

Homework #5

TA in charge: Hanhsing Tu, Room 536

RELEASE DATE: 11/27/2008

DUE DATE: 12/04/2008, 4:00 pm IN CLASS

TA SESSION: 12/03/2008, noon to 2:00 pm IN R106

Unless granted by the instructor in advance, you must turn in a hard copy of your solutions (without the source code) for all problems. For problems marked with (), please follow the guidelines on the course website and upload your source code to designated places.*

Discussions on course materials and homework solutions are encouraged. But you should write the final solutions alone and understand them fully. Books, notes, and Internet resources can be consulted, but not copied from.

Since everyone needs to write the final solutions alone, there is absolutely no need to lend your homework solutions and/or source codes to your classmates at any time. In order to maximize the level of fairness in this class, lending and borrowing homework solutions are both regarded as dishonest behaviors and will be punished according to the honesty policy.

You should write your solutions in English with the common math notations introduced in class or in the problems. We do not accept solutions written in any other languages.

5.1 Powerfulness of Neural Networks

In this problem, we will work with Neural Networks using $\text{sign}(\cdot)$ instead of $\tanh(\cdot)$ as the activation functions. In addition, we will take $+1$ to mean logic TRUE, and -1 to mean logic FALSE. Then, we will show that Neural Networks can implement any Boolean function.

- (1) (5%) Assume that we want to compute

$$f_A(x) = \text{AND}((x)_1, (x)_2, \dots, (x)_d).$$

Show a perceptron

$$g_A(x) = \text{sign}\left(\sum_{i=1}^d w_i(x)_i - \theta\right).$$

that implements $f_A(x)$.

(Note: By **SHOW**, you should write down your choice of $(\theta, w_1, w_2, \dots, w_d)$ and prove that $g_A(x) = f_A(x)$ for all $x \in \{-1, +1\}^d$)

- (2) (5%) Similar to (1), show a perceptron $g_O(x)$ that implements

$$f_O(x) = \text{OR}((x)_1, (x)_2, \dots, (x)_d).$$

- (3) (5%) Similar to (1), show a perceptron $g_N(x)$ that implements

$$f_N(x) = \text{NOT}((x)_1).$$

- (4) (10%) Show a neural network that implements

$$\text{AND}\left((x)_1, \text{NOT}((x)_2), \text{OR}\left((x)_3, (x)_4, \text{NOT}((x)_1)\right)\right).$$

- (5) (5%) Argue how you can systematically implement any Boolean function with a neural network.

- (6) (10%) Show a 3-4-1 neural network that implements $\text{XOR}((x)_1, (x)_2, (x)_3)$.

- (7) (10%) Show a 3-3-1 neural network that implements $\text{XOR}((x)_1, (x)_2, (x)_3)$.

- (8) (Bonus 10%) Show a d - $d-1$ neural network that implements $XOR((x)_1, (x)_2, (x)_3, \dots, (x)_d)$.

(Note: No partial credits will be given to the bonus question! You need to give a complete and convincing proof to get the bonus.)

5.2 Limitations of Neural Networks

- (1) (10%) Prove that $(N^D + 1)^3 \leq N^{3D+1} + 1$ for $D \geq 1$ and $N \geq 3$.
- (2) (10%) Prove that if $N \geq 3\Delta \log_2 \Delta$, $N^\Delta + 1 < 2^N$.
- (3) (20%) Consider a learning model G_3 that consists of all d -3-1 neural networks with $\text{sign}(\cdot)$ as the activation functions. Use the facts above (or not) to prove that

$$VCD(G_3) < 3 \cdot (3(d+1) + 1) \log_2(3(d+1) + 1).$$

- (4) (10%) In Problem 5.1 we proved that Neural Networks can implement any Boolean function. In this problem (including the bonus one below) we proved that $VCD(G_M)$ is bounded. Are those facts contradicting each other? Explain your opinions.
- (5) (Bonus 10%) Consider a learning model G_M that consists of all d - M -1 neural networks with $\text{sign}(\cdot)$ as the activation functions. Prove that

$$VCD(G_M) = O(M \cdot d \cdot \log(M \cdot d)).$$

(Note: No partial credits will be given to the bonus question! You need to give a complete and convincing proof to get the bonus.)

5.3 Experiment with Weight Decay in Neural Network (*)

Modify your backpropagation algorithm for d - M -1 neural network from Problem 4.4 by changing the activation function of the output neuron **only** into a linear function

$$x^{\text{out}} = \left(\sum_{i=0}^{d^{\text{in}}} w_i \cdot x_i^{\text{in}} \right).$$

In other neurons, keep using

$$x^{\text{out}} = \tanh \left(\sum_{i=0}^{d^{\text{in}}} w_i \cdot x_i^{\text{in}} \right).$$

Implement the weight decay algorithm with stochastic gradient descent on this neural network. That is, update \mathbf{w} by

$$\mathbf{w}^{(t)} \leftarrow (1 - \eta\lambda) \cdot \mathbf{w}^{(t-1)} - \eta \cdot \nabla E_n(\mathbf{w}^{(t-1)}).$$

Run the algorithm on the following set for training:

http://www.csie.ntu.edu.tw/~htlin/course/ml08fall/data/hw5_3_train.dat

and the following set for testing:

http://www.csie.ntu.edu.tw/~htlin/course/ml08fall/data/hw5_3_test.dat

Take $M = 8$, $\eta = 0.05$, and $T = 5000$. Initialize the elements of $\mathbf{w}^{(0)}$ randomly from the range $(-0.02, 0.02)$. Now, consider $\lambda = 0$ (no weight decay), 0.001, 0.01, 0.1, or 1.

- (1) (50%) Let

$$g^{(t)}(x) = \text{sign}(\text{NN}(x_n, \mathbf{w}^{(t)})).$$

For the five λ values above, plot $\nu(g^{(t)})$, and $\hat{\pi}(g^{(t)})$ as functions of t on the same figure. Compare the curves across different λ , and briefly state your findings.

(Note: You should have five separate plots with scales, legends, labels, and captions clearly marked. You need to implement your own weight decay algorithm instead of using sophisticated packages)

5.4 Experiment with Naïve Bayes (*)

(1) (50%) Implement the Naïve Bayes algorithm below.

- (a) estimate $P(+)$ and $P(-)$ from Z by $P(+) = N_+/N$, and $P(-) = N_-/N$.
- (b) discretize the region $[0, 1)$ uniformly into K bins $\{[L_k, R_k]\}_{k=1}^K$.
- (c) for each combination of $s \in \{+, -\}$, $i \in \{1, 2, \dots, d\}$, and $k \in \{1, 2, \dots, K\}$
 - Estimate $P((x)_i | s)$ by

$$P((x)_i | s) = I[(x)_i \in [L_k, R_k]] P_{ik}^{(s)},$$

$$P_{ik}^{(s)} = \frac{1}{N_s} \cdot \sum_{n=1}^N I[(x_n)_i \in [L_k, R_k] \text{ and } y_n = s].$$

(d) Return

$$g(x) = \operatorname{argmax}_{s \in \{+, -\}} P(s | x),$$

$$\text{using } P(x | s) = \prod_{i=1}^d P((x)_i | s) \text{ and the Bayes theorem.}$$

Run the algorithm on the following set for training:

http://www.csie.ntu.edu.tw/~htlin/course/ml08fall/data/hw5_4_train.dat

and the following set for testing:

http://www.csie.ntu.edu.tw/~htlin/course/ml08fall/data/hw5_4_test.dat

Assume that g_K is the decision function returned when using parameter K . For $K = 2, 5, 10, 20, 50$, report $\nu(g_K)$ and $\hat{\pi}(g_K)$. Briefly state your findings.