

# Assignment 1

**Due:** Oct 1, 2025, 11:59pm

**Groups** Groups are specified on eLearning and fixed (make sure you know your group name/number and members). You will submit reports as the same group to Gradescope. You should include (1) your group name which is specified on eLearning (e.g. “Assignment Groups01”) on the top of the report; (2) each group member’s name and netID. If we identify cases of students that do not contribute equally to the group’s work, we will be penalizing the assignment.

## 1 Overall Goal

You will build an n-gram-based language model (with variations).

**Requirements&Advices** (1) **The report is important!** The report is where you get to show that you understand not only *what* you are doing but also why and *how* you are doing it (e.g. data structure). So be clear, organized and concise; avoid vagueness and excess verbiage. Spend time doing error analysis for the model(s). This is how you understand the advantages and drawbacks of the systems you build. Using the L<sup>A</sup>T<sub>E</sub>X template below is recommended (you can also use Microsoft Word or Google Doc if preferred, the final report should be a PDF); (2) You shall not share or make public any materials of the assignments.

**Report Template:** <https://www.overleaf.com/read/rhftbfxpvrmd>

## 2 Dataset

You are given (via eLearning) a corpus of reviews of Chicago hotels. Here is an example:

After Leaving some important documents in the room, I called and asked for the lost and found department... SEVEN TIMES over the course of a week. Eventually I had enough and asked for a manager and was put on hold then disconnected. Finally , a deceptively friendly operator PROMISED me she would have someone call me back in a few minutes, It never happened. So Avoid this place because after they rip you off once they’ll keep doing it later. They cannot be trusted at their word.

In the A1\_DATASET folder, you will find two files (i.e. `train.txt`, `validation.txt`) that contain the corpus for train/validation/test set reviews respectively. Note that **each line** in these files corresponds to **a single review**.

The project will proceed generally as follows in terms of system/code development, we require that the coding language to be *python*.

1. Write code to train unsmoothed unigram and bigram language models. We will use the provided corpora of train.txt (see Section 3).
2. Implement smoothing and unknown word handling (see Section 4).
3. Implement the Perplexity calculation (see Section 5).

## 3 Unsmoothed n-grams

To start, you will write a program that computes unsmoothed unigram and bigram probabilities from the training corpus.

You are given a tokenized opinion spam corpus as input (see Section 2). You may want to do additional preprocessing, based on the design decisions you make <sup>1</sup>. Make sure to explain preprocessing decisions

---

<sup>1</sup>you can refer to text tokenization and normalization mentioned in Sec 2.4 of J&M <https://web.stanford.edu/~jurafsky/slp3/2.pdf>

you make clearly in the report. Note that you may use existing tools just for the purpose of preprocessing but you must write the code for gathering n-gram counts and computing n-gram probabilities yourself. For example, consider the simple corpus consisting of the sole sentence:

*the students like the assignment*

Part of what your program would compute for a unigram and bigram model, for example, would be the following:

$$\begin{aligned} P(\textit{the}) &= 0.4, \quad P(\textit{like}) = 0.2 \\ P(\textit{the}|\textit{like}) &= 1.0, \quad P(\textit{students}|\textit{the}) = 0.5 \end{aligned}$$

**Preprocessing** The files included are already tokenized and hence it should be straightforward to obtain the tokens by using space as the delimiter. Feel free to do any other preprocessing that you might think is important for this corpus. Do not forget to describe and explain your pre-processing choices in your report.

## 4 Smoothing and unknown words

Firstly, you should implement one or more than one methods to handle unknown words. Then You will need to implement two smoothing methods (e.g. Laplace, Add-k smoothing with different k). Teams can choose any method(s) that they want for each. **The report should make clear what methods were selected**, providing a description for any non-standard approach (e.g., an approach that was not covered in class or in the class).

## 5 Perplexity on Validation Set

Implement code to compute the perplexity of a “development set.” (“development set” is just another way to refer to the validation set — part of a dataset that is distinct from the training portion.) Compute and report the perplexity of your model (with variations) on it. Compute perplexity as follows:

$$\begin{aligned} PP &= \left( \prod_{i=1}^N \frac{1}{P(w_i|w_{i-1}, \dots, w_{i-n+1})} \right)^{1/N} \\ &= \exp \frac{1}{N} \sum_{i=1}^N -\log P(w_i|w_{i-1}, \dots, w_{i-n+1}) \end{aligned}$$

where  $N$  is the total number of tokens in the test corpus and  $P(w_i|w_{i-1}, \dots, w_{i-n+1})$  is the n-gram probability of your model. Under the second definition above, perplexity is a function of the average (per-word) log probability: use this to avoid numerical computation errors.

**If you experimented with more than one type of smoothing and unknown word handling**, you should report and compare the perplexity results of experiments among some of them.

## 6 Grading Rubrics

- Group name, student names and netids of group members
- Clear description of implementation details (include the most important code pieces) in the report (60%).
  - Unigram and bigram probability computation (15%)
  - Smoothing (15%)
  - Unknown word handling (15%)
  - Implementation of perplexity (15%)

- Eval, Analysis and Findings – including but may not be limited to: (1) report training and validation set perplexity; (2) how the smoothing strategy affects the performance on validation set (30%);
- Details of programming library usage if any (6%)
- Brief description of the contributions of each group member (2%)
- Feedback for the project (2%) – too difficult/easy? how much time did you spend on it? does the project help with your understanding of the content? any other suggestions?

## 7 Report

You should submit a short document (fewer than 5 pages) that describes your work. Your report should contain a section for each of the Sections above as well as a short workflow section that explains *how* you divided up the work for the project among group members. Important: the role of this document is to show us that you understand and are able to communicate what you did, and how and why you did it.

You might think this means you must write a lot to explain something, it may be a bad sign. Describe the general approach, the data structures used (where relevant). Use examples, identify all design choices and justify them (e.g., how did you deal with unknown words? how did you perform smoothing?). If you tried any extensions, perform measurable experiments to show whether or not your extensions had the desired effect and why?

**Code** Please upload your code to a Github Repo and put the link in the report. Please also include snippets of your code in the report, if you think it makes things clearer. Include only relevant portions of the code (such as a few important lines from a function that accomplish the task that you are describing a nearby paragraph, to help us understand your implementation). Including boilerplate code or tedious, unnecessary blocks (e.g., reading files) only makes your report cluttered and hard to follow, so avoid it.

## 8 What to Submit

- Submit your report to Gradescope (resubmissions to Gradescope REMOVES the Group structure. So be sure to re-create (add your members) your Group with every resubmission to Gradescope.)
- Submit your code on eLearning as a .zip or .tar file.

**Posting of this assignment on public or private forums, services, and other forms of distribution is not allowed. © 2025 All rights reserved. Part of this assignment was adapted from Claire Cardie.**