

Programming Project 1

This assignment is due by Monday October 2, 1:30pm (class start time).

Please submit both a hardcopy (in class) and electronically (via provide).

You can write your code in any programming language as long as the submitted code runs on *homework.eecs.tufts.edu* so that it can be tested as needed.

Overview: Unigram Model

Recall the unigram model from question 4 of assignment 1. A unigram model over a vocabulary of K words is specified by a discrete distribution with parameter μ . Under this model, the probability of the k -th word of the vocabulary appearing is given by μ_k . In this experiment we will use the unigram model to evaluate different learning methods, and to perform model selection.

Task 1: Model Training, Prediction, & Evaluation

In class we developed three methods for prediction: using the maximum likelihood estimate, using the MAP estimate, and using the predictive distribution. In this part, we evaluate the effectiveness of the three different methods. More specifically, we will use text training data to perform unigram model learning according to each method and then calculate the *perplexity* of the learned models on a held-out test data set. Perplexity is a standard effectiveness metric in probabilistic language modeling for scoring how well a model predicts a given collection of words (low perplexity values imply good performance). Perplexity is defined by

$$PP = p(w_1, w_2, \dots, w_N | \text{model})^{-\frac{1}{N}} \stackrel{\text{unigram}}{=} \exp \left(-\frac{1}{N} \sum_{i=1}^N \ln p(w_i) \right)$$

On the course webpage, you will find two files `training_data.txt` and `test_data.txt` each of size $N = 640,000$ words. We have pre-processed and “cleaned” the text so it does not require *any* further manipulation and you only have to read the space-separated strings from the corresponding ASCII text files.

For each size of training set in $[\frac{N}{128}, \frac{N}{64}, \frac{N}{16}, \frac{N}{4}, N]$, you should train a unigram model according to the three different methods discussed above (use the initial segment of the full training data). Use a Dirichlet distribution with parameter $\alpha = \alpha' \mathbf{1}$ as a prior (this is a scalar α' multiplied by a vector of ones) and set $\alpha' = 2$ for this part. Prediction equations for these models are given below.

To avoid having words in the test set that are not in your vocabulary start by building a dictionary from the entire train and test sets and use this vocabulary in the experiments. You should find $K=10000$ distinct words.

When run, your code should report the perplexities on the train set (i.e., the current portion being trained on) and test set under the three trained models for each train set size. Plot the results as a function of train set size (it is useful to plot all three methods together) and provide some observations. In your discussion of the results, please address the following:

- What happens to the test set perplexities of the different methods with respect to each other as the training set size increases? Please explain why this occurs.
- What is the obvious shortcoming of the maximum likelihood estimate for a unigram model? How do the other two approaches address this issue?
- For the full training set, how sensitive do you think the test set perplexity will be to small changes in α' ? why?

Task 2: Model Selection

Here we use the same data and dictionary as in task 1. In the previous part, we set the value of the hyperparameter, α' , manually. In this part, you will use the evidence function (question 4 assignment 1) and training data to select a value of α' . In general, direct maximization of the evidence function to determine the hyperparameter can be difficult. Here we only have one hyperparameter and can therefore use a “brute-force” grid search to select α' .

In particular, compute the log evidence at $\alpha' = 1.0, 2.0, \dots, 10.0$ for a training set of size $\frac{N}{128}$. Also, compute the perplexities on the test set (use the predictive distribution) at these same values. When run your code should output the list of evidence and perplexity values for each α .

Plot the log evidence and test set perplexity as a function of α' and discuss what you see. In your discussion of the results, please address the following:

- Is maximizing the evidence function a good method for model selection on this dataset?

Task 3: Author Identification

Can the unigram model identify authors? In this part, we apply the model to this problem. On the course web page you will find 3 additional files for this assignment. Each of them is a cleaned text version of a classic novel (thanks to the Gutenberg project).

To avoid having words in the test set that are not in your vocabulary start by building a dictionary from all 3 files and use this vocabulary in the experiments. You should find $K=18,251$ distinct words.

Train the model on `pg121.txt.clean` (use the predictive distribution with $\alpha' = 2$) and evaluate the perplexity on each of the other two texts. When run your code should output the perplexities for the two other texts. One of the test files is by the same author as the training file but the other is not. Was the model successful in this classification task?

Now repeat this evaluation but remove any word that appears less than 50 times in the training file (i.e., their count would be zero and the size of the file is similarly reduced). Was the model successful in this classification task with the reduced training file? Please report and discuss these results.

Additional Notes

- In Task 1, you'll need to handle $\ln(0) \triangleq -\infty$ as a special case in your code.
- Useful Equations:
 - Prediction using ML estimate: $p(\text{next word} = k\text{-th word of vocabulary}) = \frac{m_k}{N}$
 - Prediction using MAP estimate: $p(\text{next word} = k\text{-th word of vocabulary}) = \frac{m_k + \alpha_k - 1}{N + \alpha_0 - K}$
 - Prediction using predictive distribution: $p(\text{next word} = k\text{-th word of vocabulary}) = \frac{\frac{m_k + \alpha_k}{N + \alpha_0}}{\frac{m_k + \alpha_k}{N + \alpha_0}}$
 - Evidence: $\Pr(\text{Data}|\boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0) \prod_{k=1}^K \Gamma(\alpha_k + m_k)}{\Gamma(\alpha_0 + N) \prod_{k=1}^K \Gamma(\alpha_k)}$
 - In the above, m_k = No. of times k -th word of vocabulary appears in the document, N = Total no. of words in the document, and $\alpha_0 = \sum_{k=1}^K \alpha_k$
 - $\Gamma(x) = (x - 1)!$
- Aside from standard I/O, math, and plotting libraries, **NO** external libraries/modules should be used for this assignment. Note that since the argument for the $\Gamma()$ function are integers you can use factorials for the computation.

Submission

- You should **submit** the following items **both electronically and in hardcopy**:
 - (1) All your source code for the assignment. Please write clear code with documentation as needed. The source code should (i) run on *homework.eecs.tufts.edu*, (ii) run from the command line *without editing* with a single command (if there is more than one execution command required, include those commands in a Bash script which we can run), and (iii) output the requested results.
You can assume the data files will be available in the same directory as where the code is executed. Please use filenames as provided for the data. Please include a short README file with the code execution command.
 - (2) A PDF report on the experiments, their results, and your conclusions as requested above.
- For electronic submission, put all the files into a zip or tar archive, for example `myfile.zip` (you do not need to submit the data we give you). Please do not use another compression format such as RAR. Then submit using `provide comp136 pp1 myfile.zip`.
- Your assignment will be graded based on the clarity and correctness of the code, and presentation and discussion of the results.