

HOMEWORK 3

SVMS AND NEURAL NETWORKS

CMU 10-701: MACHINE LEARNING (FALL 2021)

piiazza.com/cmu/fall2021/10701/home

OUT: Wednesday, Oct 06, 2021

DUE: Wednesday, Oct 20, 2021, 11:59pm

START HERE: Instructions

- **Collaboration policy:** Collaboration on solving the homework is allowed, after you have thought about the problems on your own. It is also OK to get clarification (but not solutions) from books or online resources, again after you have thought about the problems on your own. There are two requirements: first, cite your collaborators fully and completely (e.g., “Jane explained to me what is asked in Question 2.1”). Second, write your solution *independently*: close the book and all of your notes, and send collaborators out of the room, so that the solution comes from you only. See the Academic Integrity Section in our course syllabus for more information: https://piiazza.com/class_profile/get_resource/ksetdrgdkob78/ksqc9bxxjt56ic
- **Late Submission Policy:** See the late submission policy here: https://piiazza.com/class_profile/get_resource/ksetdrgdkob78/ksqc9bxxjt56ic
- **Submitting your work:**
 - **Gradescope:** There will be two submission slots for this homework on Gradescope: Written and Programming.
For the written problems such as short answer, multiple choice, derivations, proofs, or plots, we will be using the written submission slot. Please use the provided template. The best way to format your homework is by using the Latex template released in the handout and writing your solutions in Latex. However 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. Each derivation/proof should be completed in the boxes provided below the question, **you should not move or change the sizes of these boxes** as Gradescope is expecting your solved homework PDF to match the template on Gradescope. If you find you need more space than the box provides you should consider cutting your solution down to its relevant parts, if you see no way to do this, please add an additional page at the end of the homework and guide us there with a ‘See page xx for the rest of the solution’.
You are also required to upload your code, which you wrote to solve the final question of this homework, to the Programming submission slot. Your code may be run by TAs so please make sure it is in a workable state.

Regrade requests can be made after the homework grades are released, however this gives the TA the opportunity to regrade your entire paper, meaning if additional mistakes are found then points will be deducted.

For multiple choice or select all that apply questions, shade in the box or circle in the template document corresponding to the correct answer(s) for each of the questions. For \LaTeX users, use \blacksquare and \bullet for shaded boxes and circles, and don't change anything else. If an answer box is included for showing work, **you must show your work!**

Problem 1: Miscellaneous [9 pts]

1. [1 pt] **Select all that apply:** In a soft-margin support vector machine, increasing the slack penalty term C causes
 - ☒ The training error to decrease
 - ☒ More robustness to outliers
 - ☒ The margin to shrink
 - ☐ The training error to increase
 - ☐ None of the above
2. [1 pt] **Select all that apply:** What is always true about the kernel matrix?
 - ☐ Symmetric
 - ☐ All positive entries
 - ☒ Invertible
 - ☒ Contains 0 as an eigenvalue
 - ☐ None of the above
3. [1 pt] **Select one:** The shortest distance of point \mathbf{z} to hyperplane $\mathbf{w}^T \mathbf{x} = 0$ is:
 - ☐ $\mathbf{w}^T \mathbf{z}$
 - ☐ $\frac{\mathbf{w}^T \mathbf{z}}{\|\mathbf{z}\|_2}$
 - ☒ $\frac{\mathbf{w}^T \mathbf{z}}{\|\mathbf{w}\|_2}$
 - ☐ $\|\mathbf{w}\|_2$
4. [1 pt] **True or False:** For any neural network, the training error must always decrease monotonically when using SGD, provided the step size is sufficiently small (assume the step size is constant).
 - ☒ True ☐ False
5. [1 pt] **True or False:** For any network, the validation set error must always decrease monotonically when using SGD and then later increase monotonically, provided the step size is sufficiently small.
 - ☒ True ☐ False

6. [1 pt] **Select all that apply:** Which of the following statements is/are true about a Convolutional Layer?

- ☒ The width and height of the layer must be equal.
- ☒ The number of parameters depends on the stride
- ☒ The number of bias terms is equal to the number of filters
- ☐ The number of parameters depends on the padding.
- ☐ None of the above

7. [1 pt] Consider the image X and filter F given below. Let X be convolved with F using no padding and a stride of 1 to produce an output Y . What is value of g in the output Y ?

$X =$

1	0	-2	3	4	1
2	9	5	6	0	-1
0	-3	1	3	4	4
6	5	2	0	6	8
-5	4	-3	1	3	-2
4	1	2	8	9	7

$F =$

-1	-1	-1
-1	8	-1
-1	-1	-1

$Y =$

a	b	c	d
e	f	g	h
i	j	k	l
m	n	o	p

0

8. [1 pt] Suppose you had an input feature map of size 6x12 and filter size 3x3, using no padding and a stride of 2, what would be the resulting output size?

2x5

9. [1 pt] **True or False:** A Neural Network with no activation function and any number of layers can always be condensed to a network with just one layer.

☒ True ☐ False

Problem 2: Soft Margin Hyperplanes [9 points]

The soft-margin primal SVM problem is:

$$\min \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \xi_i$$

subject to feasibility constraints that for all $i = 1, \dots, N$:

$$\begin{aligned} y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) &\geq 1 - \xi_i \\ \xi_i &\geq 0 \end{aligned}$$

Suppose instead the optimization objective was changed to $\frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \xi_i^p$ with $p > 1$ while the feasibility constraints are kept the same.

Hint: Before attempting this problem, we highly recommend first trying to derive the dual formulations for the hard-margin and soft-margin SVMs presented in class by working through the steps detailed below.

1. [1 pt] Incorporate the feasibility constraints into the objective function using the method of Lagrange multipliers by writing out the Lagrangian and the associated constraints.

$$\begin{aligned} \min \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \xi_i^p &\Rightarrow \begin{cases} \min \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \min C \sum_{i=1}^N \xi_i^p \end{cases} \left(\begin{array}{l} \text{with constraints} \\ y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \xi_i \\ \xi_i \geq 0 \end{array} \right) \\ \mathcal{L}(\mathbf{w}, b, \alpha, \varepsilon) &= \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \xi_i^p - \sum_{i=1}^N \alpha_i (y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1 + \xi_i) \\ &= \frac{\mathbf{w}^T \mathbf{w}}{2} + C \sum_{i=1}^N \xi_i^p - \sum_{i=1}^N \alpha_i (y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1 + \xi_i) \\ &\Rightarrow \text{objective function} \\ \min_{\mathbf{w}, b, \varepsilon} \max_{\alpha} &\frac{\mathbf{w}^T \mathbf{w}}{2} + C \sum_{i=1}^N \xi_i^p - \sum_{i=1}^N \alpha_i (y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1 + \xi_i) \end{aligned}$$

2. [3 pts] Next, take the partial derivative of the Lagrangian with respect to all of the optimization variables. Be sure to show all of your work.

$$\mathcal{L}(w, b, \alpha, \varepsilon) = \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \varepsilon_i^p - \sum_{i=1}^N \alpha_i (y^{(i)} (w^T x^{(i)} + b) - 1 + \varepsilon_i)$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial w} &= \frac{1}{2} 2w + 0 - \sum_{i=1}^N \alpha_i y^{(i)} x^{(i)} + 0 + 0 + 0 \\ &= w - \sum_{i=1}^N \alpha_i y^{(i)} x^{(i)} \end{aligned}$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial b} &= 0 + 0 - \sum_{i=1}^N 0 + \alpha_i y^{(i)} + 0 + 0 \\ &= - \sum_{i=1}^N \alpha_i y^{(i)} \end{aligned}$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \varepsilon_i} &= 0 + C p \varepsilon_i^{p-1} - \alpha_i (0 + 0 + 0 + 1) \\ &= C p \varepsilon_i^{p-1} - \alpha_i \end{aligned}$$

3. [3 pts] Finally, by setting the partial derivatives from the previous part equal to 0, derive the dual formulation of SVMs in this general case. Again, be sure to show all of your work.

$$\frac{\partial \mathcal{L}}{\partial w} = w - \sum_{i=1}^N \alpha_i y^{(i)} x^{(i)} = 0 \Rightarrow w = \sum_{i=1}^N \alpha_i y^{(i)} x^{(i)}$$

$$\frac{\partial \mathcal{L}}{\partial b} = -\sum_{i=1}^N \alpha_i y^{(i)} = 0 \Rightarrow \sum_{i=1}^N \alpha_i y^{(i)} = 0$$

$$\frac{\partial \mathcal{L}}{\partial \varepsilon_i} = CP \varepsilon_i^{p-1} - \alpha_i = 0 \Rightarrow \alpha_i = CP \varepsilon_i^{p-1}$$

$$\Rightarrow \varepsilon_i = \left(\frac{\alpha_i}{CP} \right)^{\frac{1}{p-1}}$$

Substituting w , ε_i to \mathcal{L} , the objective function

$$\Rightarrow \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \varepsilon_i^p - \sum_{i=1}^N \alpha_i (y^{(i)} (w^T x^{(i)} + b) - 1 + \varepsilon_i)$$

$$= \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \left(\frac{\alpha_i}{CP} \right)^{\frac{p}{p-1}} - \sum_{i=1}^N \alpha_i (y^{(i)} \left(\sum_{j=1}^N \alpha_j y^{(j)} x^{(j)} \right)^T x^{(i)} + b) - 1 + \varepsilon_i$$

$$= \frac{1}{2} \|w\|_2^2 + \sum_{i=1}^N \left[C \left(\frac{\alpha_i}{CP} \right)^{\frac{p}{p-1}} - CP \varepsilon_i^p \right] + \sum_{i=1}^N \alpha_i - \sum_{i=1}^N \alpha_i y^{(i)} b - \underbrace{\sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y^{(i)} y^{(j)} x^{(i)T} x^{(j)}}_{\|w\|_2^2}$$

$$= \sum_{i=1}^N \alpha_i + \sum_{i=1}^N \left[C \left(\frac{\alpha_i}{CP} \right)^{\frac{p}{p-1}} - CP \left(\frac{\alpha_i}{CP} \right)^{\frac{p}{p-1}} \right] - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y^{(i)} y^{(j)} x^{(i)T} x^{(j)}$$

$$\Rightarrow \max_{\alpha} \sum_{i=1}^N \alpha_i + \sum_{i=1}^N \left[C \left(\frac{\alpha_i}{CP} \right)^{\frac{p}{p-1}} (1-p) \right] - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y^{(i)} y^{(j)} x^{(i)T} x^{(j)}$$

$$(0 \leq \alpha_i \leq CP \varepsilon_i^{p-1})$$

4. [2 pts] How does this more general formulation ($p > 1$) compare to the standard setting ($p = 1$) discussed in lecture? Specifically, is the general formulation more or less complex in terms of the number of optimization variables and constraints? Justify your answer in a sentence or two.

Since the general formulation gives α_j a different constraint for its upper bound, it's more complex.

When $P = 1$, the constraint was $0 \leq \alpha_j \leq C$

When $P > 1$, the constraint becomes $0 \leq \alpha_j \leq CP\varepsilon_j$

Problem 3: SVM Decision Boundaries [6 pts]

Figure 1 shows SVM decision boundaries resulting from different combinations of kernels, slack penalties, and RBF bandwidths. Match each plot from Figure 1 with one of the SVM settings.

To obtain full credit, you do not need to provide any justification, but only provide the correct pairing.

RBF kernel: $K(\mathbf{x}, \mathbf{z}) = \exp(-\gamma\|\mathbf{x} - \mathbf{z}\|_2^2)$

1. Linear kernel, $C = 0.1$

e

2. Linear kernel, $C = 100$

a

3. RBF kernel, $C = 1$, $\gamma = 10$

f

4. RBF kernel, $C = 1$, $\gamma = 50$

c

5. RBF kernel, $C = 0.1$, $\gamma = 0.1$

b

6. RBF kernel, $C = 20$, $\gamma = 0.1$

d

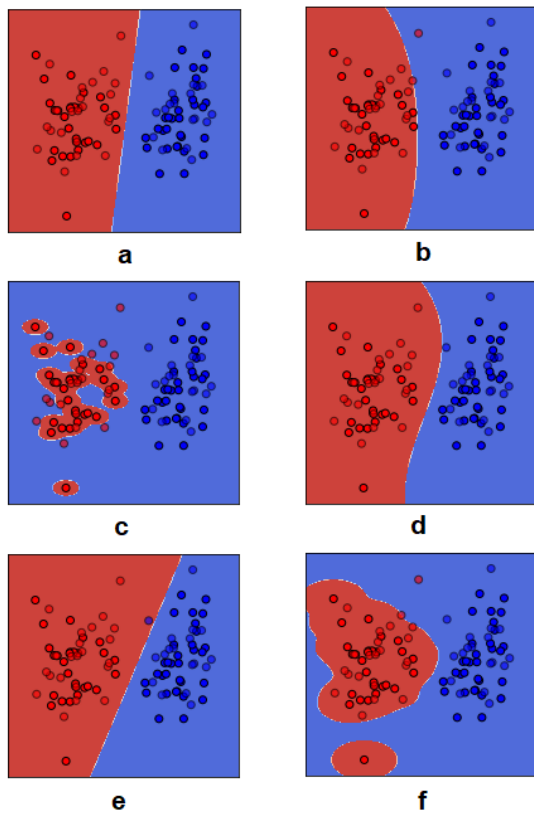


Figure 1: SVM decision boundaries learned with different parameter settings.

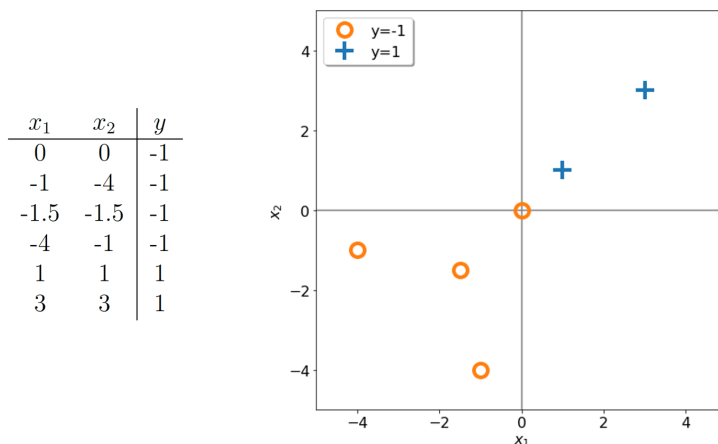
Problem 4: Visualizing SVM Kernels [10 pts]

We've seen how kernels can help us apply SVMs to produce highly nonlinear decision boundaries.

One way to way to comprehend the effect of feature mapping functions, $\phi(\mathbf{x})$, and kernels is to envision the data transformed to a higher dimensional space where a linear decision boundary can then be applied.

In this question, we'll explore a different way to look at the effect of feature mapping functions. Here, we'll visualize the function $z = f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$ with respect to the original input domain \mathbf{x} . Note the missing $\text{sign}(z)$ threshold function that would we normally use to predict $\hat{y} \in \{-1, 1\}$.

We'll use the following dataset with $\mathbf{x} \in \mathbb{R}^2$, so we can visualize $f(\mathbf{x})$ with a 3D surface plot:



For each of the following kernel functions, fit the above data using a hard-margin SVM solver and then plot surface $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$ over the domain $x_1 \times x_2 \in [-5, 5] \times [-5, 5]$. In addition, write the resulting number of support vectors for each class, $y = -1$ and $y = 1$.

You'll have to use the dual formulation of $f(x)$ to be able to evaluate this function with kernels. Recall that:

$$\mathbf{w}^T \phi(\mathbf{x}) = \sum_{i=1}^N \alpha_i y^{(i)} K(\mathbf{x}, \mathbf{x}^{(i)}) \quad (1)$$

$$b = y^{(k)} - \sum_{i=1}^N \alpha_i y^{(i)} K(\mathbf{x}^{(k)}, \mathbf{x}^{(i)}) \quad \text{for any } k \text{ where } 0 < \alpha_k < C \quad (2)$$

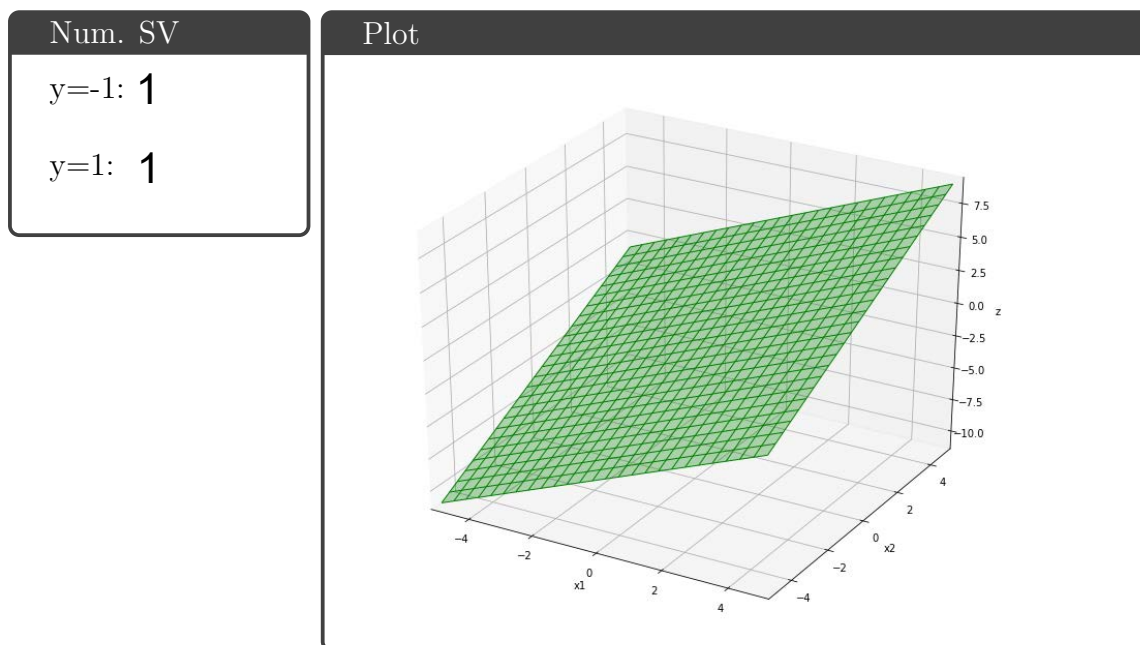
For your convenience, we have provided the following link to a Jupyter notebook with Python code to plot a surface for a function with a 2D input: [svm_surface.ipynb](#).

In this problem you MAY use an SVM solver, such as [sklearn.svm.SVC](#). Make sure to read the associated documentation to ensure you are using the solver properly. For example, you may need to set the soft-margin parameter C to infinity to use it as a hard-margin solver.

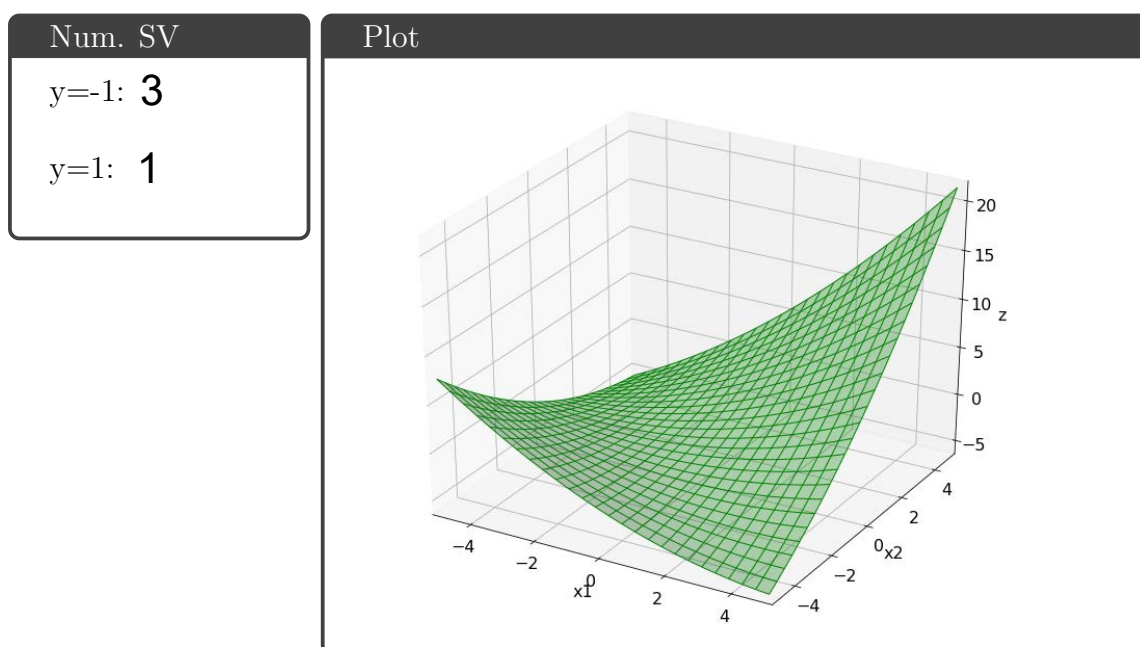
Let us know if you would like any help with Python, Jupyter notebooks, or Google Colab. You may use other programming languages as well to solve and plot. Your plot must be computer generated (definitely don't attempt to draw these by hand).

You don't need to submit your code for this problem.

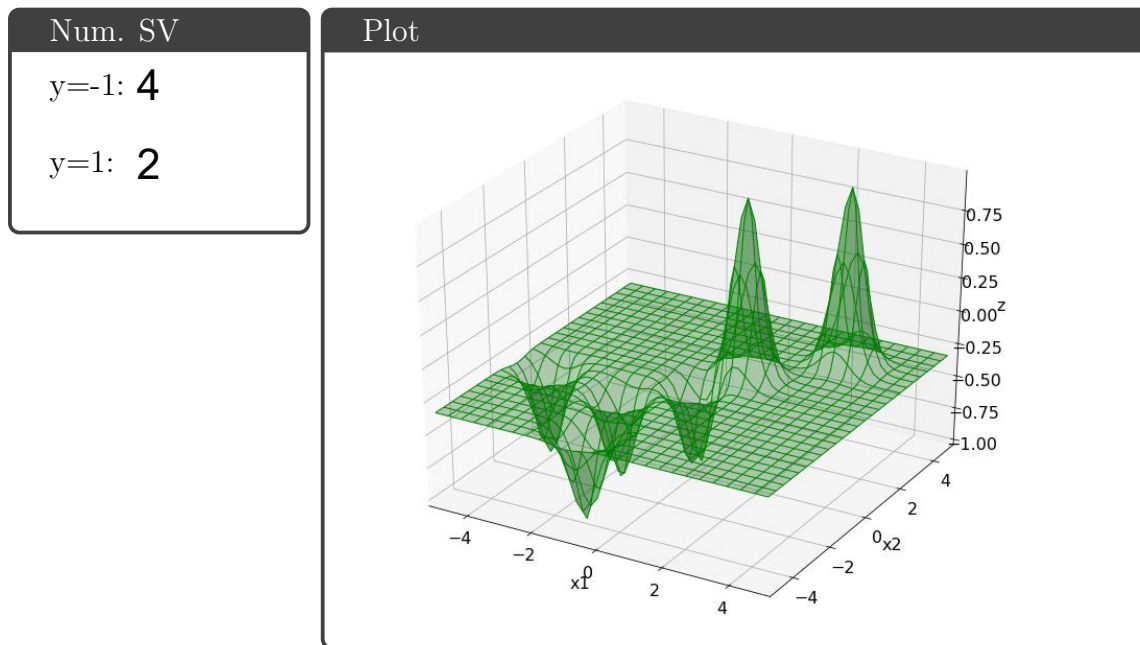
1. [2 pts] Linear kernel, $K(\mathbf{x}, \mathbf{z}) = \mathbf{x}^T \mathbf{z}$



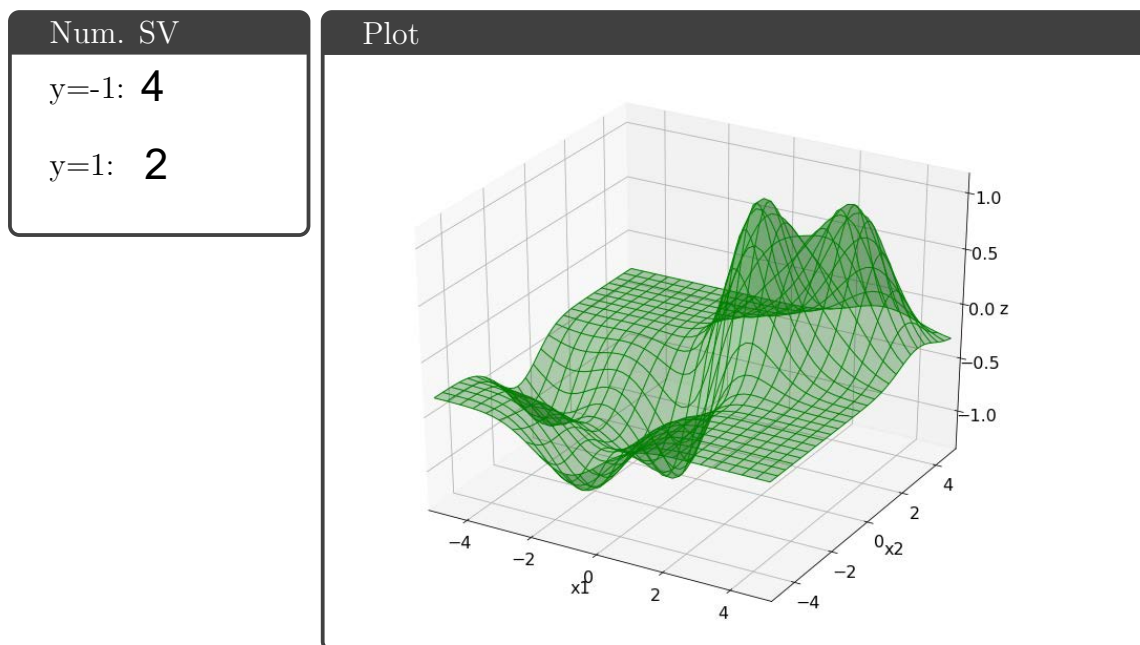
2. [2 pts] Quadratic kernel, $K(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^2$



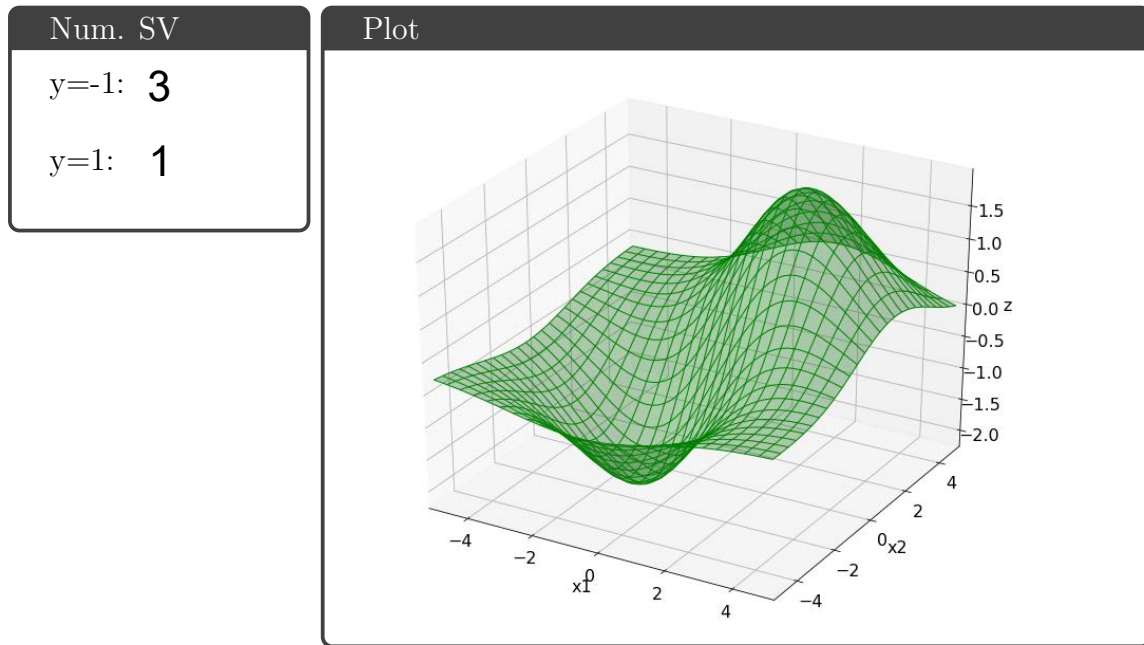
3. [2 pts] RBF kernel with $\gamma = 3$, $K(\mathbf{x}, \mathbf{z}) = \exp(-\gamma \|\mathbf{x} - \mathbf{z}\|_2^2)$



4. [2 pts] RBF kernel with $\gamma = 0.5$, $K(\mathbf{x}, \mathbf{z}) = \exp(-\gamma \|\mathbf{x} - \mathbf{z}\|_2^2)$



5. [2 pts] RBF kernel with $\gamma = 0.1$, $K(\mathbf{x}, \mathbf{z}) = \exp(-\gamma \|\mathbf{x} - \mathbf{z}\|_2^2)$



Problem 5: Neural Nets: Written Questions [30 pts]

Note: We strongly encourage you to do the written part of this homework before the programming, as it will help you gain familiarity with the calculations you will have to code up in the programming section. We suggest that for each of these problems, you write out the equation required to calculate each value in terms of the variables we created (a_j, z_j, b_k , etc.) before you calculate the numerical value.

Note: For all questions which require numerical answers, round up your final answers to four decimal places. For integers, you may drop trailing zeros.

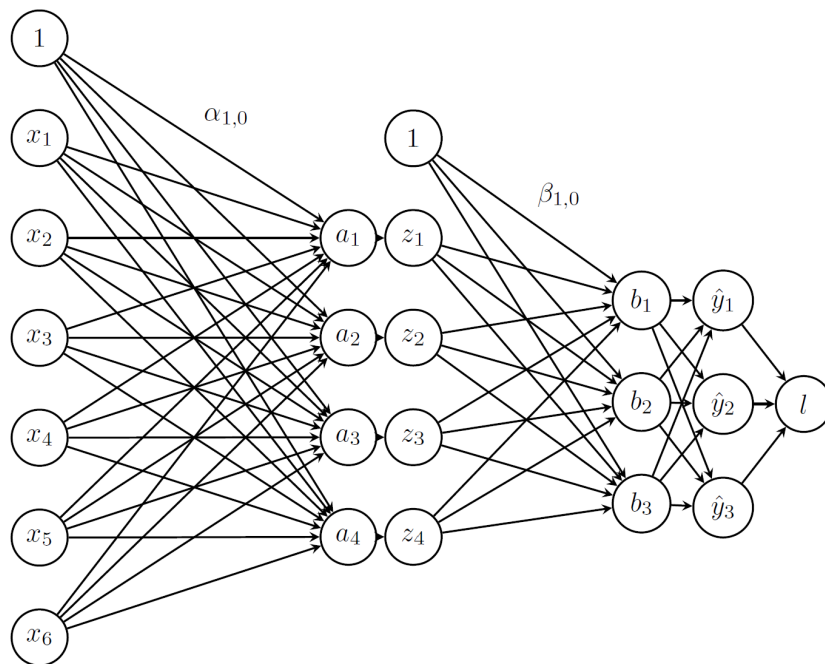


Figure 2: A One Hidden Layer Neural Network

Network Overview

Consider the neural network with one hidden layer shown in Figure 2. 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 $\hat{\mathbf{y}} = [\hat{y}_1, \hat{y}_2, \hat{y}_3]^T$ that sums to one over 3 classes. We also add a bias to the input, $x_0 = 1$ and the hidden layer $z_0 = 1$, both of which are fixed to 1.

We adopt the following notation:

1. Let α be the matrix of weights from the inputs to the hidden layer.
2. Let β be the matrix of weights from the hidden layer to the output layer.
3. Let $\alpha_{j,i}$ represent 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)

4. Let $\beta_{k,j}$ represent the weight going to the node y_k in the output layer from the node z_j in the hidden layer.
5. We will use a *sigmoid activation function* (σ) for the hidden layer and a *softmax* for the output layer.

Network Details

Equivalently, we define each of the following.

The input:

$$\mathbf{x} = [x_1, x_2, x_3, x_4, x_5, x_6]^T \quad (3)$$

Linear combination at first (hidden) layer:

$$a_j = \alpha_{j,0} + \sum_{i=1}^6 \alpha_{j,i} x_i, \quad j \in \{1, \dots, 4\} \quad (4)$$

Activation at first (hidden) layer:

$$z_j = \sigma(a_j) = \frac{1}{1 + \exp(-a_j)}, \quad j \in \{1, \dots, 4\} \quad (5)$$

Linear combination at second (output) layer:

$$b_k = \beta_{k,0} + \sum_{j=1}^4 \beta_{k,j} z_j, \quad k \in \{1, \dots, 3\} \quad (6)$$

Activation at second (output) layer:

$$\hat{y}_k = \frac{\exp(b_k)}{\sum_{l=1}^3 \exp(b_l)}, \quad k \in \{1, \dots, 3\} \quad (7)$$

Note that the linear combination equations can be written equivalently as the product of the weight matrix with the input vector. We can even fold in the bias term α_0 by thinking of $x_0 = 1$, and fold in β_0 by thinking of $z_0 = 1$.

Loss

We will use cross entropy loss, $\ell(\hat{\mathbf{y}}, \mathbf{y})$. If \mathbf{y} represents our target (true) 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 as (note that the log terms are in base e):

$$\ell(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_{k=1}^3 y_k \log(\hat{y}_k) \quad (8)$$

1. [6 pts] 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 avoid excessive loops and large run times.

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 a 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:

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

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

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

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

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

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

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

- (a) [1 pt] By examining the shapes of the initial weight matrices, how many neurons do we have in the first hidden layer of the neural network? (Not including the bias neuron)

4

- (b) [1 pt] How many output neurons will our neural network have?

3

- (c) [1 pt] What is the vector \mathbf{a} whose elements are made up of the entries a_j in equation (4). Write your answer in terms of $\boldsymbol{\alpha}$ and $\mathbf{x}^{(1)}$.

$$\boldsymbol{\alpha}\mathbf{x}^{(1)}$$

- (d) [1 pt] What is the vector \mathbf{z} whose elements are made up of the entries z_j in equation (5)? Write your answer in terms of \mathbf{a} .

$$\frac{1}{1+e^{-\mathbf{a}}}$$

- (e) [1 pt] **Select one:** We cannot take the matrix multiplication of our weights $\boldsymbol{\beta}$ and our vector \mathbf{z} since they are not compatible shapes. Which of the following would allow us to take the matrix multiplication of $\boldsymbol{\beta}$ and \mathbf{z} such that the entries of the vector $\mathbf{b} = \boldsymbol{\beta}\mathbf{z}$ are equivalent to the values of b_k in equation (5)?

- ☐ Remove the last column of $\boldsymbol{\beta}$
- ☐ 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 $\boldsymbol{\beta}$
- ☐ Append a row of 1's to be the first row of $\boldsymbol{\beta}$
- ☐ Take the transpose of $\boldsymbol{\beta}$

- (f) [1 pt] 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 .

$$\left[\frac{e^{b_1}}{e^{b_1} + e^{b_2} + e^{b_3}}, \frac{e^{b_2}}{e^{b_1} + e^{b_2} + e^{b_3}}, \frac{e^{b_3}}{e^{b_1} + e^{b_2} + e^{b_3}} \right]^T$$

2. [7 pts] We will now derive the matrix and vector forms for the backpropagation algorithm.

$$\frac{d\ell}{d\boldsymbol{\alpha}} = \begin{bmatrix} \frac{dJ}{d\alpha_{10}} & \frac{dJ}{d\alpha_{11}} & \cdots & \frac{dJ}{d\alpha_{1M}} \\ \frac{dJ}{d\alpha_{20}} & \frac{dJ}{d\alpha_{21}} & \cdots & \frac{dJ}{d\alpha_{2M}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{dJ}{d\alpha_{D0}} & \frac{dJ}{d\alpha_{D1}} & \cdots & \frac{dJ}{d\alpha_{DM}} \end{bmatrix}$$

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 $(\frac{d\ell}{d\alpha})_{2,1}$ should be equal to $\frac{d\ell}{d\alpha_{2,1}}$). Recall that ℓ is our loss function defined in equation (1.8)

- (a) [1 pt] The derivative of the softmax function with respect to b_k is as follows:

$$\frac{d\hat{y}_l}{db_k} = \hat{y}_l(\mathbb{I}[k = l] - \hat{y}_k)$$

where $\mathbb{I}[k = l]$ is an indicator function such that if $k = l$ then it returns value 1 and 0 otherwise. Using this, write the derivative $\frac{d\ell}{db_k}$ in a smart way such that you do not need this indicator function. Write your solutions in terms of $\hat{\mathbf{y}}_k, \mathbf{y}_k$.

$$\frac{d\ell}{db_k} = \hat{y}_k - y_k$$

- (b) [1 pt] What are the elements of the vector $\frac{d\ell}{db}$? (Recall that $\mathbf{y}^{(1)} = [0, 1, 0]^T$)

$$[\hat{y}_1, \hat{y}_2 - 1, \hat{y}_3]^T$$

- (c) [1 pt] What is the derivative $\frac{d\ell}{d\beta}$? Your answer should be in terms of $\frac{d\ell}{db}$ and \mathbf{Z} .

$$\frac{d\ell}{d\beta} = \frac{d\ell}{db} \frac{db}{d\beta} = \frac{d\ell}{db} \mathbf{Z}^T$$

- (d) [1 pt] Explain in one short sentence why must we go back to using the matrix β^* (The matrix β without the first column of ones) when calculating the matrix $\frac{d\ell}{d\alpha}$?

Because using β (3×5) will not match the dimension of α (4×7) when calculating $d\ell/d\alpha$

- (e) [1 pt] What is the derivative $\frac{d\ell}{dz}$? Your answer should be in terms of $\frac{d\ell}{db}$ and β^*

$$\frac{d\ell}{dz} = \frac{d\ell}{db} \frac{db}{dz} = \frac{d\ell}{db} \frac{d\beta^* z}{dz} = \beta^{*\top} \frac{d\ell}{db}$$

- (f) [1 pt] What is the derivative $\frac{d\ell}{da}$ in terms of $\frac{d\ell}{dz}$ and z

$$\frac{dl}{da} = \frac{dl}{dz} \frac{dz}{da} = \frac{dl}{dz} \circ z(1 - z)$$

$$\frac{dz}{da} = \frac{e^{-a}}{(1+e^{-a})^2} = z(1 - z)$$

- (g) [1 pt] What is the matrix $\frac{d\ell}{d\alpha}$? Your answer should be in terms of $\frac{d\ell}{da}$ and $x^{(1)}$.

$$\frac{dl}{d\alpha} = \frac{dl}{da} \frac{da}{d\alpha} = \frac{dl}{da} \frac{d\alpha x^{(1)}}{d\alpha} = \frac{dl}{da} x^{(1)\top}$$

Prediction

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

3. [8 pts] We initialize the weights as:

$$\boldsymbol{\alpha}^* = \begin{bmatrix} 2 & 1 & -1 & -1 & 0 & -2 \\ 0 & 1 & 0 & -1 & 1 & 3 \\ -1 & 2 & 1 & 3 & 1 & -1 \\ 1 & 3 & 4 & 2 & -1 & 2 \end{bmatrix}$$

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

And weights on the bias terms ($\alpha_{j,0}$ and $\beta_{j,0}$) are initialized to 1.

You are given a training example $\mathbf{x}^{(1)} = [1, 0, 1, 0, 1, 0]^T$ with label class 2, so $\mathbf{y}^{(1)} = [0, 1, 0]^T$. Using the initial weights, run the feed forward of the network over this training example (without rounding during the calculation) and then answer the following questions.

- (a) [1 pt] What is the value of a_1 ?

2

(b) [1 pt] What is the value of z_1 ?

0.8808

(c) [1 pt] What is the value of a_3 ?

2

(d) [1 pt] What is the value of z_3 ?

0.8808

(e) [1 pt] What is the value of b_2 ?

5.6290

(f) [1 pt] What is the value of \hat{y}_2 ?

0.8588

(g) [1 pt] Which class value we would predict on this training example?

2

(h) [1 pt] What is the value of the total loss on this training example?

0.1522

4. [4 pts] Now use the results of the previous question to run backpropagation over the network and update the weights. Use the learning rate $\eta = 1$.

Do your backpropagation calculations without any rounding then answer the following questions: (in your final responses round to four decimal places)

- (a) [1 pt] What is the updated value of $\beta_{2,1}$?

3.1244

- (b) [1 pt] What is the updated weight of the hidden layer bias term applied to y_1 (i.e. $\beta_{1,0}$)?

0.8682

- (c) [1 pt] What is the updated value of $\alpha_{3,4}$?

3

- (d) [1 pt] If we ran backpropagation on this example for a large number of iterations and then ran feed forward over the same example again, which class would we predict?

2

5. [5 pts] Let us now introduce regularization into our neural network. For this question, we will incorporate L2 regularization into our loss function $\ell(\hat{\mathbf{y}}, \mathbf{y})$, with the parameter λ controlling the weight given to the regularization term.

- (a) [1 pt] Write the expression for the regularized loss function of our network after adding L2 regularization (**Hint:** Remember that bias terms should not be regularized!)

$$L = - \sum_{k=1}^3 y_k \ln \hat{y}_k + \lambda \left(\sum_{j=1}^4 \sum_{i=1}^6 \alpha_{j,i}^2 + \sum_{j=1}^3 \sum_{i=1}^4 \beta_{j,i}^2 \right)$$

- (b) [1 pts] Compute the regularized loss for training example $\mathbf{x}^{(1)}$ (assume $\lambda = 0.01$ and use the weights before backpropagation)

0.9022

- (c) [1 pts] For a network which uses the regularized loss function, write the gradient update equation for $\alpha_{j,i}$. You may use $\frac{\partial \ell(\hat{\mathbf{y}}, \mathbf{y})}{\partial \alpha_{j,i}}$ to denote the gradient update w.r.t non-regularized loss and η to denote the learning rate.

$$\alpha_{j,i} \leftarrow \alpha_{j,i} - \eta \left[\frac{\partial \ell(\hat{\mathbf{y}}, \mathbf{y})}{\partial \alpha_{j,i}} + 2\lambda \alpha_{j,i} \right]$$

- (d) [2 pts] Based on your observations from previous questions, **select all statements which are true:**

- ☐ The non-regularized loss is always higher than the regularized loss
- ☒ As weights become larger, the regularized loss increases faster than non-regularized loss
- ☒ On adding regularization to the loss function, gradient updates for the network become larger
- ☒ When using large initial weights, weight values decrease more rapidly for a network which uses regularized loss
- ☐ None of the above

Problem 6: Neural Net Programming [36 pts]



Figure 3: Random Images of Each of 10 classes in Fashion-MNIST

Your goal in this assignment is to label images of fashion articles by implementing a neural network from scratch. You will implement all of the functions needed to initialize, train, evaluate, and make predictions with the network. **Important: You are not allowed to use available libraries like PyTorch or TensorFlow which will make a neural network for you, you must create this model yourself.**

The Fashion-MNIST dataset is comprised of 70,000 images of fashion articles and their respective labels. There are 60,000 training images and 10,000 test images, all of which are 28 pixels by 28 pixels. The images belong to the following 10 categories – [T-shirt, Trouser, Pullover, Dress, Coat, Sandal, Shirt, Sneaker, Bag, Ankle boot].

In this particular case, you will work with a smaller subset of Fashion-MNIST that consist of 4000 images where the pixel values in each of them ranges from 0 to 1.

Dataset format For the dataset stored under `training_data_student`, you are provided with 2 files `train.csv`, `test.csv` that contain the images and their corresponding labels (0 through 9). Each row contains $784 + 1$ columns separated by commas. Columns **1 to 784** represent the pixel values and column **785** contains the corresponding label (0 to 9). The two files contain 3000 and 1000 instances respectively.

Model Definition

Preliminaries

In this section, 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. For this particular

problem, the input vectors \mathbf{x} are of length $M = 28 \times 28$, the hidden layer \mathbf{z} consist of D hidden units, and the output layer $\hat{\mathbf{y}}$ represents a probability distribution over the $K = 10$ classes. In other words, each element \hat{y}_k of the output vector $\hat{\mathbf{y}}$ represents the probability of \mathbf{x} belonging to the class k .

Following the notation from the written section we have:

- For the output layer:

$$\hat{y}_k = \frac{\exp(b_k)}{\sum_{l=1}^K \exp(b_l)}, \quad k \in \{1, \dots, K\}, \quad (\text{softmax activation})$$

$$b_k = \beta_{k,0} + \sum_{j=1}^D \beta_{kj} z_j, \quad k \in \{1, \dots, K\}, \quad (\text{pre-activation})$$

- For the hidden layer:

$$z_j = \frac{1}{1 + \exp(-a_j)}, \quad j \in \{1, \dots, D\}, \quad (\sigma - \text{activation})$$

$$a_j = \alpha_{j,0} + \sum_{i=1}^M \alpha_{ji} x_i, \quad j \in \{1, \dots, D\}, \quad (\text{pre-activation})$$

It is possible to 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. With these considerations we have,

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

$$b_k = \sum_{j=0}^D \beta_{kj} z_j, \quad k \in \{1, \dots, K\}$$

$$z_j = \frac{1}{1 + \exp(-a_j)}, \quad j \in \{1, \dots, D\}$$

$$a_j = \sum_{i=0}^M \alpha_{ji} x_i, \quad j \in \{1, \dots, D\}$$

Objective function

Since the output corresponds to a probabilistic distribution over the K classes, the objective (cost) function we will use for training our neural network is the **average cross entropy**,

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

over the training dataset,

$$\mathcal{D} = \{(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})\}, \quad \text{for } n \in \{1, \dots, N\}$$

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

To train the network, 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**, though typically you should shuffle your data during SGD you are **not** to do so here, and instead you are to train through the dataset in the order it is given.

Implementation

To proceed, you are provided with the following guide. This is recommended but absolutely not required.

Defining layers

Just to recap, your network architecture should look like the following: **Linear** \rightarrow **Sigmoid** \rightarrow **Linear** \rightarrow **Softmax**. The size of the input is 784. The first linear layer can have 785 nodes to include bias term. The final (output) layer should have 10 nodes (one corresponding to each integer from 0 - 9). The hidden layer will have $D = 256$ nodes (excluding bias term).

Linear layer

1. Implement a forward function `linear_forward(x, weight, bias)`.
2. Implement a backward function `linear_backward(...)`. This function returns gradient of loss wrt input to the linear layer.

Sigmoid layer

1. Implement a forward function `sigmoid_forward(x)` that returns element-wise sigmoid of x . Output is of the same shape as x .
2. Implement a backward function `sigmoid_backward(...)`. This function returns gradient of loss wrt input to the sigmoid layer. The gradient of loss wrt output of sigmoid layer can be passed as an input to this function.

Softmax cross entropy loss layer

1. Implement a forward function `softmax_xeloss_forward(x, labels)`. The fusion of softmax and average cross entropy computation in one layer is common in popular deep learning frameworks.
2. Implement a backward function `softmax_xeloss_backward(...)`.

Be aware of computing issues! $\log(x)$ is problematic when $x \rightarrow 0$. Similarly $\exp(x)$ may overflow when it is huge. Think of using log to avoid some exponential calculations and dividing both numerator and denominator by a large value to avoid overflowing:

$$\frac{e^{x_i}}{\sum e^{x_j}} = \frac{e^{x_i-b}}{\sum e^{x_j-b}}$$

Pseudo-code for Training Loop

For epoch in epochs:

```
    for x, y in train_x, train_y:
        Pass x through all the forward layers and get the loss
        Pass the output through all the backward layers
        Update weights and biases
    compute the average training loss for the epoch
    compute test loss
    compute test accuracy
```

You can follow the steps mentioned above or approach it any other way.

Important Tips

- The training loss you report, should not be aggregated throughout an epoch, rather you should compute the loss on the whole training set at the end of each epoch.
- The final loss value to report when asked should be averaged across all batches.
- When performing mini-batch gradient descent, you should take the mean of the loss within a batch before applying the update.
- Do NOT shuffle the training or test data, otherwise the answers won't match.

Programming Submission

NOTE: Initial weights and biases are provided in the folder `params`. Use the following hyper-parameters:

- Learning rate = 0.01
- Number of epochs = 15
- Width of hidden layer = 256 (excluding bias)

We have provided the outputs of first five nodes of each layer and the updated weights for the first data point in the first epoch in the folder `check`. Use these to verify your implementation.

1. [1 pt] **Linear** : What is the value of a_{10} for first data point in first epoch?

-0.9639

2. [1 pt] **Sigmoid** : What is the value of z_{20} for first data point in first epoch?

0.0746

3. [1 pt] What is the predicted class for the first data point x_1 after the first forward pass in the first epoch. Remember that our y values go from 0 to 9.

0

4. [4 pts] What are the values of $\beta_{k,0}$ for $k \in \{1, \dots, K\}$ (bias terms of the second layer) at the end of third epoch? Report numbers till 4 places of decimal and in the correct order.

-1.1218, -0.6185, -0.7467, 0.6280, -0.2927,
0.1438, -0.6423, 0.1388, 2.0715, 0.4232

5. [4 pts] List the average test loss at the end of each epoch. Report numbers till 4 places of decimal.

1.5640, 1.3888, 1.2866, 1.2232, 1.1940,
1.1684, 1.1316, 1.0895, 1.0559, 1.0290,
1.0072, 0.9902, 0.9784, 0.9708, 0.9661

6. [4 pts] List the test accuracy at the end of each epoch. Report numbers till 4 places of decimal.

0.6630, 0.6980, 0.7090, 0.7220, 0.7240,
0.7320, 0.7350, 0.7440, 0.7500, 0.7590,
0.7630, 0.7640, 0.7630, 0.7700, 0.7690

7. [2 pts] Run the model for 100 epochs, report the average final training loss and the test accuracy. Report numbers till 4 places of decimal.

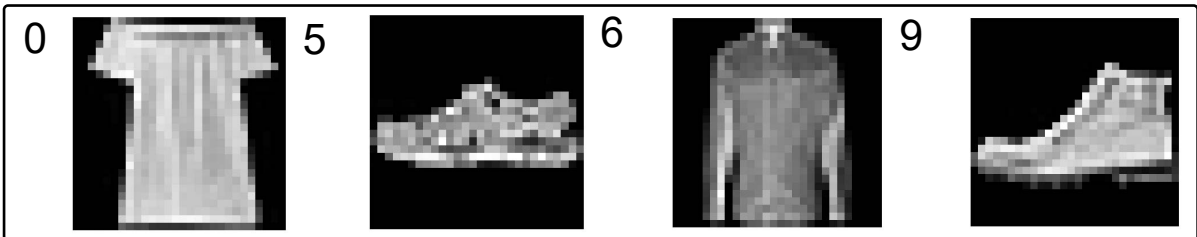
Training Loss:	0.0291
Test Accuracy:	0.7810

8. [4 pts] For the model trained for 100 epochs, display the confusion matrix for both the training and test sets. In the plot, you should have *true labels* on the y-axis and *predicted labels* on the x-axis. For example, the second entry in the fourth row should show the count for the number of times class 3 was wrongly predicted as 1.

Training										
	0	1	2	3	4	5	6	7	8	9
0	300	0	0	0	0	0	0	0	0	0
1	0	300	0	0	0	0	0	0	0	0
2	0	0	292	0	8	0	0	0	0	0
3	0	0	0	299	1	0	0	0	0	0
4	0	0	0	0	300	0	0	0	0	0
5	0	0	0	0	0	300	0	0	0	0
6	0	0	2	0	7	0	291	0	0	0
7	0	0	0	0	0	0	0	300	0	0
8	0	0	0	0	0	0	0	0	300	0
9	0	0	0	0	0	0	0	0	0	300

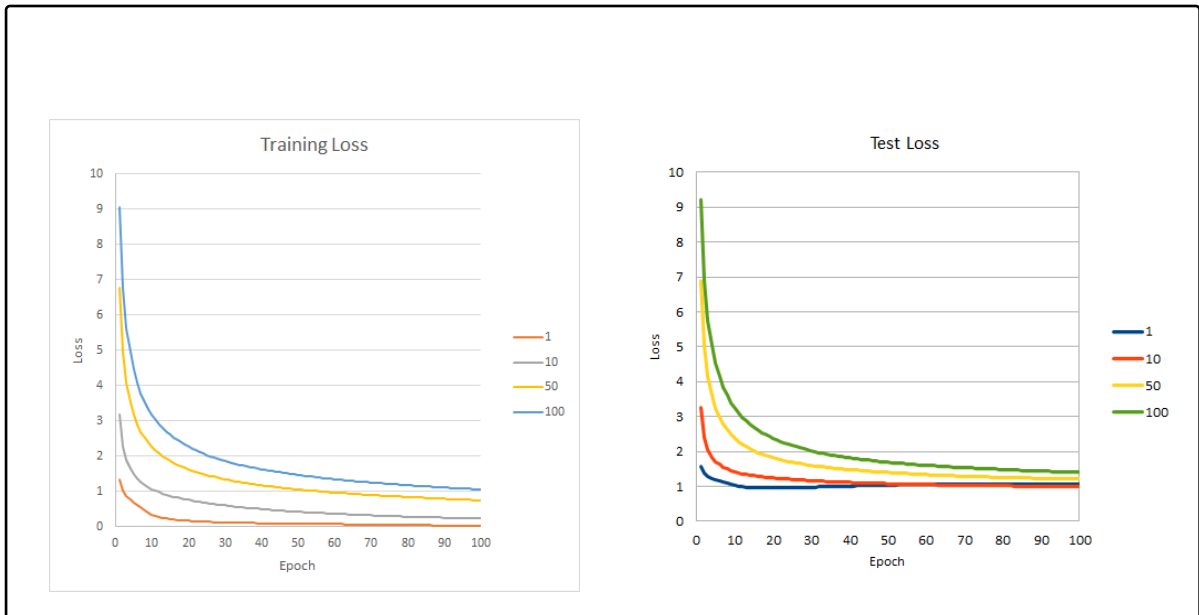
Test										
	0	1	2	3	4	5	6	7	8	9
0	81	0	3	1	1	0	11	0	3	0
1	3	88	1	5	1	0	1	0	1	0
2	3	0	61	1	26	0	9	0	0	0
3	6	2	1	78	9	0	4	0	0	0
4	1	1	4	4	83	0	7	0	0	0
5	0	0	0	0	0	90	1	7	0	2
6	22	0	7	1	22	0	43	0	5	0
7	0	0	0	0	0	7	0	85	0	8
8	1	0	1	1	4	3	3	2	84	1
9	0	0	0	0	0	3	0	8	1	88

9. [2 pts] For the categories 0 (T-Shirt), 5 (Coat), 6 (Sandal) and 9 (Ankle boot), display a sample each from the test set which was wrongly classified. (For Python - you can find a code snippet to display an image given a vector in the `utils` folder).



10. [8 pts] For the experiments until now, we were using stochastic gradient descent (SGD) i.e. batch size = 1. In this section, we will modify our code to enable training on batches of data i.e. we'll use mini-batch gradient descent. For each batch size in [1, 10, 50, 100], run your model for 100 epochs. Plot two graphs of epoch vs loss, one each for training and test loss, which have the epoch number on the x-axis and the loss value on the y-axis. Please make sure your graphs are properly labelled and include a

legend showing the color for each batch size.



11. [2 pts] Based on the plots in the previous part, what do you observe about the effect of increasing batch sizes? Also, how does this relate to the learning rate? (Explain in 1-2 lines)

When we increase the batch size, the loss values also increase. We should use a higher learning rate for a large batch size to optimize the loss function.

12. [2 pts] Based on the plots in the part 10 of the programming question, what do you think is the best next step given the trends of different batch sizes? (Explain in 1-2 lines)

We should try some numbers between 1 and 50 for batch size and higher values (i.e., 0.1) for learning rate.

13. [1 pt] Now we want to give you a chance to experiment with the model by exploring the hyper-parameters of a neural network. Produce three plots similar to that of the one requested in batch size experiment, however each of which you should have changed hyper-parameters to a different values. Tell us the details of the hyper-parameter that you changed and explain how this compares to your findings in the batch size experiment. For your reference here are some hyperparameters which you can change:

- Learning Rate
- Width of Hidden Layer
- Different weight initialization

Some examples may be:

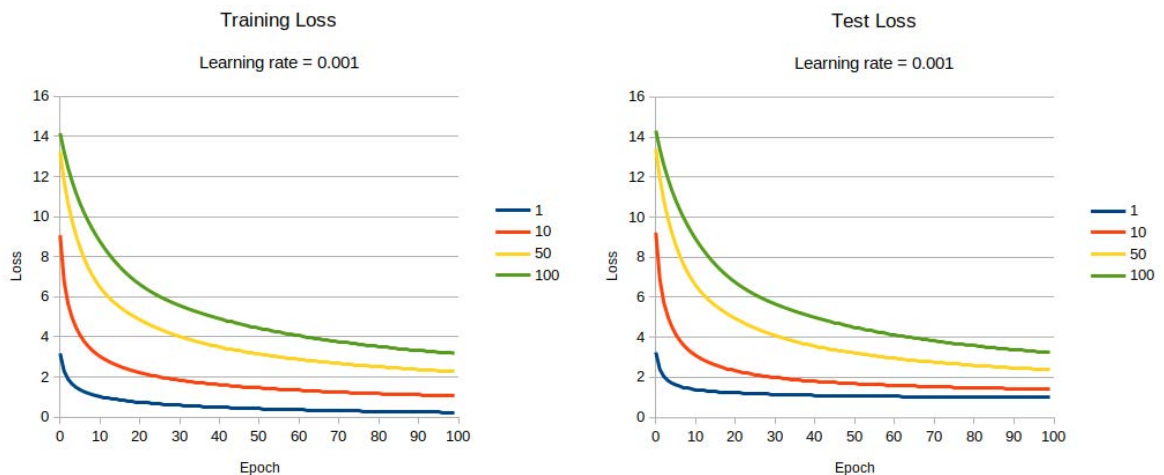
- Running your model with learning rate = 0.3, with 100 hidden nodes and weights initialized all to zeros.
- Running your model with everything the same except comparing three different numbers of hidden nodes.

This question is mostly for fun so you are welcome to do whatever you want provided you label your graphs and give some sort of explanation.

Experiment 1

Set the learning rate = 0.001

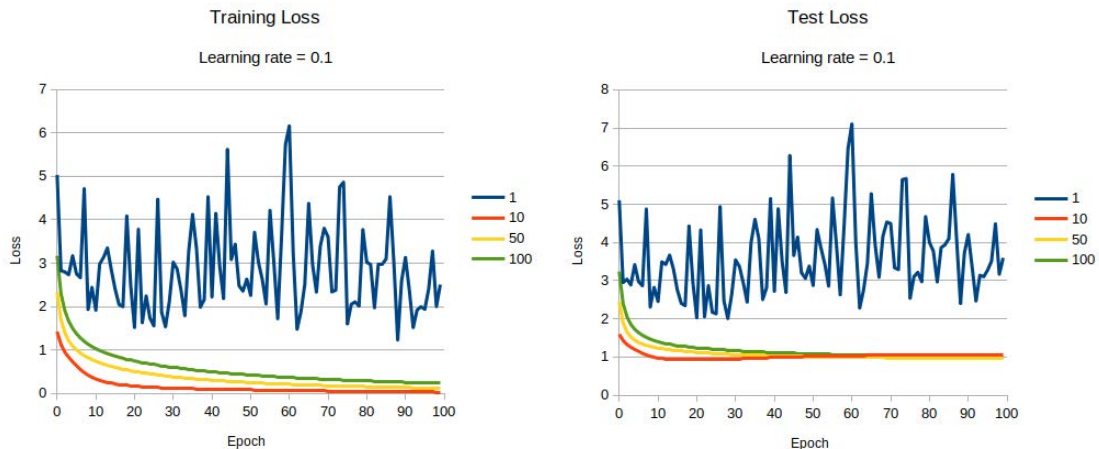
When set the learning rate to a smaller value, the loss values for training and test are both higher.



Experiment 2

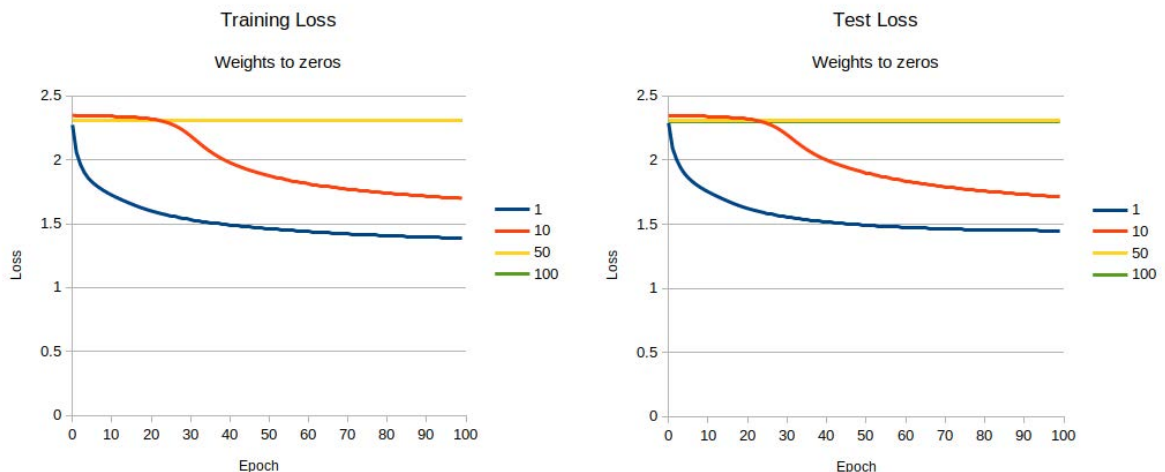
Set the learning rate = 0.1

When set the learning rate to a larger value, the loss values for large batch sizes for training and test are both lower. However, for batch size = 1, the loss is not decreasing as epoch number goes up, which means this learning rate is too large for batch size = 1.



Experiment 3

Initialize weights (α^* and β^*) to zeros and the bias terms to ones. The loss is higher for all batch sizes. For large batch size, the loss remains the same throughout 100 epochs. Thus, we will need more epochs and small batch size if we start learning from 0 weights.



1 Collaboration Questions

1. (a) Did you receive any help whatsoever from anyone in solving this assignment?

SOL: ☒ Yes / No.

- (b) If you answered 'yes', give full details (e.g. "Jane Doe explained to me what is asked in Question 3.4")

SOL:

Yajushi Khurana and Swapnil Keshari explained to me what is asked in Question 5, 6

2. (a) Did you give any help whatsoever to anyone in solving this assignment? **SOL:**

☒ Yes / No.

- (b) If you answered 'yes', give full details (e.g. "I pointed Joe Smith to section 2.3 since he didn't know how to proceed with Question 2")

SOL:

I pointed Yajushi Khurana and Swapnil Keshari to Question 5, 6

3. (a) Did you find or come across code that implements any part of this assignment?

SOL: Yes / ☒ No.

- (b) If you answered 'yes', give full details (book & page, URL & location within the page, etc.).

SOL: