

Homework Assignment #2  
Due: Friday, 11 March 2016 at 19h00 (7 PM),

## Statistical Machine Translation

TA: Patricia A. Thaine (t3araujo@cdf.utoronto.ca); Hamed Heydari (t6heydar@cdf.utoronto.ca)

### 1 Introduction

This assignment will give you experience in working with  $n$ -gram models, smoothing, and statistical machine translation through word alignment. Knowledge of French is not required.

Your tasks are to build bigram and unigram models of English and French, to smooth their probabilities using add- $\delta$  discounting, to build a word-alignment model between English and French using the IBM-1 model, and to use these models to translate French sentences into English with a decoder that we provide.

The programming language for this assignment is Matlab.

### 2 Background

#### Canadian Hansard data

The main corpus for this assignment comes from the official records (*Hansards*) of the 36<sup>th</sup> Canadian Parliament, including debates from both the House of Representatives and the Senate. This corpus is available at `/u/cs401/A2.SMT/data/Hansard/` and has been split into **Training/** and **Testing/** directories.

This data set consists of pairs of corresponding files (`*.e` is the English equivalent of the French `*.f`) in which every line is a sentence. Here, sentence alignment has already been performed for you. That is, the  $n^{\text{th}}$  sentence in one file corresponds to the  $n^{\text{th}}$  sentence in its corresponding file (e.g., line  $n$  in `fubar.e` is aligned with line  $n$  in `fubar.f`). Note that this data only consists of sentence pairs; many-to-one, many-to-many, and one-to-many alignments are not included.

Furthermore, for the purposes of this assignment we have filtered this corpus down to sentences with between approximately 4 and 15 tokens to simplify the computational requirements of alignment and decoding. We have also converted the file encodings from ISO-8859-1 to ASCII so as to be usable within the (old) version of Matlab on CDF. This involved *transliterating* the original text to remove *diacritics*, i.e., accented characters (e.g., *Chrétien* becomes *Chretien*).

To test your code, you may like to use the samples provided at `/u/cs401/A2.SMT/data/Toy`.

#### Apostrophes in strings

In a Matlab string the apostrophe (') is represented by two apostrophes as in:

```
>> in = 'The possessor''s possessions';  
in =  
The possessor's possessions
```

If you write test sentences that contain apostrophes or need to specify an apostrophe in a regular expression, for example, be sure to use two apostrophes in this way.

## Valid dictionary keys

Matlab's `struct` type can in many ways be used as a dictionary or hash, since field names of these structures can be specified dynamically with string variables. However, Matlab forbids the use of certain characters (e.g., `'='`) from being part of a valid field name. This can cause problems if you want to use a word that contains one of these forbidden characters as a field name.

To help alleviate this situation we've provided you with `/u/cs401/A2_SMT/code/convertSymbols.m` that replaces these illegal symbols with counterparts defined in `/u/cs401/A2_SMT/code/csc401_a2_defns.m`. Before these definitions can be used, you must either call that file locally, or make a reference to the relevant (pre-initialized) structure (e.g., in the file `preprocess.m`, discussed later), as in

```
global CSC401_A2_DEFNS
```

The `CSC401_A2_DEFNS` structure also defines the start- and end-of-sentence markers `SENTSTART`, and `SENTEND`, respectively. Remember to consider these markers as words when taking your  $n$ -gram counts.

## Add- $\delta$ smoothing

Recall that the maximum likelihood estimate of the probability of the current word  $w_t$  given the previous word  $w_{t-1}$  is

$$P(w_t | w_{t-1}) = \frac{\text{Count}(w_{t-1}, w_t)}{\text{Count}(w_{t-1})}. \quad (1)$$

$\text{Count}(w_{t-1}, w_t)$  refers to the number of times the word sequence  $w_{t-1}w_t$  appears in a training corpus, and  $\text{Count}(w_{t-1})$  refers to the number of times the word  $w_{t-1}$  appears in that corpus.

Laplace's method of add-1 smoothing for  $n$ -grams simulates observing otherwise unseen events by providing probability mass to those unseen events by discounting events we have seen. Although the simplest of all smoothing methods, in practice this approach does not work well because too much of the  $n$ -gram probability mass is assigned to unseen events, thereby increasing the overall entropy unacceptably.

Add- $\delta$  smoothing generalizes Laplace smoothing by adding  $\delta$  to the count of each bigram, where  $0 < \delta \leq 1$ , and normalizing accordingly. This method is generally attributed to G.J. Lidstone<sup>1</sup>. Given a known vocabulary  $\mathcal{V}$  of size  $\|\mathcal{V}\|$ , the probability of the current word  $w_t$  given the previous word  $w_{t-1}$  in this model is

$$P(w_t | w_{t-1}; \delta, \|\mathcal{V}\|) = \frac{\text{Count}(w_{t-1}, w_t) + \delta}{\text{Count}(w_{t-1}) + \delta\|\mathcal{V}\|}. \quad (2)$$

## 3 Your tasks

### 1. Preprocess input text [5 marks]

First, implement the following Matlab function:

```
preprocess(inSentence, language)
```

that returns a version of the input sentence `inSentence` that is more amenable to training. For both languages, separate sentence-final punctuation (sentences have already been determined for you), commas, colons and semicolons, parentheses, dashes between parentheses, mathematical operators (e.g., `+`, `-`, `<`, `>`, `=`), and quotation marks. When the input `language` is `'english'`, merely re-implement those rules for separating possessives and clitics that you used in the `twtt.py` function of Assignment 1. Certain contractions are required in French, often to eliminate vowel clusters. When the input `language` is `'french'`, separate the following contractions:

---

<sup>1</sup>Lidstone, G. J. (1920) Note on the general case of the Bayes-Laplace formula for inductive or *a priori* probabilities. *Transactions of the Faculty of Actuaries* 8:182–192.

Type	Modification	Example
Singular definite article ( <i>le, la</i> )	Separate leading <b>l'</b> from concatenated word	<i>l'élection</i> $\Rightarrow$ <i>l' election</i>
Single-consonant words ending in e-‘muet’ (e.g., ‘dropped’-e <i>ce, je</i> )	Separate leading consonant and apostrophe from concatenated word	<i>je t'aime</i> $\Rightarrow$ <i>je t' aime</i> , <i>j'ai</i> $\Rightarrow$ <i>j' ai</i>
<i>que</i>	Separate leading <b>qu'</b> from concatenated word	<i>qu'on</i> $\Rightarrow$ <i>qu' on</i> , <i>qu'il</i> $\Rightarrow$ <i>qu' il</i>
Conjunctions <i>puisque</i> and <i>lorsque</i>	Separate following <i>on</i> or <i>il</i>	<i>puisqu'on</i> $\Rightarrow$ <i>puisqu' on</i> , <i>lorsqu'il</i> $\Rightarrow$ <i>lorsqu' il</i>

Any words containing apostrophes not encapsulated by the above rules can be left as-is. Additionally, the following French words should not be separated: *d'abord*, *d'accord*, *d'ailleurs*, *d'habitude*.

A template of this function has been provided for you at `/u/cs401/A2_SMT/code/preprocess.m`. Make your changes to a copy of this file and submit your version.

## 2. Compute *n*-gram counts [15 marks]

Next, implement a function to simply count all unigrams and bigrams in the preprocessed training data, namely:

```
LM = lm_train(dataDir, language, fn_LM)
```

that returns a special language model structure, LM, defined below. This function trains on all of the data files in `dataDir` that end in either ‘e’ for English or ‘f’ for French (which is specified in the argument `language`) and saves the structure that it returns in the filename `fn_LM` (it should take about 5MB).

The structure returned by this function should be called ‘LM’ and must have two fields: ‘uni’ and ‘bi’, each of which holds structures which incorporate unigram and bigram counts, respectively. Remember that you can dynamically create structures in Matlab. The fieldnames (i.e. *keys*) to the ‘uni’ structure are words and the values of those fields are the total counts of those words in the training data. The fieldnames to the ‘bi’ structure are words ( $w_{t-1}$ ) and their fields are structures. The fieldnames of those structures are also words ( $w_t$ ) and the values of those fields are the total counts of ‘ $w_{t-1}w_t$ ’ in the training data.

E.g.,

```
>> LM.uni.word = 5 % the word ‘word’ appears 5 times in training
>> LM.bi.word.bird = 2 % the bigram ‘word bird’ appears twice in training
```

A template of this function has been provided for you at `/u/cs401/A2_SMT/code/lm_train.m`. Note that this template calls `preprocess`. If this procedure is giving you errors of the type `Warning: ‘...’ exceeds MATLAB’s maximum name length of 63 characters...`, that is a sign that you need to fix tokenization in your preprocessing function.

Make your changes to a copy of the `lm_train.m` template and submit your version. Train two language models, one for English and one for French, on the data at `/u/cs401/A2_SMT/data/Hansard/Training/`. You will use these models for subsequent tasks.

### 3. Compute log-likelihoods and add- $\delta$ log-likelihoods [20 marks]

Now implement a function to compute the log-likelihoods of test sentences, namely:

```
logProb = lm_prob(sentence, LM, type, delta, vocabSize) .
```

This function takes `sentence` (a previously preprocessed string) and a language model `LM` (as produced by `lm_train`). If the argument `type` is either missing or is an empty string (‘ ’), this function returns the maximum-likelihood estimate of the sentence. If the argument `type` is ‘smooth’, this function returns a  $\delta$ -smoothed estimate of the sentence. In the case of smoothing, the arguments `delta` and `vocabSize` must also be specified (where  $0 < \delta \leq 1$ ).

When computing your MLE estimate, if you encounter the situation where  $\frac{Count(w_t w_{t+1})}{Count(w_t)} = 0/0$ , then assume that the probability  $P(w_{t+1} | w_t) = 0$  or, equivalently,  $\log P(w_{t+1} | w_t) = -\infty$ . Infinity in Matlab is represented by `Inf`. Use log base 2 (i.e. `log2()`).

A template of this function has been provided for you at `/u/cs401/A2_SMT/code/lm_prob.m`. Make your changes to a copy of the `lm_prob.m` template and submit your version.

We also provide you with the function `/u/cs401/A2_SMT/code/perplexity.m`, which returns the perplexity of a test corpus given a language model. You do not need to modify this function. Using the language models learned in Task 2, compute the perplexity of the data at `/u/cs401/A2_SMT/data/Hansard/Testing/` for each language and for both the MLE and add- $\delta$  versions. Try at least 3 to 5 different values of  $\delta$  according to your judgment. Submit a report, `Task3.txt`, which summarizes your findings. Your report can additionally include experiments on the log-probabilities of individual sentences.

### 4. Implement IBM-1 [25 marks]

Now implement the IBM-1 algorithm to learn word alignments between English and French words, namely:

```
AM = align_ibm1(trainDir, numSentences, maxIter, fn_AM).
```

This function trains on the first `numSentences` read in data files from `trainDir`. The parameter `maxIter` specifies the maximum number of times the EM algorithm iterates before being terminated. This function returns a specialized alignment model structure, `AM`, in which `AM.eng_word.fre_word` holds the probability (not log probability) of the word `eng_word` aligning to `fre_word`. In this sense, `AM` is essentially the  $t$  distribution from class, e.g.,

```
>> AM.bird.oiseau = 0.8 %  $t(\text{oiseau} | \text{bird}) = 0.8$ 
```

Here, we will use a simplified version of IBM-1 in which we ignore the NULL word and we ignore alignments where an English word would align with no French word, as discussed in class. So, the probability of an alignment  $A$  of a French sentence  $F$ , given a known English sentence  $E$  is

$$P(A, F | E) = \prod_{j=1}^{len_F} t(f_j | e_{a_j})$$

where  $a_j$  is the index of the word in  $E$  which is aligned with the  $j^{th}$  word in  $F$  and  $len_F$  is the number of tokens in the French sentence. Since we are only using IBM-1, we employ the simplifying assumption that every alignment is equally likely.

**Note:** The naïve approach to initializing `AM` is to have a uniform distribution over all possible English (`e`) and French (`f`) words, i.e., `AM.e.f` =  $1/||\mathcal{V}_F||$ , where  $||\mathcal{V}_F||$  is the size of the French vocabulary. Doing so, however, will consume too much memory and computation time. Instead, you can initialize `AM.e` uniformly over only those French words that occur in corresponding French sentences. For example,

given only the training sentence pairs

<i>the house</i>	<i>la maison</i>
<i>house of commons</i>	<i>chambre des communes</i>
<i>Andromeda galaxy</i>	<i>galaxie d'Andromede</i>

, you would initialize the structure `AM.house.la = 0.2`, `AM.house.maison = 0.2`, `AM.house.chambre = 0.2`, `AM.house.des = 0.2`, `AM.house.communes = 0.2`. There would be no probability of generating *galaxie* from *house*. **Note** that you can force `AM.SENTSTART.SENTSTART = 1` and `AM.SENTEND.SENTEND = 1`.

A template of this function has been provided for you at `/u/cs401/A2_SMT/code/align_ibm1.m`. You will notice that we have suggested a general structure of empty helper functions here, but you are free to implement this function as you wish, as long as it meets with the specifications above. Make your changes to a copy of the `align_ibm1.m` template and submit your version.

## 5. Translate and evaluate the test data [20 marks]

You will now produce your own translations, obtain reference translations from IBM BlueMix, and use the latter (with additional references) to evaluate the former, with a BLEU score. This will all be done in the file `evalAlign.m` (there is a very sparse template of this file at `/u/cs401/A2_SMT/code/`).

To decode, we are providing the function

```
english = decode( french, LM, AM, lmttype, delta, vocabSize ),
```

at `/u/cs401/A2_SMT/code/decode.m`. Here, `french` is a preprocessed French sentence, `LM` and `AM` are your English language model from Task 2 and your alignment model trained from Task 4, respectively, and `lmttype`, `delta`, and `vocabSize` parameterize smoothing, as before in Task 3. You do not need to change the `decode` function, but you may (see the Bonus section, below).

For evaluation, translate the 25 French sentences in `/u/cs401/A2_SMT/data/Hansard/Testing/Task5.f` with the `decode` function and evaluate them using corresponding reference sentences, specifically:

1. `/u/cs401/A2_SMT/data/Hansard/Testing/Task5.e`, from the Hansards.
2. `/u/cs401/A2_SMT/data/Hansard/Testing/Task5.google.e`, Google's translations of the French phrases<sup>2</sup>.
3. Translations from IBM BlueMix<sup>3</sup>. For this, simply call the equivalent of the following within the `unix` command in Matlab, where FUBAR is the (unprocessed) French *sentence*:

```
curl -u "{username}":"{password}" -X POST -F "text=FUBAR"
-F "source=fr" -F "target=en"
"https://gateway.watsonplatform.net/language-translator/api/v2/translate"
```

To evaluate each translation, use the BLEU score from lecture 6-2, i.e.,

$$BLEU = BP_C \times (p_1 p_2 \dots p_n)^{(1/n)} \quad (3)$$

Repeat this task with at least four alignment models (trained on 1K, 10K, 15K, and 30K sentences, respectively) and with three values of  $n$  in the BLEU score (i.e.,  $n = 1, 2, 3$ ). You should therefore have  $25 \times 4 \times 3$  BLEU scores in your evaluation. Write a short subjective analysis of how the different references differ from each other, and whether using fewer than 3 of them might be better (or worse).

In all cases, you can use the MLE language model (i.e., specify `lmttype = ''`). Optionally, you can try additional alignment models, smoothed language model with varying  $\delta$ , or other test data from other files in `/u/cs401/A2_SMT/data/Hansard/Testing/`.

Submit your evaluation procedure, `evalAlign.m`, along with a report, `Task5.txt`, which summarizes your findings. If you make any changes to any other files, submit those files as well.

<sup>2</sup>See <https://developers.google.com/api-client-library/python/apis/translate/v2>, but be prepared to pay.

<sup>3</sup>Note: you can only decode 1M characters/month.

## Bonus [up to 15 marks]

We will give up to 15 bonus marks for innovative work going substantially beyond the minimal requirements. These marks can make up for marks lost in other sections of the assignment, but your overall mark for this assignment cannot exceed 100%.

You may decide to pursue any number of tasks of your own design related to this assignment, although you should consult with the instructor or the TA before embarking on such exploration. Certainly, the rest of the assignment takes higher priority. Some ideas:

- Read about Good-Turing smoothing in section 6.2.5 of Manning and Schütze's *Foundations of Statistical Natural Language Processing*. This method is more sophisticated than add- $\delta$  smoothing and ensures that less of the overall probability mass is assigned to unseen events. This works much better in practice than add- $\delta$ . You may implement this method of smoothing and re-run the experiments in Task 3, above. Submit your code and an associated discussion.
- Implement the IBM-2 model of word-alignment, otherwise replicating Task 4 above. Ideally, translate the test data using this model and compute the error, as you did for Task 5. How does this model compare to IBM-1? Submit your code and an associated discussion.
- We have not considered the null word when performing alignments. Re-implement the IBM-1 alignment model to include null words and the possibility that no English word aligns with a French word (or vice versa). Submit your code and an associated discussion.
- Perform substantial data analysis of the error trends observed in each method you implement. This must go well beyond the basic discussion already included in the assignment. Submit a report.
- The decoder we use here is extremely simple and incomplete. You can write your own decoder that attempts to find  $\hat{e} = \arg \max_e P(e | f)$  using a heuristic  $A^*$  search, for example. Alternatively, what happens if you weight the contributions of the alignment and the language model to the overall probability? Section 25.8 of the Jurafsky & Martin textbook offers some ideas on how to improve the decoder. Submit your code and an associated discussion, comparing the decoded results to those performed with the default decoder.

## 4 General specification

We will test your code on different training and testing documents in addition to those specified above. Where possible, do not hardwire directory names into your code. As part of grading your assignment, the grader may run your programs using test harness Matlab scripts. It is therefore important that each of your programs precisely meets all the specifications and formatting requirements, including program arguments and file names.

If a program uses a file or helper script name that is specified within the program, it must read it either from the directory in which the program is being executed, or it must read it from a subdirectory of `/u/cs401` whose path is completely specified in the program. Do **not** hardwire the absolute address of your home directory within the program; the grader does not have access to this directory.

All your programs must contain adequate internal documentation to be clear to the graders. External documentation is not required.

## 4.1 Submission requirements

This assignment is submitted electronically. You should submit:

1. All your code for `preprocess.m`, `lm_train.m`, `lm_prob.m`, `align_ibm1.m`, and `evalAlign.m`, along with any other source files to which you made changes or which are necessary to run your code in Matlab on CDF.
2. Your alignment model trained on 1k sentences from `/u/cs401/A2_SMT/data/Hansard/Training/`, dumped in file `am.mat`.
3. Your reports `Task3.txt` and `Task5.txt`.
4. Any material submitted towards a bonus mark. This should be limited to code, results, and reports as text files.
5. Your ID file as described in Assignment 1. A template of ID is available on the course web page.

**You do not need to hand in your language models or other temporary files.** The electronic submission must be made from the CDF submission site. Do not **tar** or **compress** your files, and do not place your files in subdirectories.

## 5 Using your own computer

If you want to do some or all of this assignment on your laptop or other computer, you will have to do the extra work of downloading and installing the requisite software and data. You take on the risk that your computer might not be adequate for the task. You are strongly advised to upload regular backups of your work to CDF, so that if your machine fails or proves to be inadequate, you can immediately continue working on the assignment at CDF. When you have completed the assignment, you should try your programs out on CDF to make sure that they run correctly there. **A submission that does not work on CDF will get zero marks.**

## 6 Suggestions

This assignment uses a simplified version of an alignment model which itself makes several major simplifying assumptions and, as such, the results of the decoder will not be representative of the state-of-the-art in statistical machine translation. You will generally be marked on how well you understand the underlying concepts and algorithms. This approach was chosen for this assignment in order to give you a relative reprieve in the mid-term workload. However, if you have the time you are highly encouraged to pursue bonus work as indicated above. Exploring more complex models is not only interesting, but will give you a fuller perspective on the techniques used in machine translation.

The following dates are suggestions as to how to spread out the work for this assignment. These dates may not be applicable to you personally and they are not required deadlines. However, it's a good idea to try to spread things out so you don't have to rush at the end.

**Task 1** 19 February

**Task 2** 24 February

**Task 3** 29 February

**Task 4** 7 March

**Task 5** 11 March