

## \* UNIT-5 - Uncertain knowledge:

### → Discourse Processing :

Discourse processing is a subfield of NLP that focuses on understanding the structure & meaning of text at beyond level of individual sentences,

→ The main aim is to extract meaningful information from the text which can be useful for summarization, information extraction, question answering ....

→ It involves 4 sub tasks,

i) coherence & cohesion analysis

ii) Discourse structure analysis

iii) Rhetorical Parsing

iv) Discourse - level sentiment analysis.

## i) coherence & cohesion analysis:

→ coherence refers to overall logical connection b/w sentences in Text.

→ cohesion refers to grammatical connections b/w sentences (such as: Pronouns, conjunctions etc)

e.g. John went to store. He bought some milk

He in 2<sup>nd</sup> sentence refers to John in 1<sup>st</sup> sentence.

## ii) Discourse structure analysis:

→ It involves in analyzing the hierarchical structure in the Text, identifying main Topics, supporting details, & relations

e.g. 1<sup>st</sup>: we'll discuss Problem, Then we'll

find sol<sup>n</sup>, By Discussion only we can find sol<sup>n</sup>.

Structure: Introduction: First, we'll discuss prob.  
main point: we'll find sol<sup>n</sup>  
supporting detail: By discussion only  
we can find sol<sup>n</sup>

### iii) Rhetorical Parsing:

↳ (To inform)

→ Rhetorical Parsing involves in analyzing  
structure of Text to understand  
how different parts of a Text are  
related.

Eg: "she was clever because she studys  
everyday".

because indicating its cause & effect  
of sentence.

### iv) Discourse - level sentiment Analysis:

Involves in analyzing the  
sentiment of a Text.

Eg: The battery life of mobile is great  
(Positive)

but camera quality is bad  
(negative)

## \* Key components of Discourse

processing:

- 1) cohesion
- 2) Resolution Reference
- 3) discourse cohesion & structure

↳ cohesion:

→ the grammatical connections  
between sentences

the sentences  
(or)  
features that link different  
parts of text.

→ cohesion includes:

- a) Reference words: These are like pronouns / articles which refers to entities mentioned in initial stage of text.

eg: Lucky went to bookstore, He bought a book

he in 2<sup>nd</sup> sentence refers to lucky in 1<sup>st</sup> sentence

b) conjunctions: words that connect sentences  
like: and, but, so . . .

Eg: I want to go outside, but it's raining  
but is conjunction

c) Ellipsis: words which are understood  
as necessary to repeat

Eg: "John reading a book, and  
Mary is too"

ellipse mary also reading book

d) lexical cohesion:  
creating a meaning using  
similar words

Eg: "The restaurant had a  
tasty biryani."

The dish was like by many people

Dish refers → Biryani

## a) Reference Resolution:

It is the process of determining which words in text refers to same entity.

a) Pronoun = Identify to which entity the Pronoun is referring.

Eg: After John finished work, he went outside.

he is a Pronoun referring → John.

b) Noun Phrase: Identify to which entity

the Noun is referring

Eg: The book is on table, I

read it yesterday.

Here, book & it refers same object.

c) Anaphora = Using pronoun to

refer previously mentioned entity

Eg: I bought a new laptop.

It has high resolution.

### 3) Discourse cohesion & structure:

Involves in understanding how different parts of Text related to form a sentence.

a) Discourse connectives: words which gives hints the relationships.

e.g.: "The government should education funds because increasing education is key to economic growth → {connective}

b) Discourse markers: words which helps to organise words.

e.g.: "Firstly", we need to complete syllable markers in "on the other" hand, we'll prepare for exam.

c) segmentation: dividing texts into segments.

e.g.: The cat chased the mouse. The mouse ran away.

## \* PART - 2 \*

### Overview :-

#### 1) Introduction to Language Modelling

- overview
- Applications

#### 2) Types of Language models :-

##### a) statistical language model

→ n-grams

##### b) Neural Language Models

- feedforward neural networks
- recurrent neural network
- Transformers (BERT, GPT)

#### 3) Language Model Evaluation:

- perplexity
- Accuracy

#### 4) Parameter Estimation:

- Max likelihood
- Bayesian parameter estimation.
- large-scale language model

5) language model Adaptation

6) language specific modelling

→ word order

→ character - Based models.

7) multilingual & crosslingual modelling

## ix. Introduction :

Language modeling is a crucial task in NLP involves in predicting next word in a sequence.

### Application:

Text Generation

Speech Recognition

Machine Translation.

Eg: "The cat sat on the \_\_\_\_\_".  
language model will predict next word as 'mat'.

### → Types of Model:

• statistical lang model:  
These model uses statistical probability.  
methods to estimate sequence.  
of a word

## \* N-Gram:

N gram can be defined as contiguous sequence of n items from a given text.

→ The items can be letters, words, (or) pairs

→ N grams can be; ~~to~~ times, let's say  
unigram - (considers each word as independent)  
bigram - (considers prob. of previous word)  
Trigram - (considers prob. of previous two words)

N Gram modeling:  
predicts the probability of given N-gram with sequence of previous words.  
e.g.: "the cat sat on the

- gt contains 5 words  
"The", "cat", "sat", "on", "the"  
→ To find the 6<sup>th</sup> word, find  
conditional Probability of first 5 words

$$P\left(\frac{w_6}{w_1, w_2, w_3, \dots}\right)$$

→ gt can be calculated using chain rule  
Hence, we can find 6<sup>th</sup> word.

$$P(A/B) = \frac{P(AB)}{P(B)}$$

Example:

$s_1 = "g\ am\ Happy"$

$s_2 = "g\ am\ sad"$

Predict the likelihood of  $s_1$  with  $s_2$ .

Step-1:

Let's take Bigram (2 word/pair)

$s_1 \rightarrow "g\ am"; "am\ Happy"$

$s_2 \rightarrow "g\ am"; "am\ sad"$

Step-2: count Frequency

"g am" → appeared 2 times

"am happy" → 1

Step 3: calc Probability

$$P\left(\frac{\text{am}}{g}\right) = \frac{\text{count(gam)}}{\text{count(g)}} = \frac{2}{2} = 1$$

$$P\left(\frac{\text{happy}}{\text{am}}\right) = \frac{\text{count(happy)}}{\text{count(am)}} = \frac{1}{2} = 0.5$$

$$P\left(\frac{\text{sad}}{\text{am}}\right) = \frac{\text{count(sad)}}{\text{count(am)}} = \frac{1}{2} = 0.5$$

Step - 4: predict Prob

$$P(\text{"gam happy"}) = P\left(\frac{\text{am}}{g}\right) \times P\left(\frac{\text{happy}}{\text{am}}\right)$$
$$= 1.0 \times 0.5$$
$$= 0.5$$

gm happy has 0.5 probability.

## \* Neural Language models:

### a) feed forward:

uses a simple feedforward neural network where each word is represented by embeddings.  
↓ (vector representation)

#### Architecture:

Input : Fixed no. of Previous words given

Process : word embedding given to Neural network

output : network Predicts Probability of next word.

Eg: Input: Sun rises in \_\_\_\_\_

Process : each word represented as vector

Prediction: The neural network outputs the probability of next word.

Eg: east, west, North, south.

If east has higher probability

## b) Recurrent Neural Network

designed to handle the sequences, it maintains hidden layer (state); updated for every word sequence.

Eg: I love my —

step - 1 ) "I" → update hidden state

2) "love" → update hidden state

3) "my" → update hidden state

it uses final state (hidden state)

to predict next word it might be

~~dog~~, [dog, cat]

c>

c) Transformer based (BERT, GPT);

BERT; (Bidirectional Encoder Representation from Transformer)

it uses transformer to understand meaning of words from both

~~(\*)~~ It mask some words  
wantedly (mask)

Eg: "The [mask] barked at cat".

It predicts masked word as "dog"

\*GPT: (Generative Pretrained Transformer)

It generates text by predicting next word in a sequence.

→ It takes the input prompt

and generates output.

Eg: "The weather today is"

It takes the prompt and predicts "The weather today is warm".

## \* Language model    Evaluation \*

1) Perplexity: used to measure how well a language model is predicting sequence of words.

low perplexity = better prediction of next word.

Eg: "The cat sat on the mat"

### model A

$$P(\text{The}) = 0.4$$

$$P(\text{cat}|\text{The}) = 0.2$$

$$P(\text{sat}|\text{cat}) = 0.2$$

$$P\left(\frac{\text{on}}{\text{cat sat}}\right) = 0.35$$

$$P\left(\frac{\text{mat}}{\text{cat sat on}}\right) = 0.1$$

### model B

$$P(\text{The}) = 0.5$$

$$P(\text{cat}|\text{The}) = 0.4$$

$$P(\text{sat}|\text{cat}) = 0.3$$

$$P\left(\frac{\text{on}}{\text{cat sat}}\right) = 0.1$$

$$P\left(\frac{\text{mat}}{\text{cat sat on}}\right) = 0.2$$

$$\text{Perplexity} = PPL(w) = \left[ \frac{1}{P(w)} \right]^{\frac{1}{N}}$$

model A

$$PPL(w) = \frac{1}{0.6 \times 0.2 \times 0.2 \times 0.3 \times 0.1}$$

$$= 43.84$$

model B

$$\frac{1}{0.5 \times 0.4 \times 0.3 \times 0.1 \times 0.2} = 16.67$$

↓  
low perplexity  
mean better prediction

⇒ Accuracy: measures % of correct prediction made by a model

$$\text{Accuracy} = \frac{\text{correct Prediction}}{\text{Total Prediction}}$$

eg: "g love my Dog"

Predictions; after "g" =

after g → love ✓

after g → my ✓

after g → cat X

After g → (love, my) ✓

Accuracy

$$= \frac{2}{3} \times 100$$

$$= 66 \%$$

## \* Parameter Estimations \*

1) max likelihood Estimation: used to estimate parameters of a model by likelihood function.

$$\text{max likelihood} = \arg \max P(\frac{D}{\theta}) \rightarrow \hat{\theta}_{MLE} = \arg \max P(\frac{D}{\theta})$$

eg: Training Data

g love dogs ;	g like dogs
g love cats ;	g like cats
g love coding ;	g like coding

To find Probability of next word after:

$$P(\text{love}) = \frac{\text{count(love)}}{\text{Total no. of words}}$$

$$= \frac{3}{8} = 0.375$$

50% of chance is love to appear in given context.

## 2) Bayesian Parameter Estimation

It compares prior belief on parameters & update these belief on data.

$$P\left(\frac{D}{O}\right) = \frac{P\left(\frac{D}{O}\right) \cdot P(O)}{P(D)}$$

Eg: 1) Flipping coin

$\alpha \rightarrow$  Head      } and  
 $\beta \rightarrow$  Tail      } Prior belief is  
                         $\alpha = 2 ; \beta = 2$

Now, After flipping 10 times and observes 7 heads; 3 tails

$$3) \text{ Update: } \alpha = \alpha_{\text{prior}} + \text{no.of H} = (2+7) \\ = 9$$

$$\beta = \beta_{\text{prior}} + \text{no.of T} = (2+3) \\ = 5$$

Distribution (9,5)

$$\text{Prob of head} = \frac{\alpha}{\alpha+\beta}$$

$$= \frac{9}{9+5} = 0.6$$

### 3) Large-scale Language model:

It consists of large amount of data where the prediction occurs iteratively by assigning weights on previous work.

Eg: step-1: Assigning weights

step-2: model predicts if correct  
if wrong (calculate loss)

Re-assign weight - - - again

Predict.

Eg; "The sun is shining"

step-1: Initialize weights  
(random weights)

step-2: model Predicts

"The sun is dull" X

calc loss: let's assume  
loss = 0.8

Re-Assign weights,

smoothing Technique; It is also used for Probability

Add-on smoothing:

Add something to never happen

avoid 'zero' if scenario

eg: The bird barked → Not Possible  
but, probability  $> 0$  (small non zero)

\* Language model Adaptation \*

Process of adjusting pre-trai  
model to perform

-ned language specific task.

better on a → models like GPT sometimes  
proper answers

→ models like GPT sometimes

so, Adaptation makes answers

Provide better information.

Benefits of Adaptation:

- Improved performance

- Efficiency

## Techniques of Adaptation:

### 1) Fine-tuning:

Adjust the pretrained model parameters by training on a smaller domains first, then increment further.

### 2) Prompt engineering:

Providing a proper prompt to guide the model's behavior.

Eg: "summarize about NLP" can give a full summary of NLP

### 3) Few-shot Learning:

Training on a small amount of labeled datasets.

### 4) zero shot learning:

using pretrained model without any additional training data

## \* Language - specific modelling

It is a fundamental task in NLP involves predicting next word / sequence of words in a text by word order / character.

### i) Word order:

Different languages have different word orders.

→ English have fixed word order.  
German have flexible word order.

e.g. The dog, The cat (english)  
Der Hund Der Kastze (german)

### ii) Character - Based:

Breaking the words into smaller units (characters / syllables).

e.g. "running" → "r" "u" "n" "n" "i"  
"n" "g"

## \* Multi Lingual

## Cross Lingual

### → Multi Lingual:

Focuses on building models that generates text in multiple languages.

→ used for machine translation,

Text summarization

→ eg; Translating text from English to Spanish.

### → Cross Lingual:

This involves in training model in one language to perform task in another language.

→ The trained data will be limited

& The language where it performs will be high level

low trained data → high language