# AI539 - Natural Language Processing with Deep Learning

**Homework 2 Report**
**Recurrent Neural Networks and Sentence Representations**

Author: Vy Bui
OSUID: 934370552

Instructor: Professor Stefan Lee

The School of Electrical Engineering and Computer Science
Oregon State University

**Task 1.1**
Manually find weights and biases for the univariate LSTM defined above such that the final hidden state will be greater than or equal to 0.5 for odd parity strings and less than 0.5 for even parity. The parameters you must provide values for are $w_{ix}$, $w_{ih}$, $b_i$, $w_{fx}$, $w_{fh}$, $b_f$, $w_{ox}$, $w_{oh}$, $b_o$, $w_{gx}$, $w_{gh}$, $b_g$ and are all scalars. The LSTM will take one bit of the string (0 or 1) as input $x$ at each time step. A tester is set up in `univariate_tester.py` where you can enter your weights and check performance.

*Hints: Some of the weights will likely be zero. Others may be large. Scale matters, larger weights will saturate the activation functions. Also note that* `A XOR B` *can be written as* `(A OR B) AND (A NAND B)`*. Work backwards and try to find where this sort of two-term structure could be implemented.*

**Task 2.1**
Implement the `ParityLSTM` class in `driver_parity.py`. Your model's `forward` function should process the batch of binary input strings and output a $B \times 2$ tensor $y$ where $y_{b,0}$ is the score for the $b^{th}$ element of the batch having an even parity and $y_{b,1}$ for odd parity. You may use any PyTorch-defined LSTM functions. Larger hidden state sizes will make for easier training in my experiments. Running `driver_parity.py` will train your model and output per-epoch training loss and accuracy. A correctly-implemented model should approach 100% accuracy on the training set. In your write-up for this question, describe any architectural choices you made.

*Hint: Efficiently processing batches with variable input lengths requires some bookkeeping or else the LSTM will continue to process the padding for shorter sequences along with the content of longer ones. See pack_padded_sequence and pad_packed_sequence documentation in PyTorch.*

**Task 2.2**
`driver_parity.py` also evaluates your trained model on binary sequences of length 0 to 20 and saves a corresponding plot of accuracy vs. length. Include this plot in your write-up and describe the trend you observe. Why might the model behave this way?

**Task 2.3**
We know from 1.1 that even a univariate LSTM (one with a scalar hidden state) can theoretically solve this problem. Run a few (3-4) experiments with different hidden state sizes, what is the smallest size for which you can still train to fit this dataset? Feel free to adjust any of the hyper-parameters in the optimization in the `train_model` function if you want. Describe any trends you saw in training or the generalization experiment as you reduced the model capacity.

**Task 2.4**
It has been demonstrated that vanilla RNNs have a hard time learning to classify whether a string was generated by an ERG or not. LSTMs on the other hand seem to work fine. Based on the structure of the problem and what you know about recurrent networks, why might this be the case?

**Task 3.1**
The first step for any machine learning problem is to get familiar with the dataset. Read through random samples of the dataset and summarize what topics it seems to cover. Also look at the relationship between words and part-of-speech tags – what text preprocessing would be appropriate or inappropriate for this dataset? Produce a histogram of part-of-speech tags in the dataset – is it balanced between all tags? What word-level accuracy would a simple baseline that picked the majority label achieve?

**Task 3.2**
Create a file `driver_udpos.py` that implements and trains a bidirectional LSTM model on this dataset with cross entropy loss. The BiLSTM should predict an output distribution over the POS tags for each token in a sentence. In your written report, produce a graph of training and validation loss over the course of training. Your model should be able to achieve ¿70% per-word accuracy fairly easily.

To achieve stronger performance, you will likely need to tune hyper-parameters or model architecture to achieve lower validation loss. Using pretrained word vectors will likely help as well. You may also wish to employ early-stopping – regularly saving the weights of your model during training and then selecting the saved model with the lowest validation loss. In your report, describe any impactful decisions during this process. Importantly – `DO NOT EVALUATE ON TEST DURING THIS TUNING PROCESS`.

Once you are done finetuning, evaluate on the test split of the data and report the per-word accuracy.

**Task 3.3**
Implement a function `tag_sentence(sentence, model)` that processes an input sentence (a string) into a sequence of POS tokens. This will require you to tokenize/numeralize the sentence, pass it through your network, and then print the result. Use this function to tag the following sentences:

The old man the boat.

The complex houses married and single soldiers and their families.

The man who hunts ducks out on weekends.