

Natural Language Processing – Assignment II

Reykjavik University – School of Computer Science

Fall 2023

1 NLTK: The Icelandic Gold Standard – 16 pts

In this part, you use Python and NLTK to process a Part-of-Speech tagged Icelandic corpus, the *Icelandic Gold Standard*; *MIM-GOLD*¹. A preprocessed version of MIM-GOLD containing one sentence per line (with “/” between tokens and tags) is available as the file *MIM-GOLD.sent* in the zip file linked the assignment description.²

This part is divided into the subparts described below, but you should return a single Python program, **mim_gold.py**, which includes all the code needed for carrying out these tasks.

Write code to:

1. Read in MIM-GOLD using the class *TaggedCorpusReader*, display the number of sentences, and display the individual tokens of sentence no. 100.
2. Display the number of tokens and the number of types in MIM-GOLD.
3. Display the 10 most frequent tokens in MIM-GOLD using the class *FreqDist*.
4. Display the 20 most frequent tags in MIM-GOLD using the class *FreqDist*.
5. Generate tag bigrams and use the class *ConditionalFreqDist* to print out the 10 most frequent tags that can follow the tag 'af' (this tag denotes a preposition).

Example output follows:

```
Number of sentences: 58412
```

```
Sentence no. 100:
```

```
Púttað á Listatúni í dag , laugardag , kl. 10.30 .
```

```
Number of tokens: 1000218
```

```
Number of types: 106529
```

```
The 10 most frequent tokens
```

```
. => 49066
```

```
að => 35749
```

```
og => 33813
```

```
, => 29990
```

```
í => 27622
```

```
á => 21833
```

```
er => 16604
```

```
sem => 15199
```

```
til => 9888
```

```
um => 8799
```

```
The 20 most frequent PoS tags:
```

¹See <https://clarin.is/en/resources/gold/>

²Please make sure not to distribute the MIM-GOLD corpus.

```

AF => 109899
AA => 74151
C => 68278
PL => 53212
SFG3EN => 36289
PK => 30647
SNG => 26345
TA => 22992
SFG3EP => 19773
CN => 19540
SPGHEN => 17190
PA => 13043
N----S => 12783
CT => 12664
SFG3FN => 12294
NKEN => 11411
NVEN => 11171
NVEO => 11123
NHEP => 10979
NVEP => 10281

```

The 10 most frequent PoS tags following the tag 'af':

```

NVEP => 5714
NHEP => 5361
NKEP => 4222
CN => 3903
NVEO => 3556
NKEO => 3347
NHEO => 3046
NHEPG => 2867
NVEPG => 2723
FPHEP => 2703

```

2 NLTK: PoS tagging – 16 pts

In this part, you experiment with various different types of PoS taggers in the NLTK using the Penn Treebank corpus (accessible in Python from `nltk.corpus import treebank`).

By studying chapter 5 in the NLTK book (<http://nltk.org/book/>), you should have all the necessary material to solve this part. The NLTK documentation on tagging, <http://www.nltk.org/api/nltk.tag.html> is also an important source of information.

This part is divided into the subparts described below, but you should return a single Python program, **tagging.py**, which includes all the code needed for carrying out these tasks.

Write code to:

1. Split the tagged sentences of the Penn Treebank into a training set (first 3500 sentences) and a test set (the remaining sentences), print out the total count of each set, and print the first sentence in the test set.
2. Construct four taggers trained on the training set: an instance of an *AffixTagger*, *UnigramTagger*, *BigramTagger* and a *TrigramTagger* (without any “backoff” model). Evaluate them on the test set, and print out the evaluation results.
3. Construct the latter three taggers again, but now with a backoff model, i.e. such that the trigram tagger uses a bigram tagger as backoff, which in turn uses a unigram tagger as backoff, which in turn uses the affix tagger as backoff. Print the evaluation results again.

4. You will notice that the tagging accuracy for the individual taggers without a backoff model is (much) lower than the tagging accuracy for the corresponding taggers when using a backoff model. **Explain why this is the case. In particular explain this for the case of the *BigramTagger*.**
5. Tag the test set with the default (“off-the-shelf”) tagger in the NLTK³, evaluate its accuracy and print out the result. Note that for this tagger you cannot simply call a built-in evaluation function, instead you have to write your own, which compares the results of the tagger to the gold standard.

Example output follows:

```
Number of training sentences: 3500
Number of test sentences: xxx
```

```
First sentence in test corpus:
```

```
[('About', 'IN'), ('30', 'CD'), ('%', 'NN'), ('of', 'IN'), ('Ratners', 'NNP'), (''s', 'POS'),
('profit', 'NN'), ('already', 'RB'), ('is', 'VBZ'), ('derived', 'VBN'), ('*-1', '-NONE-'),
('from', 'IN'), ('the', 'DT'), ('U.S.', 'NNP'), ('.', '.')]

```

```
Tagging accuracies:
```

```
-----
Affix tagger: xx.yy%
Unigram tagger: xx.yy%
Bigram tagger: xx.yy%
Trigram tagger: xx.yy%
```

```
Tagging accuracies with backoff:
```

```
-----
Affix tagger: xx.yy%
Unigram tagger: xx.yy%
Bigram tagger: xx.yy%
Trigram tagger: xx.yy%
```

```
Accuracy of the default tagger in NLTK: xx.yy%
```

3 Your Own Personal GPT – 32 points

Objective. Create three small character-level GPT models using Karpathy’s nanoGPT repository on a corpus of your choice. Evaluate the quality and creativity of the generated text under different parameter settings.

Background. Large-scale models like GPT-3 and GPT-4 have made headlines; however, understanding the nuances of training and generating from smaller models can provide insights into the workings of generative language models. In this part, you will use nanoGPT to train a model and experiment with its generation capabilities.

Requirements. You must have Python 3 installed. Follow the instructions on <https://www.python.org/> to install the correct version of Python 3 for your machine. The second requirement is having Pytorch installed. Follow the instructions on <https://pytorch.org/get-started/locally/> to select the proper installation command for your system. You don’t have to worry about CUDA unless you are going to be using a GPU.

Setup. Download the code for the nanoGPT repository by clicking the ‘Code’ dropdown on the Github page <https://github.com/karpathy/nanoGPT> and choosing an option that suits your needs. Make sure everything is working normally by following the instructions in the README file starting with the ‘quick start’ section. Run the

³You need to install the `numpy` module for running this tagger.

`python data/shakespeare_char/prepare.py` command in the terminal. Address any issues that may arise, like installing missing packages. Scroll down to the ‘I only have a macbook (or other cheap computer)’ section. Run the following command in the terminal `python train.py config/train_shakespeare_char.py -device=cpu -compile=False -eval_iters=20 -log_interval=1 -block_size=64 -batch_size=12 -n_layer=4 -n_head=4 -n_embd=128 -max_iters=2000 -lr_decay_iters=2000 -dropout=0.0`. Finally, generate some samples to see how the model performs using `python sample.py -out_dir=out-shakespeare-char -device=cpu`.

With the basic setup step completed, carry out the following tasks:

1. Corpus Selection (8 points)

- Choose a small corpus for training. This could be a collection of poems, short stories, news articles, or any other text your choice in any language.
- Justify your choice: Write a brief paragraph explaining why you chose this corpus and what you hope the model might learn from it.
- Unless you have beefy hardware, keep the length of the corpus under 1 million words (as counted by a standard word count program like `wc`, or those in Word or Notepad++).
- Modify the nanoGPT code slightly to have it use your data. We will use the `shakespeare_char` code for training. Open the `data/shakespeare_char/prepare.py` file. Comment out the following lines:

```
input_file_path = os.path.join(os.path.dirname(__file__), 'input.txt')
if not os.path.exists(input_file_path):
    data_url = 'https://raw.githubusercontent.com/karpathy/char-rnn/
master/data/tinyshakespeare/input.txt'
    with open(input_file_path, 'w') as f:
        f.write(requests.get(data_url).text)
```

- Immediately below these lines, add this line with the path to your text file:
`input_file_path = 'path/to/your/file.txt'`
- Run the commands in the **Setup** section again. The data files for your corpus should have been created, the model trained on your text and used to generate text based on your data.
- Note that the model will be called ‘`shakespeare_char`’ even though you’ve trained it on a completely different dataset. This is fine for the purposes of this assignment, it’s just a name.

2. Model Training with Hyperparameter Exploration (16 points)

- Train three different character-level GPT models using the method above, each with a distinct set of hyperparameters (e.g., learning rate, number of layers, number of attention heads, batch size, etc.).
- For each model:
 - Document the chosen hyperparameters and the rationale behind your choices.
 - Detail the training process, including any challenges faced and how you overcame them.
 - Generate text from the trained model and provide the output.
 - Training a new model will overwrite the previous one (unless you modify the code). Ensure you save the generated outputs and any other necessary details before moving to the next model. To save the model itself, copy the ‘`out-shakespeare-char/ckpt.pt`’ file to another folder somewhere.

3. Analysis and Comparison (8 points)

- Reflect on the entire process:
 - What did you learn about the importance of hyperparameters in training GPT models?
 - How might the size or nature of your chosen corpus have influenced the results?
 - Based on your experiments, what further tweaks or tests would you consider if you were to continue this exploration?

4 What to return

A single .zip file containing the following:

1. Two Python programs: **mim_gold.py** and **tagging.py**. Make sure you use functions (where appropriate) for specific tasks in your Python code, thus minimizing the duplication of code.
2. For part 3, address all aspects and questions posed in the task description in a PDF file:
 - Your chosen corpus and justification.
 - Documentation of the training process, hyperparameters, and generated outputs for each of the three models.
 - Your analysis and comparison of the three models.
 - Your overall reflection.

As always, please let me know if you run into issues or you need help getting set up. Try to do this as early as possible so we can address any issues.