 - 2010 : HUB

→ 2010 - 2020 : NORMALIZE (proto)

→ 2020 - 2030 : NORMALIZE

What is PageRank?The Page rank algorithm

PageRank is a Google's algorithm that decides which web page is most important.

The HITS algorithm —

function HITS(query) returns pages with hub and authority numbers

1 Pages \leftarrow EXPAND-PAGES(REVENT-PAGES(query))

for each p in Pages do

 P. AUTHORITY \leftarrow 1

 P. HUB \leftarrow 1

repeat until convergence do

 for each p in Pages do

 P. AUTHORITY $\leftarrow \sum_i \text{INLINK}(p). HUB$

 P. HUB $\leftarrow \sum_i \text{OUTLINK}(p). AUTHORITY$

 NORMALIZE(Pages)

return Pages

The PageRank for a page P is defined as:

$$PR(P) = \frac{1-d}{N} + d \sum \frac{PR(in_i)}{C(in_i)}$$

where,

N = Total Pages

d = damping factor (≈ 0.85)

i = Page linking to page P

$C(i)$ = number of outgoing links on page i