

Homework 6 report
Sanad Saha

Evaluating Entropy (2 pts)

1. Like EX2, please first evaluate the trigram entropies (base and large) on EX2 test.txt

Q: Is the large version intuitively better than the small version? If so, why? If not, why?

Answer:

I've evaluated the trigram entropy for both base and large on test.txt from EX2 using the following command.

```
cat test.txt | sed -e 's/ /_/g;s/\(.\)/\1 /g' | awk '{printf("<s> %s\n", $0)}' | carmel -sribI trigram.(base/large).wfsa.norm
```

For base trigram the entropy was 2.9148 and for the large trigram wfsa it was 2.93042.

Here the large trigram version is performing poorly than the smaller version because the large trigram model was trained on wiki-2 corpus, which contains lots of lines outside of train.txt, which by a small margin impacts the entropy on test.txt; which is related to the train.txt

2. Now write a short program to evaluate the entropy of NLMs on the same test.txt.

Hint: they should be around 2.6 and 2.0, respectively.

Q: Which version (base or large) is better? Why?

Answer:

My eval_entropy.py calculates the entropy of the test.txt using the NLMs. Using the large version we get entropy of 1.9726965712505402 and using the base version we get entropy of 2.692813611492875

I used the following command to run the eval_entropy.py

```
cat test.txt | /scratch/anaconda3/bin/python3 eval_entropy.py
```

The large version of the NLM is performing significantly better than the base version. As the large version was trained on large corpus of wiki-2; the large n-gram from NLM has both seen more cases and probabilities calculated are more accurate.

3. Is NLM better than trigram in terms of entropy? Does it make sense?

Answer:

Yes, NLM is better than trigram in terms of entropy. This makes perfect sense because NLM has a n-gram language model, which is more accurate in predicting the next letter. Generally in English language it's easier to effectively predict the next letter when you have better model than trigrams (3-gram). The higher the n, in the n-gram language model, it gives better prediction.

3 Random Generation (2 pts)

1. Random generation from n-gram models.

Use `carmel -GI 10 <your_wfsa>` to stochastically generate character sequences. Show the results. Do these results make sense?

Answer:

Using the base trigram model, the following sequences were stochastically generated. I have added (`char_seq_trigram.txt`) which contains the results in a `.txt` file. The command is given below:

`carmel -GI 10 trigram.base.wfsa.norm > gen_seq_trigram.txt`

```
pelican04 ~/Courses/CS-539/hw6-data/nlm 161% carmel -GI 10 trigram.base.wfsa.norm

Using random seed -R 1659018150
<s> andfore do to by thas expers of coffedwan itenthe recan me
urt wal ally in </s> e^-145.267935669945
<s> shos by the offevist wr he in stand re pre diands gre vou </
s> e^-114.956481700927
<s> t to impan therchob to thes is docall con ond calan any do
escot trin the the muspers hit hattty ithripsectheat of thecto
f an a my jon ductolly by for liest sawit ame fon </s> e^-328.9276132136
17
<s> atifttiones </s> 1.47113912911454e-11
<s> the grown and bled ins vereviond to prostrace sure cal pa
rm all beep int sap thoes </s> e^-165.620802753212
<s> re ripse idere aria gares band beft istag </s> e^-96.281934071179
<s> vir kicis id of and a beet wall so caftigaregamistrod yocc
ome baum this with the de bad are memorinspies the bend the wh
fround at the se is only is forke itchom to entarine ingam a
r holl cur drove ind tany an hild wit strobvicared ne all whe
a sty unwo riter ces theve </s> e^-534.064279710993
<s> the unclivings frod tho den will tragats theivanicing tolf
</s> e^-118.312238445224
<s> week find in ther doducte ralsom and cound las appoes </s> e^-99
.4168337324265
<s> relp om spocatrue does </s> 1.2206204177118e-26
</s> e^-118.312238445224
```

The results makes sense. As we have done the same experiment in EX2, we know that trigram model generates some English like words but they are not actually correct English. Small English words such as week, find, the, for, and etc. are correctly generated; but the trigram model is not good enough to generate English words having more than 4 letters, let alone English sequence.

2. Write a short code to randomly generate 10 sentences (from `<s>` to `</s>`) from NLMs.

Hint: use `random.choices()` to draw the random sample from a distribution.

Answer:

My (`gen_seq.py`) generates 10 sentences using NLMs. I've provided 2 files `char_seq_ngram_base.txt` and `char_seq_ngram_large.txt`, which contains 10 sentences generated using base, large NLM model respectively.

`random.choices()` was used to draw random samples from the distribution. The following screenshot shows some sentences generated using `gen_seq.py`.

```

pelican04 ~/Courses/CS-539/hw6-data/nlm 163% /scratch/anaconda3/bin/python3 gen_seq.py
<s> strict_and_the_the_dissiled_and_all_on_the_move_the_comple
tes_and_the_pelised_a_bad_many_had_stroldan_were_to_the_proma
ssion_and_one_experition_to_fell_bean_to_like_and_doeet_the_re
centers_in_the_promations_and_when_the_intives_with_the_stor
ies_to_the_breger_parts_pruments_to_repout_a_were_of_the_sco
mpilies_the_will_may_you_course_allorie</s>

<s> there_his_he_the_politice_and_on_the_more_for_the_new_comp
late_the_start_in_the_properican_lalled_the_comening_hith_th
e_more_decondesters_of_most_be_marke_the_not_get_in_a_prever
y_and_the_was_the_marker_from_the_the_bit_of_the_bad_many_co
nterians_to_be_politice_and_asting_the_compesed_and_simpent
he_sinder_to_the_not_who_got_strates_to_allous_seet_stard_th
e_distate_and_contertica_has_angere_to_unice_that_with_hall_
and_the_most_in_the_propered_the_and_kertion_in_the_grow_the
betioner_distrate_be_are_trest_with_betore_of_the_ploses_th
e_read_the_best_in_a_comperies_with_the_stame_of_the_polan_s
und_the_stitn_of_shick_the_conded_want_the_more_the_becens_o
f_seased_a_comen</s>

<s> whinges_to_for_the_provested_when_the_rade_in_the_mist_and
good_deant_the_strike_sonored_the_are_the_real_the_stept_th
e_indered_in_the_urer_plant_to_get_strine_to_be_are_has_man_
i_complicions</s>

```

3. Compare the results between NLMs and trigrams.

Answer:

If we look closely the generated sentences between NLM and trigrams, we can see that NLM does exceedingly better job in generating English like sentences. Although the generated sentences are not grammatically perfect, but there are far more correctly spelled bigger words than the trigram model. For example: strict, there, strike etc. words were not possible to generate for a trigram model.

4

Restoring Spaces (4 pts)

1. Recall from EX2 that you can use LMs to recover spaces:

therestcanbeatotalmessandyoucanstillreaditwithoutaproblem
thisisbecausethehumanminddoesnotreadeveryletterbyitselfbutthewordasawhole.

First, redo the trigram experiments using our provided trigrams, and using Carmel. What are the precisions, recalls, and F-scores? (use eval_space.py).

Answer:

I have redone the recover space experiment of EX2, using both of our large and base trigram models in Carmel.

Carmel returned sentences that had <s>, </s> as starting and ending symbol and _ were used instead of spaces, with each character being separated by spaces. So, I used the following command to convert the results before evaluating with the test.txt file.

```
cat space_restored_large.txt | sed -e 's/ //g;s/_/
/g;s/<s>//g;s/<\s>//g' > space_restored_test.trigram.large
```

Using the base trigram model we achieve the following scores:

Precision = 60.8% Recall = 59.8% F1-score = 60.3%

```
atukan ~/Courses/CS-539/hw6-data 172% python eval_space.py test.txt space_restored_test.trigram.base
recall= 0.598 precision= 0.608 F1= 0.603
```

And using the large trigram model the following scores were achieved:

Precision = 65.2% Recall = 62.5% F1-score = 63.8%

```
atukan ~/Courses/CS-539/hw6-data 173% python eval_space.py test.txt space_restored_test.trigram.large
recall= 0.625 precision= 0.652 F1= 0.638
```

2. Then design an algorithm to recover spaces using NLMs. Note: you can't use dynamic programming any more due to the infinite history that NLM remembers. You can use beam search instead. Describe your algorithm in English and pseudocode, and analyze the complexity.

Answer:

The algorithm to recover spaces using NLM is described below using pseudocode. For each line, we would conduct Beam search to retrieve missing spaces for that line.

First we load our language model into `h` and initialize our Beam (`beam`) which doesn't have any entry at the beginning.

Pseudocode and description:

step1:

```
for each line in the given file:
    h = load the language model
    beam = [(0, h)] (Initialize the beam)
```

Step2:

Now for each character in the given line, we have 2 decisions to make,

```
for each character in the current line:
```

```
    temp = [] (Temporary Beam)
```

Step3:

For top K-best answer in the beam in the previous step

```
for (p, h) in beam:
```

1. We can add a space and then add the current character `c`.

```
    prob1 = p + math.log(h.next_prob("_"))
    h1 = h + "_"
    prob1 += math.log(h1.next_prob(c))
```

```
h1 = h1 + c
temp.append((prob1, h1))
```

or,

2. We can add the current character *c* without adding space.

```
prob2 = p + math.log(h.next_prob(c))
h2 = h + c
temp.append((prob2, h2))
```

We calculate probability and update the new language model given these two possible decisions are made. Then we add these to our temporary beam. This temporary beam holds all possible scenarios where spaces can be added, alongside their probability.

Step4:

Now holding all possible scenarios of adding spaces, is not efficient. We can discard those with low probabilities of occurring. So, we update our original Beam with the top 20 results from the temporary beam.

```
beam = sorted(temp, reverse=True)[:beam_len]
```

After, continuing step 3 until the end of the line, we finally get our top 20 candidates lines with spaces. And as we are sorting our beam at each step, the 1st answer in the Beam is the sentence which has the best probability of occurring with spaces.

```
p, h = beam[0]
```

We retrieve that and print the history of that sentence's language model, which gives us the line with restored spaces.

```
print(("".join(h.history)).replace("_", " ").replace("<s>", ""))
```

The above algorithm has a complexity of $O(n * m * b * 2)$, where *n* is the number of lines, *m* is the number of characters per line, *b* is the beam size, and for the 2 decisions we multiply it by 2.

As, 2 is a constant factor we can write the complexity as $O(n * m * b)$

3. Implement it, and report the precision, recalls, and F-scores.

Answer:

Our implementation of the algorithm (restore_spaces.py) described in part 2, achieves the following scores [using Beam size 20]:

Recall = 82.7% Precision = 80.0% F1-score = 81.3%

```
atukan ~/Courses/CS-539/hw6-data 188% python eval_space.py test.txt spaces_restored_base.txt
recall= 0.827 precision= 0.800 F1= 0.813
```

Recall = 96.7% Precision = 95.1% F1-score = 95.9%

```
atukan ~/Courses/CS-539/hw6-data 189% python eval_space.py test.txt spaces_restored_large.txt
recall= 0.967 precision= 0.951 F1= 0.959
```

5

Restoring vowels (4 pts)

1. Redo trigram vowel recovery and report the accuracy. (use eval_vowels.py)

Answer:

I used the following command to get restore vowel using the given trigram models and carmel.

```
cat test.txt.novowels | sed -e 's/ /_/g;s/\(.\)/\1 /g' | awk
'{printf("<s> %s </s>\n", $0)}' | carmel -sribIEWk 1
trigram.large.wfsa.norm remove-vowels.fst
```

then I used the following command to remove _, <s>, </s>, and extra spaces before evaluating the restored file using eval_vowels.py

```
cat vowel_restored.trigram.base | sed -e 's/ //g;s/_/
/g;s/<s>//g;s/</s>//g' > vowel_restored_test.trigram.base
```

Using the base trigram model, we got accuracy of 42.6% and using the large trigram model 41.0% accuracy.

```
pelican02 ~/Courses/CS-539/hw6-data 198% python eval_vowels.py test.txt vowel_restored_test.trigram.large
word acc= 0.410
```

```
pelican02 ~/Courses/CS-539/hw6-data 199% python eval_vowels.py test.txt vowel_restored_test.trigram.base
word acc= 0.426
```

2. Now design an algorithm to recover spaces using NLMs.

Describe your algorithm in English and pseudocode, and analyze the complexity.

Answer:

The algorithm for recovering spaces is described below:

Pseudocode and descriptions:

Step1:

```
vowel = ['a', 'e', 'i', 'o', 'u']
```

```
for i in vowel:
    perm.append([i])
for i in vowel:
```

```

        for j in vowel:
            perm.append([i, j])
for i in vowel:
    for j in vowel:
        for k in vowel:
            perm.append([i, j, k])

```

In step1, we generate all possible combinations of vowels, upto 3 consecutive vowels and store them in a list (perm)

When we are running our algorithm, we need to search for the missing vowel combination. So, we'll be using the perm list.

step2:

```

for each line in the given file:
    h = load the language model
    beam = [(0, h)] (#Initialize the Beam)

    for each character c in the current line:
        if c == ' ': (#Replacing the spaces with _)
            c = '_'

        temp = [] (#initializing temporary beam)

```

Step3:

For top K-best answer stored in the beam, in the previous step

```

    for (p, h) in beam:

```

We have several options here,

1. We can ignore adding any vowels and just add the current character

```

        prob = p + math.log(h.next_prob(c))
        h_temp = h + c
        temp.append((prob, h_temp))

```

2. We can add vowel before the current character. We have stored all possible vowels (upto 3 consecutive) in the list (perm). So, we would add each one of those to our top K-beams from the previous step and compute the new language model and probability of occurrence. As the probabilities can be very small, we are using log to avoid underflow.

```

        for permutes in perm:
            h_temp = h + ""
            prob = p

```

```

for j in permuter:
    prob += math.log(h_temp.next_prob(j))
    h_temp += j

prob += math.log(h_temp.next_prob(c))
h_temp += c
temp.append((prob, h_temp))

```

Step4:

Now holding all possible scenarios of adding spaces, is not efficient. We can discard those with low probabilities of occurring. So, we update our original Beam with the top 20 results from the temporary beam.

```
beam = sorted(temp, reverse=True)[:beam_len]
```

After, continuing step 3 until the end of the line, we finally get our top 40 [beam_len] candidates lines with spaces. And as we are sorting our beam at each step, the 1st answer in the Beam is the sentence which has the best probability of occurring with spaces.

```
p, h = beam[0]
```

We retrieve that and print the history of that sentence's language model, which gives us the line with restored spaces.

```

print(("".join(h.history)).replace("_", " "),
      "").replace("<s>", ""))

```

The above algorithm has a complexity of $O(n * m * b * v)$,

where n is the number of lines, m is the number of characters per line, b is the beam size, and v is the number of candidate vowels.

V can change according to the number of vowel candidate, In our implementation we used maximum of 2 consecutive vowels as candidate. Because considering 3 consecutive vowels would add another 125 candidates, and it was taking a lot of time to train in pelican.

So, the complexity of the aforementioned algorithm is $O(n * m * b * v)$

3. Implement it, and report the Precisions, Recalls, and F-scores.

Answer:

The implementation of the algorithm described above for recovering vowel: restore_vowels.py is attached with this report.

For the experiment, only 30 vowel candidate was used (upto 2 consecutive vowel combinations). As the training time of the models (with 3 consecutive vowel combinations) was unexpectedly long in Pelican, I wasn't able to complete it. Nonetheless using 1, and 2 consecutive vowel candidates, give close to the intended accuracy.

I got accuracy of 52.3% using the Base NLM.

```
pelican04 ~/Courses/CS-539/hw6-data 169% python eval_vowels.py test.txt vowels_restored_base.txt  
word acc= 0.523
```

Using the Large NLM I got 76.5% accuracy.

```
pelican04 ~/Courses/CS-539/hw6-data 170% python eval_vowels.py test.txt restore_vowels_large.txt  
word acc= 0.765
```

I could have achieved 54% and 80% accuracy respectively if 3 consecutive vowel combination could be used. Because in the test.txt file there were words that have 3 consecutive vowels [e.g. ['i', 'i', 'o'], ['i', 'o', 'u'], ['e', 'o', 'u'], ['e', 'e', 'i']], which explains the small reduction in accuracy.

In this vowel, restoration problem F1 score, recall and precision are all the same, so they weren't reported individually. I've added both the recovered file with this report.