

CS 5350/6350: Machine Learning Spring 2019

Homework 5

Handed out: 9 Apr, 2019
Due date: 11:59pm, 24 Apr, 2019

- You are welcome to talk to other members of the class about the homework. I am more concerned that you understand the underlying concepts. However, you should write down your own solution. Please keep the class collaboration policy in mind.
- Feel free to discuss the homework with the instructor or the TAs.
- Your written solutions should be brief and clear. You do not need to include original problem descriptions in your solutions. You need to show your work, not just the final answer, but you do *not* need to write it in gory detail. Your assignment should be **no more than 15 pages**. Every extra page will cost a point.
- Handwritten solutions will not be accepted.
- *Your code should run on the CADE machines. You should include a shell script, `run.sh`, that will execute your code in the CADE environment. Your code should produce similar output to what you include in your report.*
You are responsible for ensuring that the grader can execute the code using only the included script. If you are using an esoteric programming language, you should make sure that its runtime is available on CADE.
- Please do not hand in binary files! We will *not* grade binary submissions.
- The homework is due by **midnight of the due date**. Please submit the homework on Canvas.

1 Paper Problems [40 points]

1. [5 points] (Warm up) Suppose we have a composite function, $z = \sigma(y_1^2 + y_2 y_3)$, where $y_1 = 3x$, $y_2 = e^{-x}$, $y_3 = \sin(x)$, and $\sigma(\cdot)$ is the sigmoid activation function. Please use the chain rule to derive $\frac{\partial z}{\partial x}$ and compute the derivative at $x = 0$.
$$\begin{aligned}\frac{\partial z}{\partial x} &= \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial x} + \frac{\partial z}{\partial y_2} \frac{\partial y_2}{\partial x} + \frac{\partial z}{\partial y_3} \frac{\partial y_3}{\partial x} \\ &= \left(\frac{2e^{-x^2 - y_2 x}}{(1 + e^{-x^2 - y_2 x})^2} \right) (3) + \left(\frac{e^{-y_2 - x^2} z}{(1 + e^{-x^2 - y_2 x})^2} \right) (-e^{-x}) + \left(\frac{e^{-y_2 - x^2} y}{(1 + e^{-x^2 - y_2 x})^2} \right) (\cos(x)) \\ &= \frac{6e^{-x^2 - y_2 x} x - e^{-x^2 - x - y_2 z} + e^{-x^2 - y_2 y} \cos(x)}{(e^{-x^2 - y_2 x} + 1)^2}\end{aligned}$$
2. [5 points] Suppose we have a three-layered feed-forward neural network in hand. The architecture and the weights are defined in Figure 1. We use the sigmoid activation function. Note that the shaded variables are the constant feature 1, i.e., $x_0 = z_0^1 =$

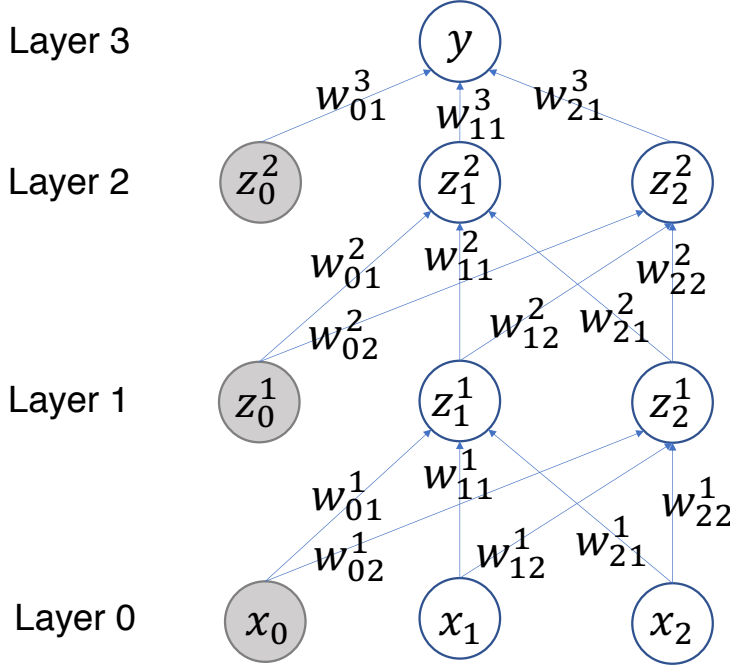


Figure 1: A three layer artificial neural network.

$z_0^2 = 1$. As we discussed in the class, they are used to account for the bias parameters. We have the values of all the edge weights in Table 1. Now, given a new input example $\mathbf{x} = [1, 1, 1]$. Please use the forward pass to compute the output y . Please list every step in your computation, namely, how you calculate the variable value in each hidden unit, and how you combine the variables in one layer to compute each variable in the next layer. Please be aware of the subtle difference in computing the variable value in the last layer (we emphasized it in the class).

$$z_1^1 = \sigma(w_{01}^1 + w_{11}^1 x_1 + w_{21}^1 x_2) = \sigma(-1 + -2 * 1 + -3 * 1) = \frac{1}{1+e^6} = 0.00247$$

$$z_2^1 = \sigma(w_{02}^1 + w_{12}^1 x_1 + w_{22}^1 x_2) = 0.99753$$

$$z_1^2 = \sigma(w_{01}^2 + w_{11}^2 z_1^1 + w_{21}^2 z_2^1) = 0.01803$$

$$z_2^2 = \sigma(w_{02}^2 + w_{12}^2 z_1^1 + w_{22}^2 z_2^1) = 0.98197$$

$$y = w_{01}^3 + w_{11}^3 z_1^2 + w_{21}^3 z_2^2 = -1 + 2 * 0.01803 + -1.5 * 0.98197 = -2.436895$$

3. [20 points] Suppose we have a training example where the input vector is $\mathbf{x} = [1, 1, 1]$ and the label $y^* = 1$. We use a square loss for the prediction,

$$L(y, y^*) = \frac{1}{2}(y - y^*)^2.$$

To make the prediction, we will use the 3 layer neural network shown in Figure 1, with the sigmoid activation function. Given the weights specified in Table 1, please use the back propagation (BP) algorithm to compute the derivative of the loss L over all the weights, $\{\frac{\partial L}{\partial w_{ij}^m}\}$. Please list every step of your BP calculation. In each step, you should

Layer	weight	value
1	w_{01}^1	-1
1	w_{02}^1	1
1	w_{11}^1	-2
1	w_{12}^1	2
1	w_{21}^1	-3
1	w_{22}^1	3
2	w_{01}^2	-1
2	w_{02}^2	1
2	w_{11}^2	-2
2	w_{12}^2	2
2	w_{21}^2	-3
2	w_{22}^2	3
3	w_{01}^3	-1
3	w_{11}^3	2
3	w_{21}^3	-1.5

Table 1: Weight values.

show how you compute and cache the new (partial) derivatives from the previous ones, and then how to calculate the partial derivative over the weights accordingly.

From the previous problem we have $y = -2.436895$

$$\frac{\partial L}{\partial w_{01}^3} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{01}^3} = (y - y^*)(1) = -2.436895 - 1 = -3.436895$$

$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{11}^3} = (y - y^*)(z_1^2) = -0.06197$$

$$\frac{\partial L}{\partial w_{21}^3} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial w_{21}^3} = (y - y^*)(z_2^2) = -3.37493$$

$$\frac{\partial L}{\partial w_{01}^2} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_1^2} \frac{\partial z_1^2}{\partial w_{01}^2} = (y - y^*)(w_{11}^3) \frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{01}^2} = (-3.436895)(2)(\sigma(s)(1 - \sigma(s)))(1) = (-3.436895)(2)(0.01803)(0.98197)(1) = -0.12170$$

$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_1^2} \frac{\partial z_1^2}{\partial w_{11}^2} = (y - y^*)(w_{11}^3) \frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{11}^2} = (-3.436895)(2)(0.01803)(0.98197)(0.00247) = -0.00030$$

$$\frac{\partial L}{\partial w_{21}^2} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_1^2} \frac{\partial z_1^2}{\partial w_{21}^2} = (y - y^*)(w_{11}^3) \frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{21}^2} = (-3.436895)(2)(0.01803)(0.98197)(0.99753) = -0.12140$$

$$\frac{\partial L}{\partial w_{02}^2} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_2^2} \frac{\partial z_2^2}{\partial w_{02}^2} = (y - y^*)(w_{21}^3) \frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{02}^2} = (-3.436895)(-1.5)(0.98197)(0.01803)(1) = 0.09127$$

$$\frac{\partial L}{\partial w_{12}^2} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_2^2} \frac{\partial z_2^2}{\partial w_{12}^2} = (y - y^*)(w_{21}^3) \frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{12}^2} = (-3.436895)(-1.5)(0.98197)(0.01803)(0.00247) = 0.00022$$

$$\frac{\partial L}{\partial w_{22}^2} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial z_2^2} \frac{\partial z_2^2}{\partial w_{22}^2} = (y - y^*)(w_{21}^3) \frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{22}^2} = (-3.436895)(-1.5)(0.98197)(0.01803)(0.99753) = 0.09105$$

$$\frac{\partial L}{\partial w_{01}^1} = \frac{\partial L}{\partial y} \left(\frac{\partial y}{\partial z_1^1} \frac{\partial z_1^1}{\partial w_{01}^1} + \frac{\partial y}{\partial z_2^1} \frac{\partial z_2^1}{\partial w_{01}^1} \right) = (y - y^*) \left((w_{11}^3) \left(\frac{\partial z_1^2}{\partial z_1^1} \frac{\partial z_1^1}{\partial w_{01}^1} \right) + (w_{21}^3) \left(\frac{\partial z_2^2}{\partial z_1^1} \frac{\partial z_1^1}{\partial w_{01}^1} \right) \right) = (y - y^*) \left((w_{11}^3) \left((w_{11}^2) \left(\frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{01}^1} \right) \right) + (w_{21}^3) \left((w_{12}^2) \left(\frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{01}^1} \right) \right) \right) = (-3.436895) \left((2) \left((-2) \left(((0.00247)(0.99753))(1) \right) + \right. \right.$$

$$(-1.5)((2)((0.00247)(0.99753))(1)))) = 0.08468$$

$$\frac{\partial L}{\partial w_{11}^1} = (y - y^*)((w_{11}^3)((w_{11}^2)(\frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{11}^1})) + (w_{21}^3)((w_{12}^2)(\frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{11}^1}))) = (-3.436895)((2)((-2)((0.00247)(0.99753))(1)))) = 0.08468$$

$$\frac{\partial L}{\partial w_{21}^1} = (y - y^*)((w_{11}^3)((w_{11}^2)(\frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{21}^1})) + (w_{21}^3)((w_{12}^2)(\frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{21}^1}))) = (-3.436895)((2)((-2)((0.00247)(0.99753))(1)))) = 0.08468$$

$$\frac{\partial L}{\partial w_{02}^1} = (y - y^*)((w_{11}^3)((w_{21}^2)(\frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{02}^1})) + (w_{21}^3)((w_{22}^2)(\frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{02}^1}))) = (-3.436895)((2)((-3)((0.99753)(0.00247))(1)))) = 0.12702$$

$$\frac{\partial L}{\partial w_{12}^1} = (y - y^*)((w_{11}^3)((w_{21}^2)(\frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{12}^1})) + (w_{21}^3)((w_{22}^2)(\frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{12}^1}))) = (-3.436895)((2)((-3)((0.99753)(0.00247))(1)))) = 0.12702$$

$$\frac{\partial L}{\partial w_{22}^1} = (y - y^*)((w_{11}^3)((w_{21}^2)(\frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{22}^1})) + (w_{21}^3)((w_{22}^2)(\frac{\partial \sigma}{\partial s} \frac{\partial s}{\partial w_{22}^1}))) = (-3.436895)((2)((-3)((0.99753)(0.00247))(1)))) = 0.12702$$

4. [10 points] Suppose we have the training dataset shown in Table 2. We want to learn a logistic regression model. We initialize all the model parameters with 0. We assume each parameter (i.e., feature weights $\{w_1, w_2, w_3\}$ and the bias w_0) comes from a standard Gaussian prior distribution,

$$p(w_i) = \mathcal{N}(w_i | 0, 1) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}w_i^2) \quad (0 \leq i \leq 3).$$

- [7 points] We want to obtain the maximum a posteriori (MAP) estimation. Please write down the objective function, namely, the log joint probability, and derive the gradient of the objective function.

Objective function:

$$\min_{\mathbf{w}} \sum_i^m \log(1 + \exp(-y_i \mathbf{w}^\top \mathbf{x}_i)) + \frac{1}{2} \mathbf{w}^\top \mathbf{w}$$

$$\text{Gradient: } \begin{bmatrix} \sum_i^m -\frac{y_i \mathbf{x}_{i0} e^{-y_i \mathbf{w}^\top \mathbf{x}_i}}{1 + e^{-y_i \mathbf{w}^\top \mathbf{x}_i}} + \mathbf{w}_0 \\ \vdots \\ \sum_i^m -\frac{y_i \mathbf{x}_{in} e^{-y_i \mathbf{w}^\top \mathbf{x}_i}}{1 + e^{-y_i \mathbf{w}^\top \mathbf{x}_i}} + \mathbf{w}_n \end{bmatrix}$$

- [3 points] We set the learning rates for the first three steps to $\{0.01, 0.005, 0.0025\}$. Please list the stochastic gradients of the objective w.r.t the model parameters for the first three steps, when using the stochastic gradient descent algorithm.

Step 0:

$$\mathbf{w} = [0, 0, 0, 0]$$

Step 1:

$$\nabla J(\mathbf{w}) = \{\text{content}...\}$$

x_1	x_2	x_3	y
0.5	-1	0.3	1
-1	-2	-2	-1
1.5	0.2	-2.5	1

Table 2: Dataset

2 Practice [62 points + 50 bonus]

1. [2 Points] Update your machine learning library. Please check in your implementation of SVM algorithms. Remember last time you created the folders “SVM”. You can commit your code into the corresponding folders now. Please also supplement README.md with concise descriptions about how to use your code to run these algorithms (how to call the command, set the parameters, etc). Please create new folders “Neural Networks” and “Logistic Regression” in the same level as these folders. *After the completion of the homework this time, please check in your implementation accordingly.*
2. [18 points] We will implement the logistic regression model with stochastic gradient descent. We will reuse the dataset “bank-note.zip” in Canvas. The features and labels are listed in the file “classification/data-desc.txt”. The training data are stored in the file “classification/train.csv”, consisting of 872 examples. The test data are stored in “classification/test.csv”, and comprise of 500 examples. In both the training and test datasets, feature values and labels are separated by commas. Set the maximum number of epochs T to 100. Don’t forget to shuffle the training examples at the start of each epoch. Use the curve of the objective function (along with the number of updates) to diagnosis the convergence. We initialize all the model parameters with 0.
 - (a) [10 points] We will first obtain the MAP estimation. In order for that, we assume each model parameter comes from a Gaussian prior distribution,

$$p(w_i) = \mathcal{N}(w_i|0, v) = \frac{1}{\sqrt{2\pi v}} \exp\left(-\frac{1}{2v} w_i^2\right)$$

where v is the variance. From the paper problem 4, you should be able to write down the objective function and derive the gradient. Try the prior variance v from $\{0.01, 0.1, 0.5, 1, 3, 5, 10, 100\}$. Use the schedule of learning rate: $\gamma_t = \frac{\gamma_0}{1 + \frac{\gamma_0}{d} t}$. Please tune γ_0 and d to ensure convergence. For each setting of variance, report your training and test error.

$\gamma_0 = 0.0000001, d = 0.000001$

(see next page)

- (b) [5 points] We will then obtain the maximum likelihood (ML) estimation. That is, we do not assume any prior over the model parameters, and just maximize the logistic likelihood of the data. Use the same learning rate schedule. Tune γ_0 and d to ensure convergence. For each setting of variance, report your training and test error.

v	Training Error	Testing Error
0.01	0.209862385321	0.214
0.1	0.202981651376	0.204
0.5	0.292431192661	0.282
1	0.279816513761	0.304
3	0.235091743119	0.232
5	0.313073394495	0.32
10	0.185779816514	0.19
100	0.209862385321	0.216

$$\gamma_0 = 0.0000001, d = 0.000001$$

v	Training Error	Testing Error
0.01	0.262614678899	0.262
0.1	0.229357798165	0.25
0.5	0.295871559633	0.286
1	0.316513761468	0.294
3	0.302752293578	0.288
5	0.292431192661	0.29
10	0.280963302752	0.276
100	0.297018348624	0.316

- (c) [3 points] How is the training and test performance of the MAP estimation compared with the ML estimation? What can you conclude? What do you think of v , as compared to the hyperparameter C in SVM?

The MAP estimation seems to perform better than the MLE in general. The training error is roughly the same, however test error appears to vary with the variance for MAP, achieving better and worse testing performance than MLE. This is likely because the variance is creating a model which generalizes better. The hyperparameter C acted as a modifier on the learning rate, while v adds variance directly to the gradient. It is easier to imagine the affect C will have on \mathbf{w} , but v can scale up and down.

3. [40 points] Now let us implement a three layer artificial neural network for classification. We will use the same dataset, “bank-note.zip”. The architecture resembles Figure 1, but we allow an arbitrary number of units in hidden layers (Layer 1 and 2). So please ensure your implementation has such flexibility. We will use the sigmoid activation function.
- (a) [20 points] Please implement the back-propagation algorithm to compute the gradient with respect to all the edge weights given one training example. For debugging, you can use the paper problem 3 and verify if your algorithm returns the same derivatives as you manually did.

- (b) [12 points] Implement the stochastic gradient descent algorithm to learn the neural network from the training data. Use the schedule of learning rate: $\gamma_t = \frac{\gamma_0}{1 + \frac{\gamma_0}{d}t}$. Initialize the edge weights with random numbers generated from the standard Gaussian distribution. We restrict the width, i.e., the number of nodes, of each hidden layer (i.e., Layer 1 & 2) to be identical. Vary the width from $\{5, 10, 25, 50, 100\}$. Please tune γ_0 and d to ensure convergence. Report the training and test error for each setting of the width.
 - (c) [5 points]. Now initialize all the weights with 0, and run your training algorithm again. What is your training and test error? What do you observe and conclude?
 - (d) [3 points]. As compared with the performance of the logistic regression and SVM, what do you conclude (empirically) about the neural network?
 - (e) [**Bonus**] [50 points] Please use tensor-flow (TF) to fulfill the neural network training and prediction. Please try two activation functions, “tanh” and “RELU”. For “tanh”, please use the “Xavier” initialization; and for “RELU”, please use the “he” initialization. You can implement these initializations by yourselves or use TF library. Vary the depth from $\{3, 5, 9\}$ and width from $\{5, 10, 25, 50, 100\}$. Please use the Adam optimizer for training. The default settings of Adam should be sufficient. Report the training and test error with each (depth, width) combination. What do you observe and conclude? Note that, we won’t provide any link or manual for you to work on this bonus problem. It is YOUR JOB to search the document and web pages, find code snippets, and test and debug with TF to ensure the correct usage of TF. This is what all machine learning practitioners do in practice.
4. [2 Points] After the completion, please upload the implementation to your Github repository immediately. How do you like your own machine learning library? *Although it is still light weighted, it is the proof of your great efforts and achievement in this class! It is an excellent start of your journey to machine learning. Wish you further success in your future endeavours!*