

CS 540: Introduction to Artificial Intelligence Homework Assignment # 7

Assigned: 3/26
Due: 4/9 before class

Hand in your homework:

If a homework has programming questions, please hand in the Java program. If a homework has written questions, please hand in a PDF file. Regardless, please zip all your files into hwX.zip where X is the homework number. Go to UW Canvas, choose your CS540 course, choose Assignment, click on Homework X: this is where you submit your zip file.

Late Policy:

All assignments are due at the beginning of class on the due date. One (1) day late, defined as a 24-hour period from the deadline (weekday or weekend), will result in 10% of the total points for the assignment deducted. So, for example, if a 100-point assignment is due on a Wednesday 9:30 a.m., and it is handed in between Wednesday 9:30 a.m. and Thursday 9:30 a.m., 10 points will be deducted. Two (2) days late, 25% off; three (3) days late, 50% off. No homework can be turned in more than three (3) days late. Written questions and program submission have the same deadline.

Assignment grading questions must be raised with the instructor within one week after the assignment is returned.

Collaboration Policy:

You are to complete this assignment individually. However, you are encouraged to discuss the general algorithms and ideas with classmates, TAs, and instructor in order to help you answer the questions. You are also welcome to give each other examples that are not on the assignment in order to demonstrate how to solve problems. But we require you to:

- not explicitly tell each other the answers
- not to copy answers or code fragments from anyone or anywhere
- not to allow your answers to be copied
- not to get any code on the Web

In those cases where you work with one or more other people on the general discussion of the assignment and surrounding topics, we suggest that you specifically record on the assignment the names of the people you were in discussion with.

Question 1: N -gram Language Modeling [60 points]

In this question you will implement a chatbot by generating random sentences from your HW1 corpus using n -gram language models with Laplace smoothing.

We have created a vocabulary file `vocabulary.txt` for you to interpret the data, though you do not need it for programming. The vocabulary is created by tokenizing the corpus, converting everything to lower case, and keeping word types that appears three times or more. There are 4699 lines in the `vocabulary.txt` file.

Download the file `corpus.txt` from the homework website. This file has one word token per line, and we have already converted the word to its index (line number) in `vocabulary.txt`. Thus you will see word indices from 1 to 4699. In addition, we have a special word type OOV (out of vocabulary) which has index 0. All word tokens that are not in the vocabulary map to OOV. For example, the first OOV in the corpus appears as

```
392    are
1512   entirely
0      undermined
12     .
```

The words on the right are provided from the original essays for readability, they are not in the corpus. The word “undermined” is not in the vocabulary, therefore it is mapped to OOV and have an index 0. OOV represents a *set* of out-of-vocabulary words such as “undermined, extra-developed, metro, customizable, optimizable” etc. But for this homework you can treat OOV as a single special word type. Therefore, the vocabulary has $v = 4700$ word types. The corpus has 228548 tokens.

Note:

Please make sure the outputs are formatted exactly as described in this document. A sample test script is provided which compares your outputs with the expected outputs. Please make sure that all the test cases pass before submitting. No points will be awarded if the output doesn’t match exactly, even if the algorithm is correctly implemented.

Write a program **Chatbot.java** with the following command line format, where the commandline input has variable length¹ and the numbers are integers:

```
$java Chatbot FLAG number1 [number2 number3 number4]
```

(Part a, 5 points) Denote the vocabulary by the v word types w_0, w_1, \dots, w_{v-1} in the order of their index (so w_0 has index 0 and represents OOV, and so on). For this homework it is important that you keep this order so that we can automatically grade your code.

You will first create a *unigram language model* with *add-one Laplace smoothing*. This is a probability distribution over the vocabulary, for word type $w \in \{0, \dots, v-1\}$ the estimated probability is

$$p_i \equiv p(w = i) = \frac{c(w = i) + 1}{n + |V|},$$

where $c(w = i)$ is the count of word type i in the corpus (i.e. how many times w_i appeared). Note you need to estimate and store p_i for all v word types, including OOV: $p(w = OOV)$ is the fraction of 0’s in the corpus. Your $p(w)$ should sum to 1 over the vocabulary, including OOV.

¹We provide code skeleton that already handles variable length input.

When FLAG=100, number1 specifies the word type index i for w_i . You should print out two numbers on two lines: $c(w = i)$ and $p(w = i)$. When printing the probabilities for this homework, keep 7 digits after the decimal point and perform rounding. For example,

```
$java Chatbot 100 0
7467
0.0320174
```

```
$java Chatbot 100 1
36
0.0001586
```

```
$java Chatbot 100 2000
3
0.0000171
```

```
$java Chatbot 100 3001
140
0.0006045
```

```
$java Chatbot 100 4699
8
0.0000386
```

(Part b, 5 points) Now you implement random sampling from a probability distribution. That is, you will generate a random word type according to its unigram probability. Here is how you do it:

1. Given the multinomial distribution parameter $(p_0, p_1, \dots, p_{v-1})$, you split the interval $[0, 1]$ into v segments. Segment 0 is $[l_0 = 0, r_0 = p_0]$. Note it is closed on the left. Segment i (for $i = 1, \dots, v - 1$) is

$$\left(l_i = \sum_{j=0}^{i-1} p_j, r_i = \sum_{j=0}^i p_j \right).$$

Note these segments are open on the left but closed on the right. Also recall that we want you to order these segments by their word type index.

2. You generate a random number r uniformly in $[0, 1]$.
3. You check which segment r falls into, and output the index of that segment.

In order to test your code in a reproducible way, we will specify the random number r from commandline. Specifically, When FLAG=200, we provide number1 and number2 (we guarantee that $number2 \geq number1$), and you should let $r = number1/number2$ (remember to use Java 'double' here so you don't get an integer zero!) instead of a random r . Your code should output three numbers on three lines: the word type index i that this r selects, l_i the left end of w_i 's interval, and r_i the right end of w_i 's interval.

```
$java Chatbot 200 32 1000
0
0.0000000
0.0320174
```

```
$java Chatbot 200 321 10000
1
0.0320174
0.0321761
```

```
$java Chatbot 200 5000 10000
2364
0.4999528
0.5144953
```

```
$java Chatbot 200 99997 100000
4699
0.9999614
1.0000000
```

(Part c, 10 points) Now you will create a *bigram language model* of the form $p(w | h)$, where both w and h (the history) are word types in the vocabulary. Fixing h , $p(w | h)$ is a multinomial distribution over word types $w = 0, \dots, v - 1$, and is estimated as follows:

$$p(w | h) = \frac{c(h, w) + 1}{(\sum_{u=0}^{v-1} c(h, u)) + |V|},$$

where $c(h, w)$ is the count of the bigram (adjacent word pair) h, w in the corpus. These counts are obtained by letting the history start at the first word position in the corpus, then gradually moving the history one position later, until finally the (history, word) pair “use up” the corpus. For bigrams, that means history stops at the 2nd to last word position in the corpus. For example, if the corpus is “cake cake cake you want cake cake” then $c(\text{cake}, \text{you}) = 1$, $c(\text{cake}, \text{cake}) = 3$, $c(\text{cake}, \text{want}) = 0$. Note it is perfectly fine to estimate $p(w = i | h = i)$ for the same word type i . It is also perfectly fine if either w or h or both are OOV.

The above discussion is for a fixed h , where $p(w | h)$ is a multinomial distribution. You will need to do so for all possible $h = 0, \dots, v - 1$, so that you will end up with v multinomial distributions. This is where the sparse storage becomes important.

When FLAG=300, number1 specifies the history word type index h , and number2 specifies the word type index w . You should print out three numbers on three lines: $c(h, w)$, $\sum_{u=0}^{v-1} c(h, u)$, and $p(w | h)$. For example,

```
$java Chatbot 300 414 2297
1054
1082
0.1824628
```

```
$java Chatbot 300 0 0
406
7467
0.0334511
```

```
$java Chatbot 300 0 1
0
7467
0.0000822
```

```
$java Chatbot 300 2110 4240
115
917
0.0206516
```

```
$java Chatbot 300 4247 0
41
1435
0.0068460
```

(Part d, 10 points) Now you will use the same function in Part b to sample from a bigram given h . That is, instead of using the unigram probability $p(w)$, we fix some h and you will generate a random word type from $p(w | h)$. The method is the same, you just need to do more bookkeeping and record the segments separately for each history h . Specifically, for history h the segments are:

$$[l_{h0} = 0, r_{h0} = p(w = 0 | h)]$$

$$\left(l_{hi} = \sum_{j=0}^{i-1} p(w = j | h), r_{hi} = \sum_{j=0}^i p(w = j | h) \right), i = 1, \dots, v - 1.$$

Again, you should use sparse storage.

When FLAG=400, we provide number1 and number2 (we guarantee that $number2 \geq number1$), number3 is the word type for history h , and you should let $r = number1/number2$ to pick the corresponding word type w from $p(w | h)$. Your code should output three numbers on three lines: the word type index i that this r selects, l_{hi} the left end of w_i 's interval conditioned on h , and r_{hi} the right end of w_i 's interval conditioned on h .

```
$java Chatbot 400 0 100 414
0
0.0000000
0.0003459
```

```
$java Chatbot 400 1 100 414
54
```

```
0.0096852
0.0100311
```

```
$java Chatbot 400 98 100 414
4584
0.9799377
0.9801107
```

```
$java Chatbot 400 81 100 4697
3807
0.8099894
0.8102019
```

```
$java Chatbot 400 15 100 4442
710
0.1499684
0.1501793
```

(Part e, 10 points) Finally you create a trigram language model of the form $p(w | h_1, h_2)$, where now the history is the pair of word types h_1, h_2 in that order. Fixing h_1, h_2 , $p(w | h_1, h_2)$ is a multinomial distribution over word types $w = 0, \dots, v - 1$, and is estimated as follows:

$$p(w | h_1, h_2) = \frac{c(h_1, h_2, w) + 1}{(\sum_{u=0}^{v-1} c(h_1, h_2, u)) + |V|},$$

where $c(h_1, h_2, w)$ is the count of the trigram (adjacent word triple) h_1, h_2, w in the corpus. For the cake corpus $c(\text{cake}, \text{cake}, \text{you}) = 1$, $c(\text{cake}, \text{cake}, \text{cake}) = 1$, $c(\text{cake}, \text{cake}, \text{want}) = 0$, $c(\text{cake}, \text{you}, \text{want}) = 1$ and for $u \neq \text{want}$ we have $c(\text{cake}, \text{you}, u) = 0$, $c(\text{want}, \text{cake}, \text{cake}) = 1$ and for $u \neq \text{cake}$ we have $c(\text{want}, \text{cake}, u) = 0$.

When FLAG=500, number1 specifies the history word type index h_1 , number2 is h_2 , and number3 is w . You should print out three numbers on three lines: $c(h_1, h_2, w)$, $\sum_{u=0}^{v-1} c(h_1, h_2, u)$, and $p(w | h_1, h_2)$. For example,

```
$java Chatbot 500 23 12 123
0
0
0.0002128
```

```
$java Chatbot 500 5 660 3425
10
402
0.0021560
```

```
$java Chatbot 500 2799 556 2364
```

```
1
3
0.0004253

$java Chatbot 500 414 2297 2364
99
1054
0.0173792

$java Chatbot 500 0 0 0
35
406
0.0070505
```

(Part f, 10 points) Now you will sample from the trigram model $p(w \mid h_1, h_2)$. When FLAG=600, we provide number1 and number2 (we guarantee that $number2 \geq number1$), number3 is h_1 and number4 is h_2 , and you should let $r = number1/number2$ to pick the corresponding word type w from $p(w \mid h_1, h_2)$. When this conditional probability is defined, your code should output three numbers on three lines: the word type index i that this r selects, $l_{h_1, h_2, i}$ the left end of w_i 's interval conditioned on h_1, h_2 , and $r_{h_1, h_2, i}$ the right end of w_i 's interval conditioned on the history. Otherwise, your code should output a single line with text *undefined*.

```
$java Chatbot 600 2 5 660 3425
1881
0.3999151
0.4001274

$java Chatbot 600 2 5 3001 104
1880
0.3998304
0.4000424

java Chatbot 600 50 100 496 4517
2340
0.4997911
0.5000000

$java Chatbot 600 33 100 2591 2473
1530
0.3298473
0.3300565

$java Chatbot 600 0 100 2297 414
0
0.0000000
```

```
0.0002128
```

```
$java Chatbot 600 0 100 496 4517
0
0.0000000
0.0002089
```

(Part g, 10 points) Now the fun begins! You will generate random sentences using your n -gram language models. But for building a chatbot, we will specify a sentence prefix s_1, s_2, \dots, s_t which are t initial words (represented by word type indices) in the sentence. Your code will complete this sentence as follows:

1. set seed for randomizer
2. Repeat:
 - (a) $h_1 = s_{t-1}, h_2 = s_t$
 - (b) generate a random word s_{t+1} from $\tilde{p}(w | h_1, h_2)$
 - (c) $t = t + 1$ // shifts the trigram history by one position in the next iteration.
3. Until the generated word is a period, or a question mark, or an exclamation mark.

Note in step 1(b) a complication arises from the sentence prefix, and we introduced a placeholder \tilde{p} :

- The sentence prefix is empty. In this case, simply let $\tilde{p}(w | h_1, h_2)$ be the unigram model $p(w)$ which does not require any history.
- The sentence prefix has only one word s_1 . In this case, let $\tilde{p}(w | h_1, h_2) = p(w | h = s_1)$ the bigram model.
- The sentence prefix $h_1 = s_{t-1}, h_2 = s_t$ as history is undefined for a trigram model. If so, let $\tilde{p}(w | h_1, h_2) = p(w | h_2)$ the bigram model.
- Otherwise, let $\tilde{p}(w | h_1, h_2) = p(w | h_1, h_2)$ the trigram model.

When `FLAG=700`, `number1=seed`, which is the seed of the random number generator; `number2=t` (which only needs to be 0, 1, or 2), and the next t numbers on the commandline specify the sentence prefix s_1, s_2, \dots, s_t . We will guarantee that s_i is not period, or a question mark, or an exclamation mark.

If `seed = -1`, you actually do not set the seed (this allows you to generate different random sentences). Otherwise you should set the seed to `seed`. To set the seed in Java, use the following code:

```
Random rng = new Random();
if (seed != -1) rng.setSeed(seed);
```

In step 1(b) each time you should generate a new random number $r \in [0, 1]$ in order to generate the random word. This should be done with


```
rng.nextDouble();
```

You should try your code multiple times with the same sentence prefix: when $seed = -1$ your code should complete the sentence in different ways; otherwise it should be the same completion.

Your code will output the completed sentence (starting at s_{t+1}), one word index per line.

Because of Laplace smoothing, the generated sentences can be too long. This is because Laplace smoothing assigns a default prior probability to all the terms even if they never appear. Hence, only partial sample outputs are shown below. For the complete outputs, please refer to the provided sample test case files.

```
$java Chatbot 700 0 0
3693
1118
2995
2587
2808
1566
1810
4628
4132
4423
...
```

```
$java Chatbot 700 1 1 523
3433
1927
976
1563
4548
28
4529
4417
4451
4404
...
```

```
$java Chatbot 700 31 2 2110 311
3434
1846
3693
3003
3508
4553
2315
985
```

2533
4521
...

(Part h, no points) For this part you do not need to develop your code any further, but you will test out the Chatbot that you have developed by actually “talking” to it.

Please download `ChatbotDriver.java` and place this file together with your `Chatbot.java` in the same directory. This driver class basically takes the user input, apply some rules to generate a prefix based on the input (or just use the input itself as the prefix), and call “`java Chatbot 700 -1 prefix`” to generate a response and visualize it as actual words (“OOV” will be displayed for OOV indices). You can compile both files and try chatting with your Chatbot. Also make sure you have the txt files in the directory as well.

Note: because of randomness, your results will differ.

```
$javac Chatbot.java ChatbotDriver.java
$java ChatbotDriver
You: What's your opinion on self-driving cars?
Chatbot: self-driving cars similar apply desired rarely describing unique behavior
proposal contributing inventions failing tolerance self-driven battle conjunction
center aspect space difficulties person population user promising 12 billion ...
You: healthcare system
Chatbot: healthcare system suggests unfair takes underlying validation corporation
investment capital falling or database professionals protocol relationships broadly
animal selection .
You: tell me a joke
Chatbot: attempts analyzed conversation amount complete accent data gone integrate
require distance parts switch else doubts enacted ct at believes happening apparent
arrive forever subjects racial x-ray screens norm 2 alter eliminating ai exercise
surgical perception consent autonomy ...
You: say anything
Chatbot: simplified household refers regulatory processing differently executive
knows carpooling face choosing unlikely sold published peer state-of-the-art cultures
healthcare start tested hired devastating boston transition questionable staple fails
collisions portability compromised true prejudices kate crisis highways ...
You: what is your idea of an ideal world?
Chatbot: my idea of an ideal world choosing individualized home/service mentions
infinite knowledge 1999 things flint year disruptions switching professional mistake
uneven compatibility direction cited developing argues eliminating brains virtually
vehicles businesses accountability attract promise involves can shut buttons ...
```

Have fun and try to make it more intelligent by modifying `ChatbotDriver.java`! (e.g. adding more rules to the `generateCommand` method.)