

Aims

This week you get to see how wonderful Perl can be for some tasks.

Assessment

Submission: give cs2041 lab07 total_words.pl count_word.pl frequency.pl log_probability.pl identify_poet.pl
[total_words.py count_word.py frequency.py log_probability.py identify_poet.py]

Deadline: either during the lab, or Monday 12 September 11:59pm (midnight)

Assessment: Make sure that you are familiar with the lab assessment criteria (lab/assessment.html).

We have covered only a small amount of Perl in lectures. In fact, to cover the whole language in detail would take a whole semester, so we're going to rely on you finding out about the language yourself in tutes, labs and assignments. A good place to start is the Perl documentation & tutorial links on the class home page For example you might find these useful:

- Perl language syntax (<http://search.cpan.org/dist/perl/pod/perlsyn.pod>)
- Perl functions (<http://search.cpan.org/dist/perl/pod/perlsub.pod>)
- Perl operators (<http://search.cpan.org/dist/perl/pod/perl.op.pod>)

Introduction

In this lab exercise you will implement a Naive Bayes Document classifier (http://en.wikipedia.org/wiki/Naive_Bayes_classifier#Document_Classification) in Perl. Naive Bayes Document classifiers are used widely for applications such as spam filtering and language recognition. You will use the approach to discover the author of a piece of poetry.

The directory /home/cs2041/public_html/lab/perl/poets/poems/ (lab/perl/poets/poems) contains text files containing poems written by 10 famous poets:

- Elizabeth Barrett Browning (lab/perl/poets/poems/Elizabeth_Barrett_Browning.txt)
- Emily Dickinson (lab/perl/poets/poems/Emily_Dickinson.txt)
- John Keats (lab/perl/poets/poems/John_Keats.txt)
- Percy Bysshe Shelley (lab/perl/poets/poems/Percy_Bysshe_Shelley.txt)
- Robert Frost (lab/perl/poets/poems/Robert_Frost.txt)
- Samuel Taylor Coleridge (lab/perl/poets/poems/Samuel_Taylor_Coleridge.txt)
- Walt Whitman (lab/perl/poets/poems/Walt_Whitman.txt)
- William Blake (lab/perl/poets/poems/William_Blake.txt)
- William Butler Yeats (lab/perl/poets/poems/William_Butler_Yeats.txt)
- William Wordsworth (lab/perl/poets/poems/William_Wordsworth.txt)

Your program will be given a poem and should determine which of these 10 famous poets is most likely to have written this poem. Your program will be based this on the probability that a particular poet will use a given word. This can be estimated from the frequency which the poet uses the word in the poems you are given. You should link the directory /home/cs2041/public_html/lab/perl/poets/poems/ into your lab07 directory. So to get started do something like this.

```
$ ln -s ~cs2041/public_html/lab/perl/poets/poems
$ ln -s ~cs2041/public_html/lab/perl/poets/poem?.txt
```

The `ln -s` command creates a symbolic link which saves disk space. If you are working on your own machine you may prefer to copy the files across and add them to your repo:

```
$ cp -r ~cs2041/public_html/lab/perl/poets/poems .
$ cp ~cs2041/public_html/lab/perl/poets/poem?.txt .
$ git add poems poem?.txt
```

Exercise: Total Words

Write a Perl script `total_words.pl` which counts the total number of words found on in its input (`STDIN`).

For the purposes of this program and the following programs we will define a word to be maximal non-empty contiguous sequences of alphabetic characters (`[a-zA-Z]`).

Any characters other than `[a-zA-Z]` separate words.

So for example the phrase "The soul's desire" contains 4 words: ("The", "soul", "s", "desire")

For example:

```
$ ./total_words.pl <poems/Walt_Whitman.txt
116941 words
$ ./total_words.pl <poems/Emily_Dickinson.txt
112882 words
$ ./total_words.pl <poems/Robert_Frost.txt
21698 words
```

Hint: if your word counts are out a little you might be counting empty strings (split can return these).

As usual:

```
$ ~cs2041/bin/autotest lab07 total_words.pl
$ git add total_words.pl
$ git commit -a -m "total_words.pl passes dryrun tests!"
```

Exercise: Count Word

Write a Perl script `count_word.pl` which counts the number of times a specified word is found on in its input (`STDIN`).

A word is as defined for the previous exercise.

The word you should count will be specified as a command line argument.

Your program should ignore the case of words.

For example:

```
$ ./count_word.pl snow <poems/Elizabeth_Barrett_Browning.txt
snow occurred 11 times
$ ./count_word.pl path <poems/Robert_Frost.txt
path occurred 3 times
$ ./count_word.pl urn <poems/John_Keats.txt
urn occurred 7 times
$ ./count_word.pl England <poems/William_Blake.txt
england occurred 34 times
```

Hint: modify the code from the last exercise.

Hint: the Perl function `uc` & `lc` convert strings to lower & uppercase respectively.

As usual:

```
$ ~cs2041/bin/autotest lab07 count_word.pl
$ git add count_word.pl
$ git commit -a -m "count_word.pl initial version - not working"
```

Exercise: Word Frequency

Write a Perl script `frequency.pl` which prints the frequency with each poet uses a word specified as argument. So if Robert Frost uses the word "snow" 30 times in the 21699 words of his poetry you are given, then its frequency is $30/21699 = 0.001382552$. For example:

```
$ ./frequency.pl snow
11/ 50080 = 0.000219649 Elizabeth Barrett Browning
56/112882 = 0.000496093 Emily Dickinson
11/ 62844 = 0.000175037 John Keats
24/ 53400 = 0.000449438 Percy Bysshe Shelley
30/ 21698 = 0.001382616 Robert Frost
19/ 38123 = 0.000498387 Samuel Taylor Coleridge
14/116941 = 0.000119718 Walt Whitman
18/ 40751 = 0.000441707 William Blake
22/121198 = 0.000181521 William Butler Yeats
25/117162 = 0.000213380 William Wordsworth
```

So of these poets Robert Frost uses the word "snow" most frequently. If you choose a word randomly from Robert Frost the probability it will be "snow" is just over 1 in a thousand (0.1%).

Make sure your Perl script produces exactly the output above (the printf format is "%4d/%6d = %.9f %s\n"). Note you should ignore case (change A-Z to a-z). You should treat as a word any sequence of alphabetic characters. You should treat non-alphabetic characters (characters other than a-z) as spaces.

Hint: use a hash table of hash tables indexed by poet and word to store the word counts.

Hint: this loop executes once for each .txt file in the directory poems .

```
foreach $file (glob "poems/*.txt") {
    print "$file\n";
}
```

Hint: reuse code from the last exercise.

```
$ ~cs2041/bin/autotest lab07 frequency.pl
$ git add frequency.pl
$ git commit -a -m "what's your frequency Kenneth?"
```

Exercise: Word Log Probability

Now suppose we have the phrase "Truth is beauty." if John Keats uses the word "truth" with frequency 0.000333885 and the word "is" with frequency 0.004944671, the word "beauty" with frequency 0.000747265. We can estimate the probability of Keats writing the phrase "Truth is beauty." as:

$$0.000333885 * 0.004944671 * 0.000747265 == 1.23369825533711e-09$$

We could similarly estimate probabilities for each of the other 9 poets, and then determine which of the 10 poets is most likely to write "Truth is beauty." (its Blake).

A sidenote: we are actually making a large simplifying assumption in calculating this probability. Its often called the *bag of words model* (http://en.wikipedia.org/wiki/Bag_of_words_model).

Multiplying probabilities like this quickly leads to very small numbers and may result in arithmetic underflow of our floating point representation. A common solution to this underflow is instead to work with the *log* of the numbers.

So instead we will calculate the the log of the probability of the phrase. You this by adding the log of the probabilities of each word. For example, you calculate the log-probability of Keats like this:

$$\log(0.000333885) + \log(0.004944671) + \log(0.000747265) == -20.5132494670232 == \log(1.23369825533711e-09)$$

Log-probabilities can be used directly to determine the most likely poet, as the poet with the highest log-probability will also have the highest probability.

Another problem is that we might be given a word that a poet has not used in the poems we have. For example:

```
$ ./frequency.pl mortality
0/ 50080 = 0.000000000 Elizabeth Barrett Browning
4/112882 = 0.000035435 Emily Dickinson
5/ 62844 = 0.000079562 John Keats
5/ 53400 = 0.000093633 Percy Bysshe Shelley
0/ 21698 = 0.000000000 Robert Frost
0/ 38123 = 0.000000000 Samuel Taylor Coleridge
1/116941 = 0.000008551 Walt Whitman
0/ 40751 = 0.000000000 William Blake
0/121198 = 0.000000000 William Butler Yeats
7/117162 = 0.000059746 William Wordsworth
```

It's not useful to assume there is zero probability that the poet would use the word, even though they haven't used it previously. You should avoid this when estimating probabilities by adding 1 to the count of occurrences of each word. So for example we'd estimate the probability of Robert Frost using the word *mortality* as $(0+1)/21699$ and the probability of John Keats using the word *mortality* as $(5+1)/62896$. This is a simple version of Additive smoothing (http://en.wikipedia.org/wiki/Additive_smoothing).

Write a perl script `log_probability.pl` which given an argument prints the estimate log of the probability that a poet would use this word. For example:

```
$ ./log_probability.pl mortality
log((0+1)/ 50080) = -10.8214 Elizabeth Barrett Browning
log((4+1)/112882) = -10.0247 Emily Dickinson
log((5+1)/ 62844) = -9.2567 John Keats
log((5+1)/ 53400) = -9.0938 Percy Bysshe Shelley
log((0+1)/ 21698) = -9.9850 Robert Frost
log((0+1)/ 38123) = -10.5486 Samuel Taylor Coleridge
log((1+1)/116941) = -10.9763 Walt Whitman
log((0+1)/ 40751) = -10.6152 William Blake
log((0+1)/121198) = -11.7052 William Butler Yeats
log((7+1)/117162) = -9.5919 William Wordsworth
```

You will only need to copy your `frequency.pl` and make a small modification. Make sure your output matches the above exactly (the printf format is `"log((%d+1)/%d) = %8.4f %s\n"`)

```
$ ~cs2041/bin/autotest lab07 log_probability.pl
$ git add log_probability.pl
$ git commit -a -m "logs are useful"
```

Exercise: Identifying the Poet

Write a Perl script `identify_poet.pl` that given 1 or more files, each containing a poem, prints the most likely poet for each poem.

In other words, for each file given as argument you should go through all (10) poets calculating the log-probability that the poet wrote that poem by summing the log-probability of that poet using each word in the file. You should print the poet with the highest log-probability.

The files `poem1.txt` (`lab/perl/poets/poem1.txt`), `poem2.txt` (`lab/perl/poets/poem2.txt`) and `poem3.txt` (`lab/perl/poets/poem3.txt`) contain famous poems by Keats, Shelley and Frost, which are not included in their poems in the poets directory.

Your program should produce exactly this output:

```
$ ./identify_poet.pl poem?.txt
poem1.txt most resembles the work of John Keats (log-probability=-2720.9)
poem2.txt most resembles the work of Percy Bysshe Shelley (log-probability=-2720.9)
poem3.txt most resembles the work of Robert Frost (log-probability=-720.9)
```

You may find it helpful to add a `-d` flag which provides debugging information (this is optional), for example:

```
$ ./identify_poet.pl -d poem2.txt
poem2.txt: log_probability of -777.0 for Percy Bysshe Shelley
poem2.txt: log_probability of -789.0 for Robert Frost
poem2.txt: log_probability of -796.3 for Samuel Taylor Coleridge
poem2.txt: log_probability of -803.7 for William Wordsworth
poem2.txt: log_probability of -804.9 for John Keats
poem2.txt: log_probability of -805.4 for William Blake
poem2.txt: log_probability of -805.4 for Elizabeth Barrett Browning
poem2.txt: log_probability of -808.2 for William Butler Yeats
poem2.txt: log_probability of -818.1 for Walt Whitman
poem2.txt: log_probability of -832.0 for Emily Dickinson
poem2.txt most resembles the work of Percy Bysshe Shelley (log-probabil:
```

```
$ ~cs2041/bin/autotest lab07 identify_poet.pl
$ git add identify_poet.pl
$ git commit -a -m "logs are useful"
```

Challenge Exercise: Poetic Python

Implement the above exercises in Python.

The example Python scripts (/~cs2041/code/python/code_examples.html) and links to external Python resources should help - but you will need more info - Google is your friend.

Hints for `total_words.py` & `count_word.py` :

This loop executes for each line of stdin:

```
import sys
for line in sys.stdin:
    print line
```

The function `re.split` or the function `re.findall` could be used to separate words.

Hints for `frequency.py`, `log_probability.py` & `identify_poet.py` :

Beware Python dicts need a slightly different approach to Perl hashes, and also Perl & Python division have different semantics.

This loop executes once for each `.txt` file in the directory `poets`.

```
import glob
for file in glob.glob("poems/*.txt"):
    print file
```

You might find `math.log`, `sorted`, `re.sub` and `collections.defaultdict` useful.

```
$ ~cs2041/bin/autotest lab07 total_words.py
$ ~cs2041/bin/autotest lab07 count_word.py
$ ~cs2041/bin/autotest lab07 frequency.py
$ ~cs2041/bin/autotest lab07 log_probability.py
$ ~cs2041/bin/autotest lab07 identify_poet.py
$ git add total_words.py count_word.py frequency.py log_probability.py identify_poet.py
$ git commit -a -m "python is poetic"
```

Testing

Remember to do your own testing as well as the autotest tests are available for this lab.

To run all tests:

```
$ ~cs2041/bin/autotest lab07
```

You can run a single test if you also pass the test label as the second argument to autotest. For example, to run just test `total_words_5` type:

```
$ ~cs2041/bin/autotest lab07 total_words_5
```

You can also tell autotest to the code you have committed to gitlab

```
$ ~cs2041/bin/autotest lab07 -gitlab
```

or a particular gitlab commit

```
$ ~cs2041/bin/autotest lab07 9bfa2c5a
```

Finalising

You must show your solutions to your tutor and be able to explain how they work. Once your tutor has discussed your answers with you, you should submit them using:

```
$ give cs2041 lab07 total_words.pl count_word.pl frequency.pl log_probability.pl identify_poet.pl [total_v
```

Whether you discuss your solutions with your tutor this week or next week, you must submit them before the above deadline.

Gitlab - More Information

I expect most students will just work in their CSE account and push work to `gitlab.cse.unsw.edu.au` from there, but you can try setting up a git repository on your home machine and pushing work to `gitlab.cse.unsw.edu.au` from there.

If you do so you'll want to use git's pull command to update the repository in your CSE account.

```
$ git pull
Unpacking objects: 100% (3/3), done.
From gitlab@gitlab.cse.unsw.EDU.AU/z5555555/16s2-comp2041-labs
 226cddf..e64fee9 master    -> origin/master
Updating 226cddf..e64fee9
Fast-forward
 total_words.pl |      1 +
 1 file changed, 1 insertion(+)
```

If ssh access doesn't work, you can also use https to access gitlab using a URL equivalent to

`https://gitlab.cse.unsw.edu.au/z5555555/16s2-comp2041-labs.git` (replace 5555555 with your student number) and use your z-id & zPass.

```
$ git remote set-url origin https://gitlab.cse.unsw.edu.au/z5555555/16s2-comp2041-labs.git
$ git push
Username for 'https://gitlab.cse.unsw.EDU.AU': z5555555
Password for 'https://z5555555@gitlab.cse.unsw.EDU.AU': zPass
```