

2. Evaluating Entropy

2.1

base: 2.9148

large: 2.93042

→ As shown above, the base has the better entropy score. I guess that if we consider the probability distribution is the same on both data sets, the large data sets has not enough impact to show the significant improvement on the entropy result.

2.2

cat test.txt | entropy.py

base: 2.693

large: 1.973

huge: 1.479

→ As the model increases, the entropy of each model is improved significantly. This is because as the model increases, the RNN model learns more errors and provides improved results with the updated probability distribution.

2.3

yes, it does make sense

3. Random Generation

3.1 Please see the below .txt files → It does not make sense.

base: ./carmel -GI 10 trigram.base.wfsa.norm > random_carmel_base.txt

large: ./carmel -GI 10 trigram.large.wfsa.norm > random_carmel_large.txt

3.2

base: python random_generation.py > random_gen_base.txt

large: python random_generation.py > random_gen_large.txt

huge: python random_generation.py > random_huge_base.txt

3.3 NLM results are more make sense than Carmel results.

4. Restoring Spaces

4.1 using carmel

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

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

`cat carmel_large_space_restored.tri | python make.py > carmel_large_space_restored.txt`

`python eval_space.py test.txt carmel_large_space_restored.txt`

base: recall= 0.598 precision= 0.608 F1= 0.603

large: recall= 0.625 precision= 0.652 F1= 0.638

4.2

pseudocode:

for char in sentence:

 for score, state in beam:

 update newscore and newstate on char

 append newscore and newstate

 update newscore and newstate on space

 append newscore and newstate

 sort newbeam on top b search

complexity: $O(n*b)$

length of sentence = n

beam search = b

base: `cat test.txt.nospaces | python restoring_spaces.py > restored_spaces.txt.base`

large: `cat test.txt.nospaces | python restoring_spaces.py > restored_spaces.txt.large`

huge: `cat test.txt.nospaces | python restoring_spaces.py > restored_spaces.txt.huge`

4.3 `python eval_space.py test.txt restored_spaces.txt.huge`

base: recall= 0.830 precision= 0.809 F1= 0.819

large: recall= 0.969 precision= 0.955 F1= 0.962

huge: recall= 0.994 precision= 0.991 F1= 0.993

5. Restoring Vowels

5.1 using carmel

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

```
base: python eval_vowels.py test.txt carmel_base_vowels_restored.txt
word acc= 0.426
large: python eval_vowels.py vowels_restored.txt carmel_large_vowels_restored.tri
word acc= 0.410
```

5.2

pseudocode:

```
for char in sentence:
    for MAX_ITERATION:
        for score, state in the previous beam search result:
            update newscore and newstate on char
            append newscore and newstate

            if before MAX_ITERATION:
                for VOWELS:
                    update newscore and newstate on vowel
                    append newscore and newstate in tmp
            append tmp to beam search result

sort newbeam on top b search
```

complexity: $O(5^k * n * b)$

```
number of vowels = 5
MAX_ITERATION = k
length of sentence = n
beam search = b
```

```
base: cat test.txt.novowels | python restoring_vowels.py > restored_vowels_base.txt
large: cat test.txt.novowels | python restoring_vowels.py > restored_vowels_large.txt
huge: cat test.txt.novowels | python restoring_vowels.py > restored_vowels_huge.txt
```

5.3

```
python eval_vowels.py test.txt restored_vowels_base.txt
base: word acc= 0.539
large: word acc= 0.789
huge: word acc= 0.931
```