

N-Gram Modelling

Markov chains - States:

- o A is a state
- o B is a state

- o Probability of $A \rightarrow B$: 50%.
- o Probability of $A \rightarrow A$: 50%.
- o Probability of $B \rightarrow A$: 50%.
- o Probability of $B \rightarrow B$: 50%.

AABA

Seq. of chain of states is called Markov chains.

N-Gram

"An n-gram is a contiguous sequence of n items from a given sample of text or speech" \rightarrow wikipedia

Items refer to states in Markov chains

Items can be "characters", "words", "sentences" etc.,

$N=2$, bigram.

$N=3$ Trigrams

characters are the states of character chains

$N=2$ "the bird is flying on the blue sky"

Bigrams = 'th', 'he', 'e', 'b', 'bi', 'ir', 'rd', 'd', 'i', etc.,

$N=3$

Trigrams = 'the', 'he', 'e b', 'bi' ... etc.,

Trigrams

next

the \rightarrow

he \rightarrow b

e b \rightarrow i

bi \rightarrow r

bir \rightarrow d

ird \rightarrow

Word Grams

(17)

"The bird is flying on the sky"

Trigram

The bird is →

bird is flying →

is flying on →

flying on the →

on the blue →

The blue sky →

Next

flying

on

The

blue

Sky

Next (Item when

large corpus)

[flying, eating, sleeping]

[on, through, on the]

The

[blue, orange]

Sky

[.]

↑

→ If we are looking at a huge corpus

Application

→ Auto Complete System

The bird is [flying, eating, sleeping]

* Select any word at random.

* The bird is flying [on, through, on]

→ The bird is flying on [the]

→ The bird is flyin on the [blue, orange]

→ The bird is flyin on the blue sky.

→ They are used in our smart phones, Google

→ Entertainment, Entertain.

<http://text-analytics101.xnlp.com/2014/11/what-are-n-grams.html>

(12)

code Ngram (character ngram model)

text = "Global warming"

n = 3

ngrams = {}

create ngrams

for i in range(len(text) - n):

gram = text[i:i+n] # text[0:3] = Glo

If gram not in ngrams.keys():

ngrams[gram] = []

ngrams[gram].append(text[i:i+n])

testing our n-gram model

CurrentGram = text[0:n]

result = CurrentGram

for i in range(100):

If CurrentGram not in ngrams.keys():

break

Possibilities = ngrams[CurrentGram]

nextitem = Possibilities[random.randrange(len(Possibilities))]

result += nextitem

CurrentGram = result[len(result)-n:len(result)]

Print(result)

Code for N-gram word

```
text = """Global warming & climate change has become a world wide concern"""
```

```
import random
```

```
n = 3
```

```
ngrams = {}
```

```
words = nltk.words(tokenize(text))
```

```
for i in range(len(words)-n):
```

```
    gram = ' '.join(words[i:i+n])
```

```
    if gram not in ngrams.keys():
```

```
        ngrams[gram] = []
```

```
        ngrams[gram].append(wordstext[i+n])
```

```
    
```

```
# testing our N-gram model
```

```
currentgram = text = ' '.join(words[0:n])
```

```
result = currentgram
```

```
for i in range(30):
```

```
    if currentgram not in ngrams.keys():
        break
```

```
    possibilities = ngrams[currentgram]
```

```
    nextitem = possibilities[random.randrange(len(possibilities))]
```

```
    result += ' ' + nextitem
```

```
    currentgram = result[len(result)-n:len(result)]
```

```
print words = nltk.word_tokenize(result)
```

```
    currentgram = ' '.join(words[len(words)-n:len(words)])
```

```
print(result)
```