

Regular exp. (cont.)

1 turkey costs 30 \$ → Positive lookahead

Is there a digit followed by a \$

$\text{id} + (? = \backslash \$)$

(30)

Positive lookahead

Is there a digit not followed by \$

$\text{id} + (? ! \backslash \$)$ → Negative lookahead

(1)

once/none

$(? < = \backslash \$) \text{id} +$ → Positive lookbehind

$(? < ! \backslash \$) \text{id} +$ → Negative lookbehind

Text Normalization

(

A set of collective tasks

Primary pre-processing tasking before feeding to model

1. Tokenization

2. Normalizing word formats

3. Segmenting sentences

- 1 "The U.S. is a big nation."
- 2 The U.S.A poster-print costs 30\$.
- 3 The driving speed unit there is m.p.h, not km."

Embedded vectors

'The',
'U.S.',
'is',
'a',
'big',
'nation'

1

2

more words
↓
more tokens

3

Tokenization: segmenting a running text into words (tokens)

Represent tokens(words) as numbers

3 diff size vectors

Rule based tokenization — prev.

Each sentence end w/ . → For 1st sentence you get 3 diff sentences

abbreviations
periods } .

white space split →

←
prob

← soln Penn Treebank standard → not all . are sentence
(PTB) breakers

postex-print — 1 single token

'30', '\$'

PTB + regex

Based on tokenization, it depends how effective
embedding will be.

Before seq2seq/transformer based models — RNN earlier

bert-based model - embedded vector
power of transformer/pretraining
knows

which words are similar

king → queen → cluster 1
student → teacher → cluster 2
country → patriot } cluster of tokens

corpus

'nation'

: similar vector value like
country

(embedding vector)

)
RNN
doesn't have this
info
only bert has

bert tokenizer / bert embedding

better representation of words ✓ → optimizing parameters

automatic token parsing techniques: (Tokenization)

1. Byte Pair Encoding (BPE)

2. unigram language modeling (ULM)

3. wordpiece → *be* *r*

2 parts:

token learner

token segmenter

[sentence piece: BPE + ULM]

corpus → { low, new, newer, wider, lowest }

PTs fails to understand what to do w/ unknown words.

BPE: corpus - 18 tokens

Divide all words based on whitespace and word frequency

count

corpus

all unique chars

← vocabulary

5

low -

→ whitespace

-, b, e, i, h,

n, o, r, s,

t, w, *e*,

2

low est -

e -, *ne*,

6

new *e* -

new, lo,

low, newer -

3

wid *e* -

low -

2

n ew -

morphemes

Task: create a vocab

List most frequent pair of symbols together

$$er = 6 + 3 = 9$$

n.e.w er _
w.i.d er _

n.e.w er _

w.i.d er _

(Learns on unseen data.
new word:
lower rate

morphemes: smallest unit of

BPE → lower $\left\{ \begin{array}{l} low \\ er \end{array} \right\}$ 2 tokens

present in learner vocab
low er _
diff token
(merged)

subwording: dividing the tokens into smaller parts

subword = morpheme

ne w er _
ne w _

BPE: K parameter \Rightarrow how many cycles you want to run

BPE(corpus, K)
 $\rightarrow 9$

Normalizing:

am is are — same word containing same meaning

case folding:

(
lower case
everything

Woodhouse
woodhouse

} same value for
embedding vector

→ Info retrieval ✓

Text/sentiment classification

& machine translation

(
not doing lower casing

US → country

us → we

Lemmaization: Determine which words are the same root

addition
additional } same root

I have some additional information.

I can do some addition on information.

same

words (additional)
 \swarrow stems (additional)
 \searrow Affixes (al)

replace each word w/ roots \rightarrow Lemmatization \checkmark
using pre-existing dictionary (better)
sing, sang, sung

Do Lemmatization before feeding corpus to language model.

Stemming : Remove affixes from words

X

~~additional~~

Changed actual meaning

(dang \rightarrow doe) \rightarrow error prone (stemming)

(dang \rightarrow do) \rightarrow not error prone (Lemmatization)

happiness \rightarrow happy
 \wedge \wedge
ih ih

Lemmatization over stemming \checkmark