

# Homework 5. Higher order functions

Due September 27th at 10pm

## Part 1: HOF warm-ups (30 points)

Each of the problems below uses at least one of the three higher-order functions that we have discussed: map, foldl, and filter.

You **must** use one of these functions for each problem below; some may require multiple higher-order functions or multiple applications of higher-order functions.

In Part 1, you may not use any built-in methods besides the ones listed for each problem, and you may not use explicit recursion. If-statements are allowed in all of the problems.

### 1. replace (5 points)

Write a function called **replace** has three arguments: a string, a second string called *before*, and a third string called *after*. The function should go through the string and replace any occurrences of *before* with *after*.

```
> (replace "pigs are pink" "i" "u")\
  "pugs are punk"
```

#### Allowed methods:

- string-append
- string->list
- string
- equal?

### 2. flatten-one (5 points)

A flattened list is a list that does not contain nested lists. In order to flatten a list, nested lists are replaced by their contents in the top-level list.

Write a function that takes a mixed list and flattens it one level. For instance, if given a list '(1 (1 2) ((1 3))) your function should return '(1 1 2 (1 3)).

Call this function **flatten-one**.

#### Allowed methods:

- append
- list?
- list

### 3. split (5 points)

Write a function that takes two strings named *text* and *sep* and returns a list of strings. Every time that *sep* appears in *text*, the preceding portion of *text* is added to the return list.

```
>(split "hello-world" "-")
`("hello" "world")
```

**Allowed methods:**

- cons
- append
- equal?
- string->list
- string

### 4. sum of squares (5 points)

Write a function that returns the sum of squares of the items in a list of numbers. Call this sum-of-squares.

For instance:

```
>(sum-of-squares '(1 2 3))
14
```

**Allowed methods:**

- +
- \*

### 5. zip (5 points)

Write a function called **zip** which takes two lists as arguments, and returns a list of lists, where each nested list at index *i* contains the *i*th item of list 1 and the *i*th item of list 2.

```
>(zip '(1 2 3) '(4 5 6))
`((1 4) (2 5) (3 6))
```

If the lists are different sizes, truncate the longer list to be the same length as the shorter one (return as soon as either list is empty).

**Allowed methods:**

- list
- take
- min

## 6. Pivot-comparison (5 points)

Write a function that takes a comparison operator, a value, and a list of numbers, and returns all items from a list of numbers that satisfy the comparator with respect to the value.

For instance, if the comparison operator is `<=` and the value is 5, the function returns all items in the list that are less than or equal to 5.

Call the function `pivot-comparison`.

**Allowed methods:** None

## Part 2: Language modeling (70 points)

A language model is a model of sentence probability. One common kind of language model is an *n-gram* model, where the probability of a sentence is estimated as the joint probability of each *n*-length sequence of words. For instance, in a unigram model, the probability of the sentence is simply the probability of each of the words, multiplied together.

1. Unigram model example:

$$p(\text{"A cat meows"}) \approx p(\text{"A"}) p(\text{"cat"}) p(\text{"meows"})$$

In a bigram model, the probability of each word in the sentence is estimated by the probability of it occurring after the word that precedes it.

2. Bigram model example:

$$p(\langle s \rangle \text{ A cat meows } \langle s \rangle) \approx p(\text{"A"} | \langle s \rangle) p(\text{"cat"} | \text{"A"}) p(\text{"meows"} | \text{"cat"}) p(\langle s \rangle | \text{"meows"})$$

Given some text, we can estimate the probability of each bigram by counting how many times it occurs relative to every other bigram that begins with the same word. This is called a *maximum likelihood estimate*.

3. Maximum likelihood estimate for the bigram "cat meows":

$$\frac{\text{count}(\text{"cat"} \text{"meows"})}{\sum_{w \in \text{vocab}} \text{count}(\text{"cat"} w)}$$

In a bigram model, we simplify this equation further. The total count of bigrams that begin with the word "cat" is equal to the total number of times that the word "cat" appears in the training data.

4. Maximum likelihood estimate for the bigram "cat meows" (simplified):

$$\frac{\text{count}(\text{"cat"} \text{"meows"})}{\text{count}(\text{"cat"})}$$

In order to estimate the probability of each bigram, you'll need to count two things: the total number of occurrences for each word, and the total number of occurrences of each bigram.

In this lab, you will write a bigram language model in Racket. Your model will be capable of generating likely English sentences.

This chapter in Jurafsky & Martin's *Speech and Language Processing* textbook is a good reference on *n*-gram language models: <https://web.stanford.edu/~jurafsky/slp3/3.pdf>.

## Data

I have provided you with two data files. The first contains the text of *Alice in Wonderland*.<sup>1</sup> I recommend reading in only the first 10,000 characters; otherwise your program will take a long time to run.

The second file is a smaller test file containing the first three paragraphs of *The Hobbit*. I recommend testing your code on the smaller file as you work; I have given sample output for the smaller data file for many of the functions that you will write.

**You should make sure that all of your code also works on the larger *Alice in Wonderland* file.**

## 1 Reading in the data

The first step is to read in data from a text file. Use the read-string function to read in 10,000 characters from the text file, and store the result in a variable.

## 2 Preprocessing the data (15 points)

In order to train your model, you will have to write some text processing functions in order to get your training data into the right format.

### 2.1 Replace

Before text data can be fed into a language model, we often want to clean up the text by removing punctuation marks. This is because we might want to include both ‘Alice:’ and ‘Alice,’ in our counts for how frequently the word ‘Alice’ occurs.

Write a function called **strip-punctuation** that uses replace to make the following substitutions:

- Replace newlines with the empty string
- Replace commas with the empty string
- Replace colons with the empty string

Hint: you might find your **string-replace** function from Part 1 useful!

Use your **strip-punctuation** function to strip out the punctuation from the text that you read in.

### 2.2 Add sentence tags

Since we are interested in modeling sentence probabilities, we need some way of identifying sentences. We are going to assume that sentences boundaries occur whenever one of three punctuation marks is present: a period, a question mark, or an exclamation mark.

We will add an end-of-sentence boundary tag to mark where the end of each sentence is. This is useful in order to properly calculate the bigram probability of the first word in the sentence: we assume that it is preceded by the special end-of-sentence tag, and calculate its probability accordingly.

Our end-of-sentence tag will be "<s>".

---

<sup>1</sup>Made available by Project Gutenberg

Write a function called **add-tags** which takes a string and adds an end-of-sentence tag after every period, exclamation point, and question mark. **Make sure to add a space between the punctuation mark and the end-of-sentence tag.**

If you apply the `add-tags` function to the Hobbit text (having stripped out the punctuation specified in Question 2.1), the beginning of the text will be as follows:

```
"In a hole in the ground there lived a hobbit. <s> Not a nasty dirty wet hole filled with the ends of worms and an oozy smell nor yet a dry bare sandy hole with nothing in it to sit down on or to eat; it was a hobbit-hole and that means comfort. <s> It had ...
```

### 2.3 Splitting into words

Write a function called **get-words** that takes a string and splits it into a list of words. Your function should also filter out any empty strings.

If you apply your `get-words` function to the tagged text you created in the previous question, the first part of the result will be as follows:

```
('("In" "a" "hole" "in" "the" "ground" "there" "lived" "a" "hobbit." "<s>" "Not" ...
```

Hint: you might find your **split** function from Part 1 useful.

## 3 Creating bigrams (10 points)

Now that we have a list of words, with sentence boundaries, it's time to turn them into bigrams.

### 3.1 Make bigrams

Write a function that takes a list of words and returns a list of bigrams. Call this **make-bigrams**.

Hint 1: You might find your **zip** function from Part 1 useful.

Hint 2: You can add an end-of-sentence token to the beginning of your list so that the first word of the first bigram is "`<s>`".

If you apply your `make-bigrams` function to the list of words we generated from the Hobbit text, the first few bigrams should be as follows:

```
('(<s> "In") ("In" "a") ("a" "hole") ("hole" "in") ("in" "the") ("the" "ground") ...
```

### 3.2 Find unique bigrams

In this step, you'll write functions that produce two lists: one of unique bigrams, and the other, of the first word in each of the unique bigrams.

First, write a function called **get-unique-bigrams**. This function will simply call the **remove-duplicates** method on your list of bigrams to get a list of unique bigrams.

Next, write a function called **get-matching-unigrams**. This function should map (`lambda (x) (first x)`) over the list of unique bigrams in order to return a list containing the first word of each bigram.

You now have a list of unique bigrams, and a list of the first word in each of those bigrams.

## 4 Calculating bigram probabilities (15 points)

### 4.1 Counting n-grams

In this step, you will write a function to get n-gram counts for a list of unique n-grams. **Your solution should generalize to different values of n.**

The input to your **get-ngram-counts** function will be two lists: a list of unique n-grams, and a list of all of the n-grams in the text. For instance, in the unigram case, your input will be a list of unique unigrams and a list of all of the words in the text.

Your function will produce a list with one item for each item in the unique n-gram list. For each unique n-gram, your function will count how many times it appears in the list of all n-grams.

For instance, if the unique list is `'("cat" "mat" "sat" "on" "the")` and the all-n-grams list is `'("the" "cat" "sat" "on" "the" "mat" "on" "the" "mat" "sat" "the" "cat")`, `get-ngram-counts` would work as shown below:

```
> (define unique (list "cat" "mat" "sat" "on" "the"))
> (define all (list "the" "cat" "sat" "on" "the" "mat" "on" "the" "mat" "sat" "the" "cat"))
> (get-ngram-counts unique all)
```

```
`(2 2 2 2 4)
```

### 4.2 Get unigram and bigram counts

Once you have written your `get-n-gram-counts` function, use it to get the unigram counts and the bigram counts for your text.

For the Hobbit text, the beginning of the unigram counts will look like:

```
‘(11 1 15 3 5 17 ...
```

You can verify these counts by looking at the original text. For instance, the fourth unigram in our list is "hole", which occurs 3 times in the original text (plus one time where it is followed by a hyphen, which we didn't remove).

The beginning of the bigram counts for the Hobbit will look like:

```
‘(1 1 1 1 4 1 1 ...
```

(Bigrams re-occur much less often. The bigram 'in the' is the first bigram in our list that occurs more than once.)

### 4.3 Estimate bigram probabilities

Now that you have a list of bigram counts, you're ready to calculate bigram probabilities. Write a function called **calc-bigram-prob** which takes a list of bigram counts and list of unigram counts, and returns a list of bigram probabilities (using the formula given in the beginning of this assignment, repeated below).

5. Maximum likelihood estimate for the bigram "cat meows" (simplified):

$$\frac{\text{count}("cat" "meows")}{\text{count}("cat")}$$

Remember: you can rely on the fact that the list of unigram counts is in the same order as the list of bigram counts. That is, our unigram counts were calculated given the list of unigrams that correspond to the first word of each bigram in the unique bigram list.

If you use your `calc-bigram-prob` function on your Hobbit data, the result should start:

```
'(1/11 1 1/15 1/3 4/5 ...
```

## 5 Setting up the language model (15 points)

We now have all the components of our language model; we just have to put them together!

### 5.1 Store bigrams and their probabilities

I have provided you with a function called **store-probabilities**. This function takes a list of bigrams and a list of their probabilities, and stores them in a hash table.

This will allow you to look up the probability of a given bigram.

Create a bigram hash table using the `store-probabilities` function, your list of unique bigrams, and your list of their probabilities.

You can look up a given bigram by calling the function **hash-ref** as shown below:

```
> (define bigram-table (store-probabilities bigrams probs))
> (hash-ref bigram-table `("hobbit" "was"))
1/2
```

As you can see, the first argument to `hash-ref` is your hash table; the second argument is your bigram, which you provide in the form of a list.

### 5.2 Get relevant bigrams

Write a function that takes a word and a list of unique bigrams and returns a list of all bigrams that begin with that word. Call this function **get-relevant-bigrams**.

For instance, when you call `get-relevant-bigrams` with "hobbit", the following list is returned:

```
'(("hobbit" "was")("hobbit" "bedrooms")("hobbit" "and"))
```

### 5.3 Get most likely next word

The final step is to write a function that returns the most likely next word and its probability, given the preceding word. Call this function **get-likely-next-word**. Its arguments will be a word, a list of unique bigrams, and a hash table.

This function should find all bigrams that begin with the word, and then return the bigram from that list that has the highest probability along with its probability. The return type should be a list whose first item is the probability and whose second item is a list representing the bigram.

For instance, when given the input word "tunnel", the function returns a list containing the bigram "tunnel" "a" along with its probability, 1/3:

```
>(get-likely-next-word "tunnel" unique-bigrams table)
` (1/3 ("tunnel" "a"))
```

## 6 Generating likely sentences (10 points)

Let's generate some sentences!

### 6.1 Generate most likely sentence

Once you have your `get-likely-next-word` function written, you can run the **generate-sentence** function that I have provided. Its arguments are a word, a list of unique bigrams, a hash table, and a maximum sentence length.

The function generates the most likely sentence that starts with the specified word according to the language model. Because the language model can predict infinitely long sentences, it curtails the sentence at the specified maximum length and returns.

For instance, given the input word "hobbit" and a maximum length of 20, the model generates:

```
"hobbit was a very comfortable tunnel a very comfortable tunnel a very comfortable tunnel a very comfortable tunnel a very"
```

**According to your language model, what is the most likely sentence generated from the *Alice in Wonderland* data when the first word is the end-of-sentence tag ("`<s>`")?**

### 6.2 Sample sentences

It's pretty boring to generate only the single most probable sentence. To *sample* a sentence based on the probabilities of its bigrams, we need to use randomness. For convenience, we will use the Racket Math library's distribution representation: **discrete-dist**.

To use distributions, you will need to add a `(require math)` statement at the top of your program (after the `#Racket` declaration).

**discrete-dist** takes two arguments: a list of events and a list of their probabilities. In our case, we will use a list of bigrams and a list of their probabilities.

**discrete-dist** returns a distribution object, which can then be sampled using the **sample** function. **sample** takes two arguments: a distribution and an integer, which represents the number of samples to be drawn.

I've given you two functions. **sample-next-word** takes a word, a list of unique bigrams, a list of bigram probabilities, and a hash table, and returns a list whose first item is the probability and the second item is a list representing the bigram. **sample-sentence** is identical to **generate-sentence**, except that it calls **sample-next-word**.

**Use these functions to generate some sentences!**

## **Extra Credit**

### **1. filter**

Write filter using foldl. Call it my-filter.

### **2. map**

Write map using foldl. Call it my-foldl.

### **3. Estimating the probability of sentences**

A language model can also be used to estimate the likelihood of a given sentence. Write a function that takes in a sentence and estimates its probability according to the language model. Call this function **estimate-sent-prob**.

### **4. Generalizing to n-grams**

Produce a version of your language model program for n-grams.