# Part-of-Speech Tagging

**(CSE 582; midterm project)**

**Teams:**
- 28 students in 6 teams (<=5 students per team); each team votes a team leader who is in charge of i) leading the project, ii) submitting the results, iii) presenting on 3/13 or 3/15**.**
- Each team has a name
- Team leader sends TA the team name and team members by 2/23.

**Training data: https://www.cnts.ua.ac.be/conll2000/chunking/train.txt.gz**
Format of training file (as the following screenshot shows): each row is for one token in the sentence; sentences are separated by an empty row. Three columns in total: token, POS tag, Chunking tag (we only use the first two columns for this midterm project)

```
a DT B-NP
substantial JJ I-NP
improvement NN I-NP
from IN B-PP
July NNP B-NP
and CC I-NP
August NNP I-NP
's POS B-NP
near-record JJ I-NP
deficits NNS I-NP
. . O

Chancellor NNP O
of IN B-PP
the DT B-NP
Exchequer NNP I-NP
Nigel NNP B-NP
Lawson NNP I-NP
's POS B-NP
restated VBN I-NP
commitment NN I-NP
to TO B-PP
a DT B-NP
firm NN I-NP
monetary JJ I-NP
```

**Dev data**: you can use a small part of training data as dev set.

**Unlabeled Test data**: will be released before 2/18

**Requirements**:
- The three algorithms you have to implement:
  - implement <u>Hidden Markov Models</u> for POS tagging
  - implement <u>Logistic Regression</u> for POS tagging
  - implement <u>Multi-layer Perceptron</u> for POS tagging

- What you can use:
  - Word embeddings
  - Features defined by you or other papers
  - Online packages such as NLTK, Pytorch, spaCy, Gensim, etc.
  - Combine above algorithms/models to get your "best model"

- What you should not use:
  - Transformer-based pretrained language models, e.g., BERT, GPT3, ChatGPT, etc.
  - Data other than the provided training data for pretraining

**What you need to submit (deadline 11:59pm on 3/12)**:

URL of your github repository, including
- **Labeled test data** by your best model: two columns (token, predicted_tag); TA will compute accuracy for each team. Filename "teamname.test.txt"
- **Code files** for the three algorithms: HMM, Logistic Regression, Multi-layer Perceptron

**Evaluation**:
- **System performance (80%)**: each team gets **your_acc/max_acc**
- **Presentation (20%)**: 20min per team. Slides quality, the work you did (how words were represented, how models were optimized, what lessons/experience you have learned, what erros/issues you found, etc.); we draw lots to decide the team order of presentation
- **Each team member gets the same score**.

Pls refer to TA, Shravya Chillamcherla (sjc6752@psu.edu), for more details on how to submit.