

CS4803/7643: Deep Learning

Fall 2020

Problem Set 3

Instructor: Dhruv Batra

TAs: Sameer Dharur, Joanne Truong, Yihao Chen, Hrishikesh Kale
Tianyu Zhan, Prabhav Chawla, Guillermo Nicolas Grande, Michael Pisen

Discussions: <https://piazza.com/gatech/fall2020/cs48037643>

Due: Wednesday October 7, 2020, 11:59pm

Instructions

1. We will be using Gradescope to collect your assignments. Please read the following instructions for submitting to Gradescope carefully!
 - Each subproblem must be submitted on a separate page. When submitting to Gradescope, make sure to mark the page(s) corresponding to each problem/sub-problem. For instance, Q3 has 3 subproblems - the solution to each must start on a new page and be marked accordingly.
 - Also remember to append your notebook PDFs from Q18 and Q19 to the solutions for the problem set, as described in the instructions.
 - For the coding problem, please use the provided `collect_submission.sh` script and upload `cs7643_hw3.zip` to the HW3 Code assignment on Gradescope. While we will not be explicitly grading your code, you are still required to submit it. Please make sure you have saved the most recent version of your Jupyter notebooks before running this script.
 - Note: This is a large class and Gradescope's assignment segmentation features are essential. Failure to follow these instructions may result in parts of your assignment not being graded. We will not entertain regrading requests for failure to follow instructions. Please check [this](#) link for additional information on submitting to Gradescope.
2. L^AT_EX'd solutions are strongly encouraged (solution template available at cc.gatech.edu/classes/AY2021/cs7643_fall/assets/sol3.tex), but scanned handwritten copies are acceptable. Hard copies are **not** accepted.
3. We generally encourage you to collaborate with other students.

You may talk to a friend, discuss the questions and potential directions for solving them. However, you need to write your own solutions and code separately, and *not* as a group activity. Please list the students you collaborated with.

This assignment introduces you to some theoretical results about learned representations in neural networks.

1 Convolution Basics

The convolution layer inside of a CNN is intrinsically an *affine transformation*: A vector is received as input and is multiplied with a matrix to produce an output (to which a bias vector is usually added before passing the result through a nonlinearity). This operation can be represented as $y = Ax$, in which A describes the affine transformation.

1. **[0.5 points]** We will first revisit the convolution layer as discussed in the class. Consider a convolution layer with a 2x2 kernel W , operated on a single input channel X , represented as:

$$W = \begin{bmatrix} w_{(0,0)}, w_{(0,1)} \\ w_{(1,0)}, w_{(1,1)} \end{bmatrix}, X = \begin{bmatrix} x_{(0,0)}, x_{(0,1)}, x_{(0,2)} \\ x_{(1,0)}, x_{(1,1)}, x_{(1,2)} \\ x_{(2,0)}, x_{(2,1)}, x_{(2,2)} \end{bmatrix} \quad (1)$$

Now let us work out a **stride-3** convolution layer, with **zero padding size of 1**. Consider ‘flattening’ the input tensor X in row-major order as:

$$X = [x_{(0,0)}, x_{(0,1)}, x_{(0,2)}, \dots, x_{(2,0)}, x_{(2,1)}, x_{(2,2)}]^\top \quad (2)$$

Write down the convolution as a matrix operation A such that: $Y = AX$. Output Y is also flattened in row-major order.

2. **[0.5 points]** Recall that transposed convolution can help us upsample the feature size spatially. Consider a transposed convolution layer with a 2x2 kernel W operated on a single input channel X , represented as:

$$W = \begin{bmatrix} w_{(0,0)}, w_{(0,1)} \\ w_{(1,0)}, w_{(1,1)} \end{bmatrix}, X = \begin{bmatrix} x_{(0,0)}, x_{(0,1)} \\ x_{(1,0)}, x_{(1,1)} \end{bmatrix} \quad (3)$$

We ‘flattened’ the input tensor in row-major order as $X = [x_{(0,0)}, x_{(0,1)}, x_{(1,0)}, x_{(1,1)}]$.

Write down the affine transformation A corresponding to a transposed convolution layer with kernel W , **stride 2**, **no padding**. Output Y is also flattened in row-major order.

3. **[1 points]** Convolution layers in most CNNs consist of multiple input and output feature maps. The collection of kernels form a 4D tensor (output channels o , input channels i , filter rows k , filter columns k), represented in short as (o, i, k, k) . For each output channel, each input channel is convolved with a distinct 2D slice of the 4D kernel and the resulting set of feature maps is summed element-wise to produce the corresponding output feature map.

There is an interesting property that a convolutional layer with kernel size $(o \times r^2, i, k, k)$ is identical to a transposed convolution layer with kernel size $(o, i, k \times r, k \times r)$. Here the word ‘identical’ means with the same input feature X , both operations will give the same output Y with only a difference in the ordering of flattened elements of Y .

Now let us prove the property in a restricted setting. Consider $o = 1, r = 2, i = 1, k = 1$. Given the same input feature X as in equation 3, write down the affine transformation for a convolutional layer with kernel size $(4, 1, 1, 1)$, and show that it is an operation identical to a transposed convolution layer with kernel size $(1, 1, 2, 2)$.

2 Logic and XOR

4. [1 points] Implement AND and OR for pairs of binary inputs using a single linear threshold neuron with weights $\mathbf{w} \in \mathbb{R}^2$, bias $b \in \mathbb{R}$, and $\mathbf{x} \in \{0, 1\}^2$:

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x} + b \geq 0 \\ 0 & \text{if } \mathbf{w}^T \mathbf{x} + b < 0 \end{cases} \quad (4)$$

That is, find \mathbf{w}_{AND} and b_{AND} such that

x_1	x_2	$f_{\text{AND}}(\mathbf{x})$
0	0	0
0	1	0
1	0	0
1	1	1

Also find \mathbf{w}_{OR} and b_{OR} such that

x_1	x_2	$f_{\text{OR}}(\mathbf{x})$
0	0	0
0	1	1
1	0	1
1	1	1

5. [1 points] Consider the XOR function

x_1	x_2	$f_{\text{XOR}}(x)$
0	0	0
0	1	1
1	0	1
1	1	0

Show that XOR can NOT be represented using a linear model with the same form as (4).

[Hint: To see why, plot the examples from above in a plane and think about drawing a linear boundary that separates them.]

3 Piecewise Linearity

Consider a specific 2 hidden layer ReLU network with inputs $x \in \mathbb{R}$, 1 dimensional outputs, and 2 neurons per hidden layer. This function is given by

$$h(x) = W^{(3)} \max\{0, W^{(2)} \max\{0, W^{(1)}x + \mathbf{b}^{(1)}\} + b^{(2)}\} + b^{(3)} \quad (5)$$

with weights:

$$W^{(1)} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \quad (6)$$

$$b^{(1)} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad (7)$$

$$W^{(2)} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \quad (8)$$

$$b^{(2)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad (9)$$

$$W^{(3)} = \begin{bmatrix} 1 & 1 \end{bmatrix} \quad (10)$$

$$b^{(3)} = 1 \quad (11)$$

An interesting property of networks with piecewise linear activations like the ReLU is that on the whole they compute piecewise linear functions. For each of the following points give the weight $W \in \mathbb{R}$ and bias $b \in \mathbb{R}$ (report the numerical values) which computes such that $Wx + b = h(x)$. Also compute the gradient $\frac{dh}{dx}$ evaluated at the given point.

6. [1 points]

$$x = 1 \quad (12)$$

7. [1 points]

$$x = -1 \quad (13)$$

8. [1 points]

$$x = -0.5 \quad (14)$$

4 Depth - Composing Linear Pieces [Extra Credit for 4803; Regular Credit for 7643]

Now we'll turn to a more recent result that highlights the *Deep* in Deep Learning. Depth (composing more functions) results in a favorable combinatorial explosion in the “number of things that a neural net can represent”. For example, to classify a cat it seems useful to first find parts of a cat: eyes, ears, tail, fur, *etc.* The function which computes a probability of cat presence should be a function of these components because this allows everything you learn about eyes to generalize to all instances of eyes instead of just a single instance. Below you will detail one formalizable sense of this combinatorial explosion for a particular class of piecewise linear networks.

Consider $y = \sigma(x) = |x|$ for scalar $x \in \mathcal{X} \subseteq \mathbb{R}$ and $y \in \mathcal{Y} \subseteq \mathbb{R}$ (Fig. 1). It has one linear region on $x < 0$ and another on $x > 0$ and the activation identifies these regions, mapping both of them to $y > 0$. More precisely, *for each linear region of the input*, $\sigma(\cdot)$ is a bijection. There is a mapping to and from the output space and the corresponding input space. However, given an output y , it's impossible to tell which linear region of the input it came from, thus $\sigma(\cdot)$ *identifies* (maps on top of each other) the two linear regions of its input. This is the crucial definition because when a function identifies multiple regions of its domain that means any subsequent computation applies to all of

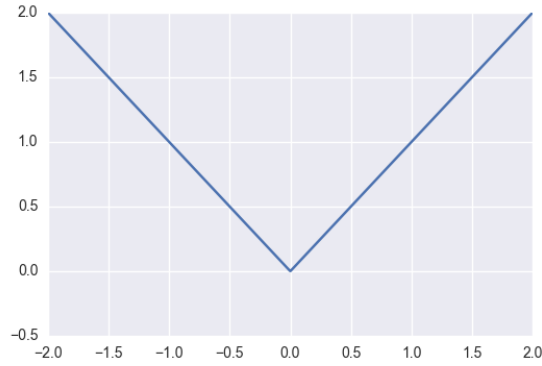


Figure 1

those regions. When these regions come from an input space like the space of images, functions which identify many regions where different images might fall (*e.g.*, slightly different images of a cat) automatically transfer what they learn about a particular cat to cats in the other regions.

More formally, we will say that $\sigma(\cdot)$ identifies a set of M disjoint input regions $\mathcal{R} = \{R_1, \dots, R_M\}$ (*e.g.*, $\mathcal{R} = \{(-1, 0), (0, 1)\}$) with $R_i \subseteq \mathcal{X}$ onto one output region $O \subseteq \mathcal{Y}$ (*e.g.*, $(0, 1)$) if for all $R_i \in \mathcal{R}$ there is a bijection from R_i to O .¹

9. **[3 points]** Start by applying the above notion of identified regions to linear regions of one layer of a particular neural net that uses absolute value functions as activations. Let $\mathbf{x} \in \mathbb{R}^d$, $\mathbf{y} \in \mathbb{R}^d$ ², and pick weights $W^{(1)} \in \mathbb{R}^{d \times d}$ and bias $\mathbf{b}^{(1)} \in \mathbb{R}^d$ as follows:

$$W_{ij}^{(1)} = \begin{cases} 2 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases} \quad (15)$$

$$b_i^{(1)} = -1 \quad (16)$$

Then one layer of a neural net with absolute value activation functions is given by

$$f_1(\mathbf{x}) = |W^{(1)}\mathbf{x} + \mathbf{b}| \quad (17)$$

Note that this is an absolute value function applied piecewise and not a norm.

How many regions of the input are identified onto $O = (0, 1)^d$ by $f_1(\cdot)$? Prove it.³

10. **[2 points]** Next consider what happens when two of these functions are composed. Suppose g identifies n_g regions of $(0, 1)^d$ onto $(0, 1)^d$ and f identifies n_f regions of $(0, 1)^d$ onto $(0, 1)^d$. How many regions of its input does $f \circ g(\cdot)$ identify onto $(0, 1)^d$?

11. **[4 points]** Finally consider a series of L layers identical to the one in question 4.1.

¹Recall that a bijection from X to Y is a function $\mu : X \rightarrow Y$ such that for all $y \in Y$ there exists a **unique** $x \in X$ with $\mu(x) = y$.

²Outputs are in some feature space, not a label space. Normally a linear classifier would be placed on top of what we are here calling \mathbf{y} .

³Absolute value activations are chosen to make the problem simpler, but a similar result holds for ReLU units. Also, O could be the positive orthant (unbounded above).

$$\mathbf{h}_1 = |W_1 \mathbf{x} + \mathbf{b}_1| \quad (18)$$

$$\mathbf{h}_2 = |W_2 \mathbf{h}_1 + \mathbf{b}_2| \quad (19)$$

$$\vdots \quad (20)$$

$$\mathbf{h}_L = |W_L \mathbf{h}_{L-1} + \mathbf{b}_L| \quad (21)$$

Let $\mathbf{x} \in (0,1)^d$ and $f(\mathbf{x}) = \mathbf{h}_L$. Note that each \mathbf{h}_i is *implicitly* a function of \mathbf{x} . Show that $f(\mathbf{x})$ identifies 2^{Ld} regions of its input.

5 Conclusion to the above questions

Now compare the number of identified regions for an L layer net to that of an $L - 1$ layer net. The L layer net can separate its input space into 2^d more linear regions than the $L - 1$ layer net. On the other hand, the number of parameters and the amount of computation time grows linearly in the number of layers. In this very particular sense (which doesn't always align well with practice) deeper is better.

To summarize this problem set, you've shown a number of results about the representation power of different neural net architectures. First, neural nets (even single neurons) can represent logical operations. Second, neural nets we use today compute piecewise linear functions of their input. Third, the representation power of neural nets increases exponentially with the number of layers. The point of the exercise was to convey intuition that removes some of the magic from neural nets representations. Specifically, neural nets can decompose problems logically, and piecewise linear functions can be surprisingly powerful.

6 Implicit Regularization of Gradient descent [Extra Credit for both 4803 and 7643, 7 points]

Classical results in statistical learning theory suggest that as the model class becomes larger, we should expect more overfitting, with everything else staying fixed (such as the training dataset size). However, this classical view appears to be at odds with the 'street wisdom' among practitioners (particularly in deep learning) that *'larger models are better'*. A gem sometimes expressed as *'add MOAR layers!'*. More formally, in deep learning it is fairly common to have more parameters in the model than number of datapoints ($d > n$). Results from classical statistics will tell you that such settings are hopeless without more assumptions. In this section and the next one, we will try to reconcile this ostensible disagreement.

Specifically, in this question we will build intuition using a simple least-squares linear regression problem, optimized with gradient descent.

Consider a linear regression problem with a d -dimensional feature vector $\mathbf{x} \in \mathbb{R}^d$ as input and $y_i \in \mathbb{R}$ as the output. The dataset consists of n training examples $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), (\mathbf{x}_3, y_3), \dots, (\mathbf{x}_n, y_n)$. Let \mathbf{X} be the $n \times d$ data matrix formed by placing the feature vectors on the rows of this matrix, *i.e.* $\mathbf{X}^T = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$. Let \mathbf{y} be a column vector with elements y_1, y_2, \dots, y_n . We will operate in the setting where $d > n$, *i.e.* there are more feature dimensions than samples. Thus, the following linear system is under-constrained / under-determined:

$$\mathbf{X}\mathbf{w} = \mathbf{y} \quad (22)$$

To avoid the degenerate case, we will assume that \mathbf{y} lies in the span of \mathbf{X} , *i.e.* this linear system has at least one solution.

Now, recall that the optimization problem in least-squares regression is the following:

$$\min_{\mathbf{w} \in \mathbf{R}^d} f(\mathbf{w}) = \sum_{i=1}^n \underbrace{(y_i - \mathbf{w}^T \mathbf{x}_i)^2}_{\text{squared error on example } i} \quad (23)$$

We will optimize (23) via gradient descent. Specifically, let us initialize $\mathbf{w}^{(0)} = 0$. And repeatedly take a step in the direction of negative gradient with a sufficiently-small constant step-size η till convergence:

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla f(\mathbf{w}^{(t)}) \quad (24)$$

Let us refer to the solution found by gradient descent at convergence as \mathbf{w}^{gd} .

12. **[4 point]** Prove that the solution found by gradient descent for least-squares is equal to the result of the following *different* optimization problem:

$$\mathbf{w}^{gd} = \underset{\mathbf{w} \in \mathbf{R}^d}{\operatorname{argmin}} \quad \|\mathbf{w}\|_2^2 \quad (25)$$

$$s.t. \quad \mathbf{X}\mathbf{w} = \mathbf{y} \quad (26)$$

13. **[1 point]** In a few words, describe what the optimization problem in (25), (26) is trying to express. Specifically, explain what the constraint (26) means/represents in terms of training error of least-squares. Explain what the objective (25) means in the context of machine learning. Putting them together, what solution is (25), (26) looking for.
14. **[1 point]** What does this tell you about gradient descent (in the context of least squares regression)? Notice that (23) is an unconstrained optimization problem. Is the solution unique? Does gradient descent appear to have a preference among the solutions?
15. **[1 point]** In class, Dhruv has been teaching you a fairly ‘clean’ separation between approximation error, estimation error, and optimization error. What does your investigation in this question tell you about that abstraction?

7 Paper Review [Extra credit for 4803, regular credit for 7643]

The paper we will study in this homework is ‘**Deep Double Descent: Where Bigger Models And More Data Hurt**’.

The paper presents a concept called ‘double descent’ (in the context of deep learning), with interesting experiments showing model performance first getting worse and then improving as we increase model size.

The paper can be viewed [here](#). The evaluation rubric for this section is as follows:

16. **[2 points]** Briefly summarize the key contributions, strengths and weaknesses of this paper.
17. **[2 points]** What is your personal takeaway from this paper? This could be expressed either in terms of relating the approaches adopted in this paper to your traditional understanding of learning parameterized models, or potential future directions of research in the area which the authors haven't addressed, or anything else that struck you as being noteworthy.

Guidelines: Please restrict your reviews to no more than 350 words (total length for answers to both the above questions).

8 Coding: Uses of Gradients With Respect to Input [20 points for both sections]

The coding part of this assignment will have you implement different ideas which all rely on gradients of some neural net output with respect to the input image. To get started click [here](#).