

Introduction

Although n gram was not invented by Google, n gram might have been introduced to general public by Google n gram viewer. N gram is an essential tool for Natural Language Processing such as Machine Translation, Speech Recognition, Word Prediction etc. However only generating n gram is not enough. Only generating it can be done efficiently and easily with existing tools such as Python NLTK, Natural Language Tool Kit. It must be stored in an efficient data structure so that further analysis can be as optimal as possible. Especially a project like Google n gram viewer must be at scale. To deal with huge data, generating and storing n gram must be efficient.

The key of generating n gram is to generate as many word/character patterns as possible out of data.

On this poster, the author introduces n gram backed by a data structure called "Patricia Trie" at scale taking advantage of multicore CPUs to challenge Google n gram viewer.

Aim

- a) To challenge Google n gram viewer in terms of number of generated n gram words
- b) To take advantage of multicore CPUs for huge data set.
- c) To be as efficient as possible from generating n gram to further analysis with n gram.

Method

- a) Implement n gram using multictores from scratch.
- b) Implement Patricia Trie from scratch.
- c) Use Patricia Trie as a storage for n gram.
- d) Gather data for n gram crawling on the Internet.
- e) Run n gram on a 36 cores with 64 GB mem machine, Two Intel(R) Xeon(R) CPU E5-2630 v3 @ 2.40GHz.
- f) Measure the performance as follows:
 - a) # of generated n gram words per second.
 - b) Total # of generated n gram words.
 - words.

n gram and Patricia Trie – Great Match In Natural Language Processing In Big Data at Scale Seiya Kawashima, The University of Chicago, Radiology Department

c) Total # of unique words used to generate n gram

struct ngrms cw *nc = create_ngrms_cw(); insert_ngrms_ngrms_cw(nc)

Maximum Likelihood Estimate: mle_ngrm(); Stupid Backoff: sbo_ngrms();

Workflow from generating n gram to analysis with n gram



Input Data: "Hello World" At character level, n gram is generated as follows: "H","He","Hel","Hello", "Hello ", "Hello W", "Hello Wo", ... At word level, n gram is generated as follows: "Hello". "Hello World"



Sample of Patricia Trie data structure



Word level
gram, thread #9
gram, thread #10
gram, thread #11
gram, thread #12
gram, thread #13
gram, thread #14
gram, thread #15
gram, thread #16

ia Trie Big	O Notation
sertion:	log(n)
earch:	log(n)
eletion:	log(n)

"2015", "abc" and "hello" are Inserted in order into Patricia Trie.

Maximum Likelihood Estimate, MLE C (W1 ... Wi) $P(W_i | W_1 \dots W_{i-1}) =$ C (W1 ... Wi - 1) Probability of W_1 followed by word up to W_{i-1}) Frequency of word up to Wi Frequency of word up to W_{i-1}

Big O Notation of Maximum Likelihood Estimate backed by Patricia Trie is log(n). Search of Frequency of word up to Wi: log(n) Search of Frequency of word up to Wi-1: log(n)

Efficient usage of Patricia Trie for n gram

Result

Total # of web pages proc # of generated n gram wo Total # of generated n grain # of generated n gram cha # of generated n gram wo # of Unique words

Conclusion

There are positive and still challenging, findings on this project. Good findings are as follows: a) It didn't break during the experiment. b) It generated a great number of n gram per

- Second.
- huge.
- Still challenging findings are as follows:

 - viewer deals with.

b) Total # of generated n gram words was not as huge as Google n gram viewer generates. Although, with these negative findings, the Author feels confident about this project to be as good as Google n gram viewer in the near future.

cessed	12,905
ords per second	3,151
am words/characters	45,608,314
aracters	27,357,970
ords	18,250,344
	611,593

c) Total # of generated n gram words/characters was

a) Data set was not as huge as Google n gram