

## Homework Assignment # 4

Assigned: 03/26/2021

Due: 04/12/2021, 11:59pm, through Canvas

Six problems, 170 points in total. Good luck!  
 Prof. Predrag Radivojac, Northeastern University

**Problem 1.** (15 points) Consider a logistic regression problem where  $\mathcal{X} = \mathbb{R}^d$  and  $\mathcal{Y} = \{-1, +1\}$ . Derive the weight update rule that maximizes the conditional likelihood assuming that a data set  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  is given.

**Problem 2.** (20 points) Consider a logistic regression problem with its initial solution obtained through the OLS regression; i.e.,  $\mathbf{w}^{(0)} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ , in the context of the code provided in class (week 6). Recall that  $\mathbf{x}$  was drawn from a mixture of two Gaussian distributions with  $\dim\{\mathbf{x}\} = 2$  (before adding a column of ones) and that  $y \in \{0, 1\}$ . You probably noticed that the initial separation line is consistently closer to the data points of class 0.

- (10 points) Why was this the case? Draw a picture (if possible) to support your argument.
- (5 points) Devise a better initial solution by modifying the standard formula  $\mathbf{w}^{(0)} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ .
- (5 points) Now again consider the case where  $y \in \{-1, +1\}$ . What is the form of the modified solution from part (b) in this case?

**Problem 3.** (40 points) Consider two classification concepts given in Figure 1, where  $x \in \mathcal{X} = [-6, 6] \times [-4, 4]$ ,  $y \in \mathcal{Y} = \{-1, +1\}$  and  $p(y|x) \in \{0, 1\}$  is defined in the drawing.

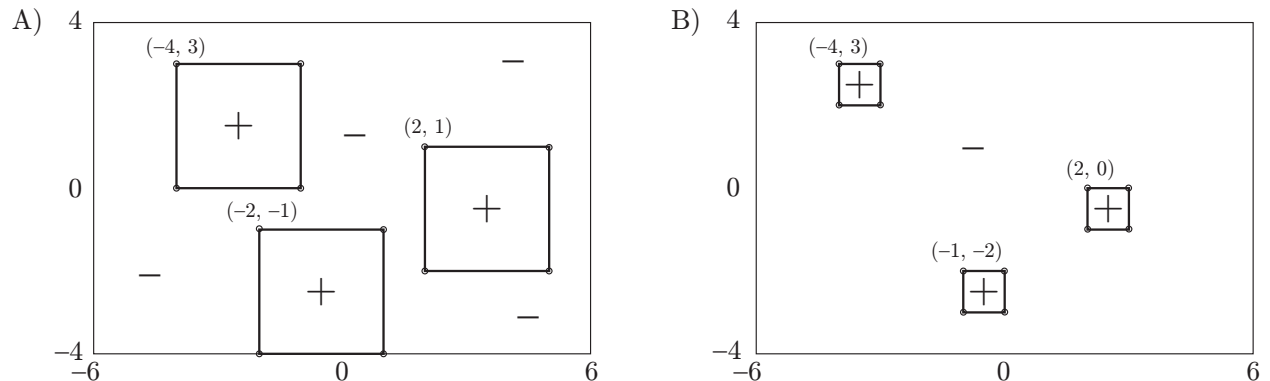


Figure 1: Two concepts where examples that fall within any of the three  $3 \times 3$  (panel A) or  $1 \times 1$  (panel B) squares are labeled positive and the remaining examples (outside each of the squares but within  $\mathcal{X}$ ) are labeled negative. The position of the point  $x = (x_1, x_2)$  in the upper left-hand corner for each square is shown in the picture. Consider horizontal axis to be  $x_1$  and vertical axis as  $x_2$ .

Your experiments in this question will rely on generating a data set of size  $n \in \{250, 1000, 10000\}$  drawn from a uniform distribution in  $\mathcal{X}$  and labeled according to the rules from Figure 1; e.g.,  $P(Y = 1|x) = 1$  if  $x$  that was randomly drawn is inside any of the three squares in either of the two panels, and  $P(Y = 1|x) = 0$  otherwise. The goal of the following two problems will be to train and evaluate classifiers created from the data generated in this way. You can use any library you want in this assignment and do programming in Python, MATLAB, R or C/C++. Your code should be easy to run for each question and sub-question below so that we can replicate your results to the maximum extent possible.

Consider single-output feed-forward neural networks with one or two hidden layers such that the number of hidden neurons in each layer is  $h_1 \in \{1, 4, 12\}$  and  $h_2 \in \{0, 3\}$ , respectively, with  $h_2 = 0$  meaning that there is no second hidden layer. Consider one of the standard objective functions as your optimization criterion and use early stopping and regularization as needed. Consider a hyperbolic tangent activation function in each neuron and the output but you are free to experiment with others if you'd like to. For each of the architectures, defined by a parameter combination  $(h_1, h_2)$ , evaluate the performance of each model using classification accuracy, balanced accuracy, and area under the ROC curve as your two performance criteria. To evaluate the performance of your models use cross-validation. However, to evaluate the performance of performance evaluation, generate another very large data set on a fine grid in  $\mathcal{X}$ . Then use the predictions from your trained model on all these points to determine the "true" performance. You can threshold your predictions in the middle of your prediction range (i.e., at 0.5 if you are predicting between 0 and 1) to determine binary predictions of your models and to then compare those with true class labels you generated on the fine grid.

Provide meaningful comments about all aspects of this exercise (performance results for different network architectures, accuracy of cross-validation, run time, etc.). The comments should not just re-state the results but rather capture trends and give reasoning as to why certain behavior was observed.

**Problem 4.** (70 points) Matrix factorization with applications. Consider an  $n \times d$  real-valued data matrix  $\mathbf{X} = (\mathbf{x}_1^T, \mathbf{x}_2^T, \dots, \mathbf{x}_n^T)$ . We will attempt to approximate this matrix using the following factorization

$$\hat{\mathbf{X}} = \mathbf{U}\mathbf{V}^T$$

where  $\mathbf{U} = (\mathbf{u}_1^T, \mathbf{u}_2^T, \dots, \mathbf{u}_n^T)$  is an  $n \times k$  and  $\mathbf{V} = (\mathbf{v}_1^T, \mathbf{v}_2^T, \dots, \mathbf{v}_d^T)$  is a  $d \times k$  matrix of real numbers, and where  $k < n, d$  is a parameter to be explored and determined. Notice that the value  $x_{ij}$  in  $\mathbf{X}$  can be approximated by  $\mathbf{u}_i^T \mathbf{v}_j$  and that  $\mathbf{x}_i^T$ , the  $i$ -th row of  $\mathbf{X}$ , can be approximated by  $\mathbf{u}_i^T \mathbf{V}^T$ , giving  $\hat{\mathbf{x}}_i = \mathbf{V}\mathbf{u}_i$ . We will formulate the matrix factorization process as the following minimization

$$\min_{\mathbf{U}, \mathbf{V}} \sum_{i,j} (x_{ij} - \mathbf{u}_i^T \mathbf{v}_j)^2 + \lambda (\sum_i \|\mathbf{u}_i\|^2 + \sum_j \|\mathbf{v}_j\|^2)$$

which minimizes the sum-of-squared-errors between real values  $x_{ij}$  and reconstructed values  $\hat{x}_{ij} = \mathbf{u}_i^T \mathbf{v}_j$ . The regularization parameter  $\lambda \geq 0$  is user-selected. This problem can be directly solved using gradient descent, but we will attempt a slightly different approach. To do this, we can see that for a fixed  $\mathbf{V}$  we can find optimal vectors  $\mathbf{u}_i$  by minimizing

$$\|\mathbf{V}\mathbf{u}_i - \mathbf{x}_i\|^2 + \lambda \cdot \|\mathbf{u}_i\|^2$$

which can be solved in a closed form using OLS regression for every  $i$ . We can equivalently express the optimization for vectors  $\mathbf{v}_j$  and find the solution for every  $j$ . Then, we can alternate these steps until convergence. This procedure is called the Alternating Least Squares (ALS) algorithm for matrix factorization. It has the following steps:

```
Initialize  $\mathbf{U}$  and  $\mathbf{V}$ 
repeat
  for  $i = 1$  to  $n$ 
```

```

     $\mathbf{u}_i = \text{formula \#1}$ 
  end
  for  $j = 1$  to  $d$ 
     $\mathbf{v}_j = \text{formula \#2}$ 
  end
until convergence

```

where you are expected to derive formula #1 and formula #2.

- (10 points) Derive the optimization steps for the ALS algorithm by finding formula #1 and formula #2 in the pseudocode listed above.
- (20 points) Consider now that some values in  $\mathbf{X}$  are missing (e.g., the rows of  $\mathbf{X}$  are users and the columns of  $\mathbf{X}$  are movie ratings, when available) and that we are interested in carrying out matrix completion using matrix factorization presented above. We would like to use the ALS algorithm, but the problem is that we must exclude all missing values from optimization. Derive now a modified ALS algorithm (formulas #1 and #2) to adapt it for matrix completion. Hint: consider adding an indicator matrix  $\mathbf{W}$  to the optimization process, where  $w_{ij} = 1$  if  $x_{ij}$  is available and  $w_{ij} = 0$  otherwise.
- (20 points) Consider now a MovieLens database available at

<http://grouplens.org/datasets/movielens/>

and find a data set most appropriate to evaluate your algorithm from the previous step; e.g., one of the 100k data sets. Now, implement the ALS algorithm on your data set and evaluate it using the mean-squared-error as the criterion of success. You can randomly remove 25% of the ratings, train a recommendation system, and then test it on the test set. You will have to make certain decisions yourselves, such as initialization of  $\mathbf{U}$  and  $\mathbf{V}$ , convergence criterion, or picking  $k$  and  $\lambda$ .

- (10 points) Describe your full experimentation process (e.g., how did you vary  $k$ ) and observations from (c). Additionally, can you provide some reasoning as to what  $k$  is and what matrices  $\mathbf{U}$  and  $\mathbf{V}$  are?
- (10 points) Compare your method against the baseline that fills in every missing movie rating value  $x_{ij}$  as an average over all users who have rated the movie  $j$ . Discuss your empirical findings.

**Problem 5.** (15 points) Prove representational equivalence of a three-layer neural network with linear activation function in all neurons and a single-layer layer neural network with the same activation function. Assume a single-output network.

**Problem 6.** (10 points) Let  $A$ ,  $B$ ,  $C$ , and  $D$  be binary input variables (features). Give decision trees to represent the following Boolean functions:

- (3 points)  $A \wedge \bar{B}$
- (3 points)  $A \vee (\bar{B} \wedge C)$
- (4 points)  $A \oplus \bar{B}$

where  $\bar{A}$  is the negation of  $A$  and  $\oplus$  is an exclusive OR operation.

### Directions and Policies

Submit a single package containing all answers, results and code. Your submission package should be compressed and named firstnamelastname.zip (e.g., predragradivojac.zip). In your package there should be a single pdf file named main.pdf that will contain answers to all questions, all figures, and all relevant results. Your solutions and answers must be typed<sup>1</sup> and make sure that you type your name and Northeastern username (email) on top of the first page of the main.pdf file. The rest of the package should contain all code that you used. The code should be properly organized in folders and subfolders, one for each question or problem. All code, if applicable, should be turned in when you submit your assignment as it may be necessary to demo your programs to the teaching assistants. Use Matlab, Python, R, Java, or C/C++. However, you are encouraged to use languages with good machine learning libraries (e.g., Matlab, Python, R), which may be handy in future assignments.

Unless there are legitimate circumstances, late assignments will be accepted up to 5 days after the due date and graded using the following rules:

on time: your score  $\times$  1

1 day late: your score  $\times$  0.9

2 days late: your score  $\times$  0.7

3 days late: your score  $\times$  0.5

4 days late: your score  $\times$  0.3

5 days late: your score  $\times$  0.1

For example, this means that if you submit 3 days late and get 80 points for your answers, your total number of points will be  $80 \times 0.5 = 40$  points.

All assignments are individual, except when collaboration is explicitly allowed. **All text must be your own or, for group assignments, by the members of the group. All sources used for problem solution must be acknowledged;** e.g., web sites, books, research papers, personal communication with people, etc. Academic honesty is taken seriously! For detailed information see Office of Student Conduct and Conflict Resolution.

---

<sup>1</sup>We recommend Latex; in particular, TexShop-MacTeX combination for a Mac and TeXnicCenter-MiKTeX combination on Windows. An easy way to start with Latex is to use the freely available Lyx. You can also use Microsoft Word or other programs that can display formulas professionally.