

[ ] = disjunction operator;  $\vee$

CHI ĐOẠN

PAGE NO. : \_\_\_\_\_

£8,830 | £4,530 2D 101150 95 Nov. 22, 200

Alphabet = {a-z, A-Z, space, 0-9, punctuations, emoji,  
descriptive ununicode blocks}

tokens  $\Rightarrow$  the Constituent Chars of a token is a subset of alphabet;

LLMs design context-sensitive languages very well. It learns contextual-embeddings.

[P] So P 2<sup>2</sup>g pattern & find 8221

$\wedge$  = negation + so  $\neg(a \rightarrow f, o \rightarrow y)$  not in  $a \rightarrow f, o \rightarrow y$

$\wedge [\wedge_{A-F0-4}] [A-Z]$

if  $a \in A$ ,  $b \in B$   $\exists x [P(x) \rightarrow x \in b]$

7A.  $\frac{1}{\pi}$  ✓

$$TABC \quad \checkmark \quad \text{programm} \quad \text{ablauf} = W /$$

$\wedge [{}^{\wedge}a-f0-4] [A-2]$  \$

So  $\text{min}_1 \text{ min}_2$  choose string accept थके.

$$A = 0 \text{ } 00 \text{ mode}$$

+ 1 = 1 02.00

$\exists x \in A$  such that  $x^2 = 9$ . So,  $A \neq \emptyset$ .

$\{n\} = \text{repeat } n \text{ times}$

$$g^{-1} = g^{-1} \text{ mod } b = [ ]$$

PAGE NO.: \_\_\_\_\_  
DATE: / /

$\alpha f_2, s_3$  can be written as  $\alpha f_2, 23 \mid 0 f_3, 33 \mid$   
 $\alpha f_4, 43 \mid \alpha f_5, 53$

Q. Write a pattern to identify the middle character in a string of length 7 where the string should start and end with a digit. The starting digit can not be zero. The middle character should be one of a, b, c, f, g, -, +.

\w = alphanumeric characters set \w+ = A

$d = \text{digit} [0-9] \times m \geq d$ , A.L. 11

\W = non-alpha numeric

\\$ = whitespace + [\$\_-A] [\$\_-C] \wedge \dots \wedge [\$\_-Z]

$\setminus S = \text{Non-white space}$

\\$ = Non-white space

Q  $x + j$  am going home  $\leftarrow$  This has 5 boundaries  
 $(S^+, x)$

$$\text{Ans} = T_{\text{end}}$$

\b = boundary

[.] is any character.

$$\frac{0.12}{P_{t+1}} \rightarrow P_t + \frac{1}{P_{t+1}}$$

PAGE NO. : \_\_\_\_\_

eg: 2 ad 3 d + v / d +

eg:  $1,23,732 \Rightarrow d+, d+, d+$   
 So to remove redundancy we use groups.

$\Rightarrow$  Order of precedence in RE

- groups
- wildord +, \*, ?,
- [ ]
- OR and laddr

Nesting of groups is also possible  $((\text{d}+)(),)+\text{d}+$

Q RE for URL

②  $(\wedge_1 [\wedge_{0+2} A-2 \cdot 0-9]) [tT]he (\wedge_{0-2} A-2) \stackrel{0+9}{\$}$

empty string or not, "the/THE" is not alpha numeric.  
(alpha numeric is digit) now = or end empty string

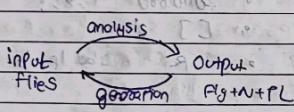
Q  $((\bar{a}-zA-z))(\bar{a}-zA-z)^+ \rightarrow 1 2$

where  $((\bar{a}-zA-z)) = 1$

$(\bar{a}-zA-z)^+ = 2 \leftarrow \text{ss} \rightarrow \text{es}$

Finite State machine can only recognize but finite state  
Transducers can generate + recognize (+1/1)

⇒ There are 2 types of tasks:  
Analysis (Recognition)  
Generation



Do pattern based decoding

Now - create a list "ending -y", "work -ing -s", etc.  
These are called Paradigms.

We store all the words in a table

work, works, walking, walked

root = walk (this is called Lemma)

Stemming = process of reducing words to their root or base form, often by removing prefixes or suffixes.

y → ies

ss → es

Morphology

inflection

Derivation

adding suffix, category (nouns etc.)  
remains the same

adding suffix, category  
changes

e.g. walk, working

Verb = Realize

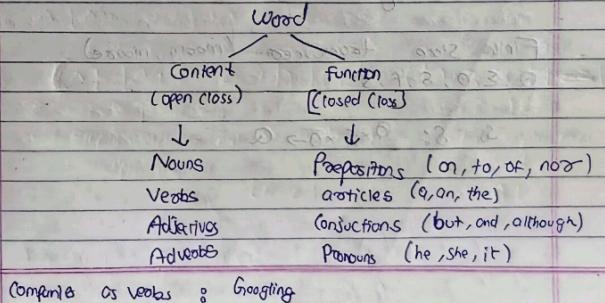
Name = Realization

posted Stemmer

⇒ Morphology is the study of word structure and how words are formed from smaller meaningful units called morphemes

Lemmatization: Converting words to their base or dictionary form (lemma) considering the context

"better" → "good"



Compound vs Verbs : Grouping

Q) Show talk boat + boatbook  
 $S_0 \xrightarrow{b} (q_1) \xrightarrow{a} (q_2) \xrightarrow{a} (q_3) \xrightarrow{a} (q_4)$   
 $\Sigma = \{b, a, !\}$   $\Sigma^+ = \text{closed property}$   
 $q_1 \xrightarrow{a} q_2$   $q_1 \xrightarrow{!} \emptyset$   
 $q_2 \xrightarrow{a} \emptyset$   
 $q_3 \xrightarrow{a} \emptyset$   
 $q_4 \xrightarrow{a} \emptyset$   
 Given a dictionary, can you design an FSA/FSM  
 if length of max word = n chars.  
 Then  $26^n$  is the size of transition table.  
 So we can use dictionary and prune the words not in dictionary.

Finite State Transducers (mealy, moore)

PAGE NO.: Signed  
DATE: 11/14/2023

Sub ( $\wedge \vee \exists \forall \exists'$ ,  $\dots$ , text)

cat, cats

$\rightarrow \text{c} \rightarrow \text{a} \rightarrow \text{t}$  (from  $\exists$ )

box, boxes

$\rightarrow \text{b} \rightarrow \text{o} \rightarrow \text{x}$  (from  $\exists$ )

for this also we need transition x-table

$\Rightarrow$  Stop words are those words which are very frequent but they do not have any "meaning" value.

We generally remove stop words

Zippf's law

freq

rank

threshold set

Normalisation / Standardization techniques:

- 1) Stopword removal
- 2) Stemming
- 3) Lemmatization
- 4) Spelling / text normalization

5) tokenization

6) Replace rare words with special symbols

To handle noisy text we do expanding

I'd go  $\Rightarrow$  I could go

Won't  $\Rightarrow$  would not

### N-Gram Language Modelling

$$P(x,y,z) = P(x) \cdot P(y|x) \cdot P(z|xy)$$

$$P(x,y,z) = P(x) \cdot P(y|x) \cdot P(z|xy)$$

$$P(w_1, w_2, \dots, w_n) = P(w_1) \cdot \prod_{i=2}^n P(w_i | w_1, \dots, w_{i-1})$$

\* N-Gram = N words in sequence

The Cat Sat on the mat.

One-gram = each word, punctuation is considered.

two-gram = 2 words or 1 consider 82

so 5

longest n-gram?

gramsize?

n-gram size?

So in trigram  $P(w_n | w_{n-2}, w_{n-1}) = \frac{C(w_{n-2}, w_{n-1}, w_n)}{C(w_{n-2}, w_{n-1})}$

If the  $w_n$  is a very new word, count 0 also Prob. 0

Smoothing.

- Add one smoothing (Laplace smoothing).  
Count of new word + 1 adds 1 to  $C(w_n)$ .

$$P(w_n) = \frac{C(w_n)}{N} \quad \text{and} \quad P(w_i | w_{i-1}) = \frac{C(w_{i-1}, w_i)}{C(w_{i-1}, v)}$$

$$= \frac{C(w_i) + 1}{N + |V|}$$

Good assumption for  $|V| = 10^4$ ,  $|V|$  is number of unique words in language i.e. vocabulary,  $N =$  dictionary size. Thus count of unique words.

problem is that the probability of those words which are already frequent that also decreases.

$$P(tot) = 0.3 = \frac{3}{100}$$

$$\text{Now } P(tot) = \frac{30 + 1}{100 + 10000} = 0.003$$

So solution is Add k Smoothing

$$P(w_n | w_{n-1}) = \frac{C(w_{n-1}, w_n) + k}{C(w_{n-1}) + k}$$

3rd type.

(odd types \*)

$$P(w_n | w_{n-1}) = \frac{C(w_{n-1}, w_n) + k}{C(w_{n-1}) + k * T(w_{n-1})}$$

here  $T(w_{n-1}) = \text{no. of unique bigrams starting with } w_{n-1}$

Now logic is bigrams etc depends on continuity & given one word to predict so logic maintains same till  $T(w_{n-1})$ .

Now to make it a probability  
smoothed count  $\leftarrow C_i^* = \frac{(C_i + 1) * N}{N + 1}$

$$\text{So } P(C_i^*) = \frac{C_i^*}{N} = \frac{C_i + 1}{N + 1}$$

$$d = \frac{C_i^*}{C_i}$$

Types = No. of unique words

PAGE NO.: / / DATE: / /

Written + Bell Smoothing

$$\sum_{w \in C(w)=0} p_i = \frac{T}{N+T} \quad \text{↳ currently unique words in text}$$

$$Z = \text{Count of all zero frequency words} = \sum_{w \in C(w)=0} 1$$

$$Z \times P_{unseen} = \frac{T}{N+T} = \frac{(N+1) - N}{N+T}$$

$$P_{unseen} = \frac{1}{N+T} \quad ; \quad Z = N-T \quad \text{for unigram}$$

$$P_{unseen} = \frac{C(w)}{N+T}$$

For bigram.

$$P_{unseen} = \frac{1}{Z(w_x)} \times \frac{T(w_x)}{N(w_x) + T(w_x)} \quad \text{which are unique}$$

↳ Total bigrams starting with w<sub>x</sub>.

$$Z(w_x) = N - T(w_x)$$

$\frac{N}{N_0}$

All these are Maximum Likelihood Estimation

PAGE NO.: / / DATE: / /

Good turing Smoothing

$$N_0 = ?$$

$$N_1 = \text{no. of words with 1 freq}$$

$$N_2 \text{ depends on } N_1, N_3 \text{ depends on } N_2$$

$$N_0^* = (\sigma+1) \cdot \frac{N(\sigma+1)}{N\sigma}$$

for unseen n-grams

$$P^*(nT) = \frac{N_1}{N} \rightarrow n\text{-grams of freq. 1}$$

N → Total no. of n-grams

Let there be a bowl and f(red) = 4, f(black) = 5, f(green) = 3  
 f(orange) = f(violet) = 2, f(yellow) = f(blue) = 1

$$\therefore N_1 = 2 \quad N_2 = 2$$

$$N_3 = N_4 = \dots = N_8 = 1$$

lets say Cyan, P(Cyan) = 0 = 0

18

so apply good turing smoothing

$$P_i = \frac{N_1}{N} = \frac{2}{18}$$

now yellow's count will be updated.

$$\sigma^* = (\sigma+1) \cdot \frac{N_2}{N\sigma+1}$$

$$\pi_1^* = \frac{(1+1) N_2}{N_1} \\ = 2 \cdot \frac{2}{2} = 2$$

$$\pi_2^* = \frac{(2+1) N_3}{N_2} \\ = 3 \cdot \frac{1}{2} = \frac{3}{2}$$

if  $\pi_1^*$  is last item of  $\pi$  set  $\pi_1^*$  as next

In bigram and Trigram we won't have all  $N_i$   
so use next available higher frequency

#### Backoff Models (interpolation)

$\Rightarrow$  Stupid backoff if  $C(w_{n-2}, w_{n-1}, w_n) = 0$   
then fall back to bigram model

if  $C(w_{n-2}, w_{n-1}) = 0$ , fall back to unigram model.

\* in fallbacks don't use smoothed models.

if nothing is there put  $1/N$ .

$$f(x) = \frac{m}{m+n} f_1(x) + \frac{n}{m+n} f_2(x) \\ m x_2 + n x_1, \quad m y_2 + n y_1 \\ \text{this is interpolation}$$

PAGE NO. :  
DATE : / /

PAGE NO. :  
DATE : / /

$$\Rightarrow P(w_i | w_{i-2}, w_{i-1}) = \lambda_1 P(w_i | w_{i-2}, w_{i-1}) \\ + \lambda_2 P(w_i | w_{i-1}) + \lambda_3 P(w_i)$$

where  $\lambda_1 + \lambda_2 + \lambda_3 = 1$ ,  $\lambda_1 \neq 0$ ,  $\lambda_2 \neq 0$ ,  $\lambda_3 \neq 0$

now Grid search will give  $\lambda_1, \lambda_2, \lambda_3$

	$\lambda_1$	$\lambda_2$	$\lambda_3$
1	0.1	0.3	0.6
2	0.2	0.2	0.6

now we use these values and apply this model on dev (val) set and find errors.

$$\text{Errors} = \text{LHS} - \text{RHS}$$

$\Rightarrow P(w_i | w_{i-2}, w_{i-1})$  are how from  
trigram and this is on voices &

so what can it correlate to?  
For this we can use smooth or unsmooth models

Deleted interpolation method (unsmooth)

### Katz backoff

$$P_{KATZ}(w_i|h) = \sum d_{w_i h} \times \hat{P}_{MLE}(w_i|h) \quad \text{where } C(w_i|h) > k \\ \alpha(\text{Smoothed Katz history}) \times P_{KATZ}(w_i|\text{Smaller history}) \quad \text{when } C(w_i|h) \leq k$$

$$d_{w_i h} = \frac{\sigma^*}{\sigma} \quad \sigma^* \text{ can be found using good} \\ \text{Tuining smoothing} \quad \text{so } d_{w_i h} = (\sigma+1) \frac{N_{w_i+1}}{N_\sigma} \times \frac{1}{C(w_i|h)}$$

$$\alpha_{KATZ} = 1 - \sum_{\substack{w_i: C(w_i, w_{i-1}, w_{i-2}) > 0}} P_{KATZ}(w_i | w_{i-2}, w_{i-1}) \\ 1 - \sum_{\substack{w_i: C(w_i, w_{i-1}) > 0}} P_{KATZ}(w_i | w_{i-1})$$

so define dual discount say and  $\alpha$  as distribution in  $\leq k$  items.

$$\sigma^* = \frac{(N+1)}{N_\sigma} \quad [ \sigma = \text{observed freq. of } n \text{ grams} ] \\ N_\sigma = \text{no. of } n \text{ grams that occur exactly } \sigma \text{ times}$$

Q11

### Kneser-Ney Smoothing

$$P(w_i|w_{i-1}) = \max_{C(w_{i-1})} (C(w_{i-1}, w_i) - d, 0) + (w_{i-1}) \times P(w_i) \\ \rightarrow d=0.75, \in [0, 1] \\ \text{bigrams starting with } w_{i-1} \leftarrow \text{Continuation} \\ \text{bigrams ending with } w_i \downarrow$$

$$\lambda(w_{i-1}) = \frac{d \times \text{no. of unique bigrams starting with } w_{i-1}}{C(w_{i-1})} \quad \text{normalization} \\ \text{Continuation} \quad \text{no. of bigrams starting with } w_{i-1} \text{ in total} \quad \text{no. of unique bigrams}$$

$$P(w_i) = \sum_{w_o} C(w_o, w_i) \quad \text{so ending with } w_i \\ \text{Continuation} \quad \text{total bigrams} \quad \text{no. of unique bigrams}$$

$$Q \quad C(\text{the, cat}) = 3, \quad C(a, cat) = 2, \quad C(\text{the, dog}) = 4, \\ C(a, dog) = 1, \quad C(cat, sleeps) = 1, \quad C(dog, sleeps) = 3 \\ C(cat) = 5$$

now let's do for trigram.

$$P_{\text{trig}}(w_i | w_{i-1}, w_{i-2}) = \max(C(w_{i-2}, w_{i-1}, i) - D, 0)$$

$$+ \lambda(w_{i-2} | w_{i-1}) \rightarrow P_{\text{trig}}(w_i | w_{i-1})$$

here  $\lambda(w_{i-2}, w_{i-1}) = \frac{d}{C(w_{i-2}, w_{i-1})} \times \text{unique trigrams following } w_{i-2}, w_{i-1}$

P की नयी फॉर्मूला ही ज्यादा चौंका देती।

Entropy = highest of uniform distribution

= measure of randomness

$$H(p) = - \sum_{x \in X} p(x_i) \log_2 p(x_i)$$

peopleility =  $2^H$

$$\text{PP}(w) = P(w_1, w_2, \dots, w_n)^{-1}$$

peopleility

$w = w_1, w_2, \dots, w_n$  sequence

normalizing peopleility complete Corpus of given.

Model can predict = a, the, red, fox, dog, .  
so  $P('a \text{ red fox .}')$

$$P(a) \times P(\text{red} | a) \times P(\text{fox} | \text{a red})$$

$$\times P(.) | \text{a red fox})$$

$$= 0.9 \times 0.27 \times 0.55 \times 0.79$$

$$= 0.469$$

$$\text{but } P_{\text{norm}}('a \text{ red fox .}') = p^n$$

$$= (0.469)^4$$

$$= 0.465$$

$$\text{peopleility} = 1 = 2.15$$

peopleility measures how well the language model will predict the next word.

$$\text{peopleility} \propto \frac{1}{\text{Probability of Sequence}}$$

Uniform distribution एवं मै peopleility अन्य जिनमें से नहीं जो नहीं अच्छी है।

likelihood = for a particular point

probability = Area in the graph over a range.

$$\begin{aligned}
 P(w) &= \frac{1}{P(w)^{1/n}} \\
 &= e^{-\frac{1}{n} \log(P(w))} \\
 &= e^{-\frac{1}{n} \sum_{i=1}^n \log(P(w_i | w_{i-1}))} \\
 &\quad \text{Chain multiplication of history} \\
 &\quad (\text{just thi koi pan model s2 prob.} \\
 &\quad \text{s12i and QM s like prev. example})
 \end{aligned}$$

$\Leftarrow$  I am going </s>  
so no need to do  $P(<s> | </s>)$

generally  
used  
formula.

bigram  
consider.

$$= P(I | <s>) \times P(am | I) \times P(going | am)$$

## Text Representation

So a sentence  $w_1, w_2, \dots, w_n$  can be represented as vector of prob.

$$(P(w_1 | s), P(w_2 | w_1), \dots, P(w_n | w_{n-1}))$$

- 1) Label Encoding (Bag of words representation)  
so let's say 10k unique words & in vocab then we give 0 to 10000 index and word at index  $i$  represent s2  $i$ .

problem: no meaning.

- 2) OHE

- 3) eg1: I am in Srinagar  
eg2: I am in Delhi

Dep. 1% Sentence can be represented in vector of 10000 dimension

and I at index 42 I, am at index 92 1, ... 0.

Dep. 2% Corpus level 10k dimension vector.

freq(word) in corpus at 0.

word  
team frequency

TF-IDF  
 $tf_i = \text{no. of times (frequency) of a team } i \text{ in a document } d.$   
 $= \text{raw count}$

but prob. of  $f_{t,d}$  is do smoothing  
 $f_{t,d} = \text{freq. of team } t \text{ in document } d$

$$tf = \log(1 + f_{t,d})$$

Normalization we do because

freq. can be very high.

$$M-1) \quad tf = \frac{f_{t,d}}{\max_{t' \in d} f_{t',d}}$$

$$M-2) \quad tf = \frac{f_{t,d}}{\max_{t' \in d} f_{t',d}}$$

M-3) Double normalization.

$$tf = 0.5 + 0.5 \times \frac{f_{t,d}}{\max_{t' \in d} f_{t',d}}$$

General formula

$$tf = k + (1-k) \times \frac{f_{t,d}}{\max_{t' \in d} f_{t',d}}$$

PAGE NO.: 1 / 1  
DATE: 1 / 1

PAGE NO.: 1 / 1  
DATE: 1 / 1

Inverse Document Frequency (IDF)

$$\text{idf weight} = \log \frac{N}{\max_{t \in D} |t|}$$

$\hookrightarrow$  no. of documents

but we take log and + 1

$$= \log(N/4) + 1$$

$$M-2) \quad idf(t) = \log \left( \frac{N}{\max_{t \in D} |t|} \right)$$

$\hookrightarrow$  max no. times this word comes in which doc

$\hookrightarrow$  max docs this word exists in

$$M-3) \quad idf = \log \left( \frac{N + 1}{\max_{t \in D} |t| + 1} \right)$$

$$\text{Finally: } TF \cdot IDF = TF \cdot IDOF$$

# How to generate new sentences using LMs?

= let's take a bigram and last 2 most probable words

& word ends at end @ </s>

problem is we can generate only 1 sentence.

So what we can do is generate 1st word a) randomly

b) based on next

use beam search  $k$  best prob.  $\leftarrow$  best probability

Ques: TF-IDF example.

- 1) This is a document.
- 2) The second document is the document.
- 3) The third document is over the first document.
- 4) The first document is there.
- 5) Is there a document?

Ans: Pre-processing : 1) Word tokenize  
2) Case same & remove stop words.

Now total no. of unique words - are 12

tf is on single document, here document is each sentence.

$$tf = \{t_i : t_i \in d\}$$

For doc 1

$$\text{tf}(this) = 1$$

$$\text{tf}(is) = 1$$

$$\text{tf}(a) = 1$$

$$\text{tf}(document) = 1$$

$$\text{tf}(.) = 1$$

? a doc. first is over second & the third this

1	0	1	1	0	1	0	0	0	0	0	1
1	0	0	2	0	1	0	1	2	0	0	0
1	0	0	2	1	1	1	0	2	0	1	0
1	0	0	1	1	1	0	0	1	1	0	0
0	1	0	1	1	0	1	0	0	0	1	0

Now for IDF

$$idf_t = \log_2 \frac{N}{df_t} \quad | \quad \begin{array}{l} \text{idf for term } t \text{ across} \\ \text{all documents} \\ \rightarrow 3 \text{rd doc has } t \text{ and } 1 \text{st} \end{array}$$

So 1 vector  $\begin{pmatrix} 1 & 0 & 1 & 1 & 0 & 1 & 0 \end{pmatrix}$

$$N = 5$$

? a doc. first  
1 2 5 2

in idf common words = low score.

$$\text{idf}_s = \log_2 \left( \frac{S}{df_s} \right) + 1$$

$$\text{idf}_s = \log_2 \left( \frac{5}{1} \right) + 1 = 2.32$$

so tf matrix  $\times$  idf matrix = tf-idf

now in testing phase, there is a new sentence so is this a document?

Testing this tf again but idf remains same and then tf  $\times$  idf (calculated earlier)

to find vocab. we can use train + test data

PAGE NO.:  
DATE: / /

Now if some new word comes in testing  
so may idf = 0 so we use.

$$idf = \log \left( \frac{N+1}{d+1} \right) + 1$$

instead of words if we want to use n-grams  
in n-gram of freq. as like, similarly in tf

e.g. The boy is sitting

bigram or example.  $\leftrightarrow$  The boy boy is is sitting  
tf of matrix

Note: if we are dealing with n-grams then we need  
n-1 start and end of sentence markers.

Also in n-grams vocab. size increases as  
 $V_1 + V_2 + V_3 + V_4$   
↳ one word freq.  $(w_1, w_2, w_3, w_4)$  of frequencies.

Also like words tf-idf can be done on  
character level that too in 2 types:

- 1) without spaces
- 2) with spaces

|V word| >> |V char|  
↳ 26 + Punctuations only

Pointwise mutual information  $\leftarrow PMI(x,y) = \log \left( \frac{P(x,y)}{P(x).P(y)} \right)$

PAGE NO.:  
DATE: / /

When we combine word and character level  
tf-idf we will get best accuracy.

\* TF-IDF doesn't take position into account (n gram thi  
tukinda) diya to see if it's collocating

### COLLOCATION

So we can maintain 2 lists, drift to find out how  
many words precedes 'I' (like you word) and other to  
find how many successors sat the word

we can do like  $V \times V$  matrix and 2 lists  
and 2 rows words which have 1 else 0  
(case of preceded 104)  $\begin{pmatrix} 1 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \\ 0 & 0 & 1 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$

Or  $V \times V$  matrix and count of words coming before  
a word.

$$P(x,y) = \frac{c(x,y)}{c(x)} \quad \begin{cases} \text{conditional} \\ \text{or distribution} \end{cases}$$

$$P(x,y) = \frac{c(x,y)}{\text{total no. of bigrams}} \quad \begin{cases} \text{no smoothing} \\ \text{or} \\ \text{smoothed} \end{cases}$$

$$P(x) = \frac{c(x)}{c(x,y)}$$

$$\text{from above eqn. or distribution}$$

main problem of n-gram = out of vocabulary (OOV)

PAGE NO.:  
DATE: / /

### Byte pair Encoding (BPE)

See the non-overlapping most frequently occurring pairs.

and pair of 2 or 3 characters which is not in vocab will replace it.

also if freq. tie then prefers original chars.

Store all the transformations in lookup table.

e.g. aaabdaaabac

→ ZabZabac ( $Z = aa$ )

→ ZydzYac ( $Y = ab$ )

→ Xdxac ( $X = ZY$ )

Stop here as no > 1 freq.

S-1: Count all word level frequencies.

h ug : 10

pug : 5

pu : 5

bun : 4

hugs : 5

S-2: Split character wise, for each word.

h u g → p u g → ...

Byte is represented as unicode generally.

PAGE NO.:  
DATE: / /

S-3: Vocabulary is of all appearing unique chars.

$$V = \{h, u, g, p, b, n, s\}$$

S-4: Find char pair frequency for all character pair

$$hu : 10$$

$$ug : 20$$

$$pu : 13$$

$$vn : 16$$

$$bu : 4$$

$$gs : 5$$

S-5: Merge the ones (char pairs) with highest frequency.

$$h \quad ug : 10$$

$$p \quad ug : 5$$

$$p \quad u \quad n : 12$$

$$b \quad u \quad n : 4$$

$$h \quad ug \quad s : 5$$

$$\text{merge} = \{ (u, g) : ug \}$$

include this char pair in vocabulary

$$V = \{h, u, g, p, b, n, s, ug\}$$

now 416 Step-4, 5

Stop: nothing left to combine

<p>PAGE NO.: 1 / 1</p> <p>DATE: 1 / 1</p> <p>if a new word out of vocabulary comes e.g. <u>thug</u> <math>\rightarrow</math> t, hug split <math>\Rightarrow</math> (longest possible merged in vocab find &amp; then split <math>\Rightarrow</math>)</p> <p><u>t, hug = unk</u> <math>\rightarrow</math> <u>hug</u> <math>\rightarrow</math> <u>unk</u> <math>\rightarrow</math> <u>unknown</u></p> <p><u>Word Piece</u></p> <p><math>\Rightarrow</math> Change - 1: we also show &amp; split <math>\Rightarrow</math> <u>hug</u> <math>\Rightarrow</math> <u>h u g</u></p> <p><u>Walking</u> <math>\rightarrow</math> <u>walk + #ing</u> <math>\rightarrow</math> <u>comes in middle</u></p> <p><u>hug</u> <math>\rightarrow</math> <u>h # u # g</u></p> <p><math>\Rightarrow</math> Change - 2: to merge we don't just see max freq. but we see score.</p> <p><u>h # u # g</u> <math>\rightarrow</math> <u>h u # g</u> <math>\rightarrow</math> <u>merge with so # original length</u></p> <p>This score penalises good frequent subwords</p> <p>Score = <math>\frac{f(A, B)}{P(A) \times P(B)}</math></p>
--

Chinese dii whitespace on EM so it is called text segmentation.

And if text splitting = english dii split 'split'

$$P(A, B) = C(A, B)$$

$$C(A) \times C(B)$$

if A, B generally comes together then score = higher  
if combination D score then score = lower

e.g. I am going to class.

I - am - go - ing - to - class.

N-gram big big tri skip-gram so we need to skip

This captures long context better

e.g. I am going home

So I  $\leftarrow$  going

am - home

going - I

home  $\leftarrow$  I am

$$\text{Total skip bigrams possible} = (n-2)(n-3) + 2(n-2)$$

Now for skip-trigrams we can do:

- 1) 3 unigrams all on dii spaces.
- 2) 1 unigram & skip us 1 bigram.
- 3) 1 bigram & skip us 1 unigram.

$$P(\text{skip trigram}) = \frac{\text{Count}(x, y, z \text{ tri in context})}{\text{Count}(x, y \text{ in context})}$$

$$P(x, y, z) = \frac{C(x, y, z)}{C(x, y)}$$

M-20  $P(x, y, z) =$  just like word-piece.

For filling the collocation table we use bigram + skip-gram  
So code dii just put 1 for all the words in the particular sentence other than this word.

So I am going home

So I, going, home  $\leftarrow$  1-1 bigram.

Now method-2 is instead of putting 1 we use TF-IDF of all and then sort them.

### Negative Sampling

So similar words in a context dii and is and for this we know these words which are in V-W  $\rightarrow$  similar words, words in context similar like am proportion dii sampling sat or dog and like probability am like 0.00001 so neural network dii this will help

So similar words one build. M most of most of us.

so tf-idf in ~~the~~ colocation table and  
 but negative sampling will give small value  
 instead.

\* also ~~to~~ each word has ~~a~~ 1 vector  
 thus  
 $I [ \_ \_ \_ ] \rightarrow \text{vector}$

so to find similarity we can use 2 methods:

- 1) Cosine Similarity (value ↑ → similar)
- 2) Euclidean distance. (value ↓ → similar)

Steps :-

- 1) Tokenization
- 2) Preprocessing techniques
  - Stopword removal
  - Lemmatization
  - Stemming
  - Normalization
- 3) PSA and FST → morphology of words
- 4) N-gram modeling → freq. distribution
- 5) Smoothing techniques
  - add 1
  - add k → add type
  - backoff → n-gram
  - Katz
  - Delayed (inter) interpolation
  - Good Turing
  - Kneser Ney

PAGE NO.:  
DATE: / /

Prediction of next word  
 Greedy Beam

PAGE NO.:  
DATE: / /

word

## 6) Text representation Schemes

TF-IDF  
 PMI

Chowdhury

— Subword

BPE

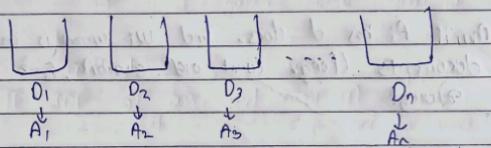
Word piece  
 Sentence piece

7) Collocation → TF-IDF

— Negative Sampling,

Naive Bayes (features one words here)

Given a new document; you have to predict the author



$$P(Y|X) = \frac{P(X|Y) \cdot P(Y)}{P(X)}$$

We don't want to model feature dependence so we remove denominators.

$$P(x_1 x_2 x_3 | Y) = P(x_1 | Y) \cdot P(x_2 | Y) \cdot P(x_3 | Y)$$

$$\text{Prior} = P(Y_k)$$

$$\text{Posterior} = P(Y_k | X)$$

$$P(D_i | A) = P(w_1 | A_i) \times P(w_2 | A_i) \times \dots \times P(w_n | A_i)$$

$D_i$  has words  $w_1$  to  $w_n$ .

So create a dictionary

$$\hat{A} = \underset{A_i}{\operatorname{argmax}} \quad P(D_i | A_i) \times P(A_i)$$

$$P(A_i) = \frac{\text{total no. of docs authored by } A_i}{\text{total no. of docs}}$$

$$P(D_i | A_i) = P(w_1 | A_i) \times P(w_2 | A_i) \times \dots \times P(w_n | A_i)$$

Now if  $A_i$  has  $d$  docs. And  $w_i$  word is in  $m$  documents (e.g. count of  $w_i$  in  $A_i$ ,  $\sum_{j=1}^m$   $1$  if  $w_i$  is in  $A_j$ )

$$\therefore P(w_i | A_i) = m$$

We can also find  $P(w_i | A_i)$  in a diff way.

$$P(w_i | A_i) = C(w_i, A_i)$$

All words written

by  $A_i$

$$P(w_n | w_{n-1}, w_{n-2}) = P(w_{n-1} | w_{n-2})$$

$$\times P(w_n | w_{n-1}, w_{n-2})$$

# Can we have  $n$ -gram features?  
= Yes

GRM replace  $w_i \rightarrow w_i, w_{i-1}$

$$P(w_1, \dots, w_m | y_k) = P(w_1 | y_k) \times P(w_2 | y_k) \times \dots \times P(w_m | y_k)$$

$$P(w_1, w_2 | y_k) \times P(w_2, w_3 | y_k) \times \dots$$

ie) Multi  $P(y_k)$  multiply all.

And  $P(w_i | w_j | y_k)$  di diff GRM & bigram & 2-gram does  
di  $w_i$  / total docs. written by that author.

problem - tf-idf or add k smoothing (tf-idf)  
or idf di  $\log\left(\frac{N+1}{d+1}\right) + 1$  for

## Sentence Classification

So we are using tf-idf as features but there are generic features

To add extra features we can add 3 columns like,  
no. of +ve words    no. of -ve words    no. of neutral words

We will have list of +ve, -ve and neutral words and  
each word in a sentence of list di find & and maintain  
the count.