# EE538 Neural Networks

## Homework 7

1. Let's consider a Bidirectional Associative Memory in Lecture 9.
   (a) Generate $S$ random binary vector pairs $(x^s, y^s)$ (for s=1,...,S) with 1024 and 512 elements for $x^s$ and $y^s$, respectively. For simplicity, each element of $x^s$ and $y^s$ are generated either -1 or +1 with a probability of 0.5. Try 3 datasets $D_{50}$, $D_{100}$, and $D_{200}$ with $S$=50, 100, and 200, respectively. (10 points)
   (b) Code a computer program to generate $y$ from $x$, and also $x$ from $y$ as
   $$y_j = \text{sgn}\left(\sum_{i=1}^{I} w_{ji} x_i\right), \quad x_i = \text{sgn}\left(\sum_{j=1}^{J} w_{ji} y_j\right), \quad w_{ji} = \sum_{m=1}^{M} x_i^m y_j^m. \text{ (10 points)}$$
   (c) For the dataset $D_{50}$, use each $x^s$ as an input vector $x$ and generate an output vector $y$ up to 10 iterations, i.e., 10 cycles of from $x$ to $y$, and then to $x$ again. Plot the average of all output errors, i.e., the number of different elements between the generated $y$ and true paired $y$, versus the epoch. (20 points)
   (d) Repeat (c) for noisy input vector $x$, which is $m$-elements different from one of the stored vector $x^{10}$. Use $m$ values of 0,1, 2, 3, 5, 7, 10, 15, 20, 30, 50, and 100, and generate 10 input vectors with random changes for each $m$ value. Make all the curves in a single plot for easy comparison. (15 points)
   (e) Repeat (d) for the dataset $D_{100}$ and $D_{200}$., and compare the results with those of (c). (15 points)

2. For the Transformer model, please answer the followings.
   (a) The Transformers usually use multi-heads, of which query and key embedding dimension is much smaller than that of the word imbedding. For example, $N_{mod} = N_{que} \times N_{head}$ and $N_{key} = N_{que}$. Discuss the advantages/disadvantages of this multi-head approach and the single-head approach with the same number of embedding dimensions. (10 points)
   (b) Please explain the idea behind the division by $\sqrt{N_{que}}$ before the Softmax operation. Hint. Read the footnotes in the Vaswani paper and judge the validity. Any other suggestion? (10 points)
   (c) Since $(q_{hp})^T k_{hp'} = (q_p)^T W_h^q (W_h^k)^T k_{p'}$, one may replace the two mapping matrices $W_h^q$ and $W_h^k$ with one matrix $W_h^{qk} = W_h^q (W_h^k)^T$. Discuss the advantages/disadvantages of these two different approaches. (10 points)