

HOMWORK 5: NEURAL NETWORKS

10-301/10-601 Introduction to Machine Learning (Summer 2022)
<https://www.cs.cmu.edu/~hchai2/courses/10601/>

OUT: Wednesday, June 22

DUE: Wednesday, July 6 at 1:00 PM

TAs: Sana, Brendon, Neural, Ayush, Boyang (Jack), Chu

Summary In this assignment, you will build an image recognition system using a neural network. In the Written component, you will walk through an on-paper example of how to implement a neural network. Then, in the Programming component, you will implement an end-to-end system that learns to perform image classification.

START HERE: Instructions

- **Collaboration Policy:** Please read the collaboration policy here: <https://www.cs.cmu.edu/~hchai2/courses/10601/#Syllabus>
- **Late Submission Policy:** See the late submission policy here: <https://www.cs.cmu.edu/~hchai2/courses/10601/#Syllabus>
- **Submitting your work:** You will use Gradescope to submit answers to all questions and code. Please follow instructions at the end of this PDF to correctly submit all your code to Gradescope.
 - **Written:** For written problems such as short answer, multiple choice, derivations, proofs, or plots, please use the provided template. Submissions can be handwritten onto the template, but should be labeled and clearly legible. If your writing is not legible, you will not be awarded marks. Alternatively, submissions can be written in LaTeX. Each derivation/proof should be completed in the boxes provided. You are responsible for ensuring that your submission contains exactly the same number of pages and the same alignment as our PDF template. If you do not follow the template, your assignment may not be graded correctly by our AI assisted grader.
 - **Programming:** You will submit your code for programming questions on the homework to Gradescope (<https://gradescope.com>). After uploading your code, our grading scripts will autograde your assignment by running your program on a virtual machine (VM). When you are developing, check that the version number of the programming language environment (Python 3.9.6) and versions of permitted libraries (numpy 1.21.2) match those used on Gradescope. You have 10 free Gradescope programming submissions. After 10 submissions, you will begin to lose points from your total programming score. We recommend debugging your implementation locally first before submitting your code to Gradescope.
- **Materials:** The data and reference output that you will need in order to complete this assignment is posted along with the writeup and template on the course website.

Instructions for Specific Problem Types

For “Select One” questions, please fill in the appropriate bubble completely:

Select One: Who taught this course?

- ☒ Henry Chai
- ☐ Marie Curie
- ☐ Noam Chomsky

If you need to change your answer, you may cross out the previous answer and bubble in the new answer:

Select One: Who taught this course?

- ☒ Henry Chai
- ☐ Marie Curie
- ☒ Noam Chomsky

For “Select all that apply” questions, please fill in all appropriate squares completely:

Select all that apply: Which are scientists?

- ☒ Stephen Hawking
- ☒ Albert Einstein
- ☒ Isaac Newton
- ☐ I don't know

Again, if you need to change your answer, you may cross out the previous answer(s) and bubble in the new answer(s):

Select all that apply: Which are scientists?

- ☒ Stephen Hawking
- ☒ Albert Einstein
- ☒ Isaac Newton
- ☒ I don't know

For questions where you must fill in a blank, please make sure your final answer is fully included in the given space. You may cross out answers or parts of answers, but the final answer must still be within the given space.

Fill in the blank: What is the course number?

10-601

10-~~6~~301

Written Questions (47 points)

1 L^AT_EX Bonus Point (1 points)

1. (1 point) **Select One:** Did you use L^AT_EX for the entire written portion of this homework?

☐ Yes

☐ No

2 Example Feed Forward and Backpropagation (28 points)

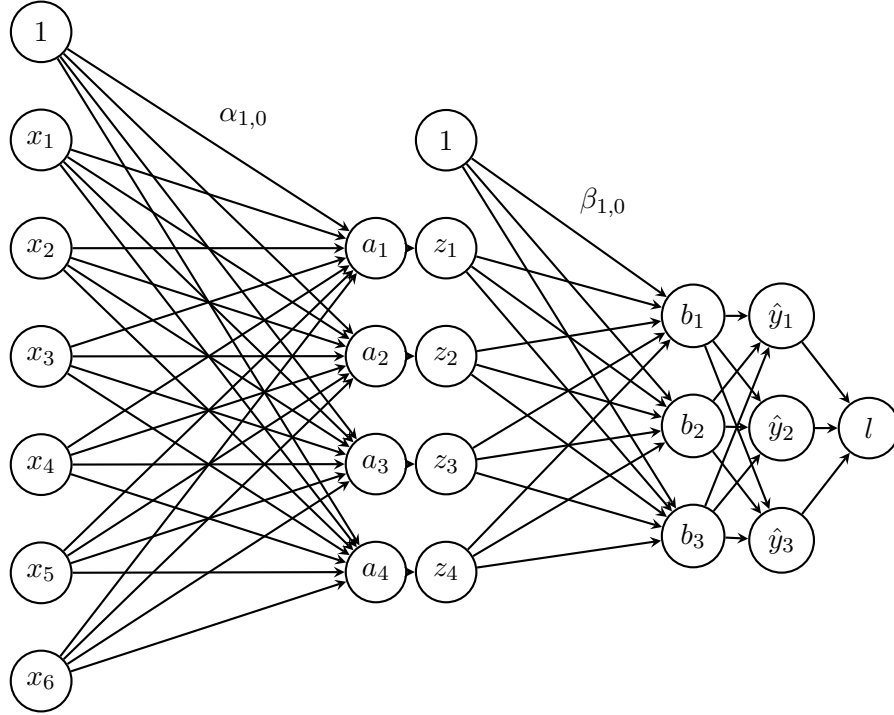


Figure 1: A One Hidden Layer Neural Network

Network Overview Consider the neural network with one hidden layer shown in Figure 1. The input layer consists of 6 features $\mathbf{x} = [x_1, \dots, x_6]^T$, the hidden layer has 4 nodes $\mathbf{z} = [z_1, \dots, z_4]^T$, and the output layer is a probability distribution $\mathbf{y} = [y_1, y_2, y_3]^T$ over 3 classes (**1-indexed** such that y_i is the probability of label i).

α^* is the matrix of weights from the inputs to the hidden layer and β^* is the matrix of weights from the hidden layer to the output layer.

$\alpha_{j,i}^*$ represents the weight going to the node z_j in the hidden layer from the node x_i in the input layer (e.g. $\alpha_{1,2}^*$ is the weight from x_2 to z_1), and β^* is defined similarly. We will use a sigmoid activation function for the hidden layer and a softmax for the output layer.

The bias vectors α_b, β_b are defined such that the j th value of α_b (which we denote $\alpha_{j,b}$) is the bias value for a_j and the k th value of β_b is the bias value for b_k .

Network Details Equivalently, we define each of the following.

The input:

$$\mathbf{x} = [x_1, x_2, x_3, x_4, x_5, x_6]^T \quad (1)$$

Linear combination at the first (hidden) layer:

$$a_j = \alpha_{j,b} + \sum_{i=1}^6 \alpha_{j,i}^* \cdot x_i, \quad \forall j \in \{1, \dots, 4\} \quad (2)$$

Activation at the first (hidden) layer:

$$z_j = \sigma(a_j) = \frac{1}{1 + \exp(-a_j)}, \quad \forall j \in \{1, \dots, 4\} \quad (3)$$

Equivalently, we can write this as vector operation where the sigmoid activation is applied individually to each element of the vector \mathbf{a} :

$$\mathbf{z} = \sigma(\mathbf{a}) \quad (4)$$

Linear combination at the second (output) layer:

$$b_k = \beta_{k,b} + \sum_{j=1}^4 \beta_{k,j}^* \cdot z_j, \quad \forall k \in \{1, \dots, 3\} \quad (5)$$

Activation at the second (output) layer:

$$\hat{y}_k = \frac{\exp(b_k)}{\sum_{l=1}^3 \exp(b_l)}, \quad \forall k \in \{1, \dots, 3\} \quad (6)$$

Loss We will use cross entropy loss, $\ell(\hat{\mathbf{y}}, \mathbf{y})$. If \mathbf{y} represents our target output, which will be a one-hot vector representing the correct class, and $\hat{\mathbf{y}}$ represents the output of the network, the loss is calculated by:

$$\ell(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_{i=1}^3 y_i \log(\hat{y}_i) \quad (7)$$

For the below questions use natural log in the equation.

Prediction When doing prediction, we will predict the argmax of the output layer. For example, if $\hat{y}_1 = 0.3, \hat{y}_2 = 0.2, \hat{y}_3 = 0.5$ we would predict class 3. If the true class from the training data was 2 we would have a one-hot vector \mathbf{y} with values $y_1 = 0, y_2 = 1, y_3 = 0$.

1. In the following questions you will derive the matrix and vector forms of the previous equations which define our neural network. These are what you should hope to program in order to keep your program under the Gradescope time-out.

When working these out it is important to keep a note of the vector and matrix dimensions in order for you to easily identify what is and isn't a valid multiplication. Suppose you are given an training example: $\mathbf{x}^{(1)} = [x_1, x_2, x_3, x_4, x_5, x_6]^T$ with **label class 2**, so $\mathbf{y}^{(1)} = [0, 1, 0]^T$. We initialize the network weights as:

$$\alpha^* = \begin{bmatrix} \alpha_{1,1} & \alpha_{1,2} & \alpha_{1,3} & \alpha_{1,4} & \alpha_{1,5} & \alpha_{1,6} \\ \alpha_{2,1} & \alpha_{2,2} & \alpha_{2,3} & \alpha_{2,4} & \alpha_{2,5} & \alpha_{2,6} \\ \alpha_{3,1} & \alpha_{3,2} & \alpha_{3,3} & \alpha_{3,4} & \alpha_{3,5} & \alpha_{3,6} \\ \alpha_{4,1} & \alpha_{4,2} & \alpha_{4,3} & \alpha_{4,4} & \alpha_{4,5} & \alpha_{4,6} \end{bmatrix}$$

$$\beta^* = \begin{bmatrix} \beta_{1,1} & \beta_{1,2} & \beta_{1,3} & \beta_{1,4} \\ \beta_{2,1} & \beta_{2,2} & \beta_{2,3} & \beta_{2,4} \\ \beta_{3,1} & \beta_{3,2} & \beta_{3,3} & \beta_{3,4} \end{bmatrix}$$

We want to also consider the bias term and the weights on the bias terms ($\alpha_{j,b}$ and $\beta_{k,b}$). To account for these we can add them as a new column to the beginning of our initial weight matrices to represent biases, (e.g. $\alpha_{1,0} = \alpha_{1,b}$).

$$\alpha = \begin{bmatrix} \alpha_{1,0} & \alpha_{1,1} & \alpha_{1,2} & \alpha_{1,3} & \alpha_{1,4} & \alpha_{1,5} & \alpha_{1,6} \\ \alpha_{2,0} & \alpha_{2,1} & \alpha_{2,2} & \alpha_{2,3} & \alpha_{2,4} & \alpha_{2,5} & \alpha_{2,6} \\ \alpha_{3,0} & \alpha_{3,1} & \alpha_{3,2} & \alpha_{3,3} & \alpha_{3,4} & \alpha_{3,5} & \alpha_{3,6} \\ \alpha_{4,0} & \alpha_{4,1} & \alpha_{4,2} & \alpha_{4,3} & \alpha_{4,4} & \alpha_{4,5} & \alpha_{4,6} \end{bmatrix}$$

$$\beta = \begin{bmatrix} \beta_{1,0} & \beta_{1,1} & \beta_{1,2} & \beta_{1,3} & \beta_{1,4} \\ \beta_{2,0} & \beta_{2,1} & \beta_{2,2} & \beta_{2,3} & \beta_{2,4} \\ \beta_{3,0} & \beta_{3,1} & \beta_{3,2} & \beta_{3,3} & \beta_{3,4} \end{bmatrix}$$

We then add a corresponding new first dimension to our input vectors, always set to 1 ($x_0^{(i)} = 1$), so our input becomes:

$$\mathbf{x}^{(1)} = [1, x_1, x_2, x_3, x_4, x_5, x_6]^T$$

- (a) (1 point) By examining the shapes of the initial weight matrices, how many neurons do we have in the first hidden layer of the neural network? Do not include the bias in your count.

Answer

- (b) (1 point) How many output neurons will our neural network have?

Answer

- (c) (1 point) What is the vector \mathbf{a} whose elements are made up of the entries a_j in Equation 2 (using $x_i^{(1)}$ in place of x_i). Write your answer in terms of α and $\mathbf{x}^{(1)}$.

Answer

- (d) (1 point) **Select all that apply:** We cannot take the matrix multiplication of our weights β and the vector $\mathbf{z} = [z_1, z_2, z_3, z_4]^T$ since they are not compatible shapes. Which of the following would allow us to take the matrix multiplication of β and \mathbf{z} such that the entries of the vector $\mathbf{b} = \beta\mathbf{z}$ are equivalent to the values of b_k in Equation 5?

- ☐ Remove the last column of β
- ☐ Remove the first row of \mathbf{z}
- ☐ Append a value of 1 to be the first entry of \mathbf{z} .
- ☐ Append an additional column of 1's to be the first column of β
- ☐ Append a row of 1's to be the first row of β
- ☐ Take the transpose of β

- (e) (1 point) What are the entries of the output vector $\hat{\mathbf{y}}$? Your answer should be written in terms of b_1, b_2, b_3 .

$\hat{\mathbf{y}}$

2. We will now derive the matrix and vector forms for the backpropagation algorithm, for example

$$\frac{\partial \ell}{\partial \alpha} = \begin{bmatrix} \frac{\partial \ell}{\partial \alpha_{10}} & \frac{\partial \ell}{\partial \alpha_{11}} & \cdots & \frac{\partial \ell}{\partial \alpha_{16}} \\ \frac{\partial \ell}{\partial \alpha_{20}} & \frac{\partial \ell}{\partial \alpha_{21}} & \cdots & \frac{\partial \ell}{\partial \alpha_{26}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \ell}{\partial \alpha_{40}} & \frac{\partial \ell}{\partial \alpha_{41}} & \cdots & \frac{\partial \ell}{\partial \alpha_{46}} \end{bmatrix}$$

The mathematics which you have to derive in this section jump significantly in difficulty; you should always be examining the shape of the matrices and vectors and making sure that you are comparing your matrix elements with calculations of individual derivatives to make sure they match (e.g., the element of the matrix $(\frac{\partial \ell}{\partial \alpha})_{2,1}$ should be equal to $\frac{\partial \ell}{\partial \alpha_{2,1}}$). Recall that ℓ is our loss function defined in Equation 7:

Note: all vectors are column vectors (i.e. an n -dimensional vector $v \in \mathbb{R}^{n \times 1}$). **Assume that all input vectors to linear layers have a bias term folded in, unless otherwise specified.** All partial derivatives should be written in [denominator layout notation](#). An example of denominator notation is that $\frac{\partial \ell}{\partial \beta} \in \mathbb{R}^{3 \times 5}$ because $\beta \in \mathbb{R}^{3 \times 5}$.

- (a) (1 point) What is the derivative $\frac{\partial \ell}{\partial \hat{y}_i}$, where $1 \leq i \leq 3$? Your answer should be in terms of y_i and \hat{y}_i . Recall that we define the loss $\ell(\hat{y}, y)$ as follows (*note: log is a natural log*):

$$\ell(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_{i=1}^3 y_i \log(\hat{y}_i) \quad (7)$$

$\partial \ell / \partial \hat{y}_i$

- (b) (1 point) Suppose we have $\hat{y} = [0.2, 0.7, 0.1]^T$ and $y = [0, 1, 0]^T$. Compute the value of $\frac{\partial \ell}{\partial b_2}$. Use your answer from part (a), as well as the fact that the derivative of the softmax function with respect to b_k is as follows:

$$\frac{\partial \hat{y}_l}{\partial b_k} = \hat{y}_l (\mathbb{I}[k = l] - \hat{y}_k) \quad (8)$$

where $\mathbb{I}[k = l]$ is an indicator function such that if $k = l$ then it returns value 1 and 0 otherwise. Show your work below.

HINT: Recall that $\frac{\partial \ell}{\partial b_k} = \sum_l \frac{\partial \ell}{\partial \hat{y}_l} \frac{\partial \hat{y}_l}{\partial b_k}$.

$\partial \ell / \partial b_2$

- (c) (3 points) Now, write the derivative $\frac{\partial \ell}{\partial b_k}$ in a smart way such that you do not need the indicator function in Equation 8. Write your solutions in terms of \hat{y}_k, y_k . Show your work below.

$\partial \ell / \partial b_k$

- (d) (2 points) What is the derivative $\frac{\partial \ell}{\partial \beta}$? Your answer should be in terms of $\frac{\partial \ell}{\partial \mathbf{b}}$ and \mathbf{z} .

You should first consider a single entry in this matrix: $\frac{\partial \ell}{\partial \beta_{kj}}$.

$\partial \ell / \partial \beta$

- (e) (1 point) Explain in one short sentence why we use the matrix β^* (the matrix β without the first column of bias values) instead of β when calculating the derivative matrix $\frac{\partial \ell}{\partial \alpha}$?

Answer

- (f) (1 point) What is the derivative $\frac{\partial \ell}{\partial \mathbf{z}}$ (**not including the bias term**)? Your answer should be in terms of $\frac{\partial \ell}{\partial \mathbf{b}}$ and β^* .

$\partial \ell / \partial \mathbf{z}$

(g) (1 point) What is the derivative $\frac{\partial \ell}{\partial a_j}$ in terms of $\frac{\partial \ell}{\partial z_j}$ and z_j ?

$\partial \ell / \partial a_j$

(h) (1 point) What is the matrix $\frac{\partial \ell}{\partial \alpha}$? Your answer should be in terms of $\frac{\partial \ell}{\partial \mathbf{a}}$ and $\mathbf{x}^{(1)}$.

$\partial \ell / \partial \alpha$

3. Now you will put these equations to use in an example with numerical values. **You should use the answers you get here to debug your code.**

You are given a training example $\mathbf{x}^{(1)} = [1, 0, 1, 1, 0, 1]^T$ with **label class 2**, so $\mathbf{y}^{(1)} = [0, 1, 0]^T$. We initialize the network weights as:

$$\boldsymbol{\alpha}^* = \begin{bmatrix} 1 & 2 & 0 & -1 & 3 & 2 \\ 2 & 3 & 1 & 0 & 1 & 1 \\ 1 & 3 & 1 & 2 & -1 & 2 \\ 0 & 1 & 2 & 0 & 0 & 3 \end{bmatrix}$$

$$\boldsymbol{\beta}^* = \begin{bmatrix} 1 & 2 & 0 & 1 \\ 1 & -1 & 3 & 2 \\ 3 & 0 & -1 & 1 \end{bmatrix}$$

We want to also consider the bias term and the weights on the bias terms ($\alpha_{j,b}$ and $\beta_{j,b}$). Let's say they are all initialized to 1. To account for this we can add a column of 1's to the beginning of our initial weight matrices.

$$\boldsymbol{\alpha} = \begin{bmatrix} 1 & 1 & 2 & 0 & -1 & 3 & 2 \\ 1 & 2 & 3 & 1 & 0 & 1 & 1 \\ 1 & 1 & 3 & 1 & 2 & -1 & 2 \\ 1 & 0 & 1 & 2 & 0 & 0 & 3 \end{bmatrix}$$

$$\boldsymbol{\beta} = \begin{bmatrix} 1 & 1 & 2 & 0 & 1 \\ 1 & 1 & -1 & 3 & 2 \\ 1 & 3 & 0 & -1 & 1 \end{bmatrix}$$

And we can set our first value of our input vectors to always be 1 ($x_0^{(i)} = 1$), so our input becomes:

$$\mathbf{x}^{(1)} = [1, 1, 0, 1, 1, 0, 1]^T$$

Using the initial weights, run the feed forward of the network over this example (rounding to 4 decimal places during the calculation) and then answer the following questions.

(a) (1 point) What is a_1 ?

a_1	Work

(b) (1 point) What is z_1 ?

z_1	Work

(c) (1 point) What is b_2 ? We have computed $z_2 = 0.9933$, $z_3 = 0.9991$, $z_4 = 0.9975$ for you.

b_2	Work

(d) (1 point) What is \hat{y}_2 ? We have computed $b_1 = 4.9367, b_3 = 3.8562$ for you.

\hat{y}_2	Work

(e) (1 point) Which class would we predict on this example? Your answer should just be an integer $\in \{1, 2, 3\}$.

Class	Work

(f) (1 point) What is the total loss on this example?

Loss	Work

4. Now use the results of the previous question to run backpropagation over the network and update the weights with **Stochastic Gradient Descent**. Use learning rate $\eta = 1$.

Do your backpropagation calculations rounding to 4 decimal places then answer the following questions:

- (a) (1 point) What is the value of $\frac{\partial \ell}{\partial \beta_{1,0}}$?

$\frac{\partial \ell}{\partial \beta_{1,0}}$	Work
<input type="text"/>	<div></div>

- (b) (1 point) What is the updated value of the weight $\beta_{1,0}$?

$\beta_{1,0}$	Work
<input type="text"/>	<div></div>

- (c) (1 point) What is the value of $\frac{\partial \ell}{\partial \alpha_{3,2}}$?

$\frac{\partial \ell}{\partial \alpha_{3,2}}$	Work
<input type="text"/>	<div></div>

(d) (1 point) What is the updated value of the weight $\alpha_{3,2}$?

$\alpha_{3,2}$	Work

(e) (2 points) What is the updated weight of the input layer bias term applied to z_2 (i.e. $\alpha_{2,0}$)?

$\alpha_{2,0}$	Work

3 Empirical Questions (18 points)

The following questions should be completed after you work through the programming portion of this assignment. **For any plotting questions, you must use curves/line graphs, title your graphs, label your axes, and provide units (if applicable) and a legend in order to receive full credit.**

For these questions, **use the small dataset**. Use the following values for the hyperparameters unless otherwise specified:

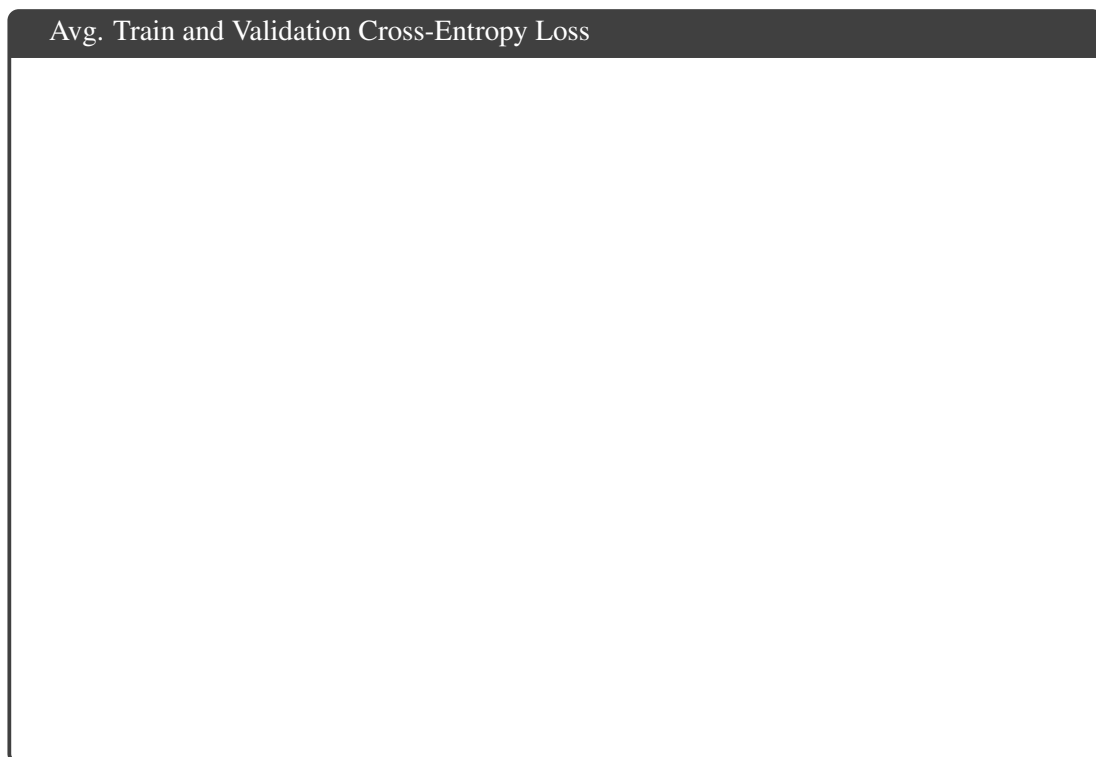
Parameter	Value
Number of Hidden Units	50
Weight Initialization	RANDOM
Learning Rate	0.01

Please submit computer-generated plots for all parts. To get full credit, your plots must be line graphs with labels for both axes with the value plotted on that axis and legends labeling every line.

1. Hidden Units


- (a) (2 points) Train a single hidden layer neural network using the hyperparameters mentioned in the table above, except for the number of hidden units which should vary among 5, 20, 50, 100, and 200. Run the optimization for 100 epochs each time.

Plot the average training cross-entropy (sum of the cross-entropy terms over the training dataset divided by the total number of training examples) of the final epoch on the y-axis vs number of hidden units on the x-axis. In the **same figure**, plot the average validation cross-entropy. The x-axis should be the number of hidden units, the y-axis should be average cross-entropy loss, and there should be one curve for validation loss and one curve for train loss.




- (b) (2 points) Examine and comment on the the plots of training and validation cross-entropy. What problem arises with too few hidden units, and why does it happen?

Answer



- (c) (2 points) In the handout folder, we provide `metrics_sgd_small.txt`, a text file with the validation cross-entropy loss values for SGD performed without AdaGrad using 100 epochs, 50 hidden units, random init, and 0.01 learning rate. In the **same figure**, plot them against your validation results for SGD **with** AdaGrad using the same set of parameters and the small dataset. The x-axis should be epoch number, the y-axis should be average cross-entropy loss, and there should be one curve for SGD without AdaGrad validation loss and one curve for SGD with AdaGrad validation loss.

Avg. Validation Cross-Entropy Loss of SGD with and without AdaGrad



- (d) (2 points) Examine and compare the two results. What do you observe?

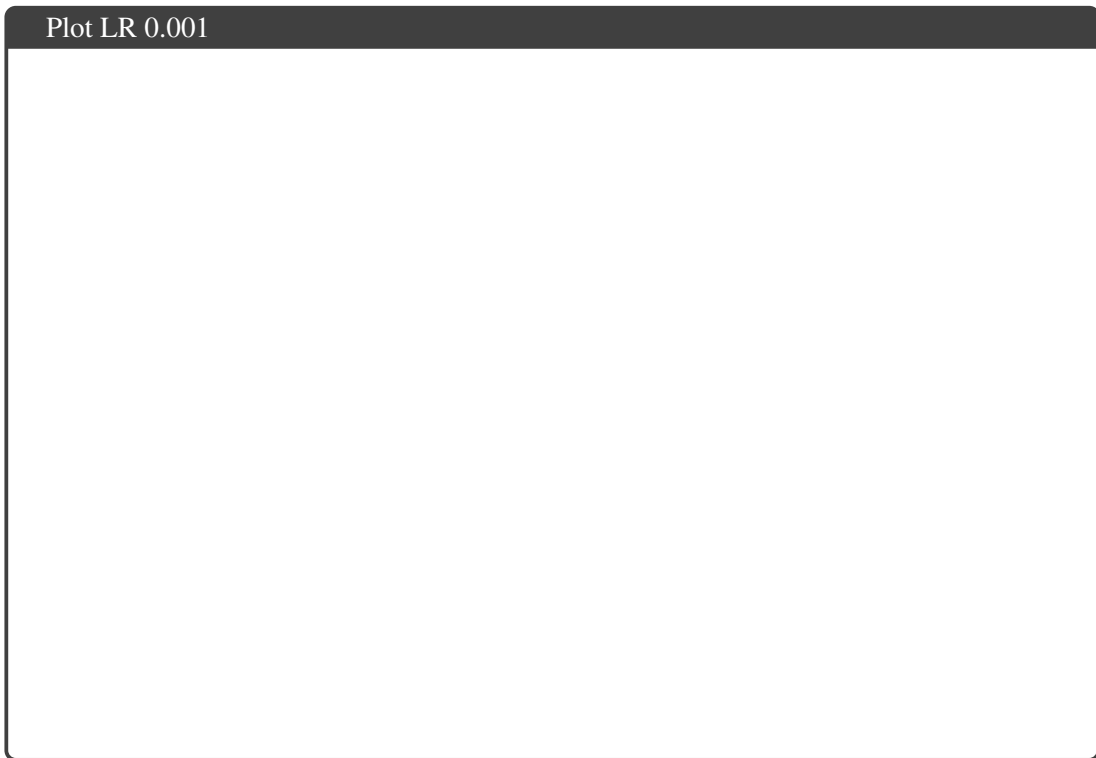
Answer

2. Learning Rate

- (a) (6 points) Train a single hidden layer neural network using the hyperparameters mentioned in the table above, except for the learning rate which should vary among 0.1, 0.01, and 0.001. Run the optimization for 100 epochs each time.

Plot the average training cross-entropy on the y-axis vs the number of epochs on the x-axis for the mentioned learning rates. In the **same figure**, plot the average validation cross-entropy loss. Make a separate figure for each learning rate. The x-axis should be epoch number, the y-axis should be average cross-entropy loss, and there should be one curve for training loss and one curve for validation loss.

Plot LR 0.1



- (b) (2 points) Examine and comment on the plots of training and validation cross-entropy. Are there any learning rates for which convergence is not achieved? Are there any learning rates that exhibit other problems? If so, describe these issues and list the learning rates that cause them.

Answer

3. Weight Initialization

- (a) (2 points) For this exercise, you can work on any data set. Initialize α and β to zero and print them out after the first few updates. For example, you may use the following command to begin:

```
$ python neuralnet.py small_train.csv small_validation.csv \  
small_train_out.txt small_validation_out.txt \  
small_metrics_out.txt 1 4 2 0.1
```

Compare the values across rows and columns in α and β . Describe the observed behavior and how this may affect model capacity and convergence.

Answer

4 Collaboration Questions

After you have completed all other components of this assignment, report your answers to these questions regarding the collaboration policy. Details of the policy can be found [here](#).

1. Did you receive any help whatsoever from anyone in solving this assignment? If so, include full details.
2. Did you give any help whatsoever to anyone in solving this assignment? If so, include full details.
3. Did you find or come across code that implements any part of this assignment? If so, include full details.

Your Answer

Programming (94 points)

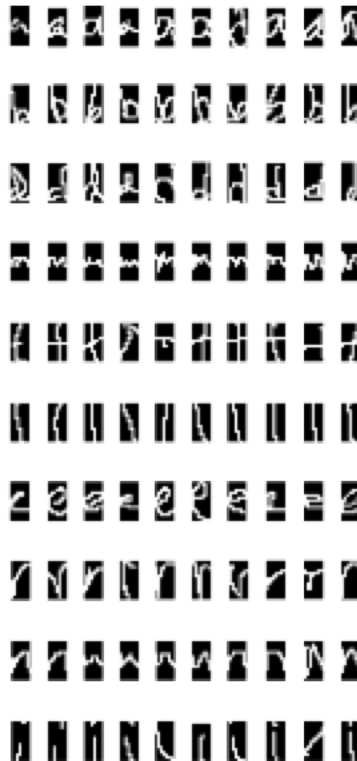


Figure 2: 10 random images of each of the 10 letters in the OCR dataset.

5 The Task

Your goal in this assignment is to implement a neural network to classify images using a single hidden layer neural network. In addition, you will implement Adagrad, a variant of stochastic gradient descent.

6 The Datasets

Datasets We will be using a subset of an Optical Character Recognition (OCR) dataset. This data includes images of all 26 handwritten letters; our subset will include only the letters “a,” “e,” “g,” “i,” “l,” “n,” “o,” “r,” “t,” and “u.” The handout contains a small dataset with 60 samples *per class* (50 for training and 10 for validation). We will also evaluate your code on a medium dataset with 600 samples per class (500 for training and 100 for validation). Figure 2 shows a random sample of 10 images of few letters from the dataset.

File Format Each dataset (small, medium, and large) consists of two csv files—train and validation. Each row contains 129 columns separated by commas. The first column contains the label and columns 2 to 129 represent the pixel values of a 16×8 image in a row major format. Label 0 corresponds to “a,” 1 to “e,” 2 to “g,” 3 to “i,” 4 to “l,” 5 to “n,” 6 to “o,” 7 to “r,” 8 to “t,” and 9 to “u.”

Because the original images are black-and-white (not grayscale), the pixel values are either 0 or 1. However, you should write your code to accept arbitrary pixel values in the range $[0, 1]$. The images in Figure 2 were produced by converting these pixel values into .png files for visualization. Observe that no feature

engineering has been done here; instead the neural network you build will *learn* features appropriate for the task of character recognition.

7 Model Definition

In this assignment, you will implement a single-hidden-layer neural network with a sigmoid activation function for the hidden layer, and a softmax on the output layer. Let the input vectors \mathbf{x} be of length M , and the hidden layer \mathbf{z} consist of D hidden units. In addition, let the output layer $\hat{\mathbf{y}}$ be a probability distribution over K classes. That is, each element \hat{y}_k of the output vector represents the probability of \mathbf{x} belonging to the class k .

We can compactly express this model by assuming that $x_0 = 1$ is a bias feature on the input and that $z_0 = 1$ is also fixed. In this way, we have two parameter matrices $\boldsymbol{\alpha} \in \mathbb{R}^{D \times (M+1)}$ and $\boldsymbol{\beta} \in \mathbb{R}^{K \times (D+1)}$. The extra 0th column of each matrix (i.e. $\boldsymbol{\alpha}_{:,0}$ and $\boldsymbol{\beta}_{:,0}$) hold the bias parameters.

$$\begin{aligned} a_j &= \sum_{m=0}^M \alpha_{j,m} x_m \\ z_j &= \frac{1}{1 + \exp(-a_j)} \\ b_k &= \sum_{j=0}^D \beta_{k,j} z_j \\ \hat{y}_k &= \frac{\exp(b_k)}{\sum_{l=1}^K \exp(b_l)} \end{aligned}$$

The objective function we will use for training the neural network is the average cross entropy over the training dataset $\mathcal{D} = \{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}$:

$$J(\boldsymbol{\alpha}, \boldsymbol{\beta}) = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_k^{(i)} \log(\hat{y}_k^{(i)}) \quad (9)$$

In Equation 9, J is a function of the model parameters $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ because $\hat{y}_k^{(i)}$ is implicitly a function of $\mathbf{x}^{(i)}$, $\boldsymbol{\alpha}$, and $\boldsymbol{\beta}$ since it is the output of the neural network applied to $\mathbf{x}^{(i)}$. $\hat{y}_k^{(i)}$ and $y_k^{(i)}$ are the k th components of $\hat{\mathbf{y}}^{(i)}$ and $\mathbf{y}^{(i)}$ respectively.

To train, you should optimize this objective function using stochastic gradient descent (SGD), where the gradient of the parameters for each training example is computed via backpropagation. You should shuffle the training points when performing SGD using the provided `shuffle` function, passing in the epoch number as a random seed. Note that SGD has a slight impact on the objective function, where we are “summing” over the current point, i :

$$J_{SGD}(\boldsymbol{\alpha}, \boldsymbol{\beta}) = -\sum_{k=1}^K y_k^{(i)} \log(\hat{y}_k^{(i)}) \quad (10)$$

Lastly, let's take a look at the Adagrad update that you will be performing. We use Adagrad (aka the Adaptive Gradient method) over standard gradient descent for two main reasons:

1. It eliminates the need to manually tune or update the learning rate. In Adagrad, the learning rate reduces as the number of iterations (below, the value of t) increases. This helps us converge closer to the optimum value.
2. It assigns a different learning rate for every parameter θ^i . As a result, updates are larger for parameters corresponding to infrequent features and smaller for more frequent ones. This makes Adagrad much more suitable to tasks where data may be sparse- such as natural language processing and image recognition.

For each parameter θ_t^i at round t , you will first compute an intermediate value s_t^i , and then use this to compute the updated θ_{t+1}^i . s_t^i will contain the element-wise sum of all the element-wise squared gradients. Therefore, s_t should have the same shape as $\frac{\partial J(\theta_t)}{\partial \theta_t^i}$. s_t should be initialized once, before the first epoch, to a zero vector. The update equations for s and θ are below.

$$s_{t+1}^i = s_t^i + \frac{\partial J(\theta_t)}{\partial \theta_t^i} \odot \frac{\partial J(\theta_t)}{\partial \theta_t^i}. \quad (11)$$

where \odot is the element-wise product operation. Then, we use s_t to scale the gradient for the update:

$$\theta_{t+1}^i = \theta_t^i - \frac{\eta}{\sqrt{s_{t+1}^i + \epsilon}} \odot \frac{\partial J(\theta_t)}{\partial \theta_t^i}. \quad (12)$$

Here, η is the learning rate, and $\epsilon = 1e-5$.

7.1 Initialization

In order to use a deep network, we must first initialize the weights and biases in the network. This is typically done with a random initialization, or initializing the weights from some other training procedure. For this assignment, we will be using two possible initialization:

RANDOM The weights are initialized randomly from a uniform distribution from -0.1 to 0.1.

The bias parameters are initialized to zero.

ZERO All weights are initialized to 0.

You must support both of these initialization schemes.

8 Implementation

Write a program `neuralnet.py` that implements an optical character recognizer using a one hidden layer neural network with sigmoid activations. Your program should learn the parameters of the model on the training data, report the cross-entropy at the end of each epoch on both train and validation data, and at the end of training write out its predictions and error rates on both datasets.

Your implementation must satisfy the following requirements:

- Use a **sigmoid** activation function on the hidden layer and **softmax** on the output layer to ensure it forms a proper probability distribution.

- Number of **hidden units** for the hidden layer should be determined by a command line flag.
- Support two different **initialization strategies**, as described in Section 7.1, selecting between them via a command line flag.
- Use stochastic gradient descent (SGD) to optimize the parameters for one hidden layer neural network. The number of **epochs** will be specified as a command line flag.
- Set the **learning rate** via a command line flag.
- Perform stochastic gradient descent updates on the training data on the data shuffled with the provided function. For each epoch, you must reshuffle the **original** file data, not the data from the previous epoch.
- In case there is a tie in the output layer \hat{y} , predict the smallest index to be the label.
- You may assume that the input data will always have the same output label space (i.e. $\{0, 1, \dots, 9\}$). Other than this, do not hard-code any aspect of the datasets into your code. We will autograde your programs on multiple data sets that include different examples.
- Do *not* use any machine learning libraries. You may use NumPy.

Implementing a neural network can be tricky: the parameters are not just a simple vector, but a collection of many parameters; computational efficiency of the model itself becomes essential; the initialization strategy dramatically impacts overall learning quality; other aspects which we will *not* change (e.g. activation function, optimization method) also have a large effect. These *tips* should help you along the way:

- Try to “vectorize” your code as much as possible—this is particularly important for Python. For example, in Python, you want to avoid for-loops and instead rely on `numpy` calls to perform operations such as matrix multiplication, transpose, subtraction, etc. over an entire `numpy` array at once. Why? Because these operations are actually implemented in fast C code, which won’t get bogged down the way a high-level scripting language like Python will.

8.1 Command Line Arguments

The autograder runs and evaluates the output from the files generated, using the following command:

```
$ python3 neuralnet.py [args...]
```

Where above `[args...]` is a placeholder for nine command-line arguments: `<train_input>` `<validation_input>` `<train_out>` `<validation_out>` `<metrics_out>` `<num_epoch>` `<hidden_units>` `<init_flag>` `<learning_rate>`. These arguments are described in detail below:

1. `<train_input>`: path to the training input `.csv` file (see Section 6)
2. `<validation_input>`: path to the validation input `.csv` file (see Section 6)
3. `<train_out>`: path to output `.txt` file to which the prediction on the *training* data should be written (see Section 8.2)
4. `<validation_out>`: path to output `.txt` file to which the prediction on the *validation* data should be written (see Section 8.2)
5. `<metrics_out>`: path of the output `.txt` file to which metrics such as train and validation error should be written (see Section 8.4)

6. `<num_epoch>`: integer specifying the number of times backpropagation loops through all of the training data (e.g., if `<num_epoch>` equals 5, then each training example will be used in backpropagation 5 times).
7. `<hidden_units>`: positive integer specifying the number of hidden units.
8. `<init_flag>`: integer taking value 1 or 2 that specifies whether to use RANDOM or ZERO initialization (see Section 7.1 and Section 7)—that is, if `init_flag==1` initialize your weights randomly from a uniform distribution over the range $[-0.1, 0.1]$ (i.e. RANDOM), if `init_flag==2` initialize all weights to zero (i.e. ZERO). For both settings, **always initialize bias terms to zero**.
9. `<learning_rate>`: float value specifying the base learning rate for SGD with Adagrad.
10. `<--debug>`: (optional argument) set the logging level, set to DEBUG to show logging

As an example, if you implemented your program in Python, the following command line would run your program with 4 hidden units on the small data provided in the handout for 2 epochs using zero initialization and a learning rate of 0.1.

```
python neuralnet.py small_train.csv small_validation.csv \
small_train_out.txt small_validation_out.txt \
small_metrics_out.txt 2 4 2 0.1
```

8.2 Output: Labels Files

Your program should write two output `.txt` files containing the predictions of your model on training data (`<train_out>`) and validation data (`<validation_out>`). Each should contain the predicted labels for each example printed on a new line. Use `\n` to create a new line.

Your labels should exactly match those of a reference implementation – this will be checked by the autograder by running your program and evaluating your output file against the reference solution.

Note: You should output your predicted labels using the same *integer* identifiers as the original training data. You should also insert an empty line (again using `'\n'`) at the end of each sequence (as is done in the input data files).

8.3 Debug Output: Logging

Note that we use the debug logging level in the starter code. If we use a higher logging level, we will log things with the default logging configuration, causing potential slowdowns when executing on an autograder.

Note also that we log NumPy matrices on separate lines from strings describing them. If we do not do this (e.g., if we call `str` on them and add them to the strings), the arrays will be turned into strings even when our logging is set to ignore debug, causing potential massive slowdowns.

8.4 Output Metrics

Generate a file where you report the following metrics:

cross entropy After each epoch, report mean cross entropy on the training data `crossentropy(train)` and validation data `crossentropy(validation)` (See Equation 9). These two cross-entropy values should be reported at the end of each epoch and prefixed by the epoch number. For example, after the second pass through the training examples, these should be prefixed by `epoch=2`. The total number of train losses you print out should equal `num_epoch`—likewise for the total number of validation losses.

error After the final epoch (i.e. when training has completed fully), report the final training error `error(train)` and validation error `error(validation)`.

A sample output is given below. It contains the train and validation losses for the first 2 epochs and the final error rate when using the command given above.

```
epoch=1 crossentropy(train): 1.9946950280285547
epoch=1 crossentropy(validation): 2.010686378337308
epoch=2 crossentropy(train): 1.912184059993547
epoch=2 crossentropy(validation): 1.944326942790059
error(train): 0.782
error(validation): 0.83
```

Take care that your output has the exact same format as shown above. There is an equal sign = between the word `epoch` and the epoch number, but no spaces. There should be a single space after the epoch number (e.g. a space after `epoch=1`), and a single space after the colon preceding the metric value (e.g. a space after `epoch=1 likelihood(train):`). Each line should be terminated by a Unix line ending `\n`.

8.5 Tiny Data Set

To help you with this assignment, we have also included a tiny data set, `tiny_train.csv` and `tiny_validation.csv`, and a reference output file `tiny_output.txt` for you to use. The tiny dataset is in a format similar to the other datasets, but it only contains two samples with five features. The reference file contains outputs from each layer of one correctly implemented neural network, for both forward and back-propagation steps. We advise you to use this set to help you debug in case your implementation doesn't produce the same results as in the written part.

For your reference, `tiny_output.txt` is generated from the following command line specifications:

```
python neuralnet.py tiny_train.csv tiny_validation.csv \
tiny_train_out.txt tiny_validation_out.txt \
tiny_metrics_out.txt 1 4 2 0.1
```

The specific output file names are not important, but be sure to keep the other arguments exactly as they are shown above.

9 Gradescope Submission

You should submit your `neuralnet.py` to Gradescope. Please do not use any other file name for your implementation. This will cause problems for the autograder to correctly detect and run your code.

10 Pseudocode

Since the network structure we will use in this homework is fairly simple, we define our neural network as a single module. Note that in most deep learning libraries, there is a module corresponding to each layer type (e.g., Linear, Sigmoid, Softmax) to allow for more flexibility.

NN Module

- 1: **procedure** FORWARD(x)
- 2: Forward pass
- 3: **procedure** BACKWARD(x, y, \hat{y})
- 4: Backward pass
- 5: **procedure** TRAIN(x, y)

- 6: Train the model with SGD
- 7: **procedure** TEST(\mathbf{x}, \mathbf{y})
- 8: Test the model and return the loss

The NN module defines a forward function $\mathbf{y} = \text{FORWARD}(\mathbf{x})$ and a backward function $\mathbf{g}_\alpha, \mathbf{g}_\beta = \text{BACKWARD}(\mathbf{x}, \mathbf{y}, \hat{\mathbf{y}})$ method. You'll want to pay close attention to the dimensions that you pass into and return from your modules.

10.1 Forward and Backward Methods

After implementing the helper functions (SIGMOID, SOFTMAX, CROSSENTROPY), we can define the methods NNFORWARD and NNBACKWARD as follows.

Algorithm 1 Forward Method

- 1: **procedure** NNFORWARD(Training example (\mathbf{x}, \mathbf{y}))
 - 2: $\mathbf{a} = \text{LINEAR}(\mathbf{x}, \boldsymbol{\alpha})$
 - 3: $\mathbf{z} = \text{SIGMOID}(\mathbf{a})$
 - 4: $\mathbf{b} = \text{LINEAR}(\mathbf{z}, \boldsymbol{\beta})$
 - 5: $\hat{\mathbf{y}} = \text{SOFTMAX}(\mathbf{b})$
 - 6: **return** $\hat{\mathbf{y}}$
-

Algorithm 2 Backward Method

- 1: **procedure** NNBACKWARD(Training example (\mathbf{x}, \mathbf{y}) , Predicted probabilities $\hat{\mathbf{y}}$)
 - 2: Place intermediate quantities $\mathbf{x}, \mathbf{a}, \mathbf{z}, \mathbf{b}, \hat{\mathbf{y}}$ in scope ▷ Hint: make use of class attributes
 - 3: $\mathbf{g}_\mathbf{b} = \text{D_CROSSENTROPY}(\mathbf{y}, \hat{\mathbf{y}})$
 - 4: $\mathbf{g}_\beta, \mathbf{g}_\mathbf{z} = \text{D_LINEAR}(\mathbf{z}, \boldsymbol{\beta}, \mathbf{g}_\mathbf{b})$
 - 5: $\mathbf{g}_\mathbf{a} = \text{D_SIGMOID}(\mathbf{a}, \mathbf{z})\mathbf{g}_\mathbf{z}$
 - 6: $\mathbf{g}_\alpha, \mathbf{g}_\mathbf{x} = \text{D_LINEAR}(\mathbf{x}, \boldsymbol{\alpha}, \mathbf{g}_\mathbf{a})$ ▷ We discard $\mathbf{g}_\mathbf{x}$
 - 7: **return** parameter gradients $\mathbf{g}_\alpha, \mathbf{g}_\beta$
-

10.2 Training Method

Consider the neural network described in Section 7 applied to the i th training example (\mathbf{x}, \mathbf{y}) where \mathbf{y} is a one-hot encoding of the true label. Our neural network outputs $\hat{\mathbf{y}} = h_{\alpha, \beta}(\mathbf{x})$, where $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ are the parameters of the first and second layers respectively and $h_{\alpha, \beta}$ is a one-hidden layer neural network with a sigmoid activation and softmax output. The loss function is negative cross-entropy $J = \ell(\hat{\mathbf{y}}, \mathbf{y}) = -\mathbf{y}^T \log(\hat{\mathbf{y}})$. $J = J_{\mathbf{x}, \mathbf{y}}(\boldsymbol{\alpha}, \boldsymbol{\beta})$ is actually a function of our training example (\mathbf{x}, \mathbf{y}) , and our model parameters $\boldsymbol{\alpha}, \boldsymbol{\beta}$ though we write just J for brevity.

In order to train our neural network, we are going to apply stochastic gradient descent. Because we want the behavior of your program to be deterministic for testing on Gradescope, we make a few simplifications: (1) you should *not* shuffle your data and (2) you will use a fixed learning rate. In the real world, you would *not* make these simplifications.

SGD proceeds as follows, where E is the number of epochs and η is the learning rate.

Algorithm 3 Training with Stochastic Gradient Descent (SGD)

```
1: procedure SGD(Training data  $\mathcal{D}_{train}$ , test data  $\mathcal{D}_t$ )
2:   Initialize parameters  $\alpha, \beta$  ▷ Use either RANDOM or ZERO from Section 7.1
3:   for  $e \in \{1, 2, \dots, E\}$  do ▷ For each epoch
4:      $\mathcal{D} = \text{SHUFFLE}(\mathcal{D}_{train}, e)$ 
5:     for  $(\mathbf{x}, \mathbf{y}) \in \mathcal{D}$  do ▷ For each training example (No shuffling)
6:       Compute neural network layers:
7:        $\mathbf{o} = \text{NNFORWARD}(\mathbf{x}, \mathbf{y}, \alpha, \beta)$ 
8:       Compute gradients via backprop:
9:       
$$\left. \begin{aligned} \mathbf{g}_\alpha &= \frac{\partial J}{\partial \alpha} \\ \mathbf{g}_\beta &= \frac{\partial J}{\partial \beta} \end{aligned} \right\} = \text{NNBACKWARD}(\mathbf{x}, \mathbf{y}, \hat{\mathbf{y}})$$

10:      Update parameters with Adagrad updates  $\mathbf{g}'_\alpha, \mathbf{g}'_\beta$ : ▷ Refer to Eq. 11
11:       $\alpha \leftarrow \alpha - \eta \mathbf{g}'_\alpha$ 
12:       $\beta \leftarrow \beta - \eta \mathbf{g}'_\beta$ 
13:      Evaluate training mean cross-entropy  $J_{\mathcal{D}}(\alpha, \beta)$ 
14:      Evaluate test mean cross-entropy  $J_{\mathcal{D}_t}(\alpha, \beta)$ 
15:   return parameters  $\alpha, \beta$ 
```

10.3 Testing Method

At test time, we output the most likely prediction for each example:

Algorithm 4 Prediction at Test Time

```
1: procedure PREDICT(Unlabeled train or test dataset  $\mathcal{D}'$ )
2:   for  $\mathbf{x} \in \mathcal{D}'$  do
3:     Compute neural network prediction  $\hat{\mathbf{y}} = h(\mathbf{x})$ 
4:     Predict the label with highest probability  $l = \text{argmax}_k \hat{y}_k$ 
```

It's also quite common to combine the Cross-Entropy and Softmax layers into one. The reason for this is the cancelation of numerous terms that result from the zeros in the cross-entropy backward calculation. (Said trick is *not* required to obtain a sufficiently fast implementation for Gradescope.)

Some additional tips: Make sure to read the autograder output carefully. The autograder for Gradescope prints out some additional information about the tests that it ran. For this programming assignment we've specially designed some buggy implementations that you might implement and will try our best to detect those and give you some more useful feedback in Gradescope's autograder. Make wise use of autograder's output for debugging your code.