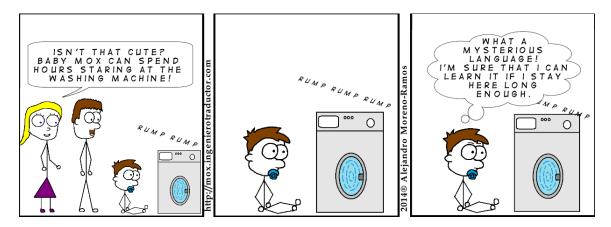
CS159 - Assignment 6

Due: Friday, April, 7 at 5pm



http://mox.ingenierotraductor.com/

In this assignment we will be implementing and experimenting with word alignment between sentence-aligned texts. This assignment has two parts to it. For the first part, you'll all be implementing a basic Model 1 word aligner. For the second part, there are a few options for how you can extend the aligner.

As with previous assignments, you may (and I would encourage you to) work with a partner. You may use whatever programming language you'd like as long as I can reasonably get it running on my computer.

Part 1: Model 1 Word Aligner

For the first part of this assignment you are to implement IBM Model 1 training algorithm that learns p(f|e) from a data set of aligned sentences and then submit a short writeup evaluating the results.

Implementation details

- Input

Create a script called wordalign.sh that takes four parameters:

word-align.sh <english_sentences> <foreign_sentences> <iterations> <probability_threshold>

See Appendix B for examples to help you.

- <english_sentences> A list of English sentences, one per line. You can assume the
 data has been preprocessed and should treat anything separated by a space as a "word".
 This will include numbers, punctuation, etc.
- <foreign_sentences> A list of foreign sentences, one per line. This file should have the same number of lines as <english_sentences> and there will be a line-wise correspondence between the files. The file has also been preprocessed, so just split on whitespace to get the "words".
- <iterations> The number of iterations of EM to run. An iteration includes one update of the model (i.e. p(e|f)).

- Output

Your program should train the model using the data provided for the number of iterations and then print out to standard output the learned probability table (i.e. p(f|e)) in the form:

```
<english_word> <foreign_word>                                                                                                                                                                                                                                                                                                                                                <
```

The output should be sorted alphabetically based on the combination of English word and foreign word (this will group them by English word) and only include entries whose probability is greater than or equal to the probability threshold. For example, the following would be the output from our example in class with 3 iterations and a probability threshold of $0.0 \ and \ NOT \ using the \ NULL \ alignment^2$:

| green | casa | 0.359999999999993 |
|-------|-------|---------------------|
| green | verde | 0.639999999999999 |
| house | casa | 0.6923076923076924 |
| house | la | 0.15384615384615385 |
| house | verde | 0.15384615384615385 |
| the | casa | 0.359999999999993 |
| the | la | 0.639999999999999 |

- Your code should include the "NULL" alignment (printed out in the probability table as NULL) that allows for foreign words to appear without having to be aligned to an English word. The good news is that for Model 1 this can be easily be done by just adding a special string

 $^{^{1}}$ I'll use the term "co-occur" throughout this document to mean when an English word e and foreign word f occur in the corresponding aligned sentences.

²Your program should include NULL, so your output of your final program won't match this, though you can remove NULL insertion and check to make sure it matches.

representing NULL as an additional word in each English sentence. The most common practice is to put it at index 0, but any location will do since word order doesn't matter. Just be careful that whatever string you choose to represent NULL does *not* occur anywhere in the actual English corpus.

- Sentences should be treated as sets of words, i.e., ignore repetitions. Since IBM model 1 ignores word order, this will simplify the calculations.
- You will need to keep track of the probability table p(f|e). In theory, this will contain all possible combinations of English and foreign words, initially with uniform probabilities. In practice, the number of words in each language can easily get into the 10K to 100K range, which makes storying all possible combinations impractical.

The key to fixing this is noticing that after the first iteration of the algorithm, many of these entries will become zero, in particular, any pair that does not co-occur in a sentence in the corpus will be set to zero after the first iteration. In your implementation you should not store all possible combinations of English and foreign words, instead only store those pairs that co-occur.

There are many ways to achieve this, but the easiest is to treat the first iteration of the EM training algorithm as a special case of the general iteration process. During the first iteration when you go through the corpus and are aggregating partial counts, rather than using entries in the word probability table (which would all be uniform anyway) just use a constant value for the probability (I used 0.01, but any constant that doesn't cause overflow is fine). You can use any constant since you will be normalizing the entries in the table at the end to be a proper probability distribution. The value of the constant will only change the value of this normalization constant, but not of the actual probabilities.

For the remaining iterations you can then just use the entries that actually are in the probability table, which will be all words that actually co-occur in the corpus.

- Your program must be able to run in a reasonable amount of time and with a reasonable amount of memory. My "basic" implementation takes about 10 seconds per EM iteration on 10K sentence pairs and can run easily with 2GB of memory, though if you need 3GB that's fined as long as you're in this same ballpark, you're in good shape. If you have problems running this on your own computer, let me know and I can help you run it on one of the CS servers.

Data

I've put some sentence-aligned data for you to play with in the usual place online at:

http://www.cs.pomona.edu/classes/cs159/assignments/assign6-starter.zip

I've included two aligned data sets:

• en-es.100k.(en|es) A parallel corpus of 100K English/Spanish sentences. This is a subset of the full Spanish/English corpus from the Europarl Corpus³, which has aligned data from

³http://www.statmt.org/europarl/

European Parliament proceedings. It has parallel corpora in 21 languages (all with English)!

• test.(en|fr) A very simple (three sentences) English/French corpus for testing your algorithm on.

Comments/Advice

- To make the algorithm extra clear, in Appendix A I have included a walkthrough of the first iteration of the example we looked at in class, following the "efficient" implementation approach. This should hopefully help you understand how we're going to implement it.
- Think about what data your program needs to store and how to store this most efficiently. For this assignment, you'll need to worry both about runtime efficiency as well as memory efficiency, though, you don't have to go overboard; just make reasonable decisions and you should be ok.
- As part of your algorithm you'll need to calculate:

$$p(f \to e) = \frac{p(f|e)}{\sum_{\hat{e} \text{ in } E} p(f|\hat{e})}$$

Make sure you calculate the denominator of this in a way that avoids repeated calculations for any given sentence pair.

- Since we're dealing with foreign text, you may have to deal with text encoding. For our basic Spanish data, you will have accented words. Most programming languages handle these things naturally, though you may have to explicitly process it using a particular encoding (in our case, UTF-8). For example, if you use the Scanner class, you will need to specify it.
- As always, develop your program incrementally and check things along the way. I'd advise first trying to get it working without the NULL alignment. Check your counts and other miscellaneous meta-data to make sure they're right and then work from there. You can then check your answer against the class notes. Once you have this working, insert the NULL feature. Think about how inserting NULL changes your answer and possibly work through it by hand.
- It's easy to get the EM algorithm "almost" right. Find ways to double-check your answers, e.g. a) work through one more example by hand (with or without the NULL) and b) do some sanity checks at different stages of the algorithm, like making sure that $\sum_f p(f|e) = 1$ for all e.

Part 2: Extending Your Word Aligner

For the second part of the assignment you need to extend your word aligner in *one additional* way listed below and then evaluate the impact of that extension. You only need to do one extension, so pick the one that seems the most interesting to you. Make sure that whatever extension you do, your program still functions correctly with the normal 4 inputs.

Extension option 1: memory efficiency

If you try and run the current version of your program on larger pieces of the en-es data set, you'll see the memory quickly climb to the point that eventually it won't function efficiently (or at all). My basic version can run the 100K sentences with 8GB of memory, but that's a lot. If you download the entire Spanish/English data set (1.9M sentences) and try it, you're not going to have a lot of success⁴.

Note, some of the things below may work better in some languages versus others. For example, if everything is an object (like in python) you can still get some savings from the tricks below, but not as much as if the language has "special" built-in data types like Java or C/C++.

There are a few ways that you can improve the memory efficiency (and to some extent, run-time efficiency) of your training algorithm:

- The biggest memory problem with the program is most likely because you left everything as a String. Strings are represented as an array of characters which is generally a byte per character (though it does depend on character encoding). In addition, Strings are classes in most languages and also incur additional memory overhead because of that.
 - A common fix for this is create an index mapping that associates each word in your corpus with a unique integer. Then, do all of your processing with integers (in Java use int not Integer) rather than Strings. First, do a preliminary pass and read through your corpus and create this vocabulary mapping. Then, start your actual algorithm and replace all your strings with ints.
- Things like lists (e.g. ArrayList) are very convenient, but they're very wasteful in memory since they reserve extra space to allow for amortized constant time adding (among other things). Arrays, on the other hand, are contiguous chunks in memory and generally don't use much more than the exact space required to store the data.
- Words that only occur a small number of times in a corpus don't provide much signal to the EM algorithm. So, rather than use all words that occur, you can filter your data set to only include words that occur X times or more (say 3-5 times). For training model 1, if all you're interested in is the p(f|e) table, you can actually just remove these words. If you're trying to do other things, it's better to replace these words with an unknown symbol (e.g. <UNK>) like we did with language modeling.

If you go this route, either include an optional 5th argument to your program that enables this or a constant within your program that needs to be set. By default, it should use all words.

Pick the first memory fix and at least one other fix to implement. Implement them and provide a short writeup with two things:

⁴I ran my basic version with 12GB of memory and it still ended up hitting the memory wall.

- 1. The output from the test example showing the same results as your basic version with your new, memory efficient version.
- 2. A table or graph comparing the memory (and run-time) usage of your basic input versus the improved version.

Extension option 2: actual alignments

For part A, all that we've calculated is the p(f|e) table, but not any actual alignments. Create a version of your program that can also calculate the alignment of any sentence pair in your training corpus. Then for your writeup, pick 10 random sentence pairs (you can pick random sentences under a certain length, to make this easier). Hand-align these sentences. Run them through your model and evaluate the precision and recall of your model. You may coordinate with other groups on generating the 10 random sentences and definitely can get help from a Spanish speaking expert (or Google translate).

If you go this route, either include an optional 5th argument to your program that enables this or a constant within your program that needs to be set. By default, it should just print out the p(f|e) table.

Extension option 3: more experimentation

Experiment and add to your writeup examining at least two of the components below. Describe your experimental setup, your results and some text interpreting the results.

- How does the number of iterations affect the performance of the algorithm? Calculate some metric that you can monitor as the number of iterations increases, e.g. sum of the difference in p(f|e) between iterations or the actual likelihood of the data under the model and plot this as the number of iterations increases. For the 10K en-es data, what would you suggest as a reasonable number of iterations.
- How does performance vary across languages? The Europarl corpus has 21 languages in it. Pick 2 others besides Spanish, train models and then compare the performance across some evaluation metric. You can use the random sampling approach from part A or some other reasonable metric. You'll need to preprocess the data in the other languages (lowercase and strip off punctuation as separate characters).
- In class I claimed that the initialization doesn't matter, i.e. initializing the p(f|e) with something besides uniform still gets you the same answer. Is this true? Verify this experimentally.
- Something of your own choosing. If you have questions about appropriateness, come talk to me.

Writeup

You should include a writeup that includes the following things, explicitly labeled:

- 1. The output of your program on the test sentences for 10 iterations and probability threshold of 0. You should also include a short (2-4 sentences) describing why the output makes sense. Describe not only whether the highest probabilities look reasonable (or not), but also other interesting characteristics like the probability associated with those that don't look good.
- 2. Run your program on the en-es data set with 10K sentence pairs (remember your command-line friend head) for 10 iterations and probability threshold of 0.3. Include 20 random entries from the output (when I say "random" I mean some computerized notion of random). Check the correctness of these values using a translation dictionary or online translation service (or even better, somebody who knows English and Spanish). Include a few sentences talking about the correctness/quality of your output. Highlight any mistakes/anomalies.
- 3. Writeup for Part 2, showing the results of your program addition. Note this will be different depending on what extension option you choose.

When you're done

When you're all done, submit via the normal submission mechanism.

Include all of your code and your writeup. Your script wordalign.sh should be in the base directory and should work from there. Please do not submit any of the data files.

Extra credit

You may implement one additional extension for 5 points of extra credit.

Commenting and code style

Your code should be commented appropriately (though you don't need to go overboard). The most important things:

- Your name (or names) and the assignment number should be at the top of each file you modify.
- Each class and method should have appropriate JavaDoc comments (doc strings if you do it in Python).
- If anything is complicated, put a short note in there to help me out if there are any issues.

This is a non-trivial assignment and it can get complicated, which makes code style and comments very important so that I can understand what you did. For this reason, you will lose points for poorly commented or poorly organized code.



Grading

| Part | | points |
|----------------|-------------|---------------|
| Part 1 | | |
| | Correctness | 20 |
| | Efficiency | 5 |
| | Writeup | 15 |
| Part 2 | | 25 |
| Comments/style | 5 | |
| extra credit | | 6 |
| total | | 70 + 5 extra |

Appendix A

In class, considered training on:

```
green house \leftrightarrow casa verde the house \leftrightarrow la casa
```

First, we can initialize the p(f|e) tables to be a constant (say 0.01) for only those words that co-occur in the corpus, so our three distributions will be:

```
p(casa|green) = 0.01
p(verde|green) = 0.01
p(casa|house) = 0.01
p(verde|house) = 0.01
p(la|house) = 0.01
p(casa|the) = 0.01
p(la|the) = 0.01
```

I'm writing them explicitly here, though you could actually just be doing this on the fly during your first iteration. Note that we do NOT normalize the distributions at this point.

Now, we go through the first iteration of the for loop which will aggregate the partial counts.

For the first sentence pair (p1):

```
p(casa -> green)_p1 = p(casa|green) / (p(casa|green) + p(casa|house)) = 1/2
p(casa -> house)_p1 = p(casa|house) / (p(casa|green) + p(casa|house)) = 1/2
p(verde -> green)_p1 = p(verde|green) / (p(verde|green) + p(verde|house)) = 1/2
p(verde -> house)_p1 = p(verde|house) / (p(verde|green) + p(verde|house)) = 1/2
```

For the second sentence pair (p2):

```
p(la \rightarrow the)_p2 = p(la|the) / (p(la|the) + p(la|house)) = 1/2
p(la \rightarrow house)_p2 = p(la|house) / (p(la|the) + p(la|house)) = 1/2
p(casa \rightarrow the)_p2 = p(casa|the) / (p(casa|the) + p(casa|house)) = 1/2
p(casa \rightarrow house)_p2 = p(casa|house) / (p(casa|the) + p(casa|house)) = 1/2
```

From these we aggregate all of the counts by iterating through all pairs of words that co-occur in each of the two sentences:

```
count(green, casa) = p(casa -> green)_p1 = 1/2
count(green, verde) = p(verde -> green)_p1 = 1/2
count(house, casa) = p(casa -> house)_p1 + p(casa -> house)_p2 = 1
```

```
count(house, verde) = p(verde -> house)_p1 = 1/2
count(house, la) = p(la -> house)_p2 = 1/2

count(the, casa) = p(casa -> the)_p2 = 1/2

count(the, la) = p(la -> the)_p2 = 1/2

count(green) = p(casa -> green)_p1 + p(verde -> green)_p1 = 1

count(house) = p(casa -> house)_p1 + p(casa -> house)_p2 + p(verde -> house)_p1 + p(la -> house)
count(the) = p(casa -> the)_p2 + p(la -> the)_p2 = 1
```

With these counts, we can then calculate the p(f—e) at the end of the first iteration:

```
p(casa|green) = count(green, casa) / count(green) = 1/2 / 1 = 1/2
p(verde|green) = count(green, verde) / count(green) = 1/2 / 1 = 1/2
p(casa|house) = count(house, casa) / count(house) = 1/2
p(verde|house) = count(house, verde) / count(house) = 1/2 / 2 = 1/4
p(la|house) = count(house, la) / count(house) = 1/2 / 2 = 1/4
p(casa|the) = count(the, casa) / count(the) = 1/2 / 1 = 1/2
p(la|the) = count(the, la) / count(the) = 1/2 / 1 = 1/2
```

If you compare these to the first iteration we calculated in the slides, you'll notice that these are the same!

Next, we would start the next iteration of EM. We would *clear* all of the counts and recalculate them again based on the new p(f—e) that we just calculated (the equations above would stay the same, but the values will change).

Appendix B

Here are two examples for scripts, one for Java programs and one for python programs.

Java

Assuming your .java files are in a package nlp.mt (i.e., in subdirectory nlp/mt/) and the name of you main class is nlp.mt.WordAligner, then you would put word-align.sh at the base directory with the following contents:

#! /bin/bash

```
javac nlp/mt/*.java
java -Xmx2G nlp.mt.WordAligner $1 $2 $3 $4
```

The first line compiles your program and then the second line calls it, passing along the 4 parameters that were passed to the script.

Python

For python, it's simpler and all you need to do is pass along the command-line parameters to your program. Assuming your .py file is in the same location as your script, it would look like:

#! /bin/bash

```
python word_aligner.py $1 $2 $3 $4
```