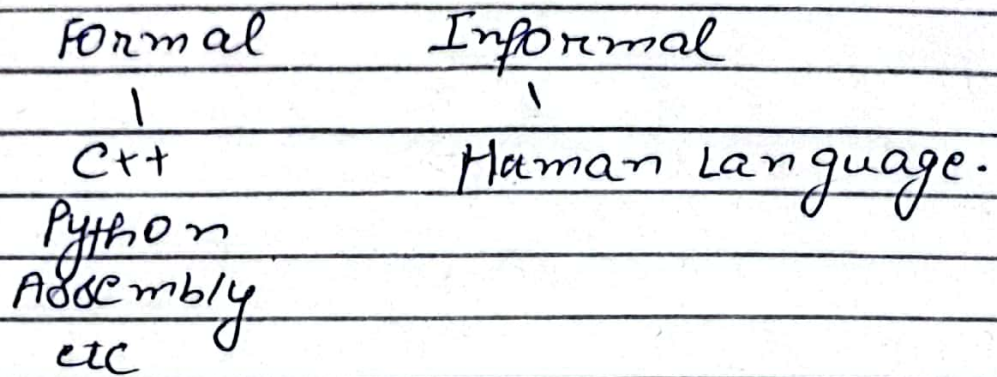


Lecture NO. 1

Language



• NLP Processing:

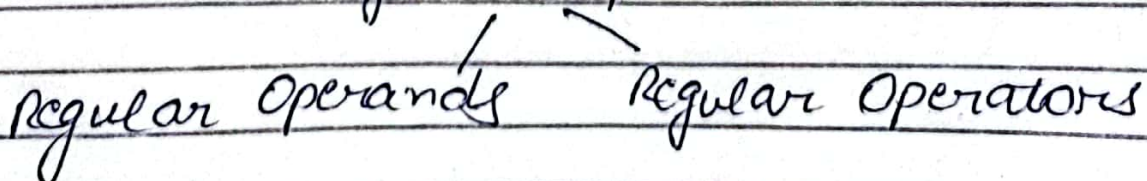
- Acquisition (Sources: Speech, Text etc)
- Representation (meaningful)
- Understanding
- Tokenization etc.

• Applications of NLP

- Machine Translation
- Sentiment Analysis
- Language Technology

- Regular expressions are used when there are infinite / unlimited sets. Limited sets are saved.

Regular Expressions



09029955

Lecture NO. 2

Regular Expression → used in string matching, finding, tokenization etc.

→ Operators:

? (0/1 occurrence) \rightarrow for specific character

* (O/N occurrences)

+ ($1/N$ occurrences)

- (for any character)

* [] (start of line)

\$ (end of line)

↳ (character 'o')

- (for any character)

[I:] ET

I:

[工][四]

$$[I, T]$$

42

-) decrease in false positive, increase in precision

→ Recall

- FP (Type I error)

set of unique elements

→ 'V' → vocabulary (unique)

→ tokens → all occurrences (redundant + unique)

→ no. of elements / size

1) $|V|$ → cardinality of set
 N → Total occurrences

→ Byte Pair Encoding → learning based
 (segmentation of words)

	Frequency	Train	Test
low_	5		new_ ✓
lowest_	2		newer_ ✓
newer_	6	a	lower_
wider_	3		low er_
			2 tokens

Vocabulary: (characters)

o, w, l, n, er, er-, ...

1) not ~~at~~ the best approach but better than tokenization based on space.

→ extrinsic model → need some real life experimentation. needs a classifier

→ ~~is~~ intrinsic model → do not need any classifier.

Lecture NO. 3+4

(grammar) Language Models

- Machine Translation
- Spell Correction
- Speech Recognition

→ use probabilities (conditional)

→ Markov Assumption

(Look at a local window)

cannot increase the

size so much

can limit the size of context

$$LW = \{0, 1, 2, 3, \dots\}$$

→ Unigram is not a good model as there is a possibility of it giving a word without any knowledge of context.

$$P(\text{water}) = \frac{C(\text{water})}{N}$$

→ Perplexity: \rightarrow intrinsic model

$$\left(\frac{P(w_1, w_2, w_3, w_4, \dots, w_n)}{P(w_1, w_2, w_3, \dots, w_n)} \right)^{\frac{1}{N}}$$

$$= \sqrt[N]{\frac{1}{P(w_1, w_2, w_3, \dots, w_n)}} = \sqrt[N]{\frac{1}{\prod P(w_i | w_{1:n-1})}}$$

product of

unaware of context
Activity-2

unigram:

$P(\text{I want to eat Turkish Pizza}) =$

$$P(I) \times P(\text{want}) \times P(\text{to}) \times P(\text{eat}) \times P(\text{Turkish}) \times P(\text{Pizza})$$

$$= \frac{7}{44} \times \frac{7}{44} \times \frac{7}{44} \times \frac{6}{44} \times \frac{1}{44} \times \frac{3}{44}$$

$$= 0.15 \times 0.15 \times 0.15 \times 0.13 \times 0.02 \times 0.06$$

$$= 5.2 \times 10^{-7} \text{ } 0.00000053$$

$P(\text{I want to eat Indian food}) =$

$$P(I) \times P(\text{want}) \times P(\text{to}) \times P(\text{eat}) \times P(\text{Indian}) \times P(\text{food})$$

$$= \frac{7}{44} \times \frac{7}{44} \times \frac{7}{44} \times \frac{6}{44} \times \frac{10}{44} \times \frac{4}{44}$$

$$= \frac{7}{44} \times 0.15 \times 0.15 \times 0.13 \times 0.23 \times 0.09$$

$$= 0.00000078$$

$$= \log \frac{7}{44} + \log \frac{7}{44} + \log \frac{7}{44} + \log \frac{6}{44} + \log \frac{10}{44} +$$

$$\log \frac{0}{44} + \log \frac{4}{44}$$

$$= -4.31$$

Bigram:

$$P(\text{I want to eat Turkish food}) =$$

$$P(\text{I} | \langle \rangle) \times P(\text{want} | \text{I}) \times P(\text{to} | \text{want}) \times \\ P(\text{eat} | \text{to}) \times P(\text{Turkish} | \text{eat}) \times P(\text{food} | \\ \text{Turkish}) \times P(\langle \rangle | \text{food})$$

- smaller the perplexity, higher the confidence of model
- Perplexity:
 - If 3 words have same probabilities then the perplexity will also be 3
 - If 10 words have same probabilities then the perplexity will be 10.

Ex:
$$3 \sqrt[3]{\frac{1}{0.7 \times 0.15 \times 0.15}} = 0.85$$

→ 0.7 → 0.15 → 0.15
 $w_1 \quad w_2 \quad w_3$

The perplexity will be less than the number of words if the probabilities are not same.

Ex:
$$3 \sqrt[3]{\frac{1}{1/3 \times 1/3 \times 1/3}} = 3 \sqrt[3]{\frac{1}{1/27}} = 3$$

- Smaller the perplexity, better the model but sometimes also depends on vocabulary.

lecture NO. 5

→ closed vocabulary:

vocab that is already defined, no more new words can be added

→ After applying \log when the prob is 0 of any word then the answer is in \log space and might be negative

$$\log(P_1 \times P_2) = \log_2 P_1 + \log_2 P_2 = x$$

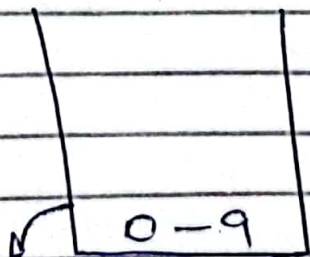
(to convert to probability 2^x space again)

we convert answer into \log space to avoid underflow.

→ Handle the context properly & efficiently.

Ex:

Perplexity
Train



prob of occurrence of each word = 0.1

$$PP = \frac{1}{10}$$

$$PP = 10$$

Test

0 1 2
3 4 6

any digit as the prob of every digit is same.

OOV \rightarrow out of vocabulary

Ex \rightarrow PP = 3 \rightarrow confused b/w 3 words.

Train		Test	
0.90 \rightarrow	0000	EX1	000
0.05 \rightarrow	1000	EX2	0011 22
0.05 \rightarrow	2000		
	1000		

max prob that it will predict 0
PP = 3
confused b/w all words

\rightarrow Laplacian Smoothing:

Disadvantages:

- \rightarrow increase in perplexity
- \rightarrow possibility that the probability of words that were not in the corpus may be greater than the prob of words that already exists in training set/corpus.

$$\frac{C(w_{i-1}, w_i) + 1}{C(w_{i-1}) + |V|}$$

\rightarrow Add k -smoothing (variant of Laplace)
(not a good approach)

$$\frac{C(w_{i-1}, w_i) + k}{C(w_{i-1}) + |V| \times k}$$

~~$C(w_{i-1}, w_i)$~~

→ Linear Interpolation:

min of unigram, bigram and trigram.

Training Data	Held Out Data	Test Data
---------------	---------------	-----------

→ Stupid Backoff
(only good for web)

- 1) trigram → available (pick that)
- 2) if trigram not available go to bigram
- 3) if bigram not available go to unigram.

→ Good Turning Smoothing:

Ex:

10 carp, 3 perch, 2 white fish,
1 trout, 1 salmon, 1 eel = 18 fish

rare

$$N_i = \frac{3}{18} \quad \left(\begin{array}{c} \text{add} \\ \text{rare} \end{array} \right)$$

$N = \text{no. of fishes (classes)}$

We will relate unseen data using rare data's frequency / probability.

$$c' = \frac{(c+1) N_{i+1}}{N_c}$$

$$C^*(\text{trout}) = 2 \times N_2 / N_1$$

$$= 2 \times 1/3$$

$$= 2/3$$

$$P_{GT}^*(\text{trout}) = C^* / N = 2/3 / 18 = 1/27$$

Ex2:

10 carp, 3 perch, 2 whitefish,
2 trout, 2 salmon, 2 eel = 21 fish

$$N_0 = 0$$

$$N_1 = 0$$

$$N_2 = 4$$

$$N_3 = 1$$

$$18/21$$

$$21$$

$$C = 2$$

$$C^*(\text{trout}) = 3 \times \frac{4}{21} = \frac{12}{21} = \frac{4}{7} = 3 \times \frac{1}{7}$$

$$P^*(\text{trout}) =$$

$$P_{GT}^* = \frac{N_2}{N} = \frac{4}{21}$$

Prob of
unseen

$$C^*(\text{trout}) = (2+1) \times \frac{N_3}{N_2}$$

$$= 3$$

$$P^*(\text{trout}) = 3/4 \times \frac{1}{21} = \frac{1}{28}$$

Regular Expressions

Morphology

- Stemming
 - Lemmatization
- } used for text normalization

→ chop off length of vocabulary and
Stemming: computational cost.

maps all variations of 1 words
on a stem

→ Stem: core meaning bearing
morphemes

EX: cat

Stem: cat

affix = s

- used in search engines
- used for info retrieval

→ In machine translation, lemmatization is better approach but increase^{on} computational cost.

11/10/00

-) Probabilistic models

parameter $\rightarrow \checkmark$

$V \rightarrow C(\text{class})$

The diagram illustrates a neural network architecture for word prediction. It consists of two input layers, C1 and C2, and an output layer. Layer C1 is labeled 'C1 → +ive' and contains a 'word prob' section and three dashed lines representing hidden units. Layer C2 is labeled 'C2 → -ive' and also contains a 'word prob' section and three dashed lines. Arrows indicate connections from C1 to C2 and from C2 to the output layer.

- "dictionary" → used to store it

Naive Bayes

First, we

-) have to convert the corpus in "bag of words"

•) $\arg\max [P(c_1|d), P(c_2|d)]$

will return the class with highest probabilities. We can add many values (classes) in it.

Activity-1

Prior Probabilities =

$$c = \frac{3}{4}$$

$$j = \frac{1}{4}$$

$$c = 0.75$$

$$j = 0.25$$

Vocabulary = Chinese, Beijing, Shanghai,
Macao, Tokyo, Japan

Vocab size = 6

$$MD_c = d_1 + d_2 + d_3$$

$$MD_c(\text{Vocab}) = 5$$

$$MD_j = d_4$$

$$MD_j(\text{Vocab}) = 3.$$

$$P(\text{China} | c) = \frac{5}{8}$$

8 \rightarrow words in class c

$$\{c: \{\text{China}: 5/8, \dots\} \dots \{j: \dots\}\}$$