**HUGH HAN, PEIYI ZHENG**

**600.465 Natural Language Processing**         **Fall 2016**
**Homework #3**         **Due:** 27 September 2016, 5:00pm

1. The log probability for each sample corpus, training on the `switchboard-small` corpus, is calculated as follows.

| | |
|---|---|
| -12111.3 | speech/sample1 |
| -7388.84 | speech/sample2 |
| -7468.29 | speech/sample3 |

To calculate the perplexity of each corpus, we can simply raise 2 to the power of each *negated* log probability. That is, we can simply do the following.

$$2^{12111.3} \quad \text{for} \quad \texttt{speech/sample1}$$
$$2^{7388.84} \quad \text{for} \quad \texttt{speech/sample2}$$
$$2^{7468.29} \quad \text{for} \quad \texttt{speech/sample3}$$

To calculate the perplexity per word, we first need to calculate the cross-entropy by dividing each negated log proability by the number of words in its respective corpus.

The number of words in each corpus is calculated as follows.

| | |
|---|---|
| 1686 | speech/sample1 |
| 978 | speech/sample2 |
| 985 | speech/sample3 |

Then we see that the perplexity per word of each file is the following.

$$2^{\frac{12111.3}{1686}} \quad \text{for} \quad \texttt{speech/sample1}$$
$$2^{\frac{7388.84}{978}} \quad \text{for} \quad \texttt{speech/sample2}$$
$$2^{\frac{7468.29}{985}} \quad \text{for} \quad \texttt{speech/sample3}$$

We obtain the following approximated perplexity-per-word values.

$$145.36 \quad \text{for} \quad \texttt{speech/sample1}$$
$$188.06 \quad \text{for} \quad \texttt{speech/sample2}$$
$$191.61 \quad \text{for} \quad \texttt{speech/sample3}$$

Training on the larger `switchboard` corpus, we get the following log probabilities for the following sample corpuses, respectively:

| | |
|---|---|
| -12561.5 | speech/sample1 |
| -7538.27 | speech/sample2 |
| -7938.95 | speech/sample3 |

Note that each of the log-probabilities grew more negative for its respective sample corpus. Thus, the each perplexity must also be smaller for its respective sample corpus. The reason for this is because if we increase the size of our training data, then we introduce a greater number of possible sentences. That is, the probability that the sentences in our specific test files appear would be much smaller, in that there is a greater variety of possibilities to choose from.

2. IMPLEMENTATION PROBLEM

3. We chose to do spam detection.

   When classifying the sample files containing genuine messages, we received the following output:

   ```
   178 looked more like gen_spam/train/gen (98.89%)
   2 looked more like gen_spam/train/spam (1.11%)
   ```

   When classifying the sample files containing spam messages, we received the following output:

   ```
   66 looked more like gen_spam/train/gen (73.33%)
   24 looked more like gen_spam/train/spam (26.67%)
   ```

   That is, we had the following error rates for the different data sets.

   $$\begin{aligned} \texttt{gen} &: \quad \text{error-rate} = 1.11\% \\ \texttt{spam} &: \quad \text{error-rate} = 73.33\% \end{aligned}$$

   (a) Each classification decision is based on probability. That is, whichever training corpus yields a higher cross-entropy with each test file is the training corpus to which that test file will be classified.

   However, we cannot necessarily deduce that the lowest cross-entropy occurs with the lowest error rate. With a low error rate, it could just be that the opposite incorrect corpus has a higher cross-entropy than the correct corpus, instead of the correct corpus having the lowest possible cross-entropy.

   Thus, we are forced to do some type of guess-and-check to find the value of the lowest cross-entropy. We proceeded by doing a human version of binary search.

   Using our binary search, we find the following two values of the approximate smallest log-probabilities.

   $$\begin{aligned} \texttt{gen} &: \quad \text{log-probability} \approx -423681 \\ \texttt{spam} &: \quad \text{log-probability} \approx -280362 \end{aligned}$$

   Next, we can count the number of words in all of the **gen** dev files and the number of words in all of the **spam** dev files using the following two commands, respectively.

   ```
   $ find gen_spam/dev/gen/ -name '*.txt' | xargs wc -w
   $ find gen_spam/dev/spam/ -name '*.txt' | xargs wc -w
   ```

We went up getting the values of 48198 and 39284 for the `gen` and `spam` dev files, respectively.

Using these values, we can then obtain the approximate smallest cross-entropies.

$$\texttt{gen} \quad : \quad \text{cross-entropy} \approx 8.79$$
$$\texttt{spam} \quad : \quad \text{cross-entropy} \approx 7.14$$

The $\lambda$ values that were used in finding these minimum cross-entropies were as follows.

$$\lambda_{\texttt{gen}} \quad = \quad 0.01357$$
$$\lambda_{\texttt{spam}} \quad = \quad 0.0083$$

(b) We can do binary search again, but let's think of how we can be a little smarter this time. We already have

$$\lambda_{\texttt{gen}} \quad = \quad 0.01357$$
$$\lambda_{\texttt{spam}} \quad = \quad 0.0083$$

so we can use them as critical points. We see that:

(i) all $\lambda < \lambda_{\texttt{spam}}$ cause increasing log-probabilities

(ii) all $\lambda > \lambda_{\texttt{gen}}$ cause increasing log-probabilities

Then the following must be true.

$$\lambda_{\texttt{spam}} \; \leq \; \lambda^* \; \leq \; \lambda_{\texttt{gen}}$$

So now we can do a smarter guess-and-check binary search, and we find the following value of $\lambda^*$.

$$\lambda^* \quad \approx \quad 0.0104$$

(c) First, let's try to classify the test files containing genuine messages.

```
347 looked more like gen_spam/train/gen (96.39%)
13 looked more like gen_spam/train/spam (3.61%)
```

Next, let's try to classify the test files containing spam messages.

```
50 looked more like gen_spam/train/gen (27.78%)
130 looked more like gen_spam/train/spam (72.22%)
```
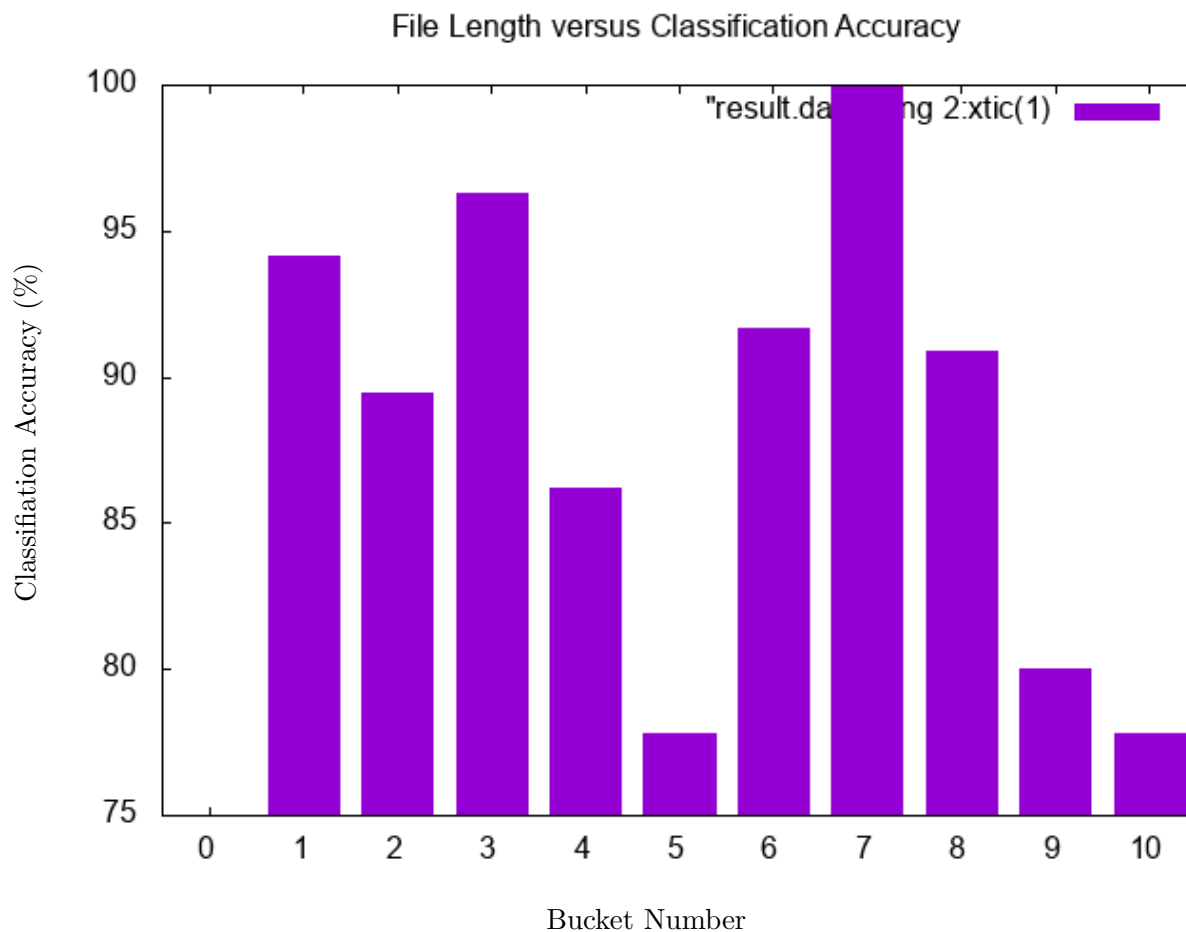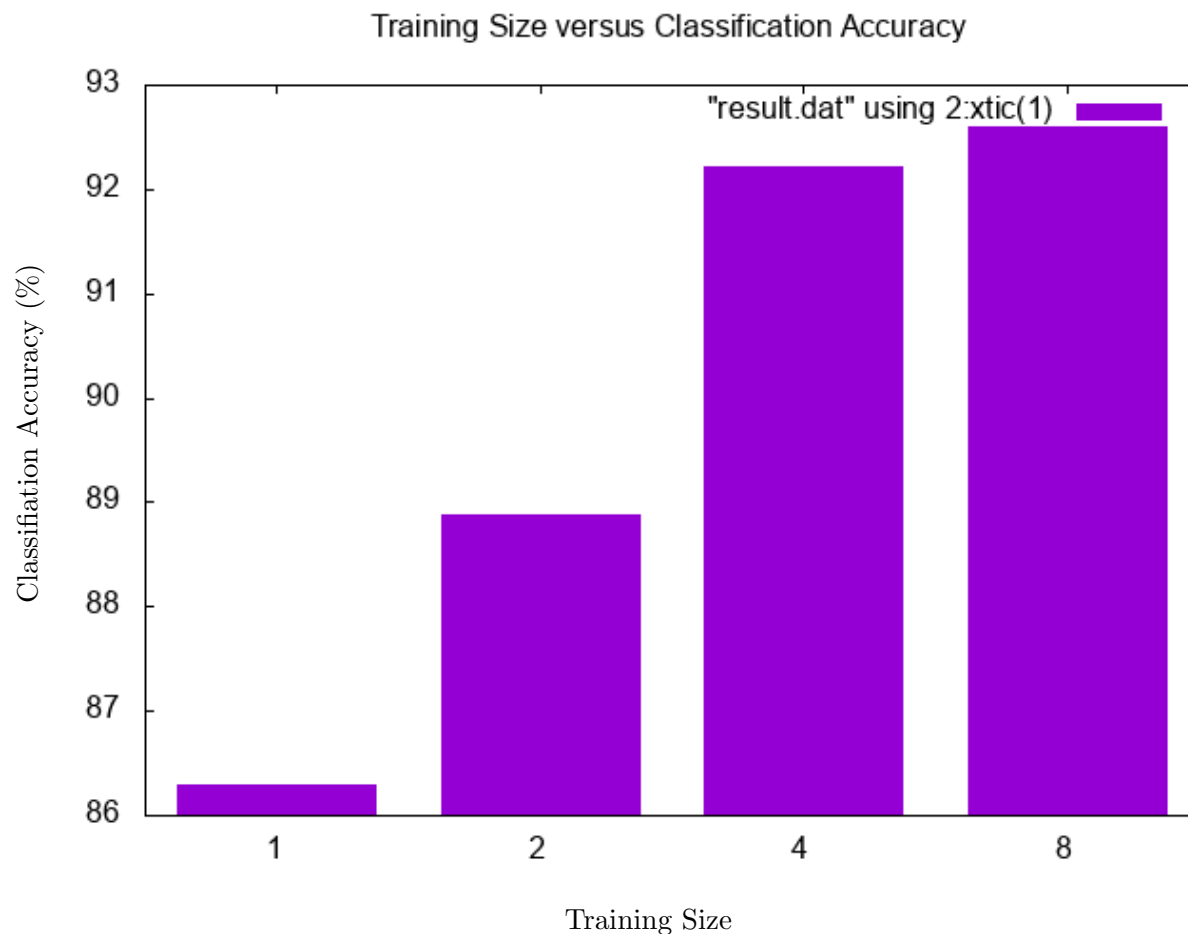
So then we get the following error rates.

$$\texttt{gen} \quad : \quad \text{error-rate} = 3.61\%$$
$$\texttt{spam} \quad : \quad \text{error-rate} = 27.78\%$$

(d) For this problem, we can all of the possible file lengths into 10 buckets of file length intervals. For instance, the first bucket would contain the classification results from the files with lengths in the range [0, 40), the second bucket would contain those in the range [40, 80), and so on. We can reserve our last bucket as a "special" bucket. That is, this special bucket will contain the classificatioon results of all of the files with lengths greater than 400. Now note that the first bucket can serve as a baseline, which represents the lowest accuracy of 75%. The reason for this is that small files containing less information, so they could be more difficult to classify.

The plot is shown below.



File Length versus Classification Accuracy

(e) It is easy to see that as the training data size increases, the classification accuracy also increases. But increasing the size of the training set does not necessarily solve all classification accuracy issues. At some point, it will level off. For example, it was noted that when we switched from `gen-times4` to `gen-times8`, the classification accuracy was improved only slightly.



Training Size versus Classification Accuracy

4. (a) Let's say that we take $V = 19{,}999$.

    i. Using the UNIFORM estimate, if we see a word that is out of the vocabulary, we would have no choice but to assign that OOV word a probability of 0, due to it being the 20,000th word of a vocabulary of size 19,999.
The reason for this is that the probabilities of all other words must sum to 1, and we cannot have a probability greater than 1.

$$\sum_{i=1}^{V} \frac{1}{V} \;=\; 1, \;\; V = 19{,}999$$

From the above equation, it is clear that we cannot assign OOV a probability of any value other than 0, if we are not to exceed a total probability of 1.
However, assigning OOV a probability of 0 is a bad estimation, because no word should have a probability of 0, regardless of whether that word was observed in the training data.

    ii. Using the ADDL estimate, a somewhat similar situation occurs. The ADDL probability is defined as follows.

$$\hat{p}(z \mid xy) \;=\; \frac{c(xyz) + \lambda}{c(xy) + \lambda V}$$

Now if we were to sum the probabilities $\hat{p}(z \mid xy)$ for all $z \in$ vocabulary, we would take a summation over 20,000 elements. However, if $V = 19{,}999$, we would end up getting a sum that is greater than 1, which canot be true due to the laws of probability. Hence, taking $V = 19{,}999$ is problematic.

(b) If we let $\lambda = 0$, then we are not modifying our probability at all. That is, we get the following from our probabiliy estimate.

$$\begin{aligned}
\hat{p}(z \mid xy) &= \frac{c(xyz) + \lambda}{c(xy) + \lambda V} \\
&= \frac{c(xyz) + 0}{c(xy) + 0} \\
&= \frac{c(xyz)}{c(xy)}
\end{aligned}$$

This ends up being the naive estimate of $\hat{p}(z \mid xy)$. If no smoothing occurs, then the probabilitiy given to any trigram that did not appear in the training corpus is 0, which is not well represntative of natural language.

(c)   i. First, let's think about the case when $c(xyz) = c(xyz') = 0$.

We see that

$$\hat{p}(z \mid xy) \;\; = \;\; \frac{0 + \lambda V \cdot \hat{p}(z \mid y)}{c(xy) + \lambda V}$$

Then the difference between $\hat{p}(z \mid xy)$ and $\hat{p}(z' \mid xy)$ is the possible differing values of $\hat{p}(z \mid y)$ and $\hat{p}(z' \mid y)$. And we cannot guarantee that $\hat{p}(z \mid y) = \hat{p}(z' \mid y)$.

For example, we might see neither of the phrases "kiss the tree" nor "kiss the Jupiter". However, when we backoff, it is possible that we would have seen both of the phrases "the tree" and "the Jupiter" (which are both much more common).

And in English, it is much more likely that the phrase "the tree" is more common than "the Jupiter", which demonstrates a case in which $\hat{p}(z \mid y) \neq \hat{p}(z' \mid y)$.

$\implies \hat{p}(z \mid xy) \neq \hat{p}(z' \mid xy)$.

  ii. Now let's think about the case when $c(xyz) = c(xyz') = 1$.

We see that

$$\hat{p}(z \mid xy) \;\; = \;\; \frac{1 + \lambda V \cdot \hat{p}(z \mid y)}{c(xy) + \lambda V}$$

Again, the difference between $\hat{p}(z \mid xy)$ and $\hat{p}(z' \mid xy)$ is the possible differing values of $\hat{p}(z \mid y)$ and $\hat{p}(z' \mid y)$. And again, we cannot guarantee that $\hat{p}(z \mid y) = \hat{p}(z' \mid y)$.

For example, consider the phrases "kick the Spongebob" and "kick the man", both of which have been counted exaclty once in the training corpora. But when we backoff, it is possible that we would have seen the "the man" much more frequently than "the Spongebob". Again, this is another case in which $\hat{p}(z \mid y) \neq \hat{p}(z' \mid y)$.

$\implies \hat{p}(z \mid xy) \neq \hat{p}(z' \mid xy)$.

(d) Again, we can write out the formula for the probability of a particular trigram using BACKOFF_ADDL smoothing.

$$\hat{p}(z \mid xy) \;\; = \;\; \frac{c(xyz) + \lambda V \cdot \hat{p}(z \mid y)}{c(xy) + \lambda V}$$

If we were to increase the value of $\lambda$, we would be smoothing at a greater rate. That is, we would be more heavily weighting the value of $\hat{p}(z \mid y)$ into our calculation of $\hat{p}(z \mid xy)$. So increasing the value of $\lambda$ would push our probabilities $\hat{p}(z \mid xy)$ closer to the probabilities $\hat{p}(z \mid y)$.

Now note that $\hat{p}(z \mid y)$ backs off onto $\hat{p}(z)$, which is calculated using the naive estimate. Also ote that the calculation of $\hat{p}(z \mid y)$ is *also* influenced by $\lambda$ (i.e. larger values of $\lambda$ push it towards $\hat{p}(z)$). So if we increase $\lambda$ to be large enough, we would eventually push the probabilities $\hat{p}(z \mid xy)$ towards $1/V$, or the uniform distribution.

5. (a) IMPLEMENTATION PROBLEM

   (b) Using $\lambda^* = 0.0104$, as calculated in problem **3**(b), we can calculate the cross-entropies for the switchboard corpora and analyze the text categorization error rates for `gen`/`spam`. The following generalizations are made with the assumption that $\lambda^* = 0.0104$.

   **cross-entropies for the switchboard corpora**

   In problem **1**, we calculated the log-probabilities using ADDL using a $\lambda$ value of 0.01. Because $\lambda^* = 0.0104$, we need to recalculate the log-probabilities.

   Below contains log-probabilities using ADDL for switchboard corpora.

   ```
   -12572.8   speech/sample1
   -7544.45   speech/sample2
   -7944.22   speech/sample3
   ```

   Below contains log-probabilities using BACKOFF_ADDL for switchboard corpora.

   ```
   -10040.1   speech/sample1
   -6039.31   speech/sample2
   -6413.95   speech/sample3
   ```

   Now to calculate the cross-entropies of the sample files, we must count the number of words in each sample file. We can do so using the following command.

   ```
   $ find speech -name 'sample*' | xargs wc -w
   ```

   ```
   1686   speech/sample1
   978    speech/sample2
   985    speech/sample3
   ```

   Then we can find the cross-entropies of each sample file by dividing the negative log-probability of each sample file by its corresponding word count.

   Below contains the cross-entropies of the switchboard corpora for each of the three samples, using ADDL and BACKOFF_ADDL smoothing.

   |                  | ADDL  | BACKOFF_ADDL |
   |------------------|-------|--------------|
   | speech/sample1   | 7.457 | 5.955        |
   | speech/sample2   | 7.714 | 6.175        |
   | speech/sample3   | 8.065 | 6.512        |

   For each of the three speech samples, we see that switching from ADDL smoothing to BACKOFF_ADDL smoothing decreases the cross-entropy greatly, against the training switchboard corpus.

   **text categorization error rates for gen/spam**

   Below contains output using BACKOFF_ADDL for `gen` test files.

   ```
   337 looked more like gen_spam/train/gen (93.61%)
   23 looked more like gen_spam/train/spam (6.39%)
   ```

Below contains output using BACKOFF_ADDL for `spam` test files.

```
17 looked more like gen_spam/train/gen (9.44%)
163 looked more like gen_spam/train/spam (90.56%)
```

From these outputs, we can see that the error rates using BACKOFF_ADDL were *significantly* lower than the error rates using ADDL as in **3**(c) for `spam` test files.

At the same time, the error rates using BACKOFF_ADDL were *slightly* higher than the error rates using ADDL as in **3**(c) for `gen` test files.

We can conclude that switching from ADDL to BACKOFF_ADDL causes the error rates for classifying spam files to decrease extraordinarily, but at the cost of slightly increasing error rates for classifying genuine files. It seems to balance out the classification correctness rates a quite a bit.

(c) *Extra credit:* We were able to find $\lambda_1 = 0.1$, which works a bit better than $\lambda^*$ in a sense. By using $\lambda_1$, our classification results on test data are exactly the same, but the cross-entropies are much smaller.

6. (a)
   (b)
   (c)
   (d)
   (e) *Extra credit.*
   (f)
   (g)  i.
       ii.
       iii.
       iv.
       v.
       vi.
       vii.

7. Let $w$ denote an arbitrary word. When we want to do the classification, we may want to consider the probability $P(w \text{ is spam} \mid w)$. According to Bayes's Theorem, we can write the following equation.

$$P(w \text{ is spam} \mid w) = \frac{P(w \mid w \text{ is spam}) \cdot P(\text{spam})}{P(w)}$$

$$= \frac{P(w \mid w \text{ is spam}) \cdot P(\text{spam})}{P(w \mid w \text{ is spam}) \cdot P(\text{spam}) + P(w \mid w \text{ is gen}) \cdot P(\text{gen})}$$

$P(w \text{ is spam} \mid w)$ and $P(w \mid w \text{ is gen})$ are the probabilities generated separately from the two models we train using the `spam` and `gen` corpora.

Now we know the following.

$$P(\text{spam}) = \frac{1}{3}$$

Then the following naturally follows.

$$P(\texttt{gen}) \;=\; 1 - P(\texttt{spam})$$
$$=\; \frac{2}{3}$$

We can rewrite then rewrite the following.

$$P(w \text{ is } \texttt{spam} \mid w) \;=\; \frac{1}{1 + \frac{2P(w \mid w \text{ is } \texttt{gen})}{P(w \mid w \text{ is } \texttt{spam})}}$$

By doing this, we can calculate the probability $P(w \text{ is } \texttt{spam} \mid w)$. Obviously we don't need *priori* when training the model.

In order to implement this method, we need to

i. train the **gen** model and store the probabilities for all test files,

ii. train the **spam** model and store the probabilities for all test files,

iii. iterate through all $P(w \mid w \text{ is } \texttt{gen})$ and $P(w \mid w \text{ is } \texttt{spam})$, so that we can calculate $P(w \text{ is } \texttt{spam} \mid w)$.

*Extra credit:* We implemented this change in `problem7.py`. We tested our implementation using an extremely small value for $\lambda$ (which was `add0.00001` in our case). From this specific test, the program successfully classified approximately $33.3\%$ of the test data to be spam. That is, the result is extremely close to *priori*.

8. (a) Using Bayes Theorem, we can express the probability $P(\vec{w} \mid U)$ in the following equation.

$$P(\vec{w} \mid U) \;=\; \frac{P(U \mid \vec{w}) * P(\vec{w})}{P(U)}$$

Now in order to maximize $P(\vec{w} \mid U)$, we can think about maximizing $log_2(P(\vec{w} \mid U))$. Using the rules of logarithms, note that the following statement is true.

$$log_2(P(\vec{w} \mid U)) \;=\; log_2(P(U \mid \vec{w})) + log_2(P(\vec{w})) - log_2(U)$$

Since $log_2(U)$ remains constant for all of the possible $\vec{w}$ being considered, we can ignore the final term on the right side of the equation when maximizing $log_2(P(\vec{w} \mid U))$.

That is, we only need to worry about maximizing the two other terms, $log_2(P(U \mid \vec{w})) + log_2(P(\vec{w}))$. Luckiily, $log_2(P(U \mid \vec{w}))$ is provided to us in the input files. Then all we need to do is calculate $log_2(P(\vec{w}))$ using our trigram model, which will measure the extent to which it looks like English.

So we iterate all 9 candidates and compute $log_2(P(U \mid \vec{w})) + log_2(P(\vec{w}))$, and can thus find the candidate with the highest probability $P(\vec{w} \mid U)$.

(b) IMPLEMENTATION PROBLEM

(c) Some trigrams in unrestricted data are very uncommon due to the existence of words like "uh", "um", or "huh" intertwined inside of the sentences.

Therefore in order to overcome this problem we need to utilize backing off. In pratice, we set the value of $\lambda$ to be 0.01 according to our experiments on development data. We can see the overall error rates in the table below.

|  | test/easy | test/unrestricted |
|---|---|---|
| 3-gram model | 0.141 | 0.380 |
| 2-gram model | 0.159 | 0.396 |
| 1-gram model | 0.210 | 0.408 |

9. *Extra credit.*

10. *Extra credit.*

    (a)

    (b)

    (c)

    (d)

    (e)   i.

         ii.

        iii.

    (f)

    (g)

11. *Extra credit.*