

# CSCI – 3412 ALGORITHMS

## SPRING 2013

### PROBLEM SET -2 PROGRAMMING

#### **PART 1**

According to instructions given in Problem Set -2, the program should find the stop words from the file and eliminate them while counting the n-grams. So, for that my program reads the .csv file which contains the stop words and read each line in the loop of counting the n-grams. Basically, while reading the main file (shakespear.txt) for its each and every line, this program will compare with each of the comma separated words from the .csv file each time. If it matches, then it just ignores it and continue for the next one. Otherwise it will store the n-grams.

At the end of the loop, this program has used the Hash Map to store n-grams. Then it uses TreeMap for sorted array.

Here is the Result of this program:

```

C:\windows\system32\cmd.exe
D:\SPRING 2013_DP\ALGORITHM\ASSIGNMENT 2>java Ngrams
----- Part 2 -----

----- For 1-grams -----
There are UNIQUE : 26853 grams

(" shall "> [3603];
(" good "> [2817];
(" now "> [2798];
(" lord "> [2717];
(" come "> [2565];
(" sir "> [2541];
(" well "> [2519];
(" more "> [2291];
(" here "> [2149];
(" ill "> [2002];

----- For 2-grams -----
There are UNIQUE : 237428 grams

(" c " " sar "> [346];
(" good " " lord "> [244];
(" here " " comes "> [189];
(" sir " " john "> [152];
(" come " " come "> [135];
(" fare " " well "> [117];
(" good " " morrow "> [113];
(" come " " hither "> [112];
(" mine " " eyes "> [106];
(" shall " " see "> [96];

----- For 3-grams -----
There are UNIQUE : 206393 grams

" sir " " john " " falstaff "> [19];
" c " " sar " " shall "> [14];
" look " " here " " comes "> [14];
" three " " thousand " " ducats "> [13];
" ha " " ha " " ha "> [12];
" one " " word " " more "> [12];
" now " " whats " " matter "> [12];
" fie " " fie " " fie "> [11];
" mistress " " anne " " page "> [11];
" julius " " c " " sar "> [11];

----- For 4-grams -----
There are UNIQUE : 133870 grams

" hey " " ho " " wind " " rain ">[6];
" chooseth " " shall " " much " " deserves ">[5];
" kill " " kill " " kill " " kill ">[5];
" come " " hither " " come " " hither ">[5];
" talk " " young " " master " " launcelot ">[4];
" arrest " " high " " treason " " name ">[4];
" hey " " ho " " hey " " nonino ">[4];
" here " " sir " " here " " sir ">[4];
" east " " west " " north " " south ">[4];
" william " " de " " la " " pole ">[4];

```

## PART 2

In this section, for 2 grams ( $n=2$ ), it generates all the 2-grams and creates new file (.dot file), and puts all the values to this newly created file in the following manner:

digraph words {

" shall " -> " find " [weight=90];

```
"mine " -> "honour " [weight=79];  
"whats " -> "matter " [weight=76];  
"give " -> "leave " [weight=75];  
"good " -> "sir " [weight=74];  
  
.  
.  
.  
}
```

Now, this program only executes the graphviz command and gives this file as an input and generates pdf file which gives the Directed graph for the given 2-grams.

Pdf file is attached with the program.

Let me know in case of any changes.

Thanks,

Darpan Shah.