

Programming Project Evaluation Form

Student Name:	Project:	Date:
Correctness Criteria: <ul style="list-style-type: none"> Reads normal input data Handles incorrect input data Calculates correct results Outputs results properly Prints appropriate error messages 	Score: _____ / 50	
Design Criteria: <ul style="list-style-type: none"> Problem decomposition Choice of data structures Choice of algorithms Program efficiency (space/time) 	Score: _____ / 20	
Documentation Criteria: <ul style="list-style-type: none"> Comments for classes and methods Comments describing algorithms Comments describing data structures Description of program design Description of test results Description of known problems 	Score: _____ / 10	
Style Criteria: <ul style="list-style-type: none"> Method and variable names Program indenting Use of white space Ordering of methods Easy to read code 	Score: _____ / 10	
Testing Criteria: <ul style="list-style-type: none"> Normal input data Incorrect input data Special cases for data structures Special cases for algorithms 	Score: _____ / 10	
Grader Comments: 		
Total: _____ / 100		

Installation

Programming language: Python

Source code location: /home/sbillah/nlp/

Documentation:

```
sbillah@turing:~/nlp$ ./NLPEngine.py -h
```

```
usage: NLPEngine.py [-h] [-o DST_DIR] [-s SRC_DIR] [-nlp BIGRAMS]
```

This is Syed's NLP program for HW01

optional arguments:

-h, --help show this help message and exit

-o DST_DIR, --output DST_DIR

 The directory where output goes

-s SRC_DIR, --source SRC_DIR

 The directory where raw files reside

-nlp BIGRAMS, --count_bigrams BIGRAMS

 count the bigrams in src directory and write to
 the dst folder in descending order

Here is a complete command line example:

```
sbillah@turing:~/nlp$ ./NLPEngine.py -s  
/home/sgauch/public_html/5013IR/files/ -o parsed/ -nlp bigram_counts
```

Done!

Algorithm

I use an in-memory Hash-table to count all bigrams. The pseudo code is given below:

```
initialize hash-table ht<(tuple), int>
foreach file f in input_directory:
    plain_text = html_parser(f.read())
    tokens = tokenizer(plain_text)

    i = 0
    foreach token t in tokens:
        if i>0:
            ht[(t[i-1],t[i])] += 1
        i++
sort ht
write ht to file
```

Time Complexity:

N = num_files

M= avg num_of_words_per_file

Bigram generation complexity: $O(N*M)$

Hash-table sorting complexity: $O(N*M*\log(N*M))$

Total complexity: $O(N*M) + O(N*M*\log(N*M)) = O(N*M*\log(N*M))$

Parser Configuration

Here is the configuration of my html parser and tokenizer:

str_src_dir	/home/sgauch/public_html/5013IR/files/
str_dst_dir	parsed/
str_doc_id_file_name	bigram.txt
min_token_freq	3
max_token_freq	1000
min_token_len	3
max_token_len	12
str_stop_list	Stoplist from this link: http://www.csce.uark.edu/~sgauch/5013IR/S12/index.html

Runtime & Memory usage

Input size (# files)	Run time (sec)	Memory size (MB)	Total bigrams
100	5.60	110	55,013
200	19.21	210	116,697
300	30.31	277	166,179
505	54.64	440	272,646

Top 50 bigrams

risks	jul	607
net	alter	344
alter	dynip	340
health	care	215
paper	title	208
com	interramp	204
net	sunbelt	189
edu	psu	177
net	mci	152
edu	nodak	144
edu	uiuc	142
critical	analysis	142
mass	media	141
edu	umich	137
net	idt	133
edu	umn	132
mil	navy	130
political	science	129
rights	reserved	127
edu	arizona	120
human	rights	117
edu	indiana	115
social	security	111
hogy	nem	111
edu	utexas	111

nemzet	magyar	110
los	angeles	109
horn	gyula	109
send	comments	108
mci	campus	108
black	studies	107
world	war	106
home	page	106
book	report	106
term	papers	105
http	www	105
written	price	104
urban	studies	104
termpaper	com	104
term	paper	104
subject	index	104
sports	recreation	104
specific	paper	104
paper	written	104
paper	click	104
description	paper	104
copyright	asm	104
comments	termpaper	104
comments	comments	104
cold	surges	104

Bottom 20 bigrams

abacs	kiskun	1
ababa	response	1
aaU	zoo	1
aaU	psy	1
aaU	hum	1
aas	nearly	1
aarp	national	1
aarp	american	1
aaron	word	1
aaron	netland	1

aaron	moshiashwili	1
aaron	jon	1
aaron	happened	1
aaron	comparison	1
aalen	image	1
aaemassago	sem	1
aeliberalis	demokratiato	1
aachen	rad	1
aachen	oph	1
aaa	passed	1

NLP Engine

```
1#!/usr/bin/env python
2
3'''
4Created on Feb 23, 2012
5
6@author: Masum
7'''
8
9from myparser import BiGram
10from myparser import src_dir, dst_dir
11from mycollection import argparse
12
13class NLP Engine():
14    config= {}
15
16    def __init__(self):
17        self.read_config()
18
19    def read_config(self, name="config.txt"):
20
21    def build_bigram_index(self):
22        _ = BiGram(self.config)
23
24import argparse
25if __name__=='__main__':
26    args = argparse.ArgumentParser(description="This is Syed's NLP program")
27    args.add_argument("-o", "--output", dest="dst_dir",
28                      help="The directory where output goes", default="")
29    args.add_argument("-s", "--source", dest="src_dir",
30                      help="The directory where raw files reside", default="")
31
32    args.add_argument("-nlp", "--count_bigrams", dest="bigrams",
33                      help="Re-index of all files - and overwrite existing dict and post files,
34                           i.e. -i all", default="")
35
36
37    args = args.parse_args()
38    src_dir = args.src_dir
39    dst_dir = args.dst_dir
40
41    if not args.bigrams:
42        print 'Invalid or No arguments. Try with -h for help!!'
43    else:
44        nlp = NLP Engine()
45        if args.bigrams: nlp.build_bigram_index()
46        print 'Done!'
47
48    nlp = SearchEngine()
49    nlp.search_query('.8 susan .1 uark .1 edu', True )
50    print 'Done!'
51
52
53
```

shared

```
1 '''
2 Created on Mar 11, 2012
3
4 @author: Masum
5 '''
6
7 import re
8 from collections import defaultdict
9
10 #===== global regex =====
11 is_word = re.compile("[a-zA-Z_]+$")
12 delim = re.compile("\t|\.|-|:|@|\\|/|,|;|\"|!|\\*|"+")
13
14 #=====path =====
15 dst_dir = ''
16 src_dir = ''
17
18 #===== global hash variables =====
19 doc_id= defaultdict(int)
20 term_count= defaultdict(int)
21
22 def parse(line, config):
23     line =line.strip()
24     if not line: return []
25
26     tokens= []
27     for word in delim.split(line):
28         word = word.strip("~`$%^&()[|<>=+-_/'").lower()
29         if word in config['str_stop_list']: continue
30         if config['min_token_len']<= len(word) <= config['max_token_len'] and
is_word.match(word):
31             tokens.append(word)
32     return tokens
```

BiGram

```
1 '''
2 Created on Feb 4, 2013
3
4 @author: Masum
5 '''
6 import os
7
8 from BiGramHTMLParser import BiGramHTMLParser
9 from collections import defaultdict
10
11 class BiGram():
12     def __init__(self, config):
13         self.config = config
14         self.ht = defaultdict(int)
15         self.htmlparser = BiGramHTMLParser(self.config, self.ht)
16         self.start_batch_processing()
17         self.write_file_map()
18
19     def start_batch_processing(self):
20         file_id=0
21
22         for in_file in os.listdir(self.config['str_src_dir']):
23             #if in_file not in ['medium.html','simple.html']: continue #for testing
24             with open(self.config['str_src_dir']+ in_file, 'r') as f:
25                 self.htmlparser.feed(f.read(),file_id)
26                 file_id += 1
27
28     def write_file_map(self):
29         #writing bigram file to a file named under document id
30         with open(self.config['str_dst_dir']+ self.config['str_doc_id_file_name'],'wb+') as f:
31             for words, count in sorted(self.ht.iteritems(), key=lambda (k,v): (v,k), reverse=
32 True):
33                 f.write(words[0]+" "+words[1]+" "+str(count)+"\n")
```

BiGramHTMLParser

```
1 '''
2 Created on Mar 11, 2012
3
4 @author: Masum
5 '''
6 import sys
7 from HTMLParser import HTMLParser
8 from shared import parse
9
10
11 class BiGramHTMLParser(HTMLParser):
12     text, N, newline = '', 0, {'br': '\n', 'BR': '\n'}
13
14     def __init__(self, config, ht):
15         HTMLParser.__init__(self)
16         self.config = config
17         self.ht = ht
18
19
20     def handle_data(self, raw):
21         self.text = self.text+ raw+ self.newline.get(self.lasttag, '')
22
23     #format is (token =>value, here val)
24     def feed(self, data, did):
25         #extracting html text
26         try:
27             HTMLParser.feed(self, data)
28         except:
29             sys.exc_clear()
30
31         #tokenizing extracted data
32         for line in self.text.splitlines():
33             tokens= parse(line, self.config)
34             if not tokens: continue
35             for i in xrange(len(tokens)):
36                 if i > 0 :
37                     self.ht[(tokens[i], tokens[i-1])] +=1;
38                     self.N +=1
39
40         #cleaning for next feed
41         self.text = ''; self.N=0
42         HTMLParser.reset(self)
43
44
```