

# Rust-Python HMM

For this assignment, I decided to focus on the model computational performance. Specifically, I implemented the assignment also in [Rust](#) to see what the performance gain will be. The reasons are that (1) I want to learn Rust but every class is Python + PyTorch, and (2) current NLP research is all done by prototyping in Python, yet there is a virtue in the experience of programming usable deployable and fast solutions. Note: I don't want to be mean to Python; it certainly has its place and advantages.

I understand that the task was to program this in Python, which I hope I fulfilled. Yet I also hope that you will find this comparison interesting.

The Rust code is twice as large\* and took much longer to complete. The benefits are, however, that once it compiled, I was convinced of its functionality, which was not the case with Python. \*Running `grep -r -E "\s*[\{\}]\s*$" rust/src/ | wc -l` reveals that more than 100 lines are just opening or closing brackets, so the code size is not that significant. Also, strong typing allows for more finer tooling. An example would be clippy, which helped me discover multiple bad design patterns.

## Project structure

```
data/                # not supplied, paste the files here for reproducibility
- de-{eval,train}.tt
- de-test.t
data_measured/
- {r,p}-de-eval{,-smooth}.tt # model outputs
- time-{1,2,3}              # measured results for graphs
- time-{1,2,3}.png          # exported graphs
meta/                  # scripts for measuring performance and accuracy
- graph.py                # produce graphs given logs time-{1,2,3} in data_measured
- run_times.py             # measure performance from r-build-time and p-run-time recipes
- eval.py                 # computes metrics and prints table
rust/                  # Rust source code
python/                # Python source code
Makefile               # Makefile for common recipes
```

## Makefile and reproducing results

`make r-print-eval` trains two models on the data and outputs the CONLL-U file to `data_measured/p-de-eval.tt` and `data_measured/p-de-eval-smooth.tt`, similarly `make r-print-eval` produces `data_measured/r-de-eval.tt` and `data_measured/r-de-eval-smooth.tt` (assuming stable Rust compiler in path). The semantics of the rest of the command line arguments is intuitive from the Makefile: `print_acc` self-reports the accuracy on anything it computes (`comp_test`, `comp_train` OR `comp_eval`).

File paths are relative hardcoded because there are no plans to make this portable and there were already too many switches. Both versions assume that they are run from the top-level directory (the directory the `README.md` is in). If `print_pred` is present, the program outputs predictions to stdout. Progress is output to stderr.

## Correctness

Even though the Viterbi algorithm should be mostly deterministic, there is a big issue with number representation and rounding. There appears to be a big difference in accuracy based on the underlying numeric type used (f32 vs f64). All parameters were multiplied by 1500 in both versions because this maximized the performance (possibly

striking the sweet spot between diminishing and exploding values). In trellis computation, the layers are all normalized to sum to 1 after every step. Making the normalization to sum to something other than 1 did not affect the performance.

Despite my best efforts, the two versions produce slightly different results. This may be due to different corner-case numeric handling in the two systems.

Unseen tokens were dealt with by substituting the emission probability with 1, thus relying on the surrounding transition probabilities.

I tried to use the same algorithmic steps in both solutions so that they are comparable. It is, however, still possible, that I mistakenly used some other data structure, assuming it was the same.

## Log space

Another solution to the issue of storing very small probabilities would be to work in log space. One of the issues is that it no longer supports the computation of cumulative probability (because the probabilities there are summed) and also it had a negative effect on performance relative to the current solution: for (train, eval) accuracy, the new results in Rust were (89.16%, 78.96%) and in Python (66.67%, 66.98%).

## Code structure

Structures in both versions follow the same naming scheme. The programs function as follows:

1. Train Loader is created, which also creates a Mapper objects (see Note)
2. HMM Model parameters are estimated from the training data.
3. Eval or Test Loader is created, reusing Training Mapper.
4. Based on the arguments, datasets are evaluated (comp\_test, comp\_train OR comp\_eval).

The HMM class contains code for initialization and Viterbi and can be used generically. HMMTag inherits from this class and adds specific functions for initialization from Loader and evaluation. Both implementations start with `main.{rs,py}`.

## Performance Graphs

The performance was measured with respect to changing training data size (steps of 10000 tokens). The task was (1) train, (2) train + evaluate on eval, (3) train + evaluate on train and eval. Accuracy of these models was also measured. The measured times are without writing to files. Rust version is compiled with the `--release` flag and Python is run with `-0`. Both versions use aforementioned smoothing.

Figure 1 shows simply that in training, the Rust implementation seems to be faster by the factor of  $\sim 8$ .

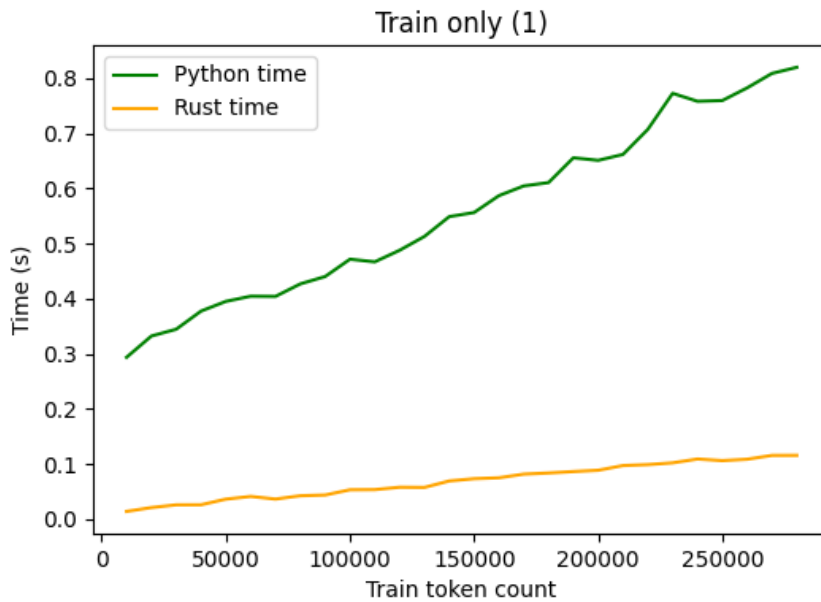
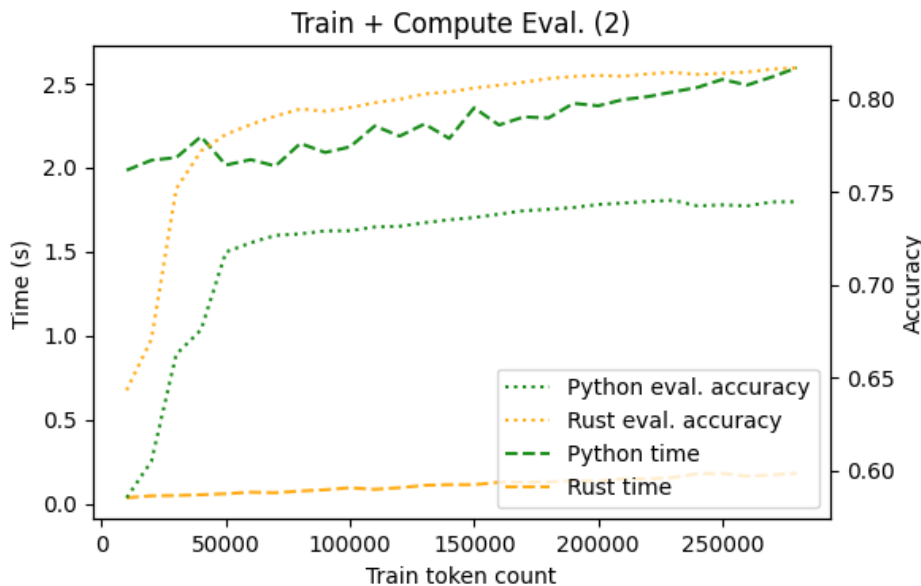
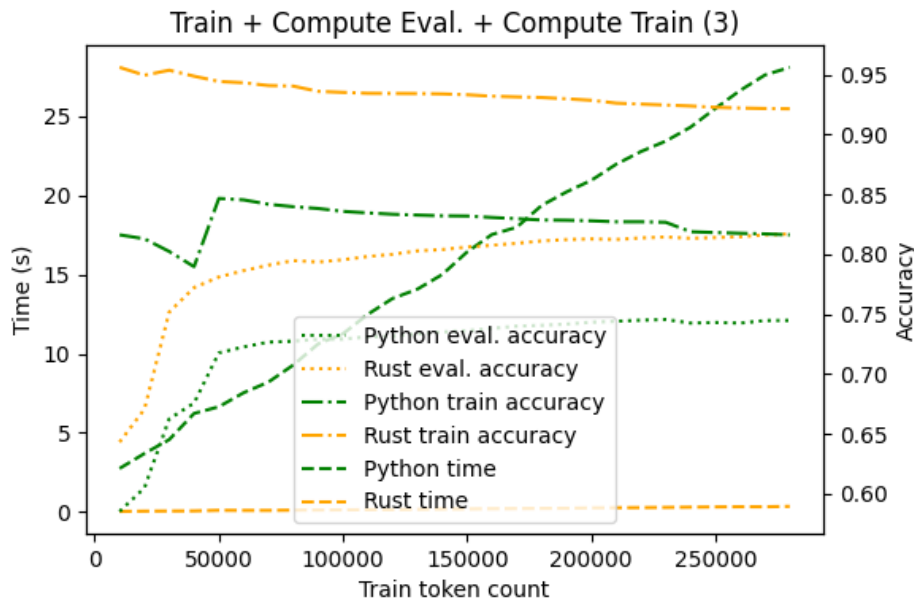


Figure 2 also shows that the Rust implementation is more stable (possibly because of the lack of runtime). We also see that there seems to be diminishing return in performance after we pass 50k train tokens. Python ends at 2.56s and Rust on 0.18s (factor of  $\sim 14$ ).



Evaluating the whole data proved to be the most challenging task. This is shown in Figure 3. While Python ends at 28.08s, for Rust it is 0.36s (factor of  $\sim 78$ ). The training accuracy is also decreasing (apart from the spike in Python) because the capacity of the model is getting shared with a larger amount of examples. Train accuracies were 92.16% and 81.63% for Rust and Python, respectively. Evaluation accuracies were 81.70% and 74.48%.



## Note on Performance

I did not try to especially optimize algorithmic performance. For example, the trellis is allocated and cleared for every sentence in the data. This could be done much more efficiently by creating one static one (the size of the longest sentence) and reusing that for the computation. It does not need to be cleared, because every cell is first written to and only then read.

One of the biggest performance boosts was gained by creating a hashmap mapping from string (both for words and for tags), convert everything to numbers (Rust version uses 8bytes, which is unnecessary), manipulate just these numbers and only when printing revert back. This is done by the `Loader` and `Mapper` classes in both versions.

Also, both versions contain code for computing sequence observation probability in trellis (`sum` instead of `max`), but is turned off in both versions. The Rust version gets an unfair advantage in this because it is removed compile-time, while in Python, the interpreter has a bit more work to do.

## Additional

### Smoothing

I also experimented with rudimentary smoothing. This can be done easily by changing the initial probabilities in the constructor (`class HMM`) to some parameter `alpha` instead of zeroes. Since probabilities are scaled up by the factor of 4096, it makes sense to use higher values.

Interestingly enough, the performance increased by tinkering with start and transition probabilities and not emission probabilities. Furthermore, setting initial transition count to a negative number -32 and -16 resulted in the best results (I did not employ grid-search, so there surely exists a better set of parameters. The resulting (train, eval) accuracies were (92.47%, 82.14%) and (81.63%, 74.48%) for Rust and Python respectively. This is an improvement of (+0.32%, +0.45%) and (+1.93%, +1.13%) for Rust and Python. The resulting inferences are stored in `data_measured/{p,r}-de-eval{-smooth}.tt`.

## Ice cream

The Rust code also contains the toy ice-cream X weather example. It can be run from the rust directory with `cargo test -- --nocapture`.

## Unknown word handling by subwords

This is an idea beyond the scope of this homework, but I would nevertheless like to see it implemented (and especially to see the performance) or any comments that show the caveats of this approach.

In order to better handle unknown word handling, all tokens could be split into subword units, e.g. by Byte Pair Encoding. This would allow the splitting to be trained not only on annotated data but also on unannotated. The HMM parameters could be then estimated as follows:

Assume the sequence SENT A-B C-B (BPE compound A-B at the beginning of the sentence, followed by C-B). Since individual subwords have the same POS tags, the starting and transition probabilities can be computed in an almost normal way: both A and B are starting and both A, B are followed by C, D (4 transitions). Furthermore, emission probabilities can also remain unchanged. This is counterintuitive because it will lead to affixes with POS tags as the same word (e.g. un-do-able -> (un, ADJ), (do, ADJ), (able, ADJ).) To avoid this, I would suggest early stopping of the BPE algorithm.

Further assume, that we trained two sets of HMM parameters: in the standard way ( $E, T, P$ ) and also with subword units ( $E', T', P'$ ). The main difference would be in inference. If the next token to be processed is present in the training data, the standard parameters and approach would be used. If it is, however, not in the training data, it is split to subwords:  $c = A_1A_2...A_n$ . The starting and transition probability would be estimated from  $P'$  and  $T'$ . Emission probability would then be the average of the parameters for individual subwords:  $E''(c, s) = [E'(A_1, s) + E'(A_2, s) + ... + E'(A_n, s)]/n$ .

The emission probability function can be extended to convex interpolate between the standard and subword version:  $E'''(c, s) = a * E''(c, s) + (1-a) * E(UNK, s)$ . Here  $a$  is a parameter, which can be estimated from held-out data.

## Unknown word handling by Stemming

Completely another approach would be some sort of sensitive stemming, which would remove affixes that do not change the part of speech. This would reduce the amount of word forms while preserving correctness in the annotation.

## Eval Results

This sections lists results of models in descending order. The file `eval.py` is included, because I changed it to produce markable tables.

## Rust + scaling, smoothing

Highest results from `r-eval-smooth.tt`, accuracy: 82.14%.

Tag	Prec.	Recall	F1 score
DET	0.8556	0.9430	0.8972
ADV	0.4454	0.8374	0.5815

<b>Tag</b>	<b>Prec.</b>	<b>Recall</b>	<b>F1 score</b>
NOUN	0.8882	0.7526	0.8148
VERB	0.9609	0.8213	0.8856
ADP	0.9296	0.8316	0.8779
.	0.9921	0.9979	0.9950
CONJ	0.8953	0.8254	0.8590
PRON	0.8998	0.7258	0.8035
ADJ	0.7517	0.6264	0.6834
NUM	0.3877	0.7222	0.5045
PRT	0.7359	0.8078	0.7702
X	0.1538	0.0909	0.1143

## **Rust + scaling**

File `r-eval.tt`, accuracy 81.69%.

<b>Tag</b>	<b>Prec.</b>	<b>Recall</b>	<b>F1 score</b>
DET	0.8463	0.9496	0.8950
ADV	0.4433	0.8335	0.5788
NOUN	0.8896	0.7497	0.8137
VERB	0.9609	0.8203	0.8851
ADP	0.9294	0.8040	0.8622
.	0.9933	0.9937	0.9935
CONJ	0.9013	0.8254	0.8617
PRON	0.9140	0.7110	0.7998
ADJ	0.7489	0.6255	0.6816
NUM	0.3816	0.7222	0.4994
PRT	0.6410	0.8143	0.7174
X	0.1714	0.2727	0.2105

## **Python + scaling, smoothing**

File `p-eval-smooth.tt`, accuracy: 74.48%.

<b>Tag</b>	<b>Prec.</b>	<b>Recall</b>	<b>F1 score</b>
DET	0.8487	0.2865	0.4284
NUM	0.1407	0.9593	0.2454
NOUN	0.9110	0.7230	0.8062
VERB	0.9424	0.8315	0.8835
ADP	0.9106	0.9030	0.9067
.	0.9925	0.9975	0.9950
CONJ	0.5172	0.9433	0.6681
PRON	0.5374	0.8755	0.6660
ADV	0.7011	0.5201	0.5972
ADJ	0.7225	0.5287	0.6106
PRT	0.7126	0.5896	0.6453
X	0.5000	0.0909	0.1538

## Python + scaling

File p-eval.tt, accuracy: 73.35%.

<b>Tag</b>	<b>Prec.</b>	<b>Recall</b>	<b>F1 score</b>
DET	0.8487	0.2865	0.4284
X	0.0059	0.5000	0.0117
NOUN	0.9121	0.7227	0.8064
VERB	0.9405	0.8121	0.8716
ADP	0.9140	0.9010	0.9074
.	0.9936	0.9836	0.9886
CONJ	0.5159	0.9204	0.6612
PRON	0.6296	0.8216	0.7129
ADV	0.7034	0.5185	0.5970
ADJ	0.7203	0.5278	0.6092
NUM	0.4717	0.8037	0.5945
PRT	0.7109	0.5928	0.6465

## Python

Because in the first homework I was penalized for not sticking literally to the assignment (plotting top 100 most frequent words instead of all, because in the later case nothing could be inferred - a deliberate decision), I do not trust that I would not be penalized also for not providing vanilla HMM POS tagger implementation in Python. Since the code would be too convoluted to parametrize also for vanilla implementation and I do not want to regress the performance of the main project, the code (duplicated with slight changes) is located in `python/vanilla/`. The resulting inference is in `data_measured/p-de-eval-vanilla.tt`. Accuracy: 65.63%.

<b>Tag</b>	<b>Prec.</b>	<b>Recall</b>	<b>F1 score</b>
DET	0.9191	0.2487	0.3914
NUM	0.1213	0.8778	0.2131
NOUN	0.9163	0.6168	0.7373
VERB	0.9395	0.7443	0.8306
ADP	0.9152	0.8111	0.8600
.	0.9946	0.8471	0.9149
CONJ	0.4883	0.8637	0.6239
PRON	0.5292	0.8420	0.6499
ADV	0.2546	0.4594	0.3276
ADJ	0.7216	0.4569	0.5595
PRT	0.6716	0.5863	0.6261
X	0.0625	0.0909	0.0741

I also did not provide any comment for the `alpha=0.9` parameter in `matplotlib` function call in `graph.py`, because that seems just absurd and commenting every other line reduces readability.