

Homework #2: CMPT-825

Anoop Sarkar – anoop@cs.sfu.ca

For the programming questions use a `makefile` such that `make` compiles all your programs, and `make test` runs each program, and supplies the necessary input files. For the written questions, please submit either a \LaTeX or ASCII file.

(1) Grammar Development

Questions in English syntax have an interesting property. When constructing a question if we use a *wh*-word we have to front this *wh*-word to the start of the sentence. Consider the question variant of the transitive verb *likes*, where *Calvin* is the subject and *ice-cream* is the object of *likes*. The following sentences are an example of how the *wh*-word has to be fronted in English. The non-fronted variant is often called an *echo* question which means something quite different.

Calvin often likes ice-cream
*does Calvin often like what
what does Calvin often like

The following CFG attempts to capture *wh*-fronting in English questions (the CFG has start symbol *S*). Where the object usually appears we use a non-terminal *NP/NP* which derives the empty string. The new non-terminals *VP/NP* and *S/NP* are used instead of *VP* and *S*. These are simply new non-terminals without any special properties.

$S \rightarrow NP\ VP$	$VP \rightarrow V_{tms}\ NP$
$S \rightarrow WH\ does\ S/NP$	$WH \rightarrow what \mid who$
$S/NP \rightarrow NP\ VP/NP$	$NP \rightarrow Calvin \mid Hobbes \mid ice-cream$
$VP/NP \rightarrow V_{inf}\ NP/NP$	$V_{tms} \rightarrow likes \mid hates$
$NP/NP \rightarrow \epsilon$	$V_{inf} \rightarrow like \mid hate$
$VP/NP \rightarrow often\ VP/NP$	$VP \rightarrow often\ VP$

- a. Using the CFG above, provide parse trees for the following sentences:

Calvin often likes ice-cream
what does Calvin often like

- b. Extend the CFG given above by adding exactly **two additional binary rules** (at most two non-terminals in the rhs of the rule) to deal with the following type of question:

who often likes ice-cream

The word *who* must be derived using the $WH \rightarrow who$ rule and the subject of *likes* should be derived by the NP/NP rule.

- c. The syntax of coordination in English which combines constituents using conjunctions like *and*, *or*, etc. interacts with the syntax of questions. Extend the above CFG by adding exactly **five additional rules** so that the resulting CFG will correctly handle the following examples of coordination and questions in English. Your CFG should *accept* the sentences listed below that are *not* marked with * and it *does not* accept sentences listed below that *are* marked with *.

Calvin likes ice-cream and Hobbes and Calvin
 Calvin likes ice-cream and hates Hobbes
 Calvin likes ice-cream and hates Hobbes and likes ice-cream
 what does Calvin like and Hobbes hate
 what does Calvin like and Hobbes hate and Calvin hate
 what does Calvin like and hate
 what does Calvin like and hate and Hobbes like and hate
 Hobbes likes ice-cream and what does Calvin like and hate
 * what does Calvin like and Hobbes hates beans
 * what does Calvin likes ice-cream and Hobbes hate
 * what does Calvin likes and hates

- d. Consider a deterministic shift-reduce parser without any lookahead that uses the completed grammar. Explain briefly why such a parser could be led down the garden-path for the input: *Calvin likes Hobbes and Calvin hates ice-cream*.

(2) **Part-of-speech Tagging:**

Consider the task of assigning the most likely part of speech tag to each word in an input sentence. We want to get the best (or most likely) tag sequence as defined by the equation:

$$T^* = \operatorname{argmax}_{t_0, \dots, t_n} P(t_0, \dots, t_n \mid w_0, \dots, w_n)$$

- a. Write down the equation for computing the probability $P(t_0, \dots, t_n \mid w_0, \dots, w_n)$ using Bayes Rule and a trigram probability model over part of speech tags.
 b. We realize that we can get better tagging accuracy if we can condition the current tag on the previous tag and the next tag, i.e. if we can use $P(t_i \mid t_{i-1}, t_{i+1})$. Thus, we define the best (or most likely) tag sequence as follows:

$$\begin{aligned} T^* &= \operatorname{argmax}_{t_0, \dots, t_n} P(t_0, \dots, t_n \mid w_0, \dots, w_n) \\ &\approx \operatorname{argmax}_{t_0, \dots, t_n} \prod_{i=0}^n P(w_i \mid t_i) \times P(t_i \mid t_{i-1}, t_{i+1}) \text{ where } t_{-1} = t_{n+1} = \text{none} \end{aligned}$$

Explain why the Viterbi algorithm cannot be directly used to find T^* for the above equation.

- c. BestScore is the score for the maximum probability tag sequence for a given input word sequence.

$$\text{BestScore} = \max_{t_0, \dots, t_n} P(t_0, \dots, t_n \mid w_0, \dots, w_n)$$

It is a bit simpler to compute than Viterbi since it does not compute the best sequence of tags (no back pointer is required). For the standard trigram model $P(t_i \mid t_{i-2}, t_{i-1})$:

$$\text{BestScore} = \max_{t_0, \dots, t_n} \prod_{i=0}^n P(w_i \mid t_i) \times P(t_i \mid t_{i-2}, t_{i-1})$$

Assuming that $t_{-1} = t_{-2} = t_{n+1} = \text{none}$, we can compute BestScore recursively from left to right as follows:

$$\begin{aligned} \text{BestScore}[i+1, t_{i+1}, t_i] &= \max_{t_{i-1}} (\text{BestScore}[i, t_{i-1}, t_i] \times P(w_{i+1} \mid t_{i+1}) \times P(t_{i+1} \mid t_i, t_{i-1})) \\ &\quad \text{for all } -1 \leq i \leq n \\ \text{BestScore} &= \max_{\langle t, \text{none} \rangle} \text{BestScore}[n+1, \langle t, \text{none} \rangle] \end{aligned}$$

This algorithm for computing BestScore is simply the recursive forward algorithm for HMMs but with the sum replaced by max.

Provide an algorithm in order to compute BestScore for the improved trigram model $P(t_i | t_{i-1}, t_{i+1})$:

$$\text{BestScore} = \max_{t_0, \dots, t_{n-1}} \prod_{i=0}^{n+1} P(w_i | t_i) \times P(t_i | t_{i-1}, t_{i+1})$$

where $t_{-1} = w_{n+1} = t_{n+2} = t_{n+1} = t_n = \text{none}$, and assume that: $P(w_{n+1} = \text{none} | t_{n+1} = \text{none}) = 1$; $P(t_0 | t_{-1} = \text{none}, t_1) \approx P(t_0 | t_1)$; $P(t_{n+1} | t_n, t_{n+2}) = 1$ and $P(t_n | t_{n-1}, t_{n+1} = \text{none}) \approx P(t_n | t_{n-1})$

You can provide either pseudo code, a recursive definition of the algorithm, or a recurrence relation.

Hint: The first step would be to extend the recursive BestScore algorithm given above to read the input from right to left.

- d. Provide an algorithm in order to compute BestScore for the even more improved model:

$$\text{BestScore} = \max_{t_0, \dots, t_{n-1}} \prod_{i=0}^{n+1} P(w_i | t_{i-1}, t_i, t_{i+1}) \times P(t_i | t_{i-2}, t_{i-1}, t_{i+1})$$

(3) Hidden Markov models

Suppose you are given a large amount of text in a language you cannot read (and perhaps you are not even sure that it *is* a language). Faced with a total lack of knowledge about the language and the script used to transcribe it, you try to see if you can discover some very basic knowledge using unsupervised learning. Perhaps it is worth trying to see if the letters in the script can be clustered into meaningful groups.

Hidden Markov models (HMMs) are particularly suited for such a task. The hidden states correspond to the meaningful clusters we hope to find, and the observations are the characters of the text. The text provided to you is taken from the Brown corpus (which is in English, but we will try to “decipher” it anyway). Let us take an HMM with two hidden states, and the observations are to be taken from the list of 26 lowercase English characters plus one extra character for a single space character that separates the words in the text.

We will use the software `umdhmm` which implements the Baum-Welch algorithm for unsupervised training of HMMs. It has been installed in `~anoop/cmpt825/sw`. The web page for this software is:

<http://www.kanungo.com/software/software.html>

You can use the `esthmm` program to estimate the parameters of an HMM based on maximizing the likelihood of a training data sequence. The parameters of an HMM with n states and m observation symbols are: the probability of moving from one state to another (an n by n matrix A), the probability of producing a symbol from a state (an n by m matrix B), and the probability of starting a sequence from a state (a vector of size n).

In the notation used by `umdhmm` the parameters for $m = 2$ observation symbols and $n = 3$ states is:

```
M= 2
N= 3
A:
0.333 0.333 0.333
0.333 0.333 0.333
0.333 0.333 0.333
B:
0.5    0.5
0.75   0.25
0.25   0.75
pi:
0.333 0.333 0.333
```

While you do not need to implement the unsupervised training of an HMM for this question, you do need to interpret the HMM that is learned.

You are provided the training data from the Brown corpus and it is already in the correct format for the `esthmm` program. The file name for the training data is `brownchars.sf.txt` and it contains the first 50000 characters from the science fiction section of the Brown corpus. In this file, the numbers 1 to 26 are an index that specify the ASCII characters 'a' to 'z'. The number 27 is mapped to the space character, and in this data a single space character is used to separate words.

- a. Use the training corpus `brownchars.sf.txt` to train an HMM using `esthmm`. Provide the HMM that is learned. Try different initializations and see if the learned HMM is different.
- b. With careful initialization of the HMM, it is possible to train the HMM to automatically distinguish vowels (let us coarsely define these to be the ASCII characters *a, e, i, o, u*) and consonants (everything else) with the space character being in neither set. Provide a graph that clearly shows (visually) that the trained HMM has learned to distinguish one group (vowels) from the other (consonants).

(4) **Beam Search Thresholding**

This question will introduce you to the basics of statistical parsing. The task will be to speed up an existing parser written in Python in a sequence of steps listed below.

A recurring theme in statistical parsing is the need to limit the search space of the parser. Many productions that a hand-constructed context-free grammar (CFG) would rule out as impossible, a learned probabilistic context-free grammar (PCFG) might permit with low but non-zero probability. Techniques to limit the search space might still guarantee that the optimal parse is found (such as A* search). Or they might do more aggressive pruning for even more speed, but sometimes fail to find the most probable parse. The Beam Search Thresholding you will implement is an example of this second, more aggressive pruning technique.

You will be taking an existing probabilistic context-free CKY parser and adding beam search thresholding to improve its running time. The CKY algorithm without pruning is exhaustive, guaranteed to find every possible parse of a sentence. CKY for PCFGs (also called Viterbi for PCFGs) are guaranteed to find the parse with highest probability according to the grammar.

Beam search thresholding is a particular form of pruning which ignores items stored in the parser chart that have a very low probability, relative to other items in the same cell. Pruning these items is like making a bet that these items are so unlikely, they will not be in the most probable parse, even if they possibly might combine with other high-probability items for larger spans. The ratio in probabilities from the most probable item in a cell to the lowest probability permitted for a “competing” item is the Beam Width, and this parameter can be tuned to trade off speed (a narrower beam can prune more items) vs. accuracy (pruning more items risks pruning the most probable parse).

The hand-annotated (also called gold) reference parses, and the output of the provided parser, use the standard nested parentheses to represent the parse structure. The evaluation of the quality of parses uses the `evalb` tool written by Satoshi Sekine and Michael Collins, which reports the precision and recall of the brackets (the parse tree nodes and their spans) relative to the reference parse.

Here’s an example of what the evaluation is:

```
candidate: (S (A (P this) (Q is)) (A (R a) (T test)))
gold:      (S (A (P this)) (B (Q is) (A (R a) (T test))))
```

In order to evaluate this, we list all the labeled brackets with their spans, skipping the brackets of the form (A a), i.e. those brackets that have a span of 1 (note that evaluation of spans of length 1 would be equivalent to evaluating the part of speech tagging accuracy):

Candidate	Gold
(0,4,S)	(0,4,S)
(0,2,A)	(0,1,A)
(2,4,A)	(1,4,B)
	(2,4,A)

Precision is defined as $\frac{\#correct}{\#proposed} = \frac{2}{3}$ and Recall as $\frac{\#correct}{\#in\ gold} = \frac{2}{4}$.

The F-score (or balanced F-score) is defined as harmonic mean of Precision and Recall, defined as:

$$F = 2 \cdot \frac{P \cdot R}{P + R}$$

Use the `evalb` program to check the values of Precision and Recall for the above example. It is installed in the directory `EVALB` inside the `~anoop/cmpt825/sw` directory.

Here are the data files you will be using for the experiments:

- Grammar: `atis3.gr`
- Part of speech tags to be used for previously unseen words: `atis3.unseen`
- Sentences to be parsed: `atis3.test.sents`
- "Gold" reference parses for the input sentences: `atis3.gold`
- Pretty printing of parse trees: `indentrees.pl`

Notes on the grammar: The grammar was extracted from the Penn Treebank of Air Travel Information System (ATIS) articles with no smoothing of the count statistics. The grammar was converted into a form that is (almost) in Chomsky Normal Form. The implementation has some additional code to handle unary rules (rules of the form: $A \rightarrow B$), since the grammar used to annotate the treebank include many nodes with a single non-terminal as a (Right-Hand Side) RHS.

Optionally copy the above files to your directory. Copy the Python program `ckyParser-nopruning.py` as the file `ckyParser.py`. You will modify this Python program to complete the homework.

`ckyParser.py` allows you to set many options on the command line. The main ones are:

- v goes into verbose mode
- h prints a message on how to run the parser and quits
- g provides a grammar file to the parser (use `atis3.gr`)
- s provides a start symbol (TOP by default)
- u file contains tags used for unseen words (use `atis3.unseen`)

a. Run the provided parser on the input file:

```
python2.6 ckyParser.py -g atis3.gr -u atis3.unseen < atis3.test.sents > out
```

The code provided to you does the job of parsing and recovering the most likely parse tree using the CKY algorithm for probabilistic context-free grammars. However, it runs without any pruning and it is relatively slow even on these short ATIS sentences. The above command creates the output file `out` containing the parses for the input `atis3.test.sents`. Compare the output parses to the gold reference parses in `atis3.gold` using the program `indentrees.pl`:

```
cat out | perl indentrees.pl | more
cat atis3.gold | perl indentrees.pl | more
```

Check the accuracy of the parser output (put `evalb` in your PATH):

```
evalb -p COLLINS.prm atis3.gold out
```

The Bracketing Recall and Precision figures reported by `evalb` are the best the CKY parser can do with this training data. We will try to maintain this accuracy while trying to speed up the time taken by the parser.

- b. Extend the function `prune` in the `ckyParser` class so that after the parser outputs its best parse for a sentence, you can output the mean, max, and min number of items per chart cell for this sentence. Collect these statistics by modifying the `prune` function. Print these values to `sys.stderr` using the format:

```
# Items/cell mean: X.X max: N min: N Total items: N cells: N
```

- c. Extend the function `prune` further to also compute the maximum and minimum log-probabilities of the items in each cell, and the log-ratio of max:min probabilities in that cell. Output another line to `stderr` after the `Items/cell` output in the following format:

```
# Log max:min probability log-ratios mean: X.X max: X.X min: X.X
```

Collect these statistics by modifying the `prune` function. (Note: each non-empty cell has a single max:min ratio, and your output should show the statistics for this value across all non-empty cells.)

- d. Modify the function `prune` in the file `ckyParser.py`. Remove all items in each cell i, j whose probability is less than that of the most probable item in the same cell multiplied by β . The value of β is stored in the variable `beam` which can be assigned the value `log(value)` using the command line option `-b value` to the parser (the `-p` option enables pruning):

```
python2.6 ckyParser.py -p -b 0.0001
```

For example, let us assume that the value of the beam, $\beta = 0.0001$. Let us assume that the chart entry for the span $2, 4$ already contains the following non-terminals: NP, VP, PP. The list of non-terminals for each span i, j is stored in `chart[(i, j)]`. The log probability that a given non-terminal X derives the input from i to j is obtained by calling the function `chartGetLogProb(i, j, X)`.

In this example, let us take $P(\text{NP}, 2, 4) = 0.12$, $P(\text{VP}, 2, 4) = 0.03$, $P(\text{PP}, 2, 4) = 0.000011$. Beam thresholding would imply that we prune away $P(\text{PP}, 2, 4)$ since it is less than the maximum probability for the same span times the beam value: $P(\text{NP}, 2, 4) \times \beta$.

$P(\text{VP}, 2, 4)$ is not pruned since it sneaks by under the beam threshold.

Important: Remember that unlike the above example the values in the chart and the variable `beam` are stored as log probabilities.

Measure accuracy (using `evalb` as above) using beam thresholding with several different beam widths. Include the beam values 10^{-4} and 10^{-5} . You will observe that the F-score drops by a small amount but parsing times are reduced.

It's important to use the `-p` option which enables pruning in the `prune` function and to provide the beam value using the `-b` option. The source code that takes the values from the command line and assigns them to the right variables in the parser is provided to you already as part of `ckyParser.py`. You can simply use the variable `self.beam` in your code and assume that it has the right value taken from the command line.

- e. A problem with pruning using beam thresholding as we have done in the previous step is that some non-terminals can get the highest probability for a span, but when combined with other non-terminals in the right-hand side of a CFG rule, it might lead to a very improbable derivation. For example, $P(\text{FRAG}, 2, 4)$ might be very high but the rule $\text{VP} \rightarrow \text{FRAG PP}$ might be very improbable leading to search errors.

A solution to this problem is to use a rough estimate as to how likely a subtree will lead to full parse tree. This rough estimate can be as simple as the probability of a non-terminal independent of any context. This is the *prior* probability of a non-terminal. In the example above $P_{\text{prior}}(\text{FRAG})$ would be multiplied with the probability $P(\text{FRAG}, 2, 4)$ before we do beam thresholding as in the previous step. Multiply the prior probability for each non-terminal in the pruning process. The `getPrior` method in the `pcfg` class returns the prior probability for a non-terminal. Within the `prune` function the prior probability for non-terminal X can be obtained by calling `self.gram.getPrior(X)`.

Provide the accuracy of the parsing with prior probabilities combined with beam values of 10^{-4} and 10^{-5} .

- f. Plot the change in parsing time, comparing the parser without pruning and the parser with pruning at the various beam widths. Provide this graph which should show that your pruning strategy succeeds in speeding up parse times as you tighten the beam threshold value. Plot the output with and without the use of priors to check that adding the prior does not significantly slow the parser. Your output graphs should look similar to the following graphs (notice the use of error bars to show the variation in parsing times for sentences of the same length):

