# Report for Assignment - 1

This report was prepared by four team members, each contributing a specific role in the analysis of the given datasets and the development of our HMM–based Part-of-Speech tagging system:

1. Deepansh Pandey (Roll No. MT2024038)
2. Eshwar Chawda (Roll No. MT2024047)
3. Himanshu Shrivastava (Roll No. MT2024059)
4. Shivam Pandey (Roll No. MT2024145)

Our primary reference files are **TRAIN.csv** and **TEST.csv** which contain tokenized sentences labeled with POS tags. The goal of our project is to use Viterbi algorithm to train, validate, and evaluate an HMM for accurate POS classification, using these datasets as the main source of training and testing data. The sections that follow describe our methods, findings, and concluding remarks. We have attached our Python Notebook with this report zip.

## 1. Preprocessing

We began by collecting and examining our tagged text data. Our primary goal was to understand how many tokens (words) and how many unique tags (NOUN, VERB, ADJ, etc.) we had in total. This led us to perform frequency counts for each tag, which quickly showed that certain classes—particularly NOUN and VERB—accounted for a large fraction of the data, while tags like INTJ or PRT were underrepresented.

Given that some tags were much rarer than others, we checked for potential imbalances that could affect model performance. We also looked for out-of-vocabulary words in the test set that weren't present in the training data. Identifying these OOV words helped us anticipate problems the model might face in correctly tagging unseen tokens.

We took note of any tokens that can appear in multiple parts of speech—for example, "like," which can be a verb ("I like apples"), a preposition ("something like a tree"), or even a filler word in informal speech. Observations like these would prove relevant later, since our approach must handle a certain level of linguistic ambiguity.

The data also contained many stopwords that could be removed, as well as various words that could be converted to lowercase for consistency.

# 2. Viterbi

Once the dataset was ready for modeling, we implemented the Viterbi algorithm to perform sequence decoding. We first established dynamic programming matrices: one to store the maximum log-likelihood for each state (i.e., each tag) at each position in the sentence, and another for backpointers to reconstruct the most likely tag sequence.

We made sure to handle unknown words by assigning them a low default probability, preventing numerical underflow with small constants or log-space arithmetic. Because Viterbi accumulates probabilities across a sentence, any small misestimation early on could lead to large downstream errors.

During development, we ran the Viterbi procedure on a few example sentences to confirm that it tagged common words correctly—particularly boundary cases like the first word in the sentence, where an initial probability distribution is used instead of a transition from a previous tag. After verifying these steps, we integrated the code to evaluate our entire test set.

# 3. HMM

We used a standard Hidden Markov Model (HMM) framework to estimate the probability of each tag sequence. First, we derived emission probabilities from the training set by counting how often a given tag produces a given word. Then, we computed transition probabilities to capture which tag is most likely to follow another (for instance, ADJ often precedes NOUN).

To initialize the model, we incremented all counts by one (Laplace smoothing) to avoid zero probabilities for any unseen (tag, word) or (tag, tag) pair. This provided a more robust set of estimates, especially for tags that appeared infrequently.

Finally, we stored the resulting probabilities as log-probabilities to avoid numerical instability— important because multiplying many small probabilities over long sentences can quickly approach floating-point limits. Armed with these emission and transition distributions, we handed the data off to be decoded by the Viterbi procedure.We used a standard Hidden Markov Model (HMM) framework to estimate the probability of each tag sequence. First, we derived emission probabilities from the training set by counting how often a given tag produces a given word. Then, we computed transition probabilities to capture which tag is most likely to follow another (for instance, ADJ often precedes NOUN).

To initialize the model, we incremented all counts by one (Laplace smoothing) to avoid zero probabilities for any unseen (tag, word) or (tag, tag) pair. This provided a more robust set of estimates, especially for tags that appeared infrequently.

Finally, we stored the resulting probabilities as log-probabilities to avoid numerical instability—important because multiplying many small probabilities over long sentences can quickly approach floating-point limits. Armed with these emission and transition distributions, we handed the data off to be decoded by the Viterbi procedure.

# 4. Confusion Matrix

After we obtained predicted tags on the test set, we generated the confusion matrix to visualize the model's performance. On the rows, we placed the actual (ground truth) tags, and on the columns, we placed the predicted tags. Because the matrix was normalized, each row sums to 1, making it easier to compare how well each tag was correctly identified versus how it was misclassified.

Many tags, particularly NOUN, VERB, and ADJ, had strong diagonal scores near 0.9 or higher. This indicates that the HMM + Viterbi pipeline handled common word classes effectively. However, the matrix revealed potential confusion between PROPN (proper nouns) and NOUN—likely due to capitalization inconsistencies or ambiguous contexts. There were also occasional mix-ups between AUX and VERB, which makes sense because words like "have," "do," or modal auxiliaries can appear in full verb contexts.

By analyzing the largest off-diagonal cells, we discovered that rare tags—INTJ, PRT, X—were more prone to errors. This finding led us to suggest targeted data augmentation and additional linguistic rules to better handle ambiguous or low-frequency categories.

# 5. Conclusion

Overall, the **HMM** approach with **Viterbi decoding** proved quite robust for most frequent tags, achieving high accuracy on classes like NOUN, VERB, and ADJ. The **confusion matrix** confirmed that while common categories are well-modeled, the system often confuses certain classes—particularly those that are inherently ambiguous or underrepresented in the training corpus.

# Results

Initial accuracy: 68.78%

Initial Dataset: TRAIN (1056973 rows) & TEST (1318475 rows)
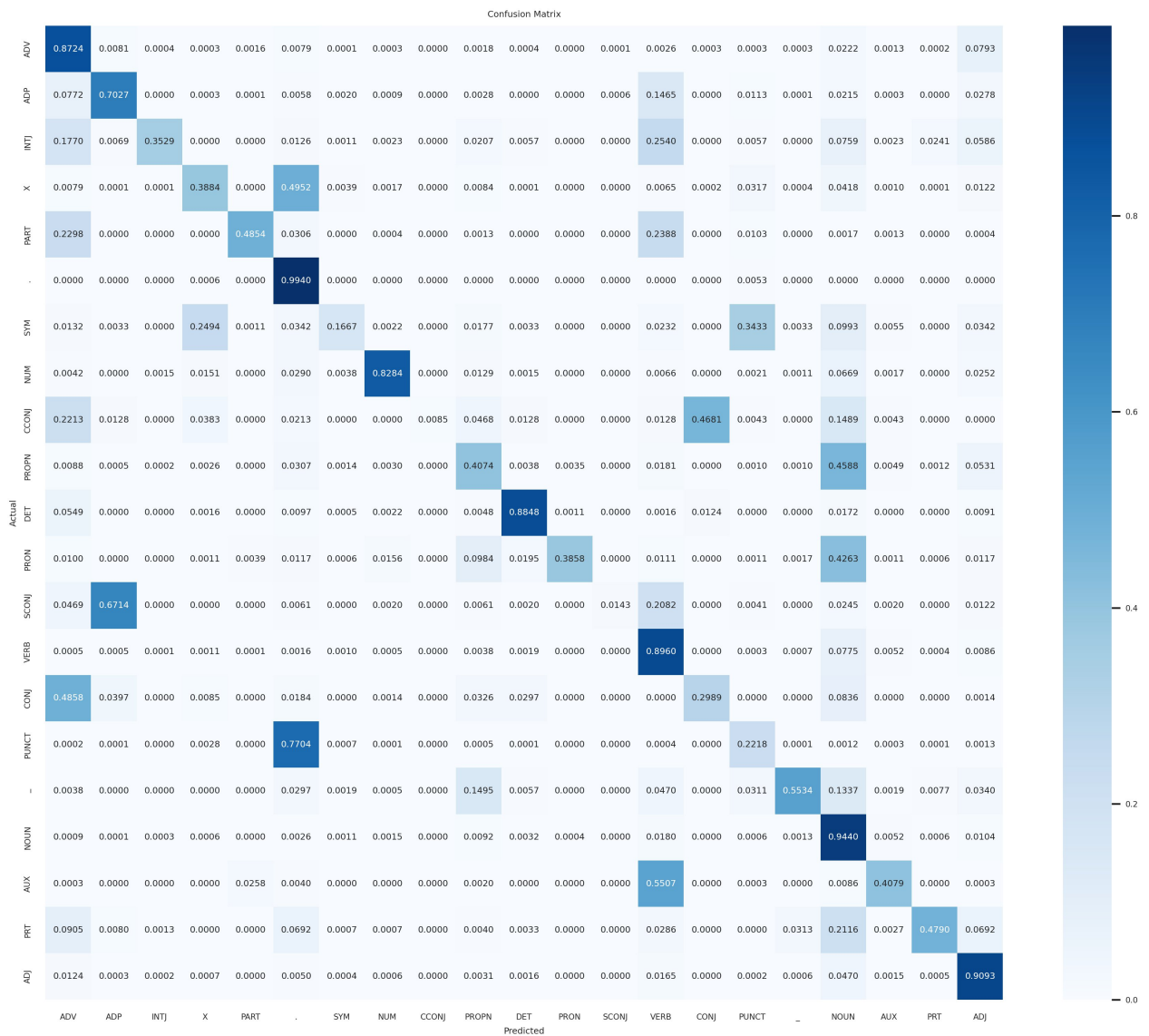
Initial Train-Test time: 8.5 Minutes

After Lowercasing and Stopword Removal.

Final Accuracy: 85.89%

Final Dataset: TRAIN (633439 rows) & TEST (790497 rows)

Final Train-Test time: 5 Minutes

Confusion Matrix:

In these results, the confusion matrix illustrates how often each actual POS tag (rows) is classified as each predicted tag (columns). Most common tags—such as NOUN, ADJ, and VERB—show high diagonal values, indicating strong model accuracy. Rarer or ambiguous tags, including PRT and INTJ, exhibit more confusion, likely due to fewer examples or overlapping linguistic functions. Notable misclassifications include occasional mix-ups between PROPN and NOUN, as well as AUX and VERB, which arise from shared grammatical contexts. Overall, the matrix confirms robust performance on frequent classes while highlighting improvement opportunities for less frequent or context-dependent tags.