

Objectives

- Proficiency at text processing in Perl.
- Understanding multi-dimensional hashes.
- Explore a simple learning algorithm.

Preparation

Before the lab you should re-read the relevant lecture slides and their accompanying examples.

Getting Started

Create a new directory for this lab called **lab07** by typing:

```
$ mkdir lab07
```

Change to this directory by typing:

```
$ cd lab07
```

Exercise: How many words in standard input?

In these exercises you will work with a dataset containing sing lyrics.

This [zip file](#) contains the lyrics of the songs of 10 well known artists.

```
$ wget -q https://cgi.cse.unsw.edu.au/~cs2041/18s2//activities/total_words/lyrics.zip
$ unzip lyrics.zip
Archive:  lyrics.zip
  creating: lyrics/
  inflating: lyrics/David_Bowie.txt
  inflating: lyrics/Adele.txt
  inflating: lyrics/Metallica.txt
  inflating: lyrics/Rage_Against_The_Machine.txt
  inflating: lyrics/Taylor_Swift.txt
  inflating: lyrics/Keith_Urban.txt
  inflating: lyrics/Ed_Sheeran.txt
  inflating: lyrics/Justin_Bieber.txt
  inflating: lyrics/Rihanna.txt
  inflating: lyrics/Leonard_Cohen.txt
  inflating: song0.txt
  inflating: song1.txt
  inflating: song2.txt
  inflating: song3.txt
  inflating: song4.txt
```

The lyrics for each song have been re-ordered to avoid copyright concerns.
The zip file also contains lyrics from 5 songs where we don't know the artists.

```
$ cat song0.txt
I've made up my mind,  Don't need to think it over,
If I'm wrong I am right,
Don't need to look no further,
This ain't lust,
I know this is love but,

If I tell the world,
I'll never say enough,
Cause it was not said to you,
And that's exactly what I need to do,
If I'm in love with you,
$ cat song1.txt
Come Mr. DJ song pon de replay
Come Mr. DJ won't you turn the music up
All the gal pon the dance floor wantin' some more what
Come Mr. DJ won't you turn the music up
$ cat song2.txt
And they say
She's in the class A team
Stuck in her daydream
Been this way since eighteen
But lately her face seems
Slowly sinking, wasting
Crumbling like pastries
```

They are each from one of the artists in the dataset but they are not from a song in the dataset.

To start on this analysis write a Perl script `total_words.pl` which counts the total number of words found 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 <lyrics/Justin_Bieber.txt
46589 words
$ ./total_words.pl <lyrics/Metallica.txt
38096 words
$ ./total_words.pl <lyrics/Rihanna.txt
53157 words
```

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

When you think your program is working you can use `autotest` to run some simple automated tests:

```
$ 2041 autotest total_words
```

When you are finished working on this exercise you must submit your work by running **give**:

```
$ give cs2041 lab07_total_words total_words.pl
```

Sample solution for `total_words.pl`

```
#!/usr/bin/perl -w
# count words in STDIN
# written by andrewt@unsw.edu.au for COMP[29]041

$count = 0;
while ($line = <STDIN>) {
    my @words = split /[^\a-zA-Z]+/, $line;
    foreach $word (@words) {
        $count++ if $word ne '';
    }
}
print "$count words\n"
```

Alternative solution for `total_words.pl`

```
#!/usr/bin/perl -w

# count words in STDIN
# written by andrewt@unsw.edu.au for COMP[29]041
# more concise but less readable solution

$count = 0;
while (<STDIN>) {
    $count++ foreach /[a-z]+/ig
}
print "$count words\n"
```

Alternative solution for `total_words.pl`

```
#!/usr/bin/perl -w

# count words in STDIN
# written by andrewt@unsw.edu.au for COMP[29]041
# very concise but less readable code

@words = join(' ', <STDIN>) =~ /[a-z]+/ig;
print @words+0, " words\n";
```

Exercise: How many times does a word occur in standard input

Write a Perl script `count_word.pl` which counts the number of times a specified word is found 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 death <lyrics/Metallica.txt
death occurred 69 times
$ ./count_word.pl death <lyrics/Justin_Bieber.txt
death occurred 0 times
$ ./count_word.pl love <lyrics/Ed_Sheeran.txt
love occurred 218 times
$ ./count_word.pl love <lyrics/Rage_Against_The_Machine.txt
love occurred 4 times
```

Hint: modify the code from the last exercise.

Hint: the Perl functions `uc` & `lc` convert strings to lowercase & uppercase respectively.

When you think your program is working you can use `autotest` to run some simple automated tests:

```
$ 2041 autotest count_word
```

When you are finished working on this exercise you must submit your work by running **give**:

```
$ give cs2041 lab07_count_word count_word.pl
```

Sample solution for count_word.pl

```
#!/usr/bin/perl -w

# count occurrences in STDIN of specified WORD
# written by andrewt@unsw.edu.au for COMP[29]041

die "Usage: $0 <word>\n" if @ARGV != 1;

$specified_word = lc $ARGV[0];

$count = 0;
while ($line = <STDIN>) {
    $line = lc $line;
    my @words = split /[^\a-z]+/i, $line;
    foreach $word (@words) {
        $count++ if $word eq $specified_word;
    }
}

print "$specified_word occurred $count times\n"
```

Exercise: Do you use that word often?

Write a Perl script `frequency.pl` which prints the frequency with which each artist uses a word specified as an argument. So if Justin Bieber uses the word **"love"** 493 times in the 46583 words of his songs, then its frequency is $493/46583 = 0.0105832599875491$.

For example:

```
$ ./frequency.pl love
165/ 16359 = 0.010086191 Adele
189/ 34080 = 0.005545775 David Bowie
218/ 18207 = 0.011973417 Ed Sheeran
493/ 46589 = 0.010581897 Justin Bieber
217/ 27016 = 0.008032277 Keith Urban
212/ 26192 = 0.008094075 Leonard Cohen
57/ 38096 = 0.001496220 Metallica
4/ 18985 = 0.000210693 Rage Against The Machine
494/ 53157 = 0.009293226 Rihanna
89/ 26188 = 0.003398503 Taylor Swift
```

So of these artists, Ed Sheeran uses the word **"love"** most frequently. If you choose a word a randomly from an Ed Sheeran song the probability it will be "love" is just over 1 in a hundred (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 artist and word to store the word counts.

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

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

Hint: reuse code from the last exercise.

When you think your program is working you can use `autotest` to run some simple automated tests:

```
$ 2041 autotest frequency
```

When you are finished working on this exercise you must submit your work by running **give**:

```
$ give cs2041 lab07_frequency frequency.pl
```

Sample solution for `frequency.pl`

```
#!/usr/bin/perl -w

# print frequencies of specified words in lyrics files
# written by andrewt@unsw.edu.au for COMP[29]041

foreach $file (glob "lyrics/*") {
    my $artist = $file;
    $artist =~ s/.*\\//;
    $artist =~ s/\\.txt$//;
    $artist =~ s/_/ /g;
    open my $f, '<', $file or die "can not open $file: $!";
    while ($line = <$f>) {
        $line = lc $line;
        foreach $word ($line =~ /[a-z]+/g) {
            $frequency{$artist}{$word}++;
            $n_words{$artist}++;
        }
    }
    close $f;
}

foreach $word (@ARGV) {
    $word = lc $word;
    foreach $artist (sort keys %frequency) {
        my $f = $frequency{$artist}{$word}||0;
        my $n = $n_words{$artist};
        printf "%4d/%6d = %.9f %s\n", $f, $n, $f/$n, $artist;
    }
}
```

Exercise: When numbers get very small, logarithms are your friend

Now suppose we have the song line "**truth is beauty**". Given that David Bowie uses the word "**truth**" with frequency 0.000146727 and the word "**is**" with frequency 0.005898407, the word "**beauty**" with frequency 0.000264108; we can estimate the probability of Bowie writing the phrase "**truth is beauty**" as:

$$0.000146727 * 0.005898407 * 0.000264108 = 2.28573738067596e-10$$

We could similarly estimate probabilities for each of the other 9 artists, and then determine which of the 10 artists is most likely to sing "**truth is beauty**" (it's Leonard Cohen).

A sidenote: we are actually making a large simplifying assumption in calculating this probability. It is often called the [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 do this by adding the log of the probabilities of each word. For example, you calculate the log-probability of Bowie singing the phrase "**Truth is beauty.**" like this:

$$\log(0.000146727) + \log(0.005898407) + \log(0.000264108) = -22.1991622527613 = \log(2.28573738067596e-10)$$

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

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

```
$ ./frequency.pl fear
 2/ 16359 = 0.000122257 Adele
13/ 34080 = 0.000381455 David Bowie
 0/ 18207 = 0.000000000 Ed Sheeran
10/ 46589 = 0.000214643 Justin Bieber
 0/ 27016 = 0.000000000 Keith Urban
 4/ 26192 = 0.000152718 Leonard Cohen
39/ 38096 = 0.001023730 Metallica
26/ 18985 = 0.001369502 Rage Against The Machine
 3/ 53157 = 0.000056437 Rihanna
 3/ 26188 = 0.000114556 Taylor Swift
```

It is not useful to assume there is zero probability that **Ed Sheeran** would use the word **fear** in a song even though he hasn'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 Ed Sheeran using the word **fear** as $(0+1)/18205$ and the probability of Metallica using the word **fear** as $(39+1)/38082$. This is a simple version of [Additive smoothing](#).

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

```
$ ./log_probability.pl fear
log((2+1)/ 16359) = -8.6039 Adele
log((13+1)/ 34080) = -7.7974 David Bowie
log((0+1)/ 18207) = -9.8096 Ed Sheeran
log((10+1)/ 46589) = -8.3512 Justin Bieber
log((0+1)/ 27016) = -10.2042 Keith Urban
log((4+1)/ 26192) = -8.5638 Leonard Cohen
log((39+1)/ 38096) = -6.8590 Metallica
log((26+1)/ 18985) = -6.5556 Rage Against The Machine
log((3+1)/ 53157) = -9.4947 Rihanna
log((3+1)/ 26188) = -8.7868 Taylor Swift
```

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)/%6d) = %8.4f %s\n**")

When you think your program is working you can use **autotest** to run some simple automated tests:

```
$ 2041 autotest log_probability
```

When you are finished working on this exercise you must submit your work by running **give**:

```
$ give cs2041 lab07_log_probability log_probability.pl
```

Sample solution for **log_probability.pl**


```
#!/usr/bin/perl -w

# calculate log probability of an artist using a word
# written by andrewt@unsw.edu.au for COMP[29]041
# more concise but less readable solution

foreach $file (glob "lyrics/*") {
    my $artist = $file;
    $artist =~ s/.*\\//;
    $artist =~ s/\\.txt$//;
    $artist =~ s/_/ /g;
    open my $f, '<', $file or die "can not open $file: $!";
    while ($line = <$f>) {
        $line = lc $line;
        foreach $word ($line =~ /[a-z]+/g) {
            $frequency{$artist}{$word}++;
            $n_words{$artist}++;
        }
    }
    close $f;
}

foreach $word (@ARGV) {
    $word = lc $word;
    foreach $artist (sort keys %frequency) {
        my $f = $frequency{$artist}{$word}||0;
        my $n = $n_words{$artist};
        printf "log((%d+1)/%6d) = %8.4f %s\n", $f, $n, log(($f+1)/$n), $artist;
    }
}
```

Exercise: Who sang those words?

Write a Perl script `identify_artist.pl` that given 1 or more files, each containing part of song), prints the most likely artist to have sung those words.

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

Your program should produce exactly this output:

```
$ ./identify_artist.pl song?.txt
song0.txt most resembles the work of Adele (log-probability=-352.4)
song1.txt most resembles the work of Rihanna (log-probability=-254.9)
song2.txt most resembles the work of Ed Sheeran (log-probability=-206.6)
song3.txt most resembles the work of Justin Bieber (log-probability=-1089.8)
song4.txt most resembles the work of Leonard Cohen (log-probability=-493.8)
```

Hint: only read each file once. Store the data in a (2-dimensional) hash. If you read the files many times your program will be very slow and exceed autotest time limits.

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


```
$ ./identify_artist.pl -d song2.txt
song2.txt: log_probability of -206.6 for Ed Sheeran
song2.txt: log_probability of -210.8 for Adele
song2.txt: log_probability of -211.5 for Taylor Swift
song2.txt: log_probability of -211.7 for Keith Urban
song2.txt: log_probability of -215.0 for Leonard Cohen
song2.txt: log_probability of -215.4 for Rage Against The Machine
song2.txt: log_probability of -215.7 for David Bowie
song2.txt: log_probability of -217.2 for Justin Bieber
song2.txt: log_probability of -222.2 for Metallica
song2.txt: log_probability of -223.4 for Rihanna
song2.txt most resembles the work of Ed Sheeran (log-probability=-206.6)
```

Hint: you may like to create simpler input to use in debugging, for example:

```
$ echo Andrew Rocks >andrew_rocks.txt
```

```
$ ./log_probability.pl Andrew
log((0+1)/ 16359) = -9.7025 Adele
log((0+1)/ 34080) = -10.4365 David Bowie
log((0+1)/ 18207) = -9.8096 Ed Sheeran
log((0+1)/ 46589) = -10.7491 Justin Bieber
log((0+1)/ 27016) = -10.2042 Keith Urban
log((0+1)/ 26192) = -10.1732 Leonard Cohen
log((0+1)/ 38096) = -10.5479 Metallica
log((0+1)/ 18985) = -9.8514 Rage Against The Machine
log((0+1)/ 53157) = -10.8810 Rihanna
log((0+1)/ 26188) = -10.1731 Taylor Swift
$ ./log_probability.pl Rocks
log((0+1)/ 16359) = -9.7025 Adele
log((10+1)/ 34080) = -8.0386 David Bowie
log((0+1)/ 18207) = -9.8096 Ed Sheeran
log((1+1)/ 46589) = -10.0560 Justin Bieber
log((0+1)/ 27016) = -10.2042 Keith Urban
log((0+1)/ 26192) = -10.1732 Leonard Cohen
log((1+1)/ 38096) = -9.8547 Metallica
log((0+1)/ 18985) = -9.8514 Rage Against The Machine
log((2+1)/ 53157) = -9.7824 Rihanna
log((6+1)/ 26188) = -8.2271 Taylor Swift
$ ./identify_artist.pl -d andrew_rocks.txt
andrew_rocks.txt: log_probability of -18.4 for Taylor Swift
andrew_rocks.txt: log_probability of -18.5 for David Bowie
andrew_rocks.txt: log_probability of -10.4 for Adele
```

Hint: if a word appears multiple times its log-probability needs to be summed multiple times.

```
$ echo echo echo >echo.txt
```

```
$ cat echo.txt
echo echo
$ ./log_probability.pl echo
log((0+1)/ 16359) = -9.7025 Adele
log((0+1)/ 34080) = -10.4365 David Bowie
log((0+1)/ 18207) = -9.8096 Ed Sheeran
log((0+1)/ 46589) = -10.7491 Justin Bieber
log((0+1)/ 27016) = -10.2042 Keith Urban
log((0+1)/ 26192) = -10.1732 Leonard Cohen
log((0+1)/ 38096) = -10.5479 Metallica
log((14+1)/ 18985) = -7.1434 Rage Against The Machine
log((0+1)/ 53157) = -10.8810 Rihanna
log((1+1)/ 26188) = -9.4799 Taylor Swift
$ ./identify_artist.pl -d echo.txt
echo.txt: log_probability of -14.3 for Rage Against The Machine
echo.txt: log_probability of -19.0 for Taylor Swift
echo.txt: log_probability of -19.4 for Adele
echo.txt: log_probability of -19.6 for Ed Sheeran
echo.txt: log_probability of -20.3 for Leonard Cohen
echo.txt: log_probability of -20.4 for Keith Urban
echo.txt: log_probability of -20.9 for David Bowie
echo.txt: log_probability of -21.1 for Metallica
echo.txt: log_probability of -21.5 for Justin Bieber
echo.txt: log_probability of -21.8 for Rihanna
echo.txt most resembles the work of Rage Against The Machine (log-p
probability: -14.3)
```

When you think your program is working you can use **autotest** to run some simple automated tests:

```
$ 2041 autotest identify_artist
```

When you are finished working on this exercise you must submit your work by running **give**:

```
$ give cs2041 lab07_identify_artist identify_artist.pl
```

Sample solution for `identify_artist.pl`

```
#!/usr/bin/perl -w

# identify artists most likely to have sung lyrics using bag of words model
# written by andrewt@unsw.edu.au for COMP[29]041

foreach $file (glob "lyrics/*") {
    my $artist = $file;
    $artist =~ s/.*\\//;
    $artist =~ s/\\.txt$//;
    $artist =~ s/_/ /g;
    open my $f, '<', $file or die "can not open $file: $!";
    while (<$f>) {
        tr/A-Z/a-z/;
        foreach (/[a-z]+/g) {
            $frequency{$artist}{$_}++;
            $n_words{$artist}++;
        }
    }
    close $f;
}

@artists = keys %frequency;

$debug = 0;
foreach $file (@ARGV) {
    my %log_probability;
    if ($file eq '-d') {
        $debug = 1;
        next;
    }

    open my $f, '<', $file or die "can not open $file: $!";
    while (<$f>) {
        tr/A-Z/a-z/;
        foreach $word (/[a-z]+/g) {
            foreach $artist (@artists) {
                $log_probability{$artist} += log(((frequency{$artist}{$word}||0) + 1)/$n_words{$artist});
            }
        }
    }
    close $f;

    @sorted_artists = sort {$log_probability{$b} <=> $log_probability{$a}} @artists;

    if ($debug) {
        foreach $artist (@sorted_artists) {
            printf "%s: log_probability of %.1f for %s\n", $file, $log_probability{$artist}, $artist;
        }
    }

    printf "%s most resembles the work of %s (log-probability=%.1f)\n", $file, $sorted_artists[0], $log_probab
}
```

Submission

When you are finished each exercises make sure you submit your work by running **give**. You can run **give** multiple times. Only your last submission will be marked.

Don't submit any exercises you haven't attempted.

If you are working at home, you may find it more convenient to upload your work via [give's web interface](#).

Remember you have until **Saturday 5 January 19:00** to submit your work.

You cannot obtain marks by e-mailing lab work to tutors or lecturers.

You check the files you have submitted [here](#)

Automarking will be run by the lecturer several days after the submission deadline for the test, using test cases that you haven't seen: different to the test cases `autotest` runs for you.

(Hint: do your own testing as well as running `autotest`)

After automarking is run by the lecturer you can [view it here](#) the resulting mark will also be available via [via give's web interface](#)

Lab Marks

When all components of a lab are automarked you should be able to view the the marks [via give's web interface](#) or by running this command on a CSE machine:

```
$ 2041 classrun -sturec
```

The lab exercises for each week are worth in total 1 mark.

The best 10 of your 11 lab marks will be summed to give you a mark out of 9. If their sum exceeds 9 - your mark will be capped at 9.

- You can obtain full marks for the labs without completing challenge exercises
- You can miss 1 lab without affecting your mark.

COMP[29]041 18s2: Software Construction is brought to you by
the [School of Computer Science and Engineering](#) at the [University of New South Wales](#), Sydney.

For all enquiries, please email the class account at cs2041@cse.unsw.edu.au

CRICOS Provider 00098G