

# Homework 1

**Version:** 1.1

**Version Release Date:** 2022-01-20

**Changes by Version:**

- (v1.1)
  - Q1.3: Change  $\mathbf{W}^{(3)} \in \mathbb{R}^4$  to  $\mathbf{W}^{(3)} \in \mathbb{R}^{1 \times 4}$ ,
  - Q2.2.2: Change  $\|\mathbf{x}\| \|\bar{\mathbf{z}}\|$  to  $\|\mathbf{x}\|_2^2 \|\bar{\mathbf{z}}\|_2^2$ ,
  - Q2.2.2: Fix missing transpose,
  - Q3.2.2: Change  $\frac{1}{n} \|(X(XX^\top)^{-1}X^\top - I)\epsilon\|^2$  to  $\frac{1}{n} \|(X(X^\top X)^{-1}X^\top - I)\epsilon\|^2$ ,
  - Q3.2.2: Change “in terms of  $n$  and  $\sigma$ ” to “in terms of  $n$ ,  $d$  and  $\sigma$ ”

**Deadline:** Friday, Jan. 28, at 11:59pm.

**Submission:** You must submit your solutions as a PDF file through MarkUs<sup>1</sup>. You can produce the file however you like (e.g. LaTeX, Microsoft Word, scanner), as long as it is readable.

See the syllabus on the course website<sup>2</sup> for detailed policies. You may ask questions about the assignment on Piazza<sup>3</sup>. *Note that 10% of the homework mark (worth 1 pt) may be removed for a lack of neatness.*

You may notice that some questions are worth 0 pt, which means we will not mark them in this Homework. Feel free to skip them if you are busy. However, you are expected to see some of them in the midterm. So, we won't release the solution for those questions.

The teaching assistants for this assignment are Yongchao Zhou and Shervin Mehryar. Send your email with subject “[CSC413] HW1 ...” to [csc413-2022-01-tas@cs.toronto.edu](mailto:csc413-2022-01-tas@cs.toronto.edu) or post on Piazza with the tag `hw1`.

## 1 Hard-Coding Networks

Can we use neural networks to tackle coding problems? Yes! In this question, you will build a neural network to find the  $k^{th}$  smallest number from a list using two different approaches: sorting and counting (Optional). You will start by constructing a two-layer perceptron “Sort\_2” to sort two numbers and then use it as a building block to perform your favorite sorting algorithm (e.g., Bubble Sort, Merge Sort). Finally, you will output the  $k^{th}$  element from the sorted list as the final answer.

Note: Before doing this problem, you need to have a basic understanding of the key components of neural networks (e.g., weights, activation functions). The reading on multilayer perceptrons located at <https://uoft-csc413.github.io/2022/assets/readings/L02a.pdf> may be useful.

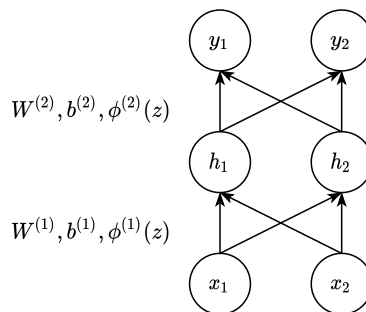
<sup>1</sup><https://markus.teach.cs.toronto.edu/2022-01/>

<sup>2</sup><https://uoft-csc413.github.io/2022/assets/misc/syllabus.pdf>

<sup>3</sup><https://piazza.com/class/ky8yug7i9cl3b1>

### 1.1 Sort two numbers [1pt]

In this problem, you need to find a set of weights and bias for a two-layer perceptron “Sort\_2” that sorts two numbers. The network takes a pair of numbers  $(x_1, x_2)$  as input and output a sorted pair  $(y_1, y_2)$ , where  $y_1 \leq y_2$ . You may assume the two numbers are distinct and positive for simplicity. You will use the following architecture:



Please specify the weights and activation functions for your network. Your answer should include:

- Two weight matrices:  $\mathbf{W}^{(1)}, \mathbf{W}^{(2)} \in \mathbb{R}^{2 \times 2}$
- Two bias vector:  $\mathbf{b}^{(1)}, \mathbf{b}^{(2)} \in \mathbb{R}^2$
- Two activation functions:  $\phi^{(1)}(z), \phi^{(2)}(z)$

You do not need to show your work.

*Hints: Sorting two numbers is equivalent to finding the min and max of two numbers.*

$$\max(x_1, x_2) = \frac{1}{2}(x_1 + x_2) + \frac{1}{2}|x_1 - x_2|, \quad \min(x_1, x_2) = \frac{1}{2}(x_1 + x_2) - \frac{1}{2}|x_1 - x_2|$$

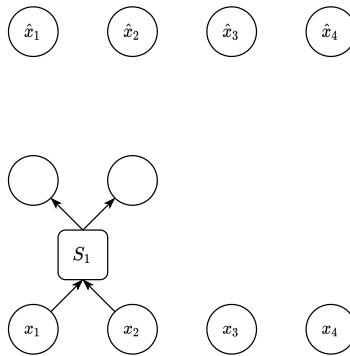
**Answer**

$$\phi^{(1)}(z) = |z|, \quad \phi^{(2)}(z) = z$$

$$\mathbf{W}^{(1)} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}, \mathbf{b}^{(1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \mathbf{W}^{(2)} = \begin{bmatrix} 1/2 & -1/2 \\ 1/2 & 1/2 \end{bmatrix}, \text{ and } \mathbf{b}^{(2)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

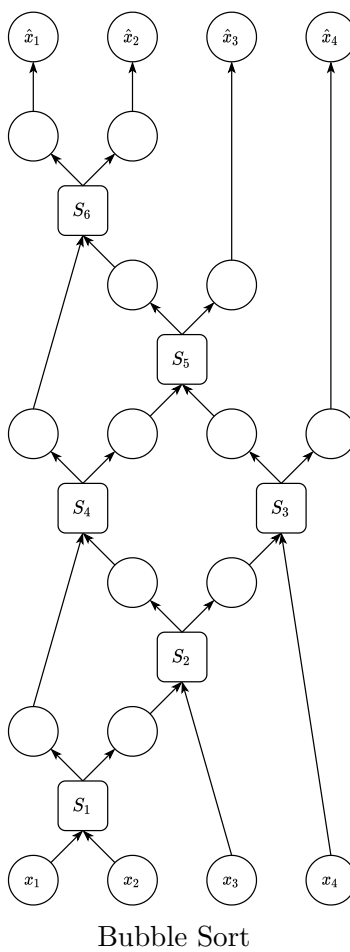
### 1.2 Perform Sort [1.5pt]

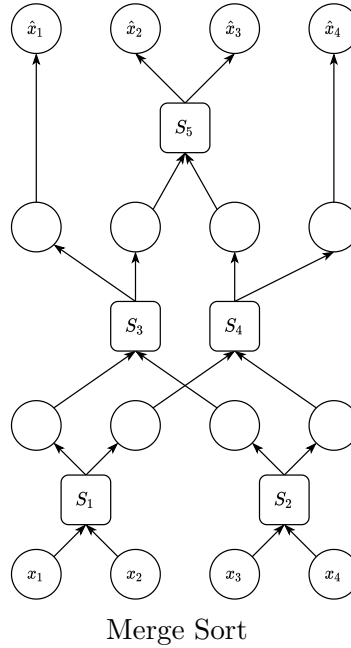
Draw a computation graph to show how to implement a sorting function  $\hat{f} : \mathbb{R}^4 \rightarrow \mathbb{R}^4$  where  $\hat{f}(x_1, x_2, x_3, x_4) = (\hat{x}_1, \hat{x}_2, \hat{x}_3, \hat{x}_4)$  where  $(\hat{x}_1, \hat{x}_2, \hat{x}_3, \hat{x}_4)$  is  $(x_1, x_2, x_3, x_4)$  in sorted order. Let us assume  $\hat{x}_1 \leq \hat{x}_2 \leq \hat{x}_3 \leq \hat{x}_4$  and  $x_1, x_2, x_3, x_4$  are positive and distinct. Implement  $\hat{f}$  using your favourite sorting algorithms (e.g. Bubble Sort, Merge Sort). Let us denote the “Sort\_2” module as  $S$ , please complete the following computation graph. Your answer does not need to give the label for intermediate nodes, but make sure to index the “Sort\_2” module.



*Hints: Bubble Sort needs 6 “Sort\_2” blocks, while Merge Sort needs 5 “Sort\_2” blocks.*

**Answer**





### 1.3 Find the $k^{th}$ smallest number [0pt]

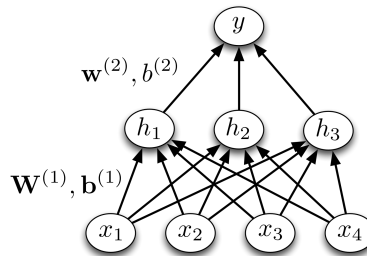
Based on your sorting network, you may want to add a new layer to output your final result ( $k^{th}$  smallest number). Please give the weight  $\mathbf{W}^{(3)}$  for this output layer when  $k = 3$ .

*Hints:*  $\mathbf{W}^{(3)} \in \mathbb{R}^{1 \times 4}$ .

### 1.4 Counting Network [0pt]

The idea of using a counting network to find the  $k^{th}$  smallest number is to build a neural network that can determine the rank of each number and output the number with the correct rank. Specifically, the counting network will count how many elements in a list are less than a value of interest. And you will apply the counting network to all numbers in the given list to determine their rank. Finally, you will use another layer to output the number with the correct rank.

The counting network has the following architecture, where  $y$  is the rank of  $x_1$  in a list containing  $x_1, x_2, x_3, x_4$ .



Please specify the weights and activation functions for your counting network. Draw a diagram to show how you will use the counting network and give a set of weights and biases for the final layer to find the  $k^{th}$  smallest number. In other words, repeat the process of sections 1.1, 1.2, 1.3 using the counting idea.

*Hints: You may find the following two activation functions useful.*

1) *Hard threshold activation function:*

$$\phi(z) = \mathbb{I}(z \geq 0) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

2) *Indicator activation function:*

$$\phi(z) = \mathbb{I}(z \in [-1, 1]) = \begin{cases} 1 & \text{if } z \in [-1, 1] \\ 0 & \text{otherwise} \end{cases}$$

## 2 Backpropagation

This question helps you to understand the underlying mechanism of back-propagation. You need to have a clear understanding of what happens during the forward pass and backward pass and be able to reason about the time complexity and space complexity of your neural network. Moreover, you will learn a commonly used trick to compute the gradient norm efficiently without explicitly writing down the whole Jacobian matrix.

Note: The reading on backpropagation located at <https://uoft-csc413.github.io/2022/assets/readings/L02b.pdf> may be useful for this question.

### 2.1 Automatic Differentiation

Consider a neural network defined with the following procedure:

$$\begin{aligned} \mathbf{z}_1 &= \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)} \\ \mathbf{h}_1 &= \text{ReLU}(\mathbf{z}_1) \\ \mathbf{z}_2 &= \mathbf{W}^{(2)}\mathbf{x} + \mathbf{b}^{(2)} \\ \mathbf{h}_2 &= \sigma(\mathbf{z}_2) \\ \mathbf{g} &= \mathbf{h}_1 \circ \mathbf{h}_2 \\ \mathbf{y} &= \mathbf{W}^{(3)}\mathbf{g} + \mathbf{W}^{(4)}\mathbf{x}, \\ \mathbf{y}' &= \text{softmax}(\mathbf{y}) \\ \mathcal{S} &= \sum_{k=1}^N \mathbb{I}(t = k) \log(\mathbf{y}'_k) \\ \mathcal{J} &= -\mathcal{S} \end{aligned}$$

for input  $\mathbf{x}$  with class label  $t$  where  $\text{ReLU}(\mathbf{z}) = \max(\mathbf{z}, 0)$  denotes the ReLU activation function,  $\sigma(\mathbf{z}) = \frac{1}{1+e^{-\mathbf{z}}}$  denotes the Sigmoid activation function, both applied elementwise, and  $\text{softmax}(\mathbf{y}) = \frac{\exp(\mathbf{y})}{\sum_{i=1}^N \exp(\mathbf{y}_i)}$ . Here,  $\circ$  denotes element-wise multiplication.

#### 2.1.1 Computational Graph [0pt]

Draw the computation graph relating  $\mathbf{x}$ ,  $t$ ,  $\mathbf{z}_1$ ,  $\mathbf{z}_2$ ,  $\mathbf{h}_1$ ,  $\mathbf{h}_2$ ,  $\mathbf{g}$ ,  $\mathbf{y}$ ,  $\mathbf{y}'$ ,  $\mathcal{S}$  and  $\mathcal{J}$ .

**2.1.2 Backward Pass [1pt]**

Derive the backprop equations for computing  $\bar{\mathbf{x}} = \frac{\partial \mathcal{J}}{\partial \mathbf{x}}^\top$ , one variable at a time, similar to the vectorized backward pass derived in Lec 2.

*Hints: Be careful about the transpose and shape! Assume all vectors (including error vector) are column vector and all Jacobian matrices adopt numerator-layout notation<sup>4</sup>. You can use  $\text{softmax}'(\mathbf{y})$  for the Jacobian matrix of softmax.*

---

<sup>4</sup>Numerator-layout notation: [https://en.wikipedia.org/wiki/Matrix\\_calculus#Numerator-layout\\_notation](https://en.wikipedia.org/wiki/Matrix_calculus#Numerator-layout_notation)

**Answer**

$$\begin{aligned}
\overline{\mathcal{J}} &= 1 \\
\overline{\mathcal{S}} &= \frac{\partial \mathcal{J}^\top}{\partial \mathcal{S}} \overline{\mathcal{J}} \\
&= -\overline{\mathcal{J}} \\
\overline{\mathbf{y}'} &= \frac{\partial \mathcal{S}^\top}{\partial \mathbf{y}'} \overline{\mathcal{S}} \\
&= [0, \dots, 0, 1, 0, \dots, 0]^\top \circ [\dots, 0, \frac{1}{\mathbf{y}'_k}, 0, \dots, 0]^\top \overline{\mathcal{S}} \\
&= [\dots, 0, \frac{1}{\mathbf{y}'_k}, 0, \dots, 0]^\top \overline{\mathcal{S}} \\
&\text{where the } \frac{1}{\mathbf{y}'_k} \text{ is at the } k^{th} \text{ index.} \\
\overline{\mathbf{y}} &= \frac{\partial \mathbf{y}'^\top}{\partial \mathbf{y}} \overline{\mathbf{y}'} \\
&= \text{softmax}'(\mathbf{y})^\top \overline{\mathbf{y}'} \\
\overline{\mathbf{g}} &= \frac{\partial \mathbf{y}^\top}{\partial \mathbf{g}} \overline{\mathbf{y}} \\
&= [\mathbf{W}^{(3)}]^\top \overline{\mathbf{y}} \\
\overline{\mathbf{h}_1} &= \frac{\partial \mathbf{g}^\top}{\partial \mathbf{h}_1} \overline{\mathbf{g}} \\
&= \mathbf{h}_2 \circ \overline{\mathbf{g}} \\
\overline{\mathbf{h}_2} &= \frac{\partial \mathbf{g}^\top}{\partial \mathbf{h}_2} \overline{\mathbf{g}} \\
&= \mathbf{h}_1 \circ \overline{\mathbf{g}} \\
\overline{\mathbf{z}_1} &= \frac{\partial \mathbf{h}_1^\top}{\partial \mathbf{z}_1} \overline{\mathbf{h}_1} \\
&= \overline{\mathbf{h}_1} \circ \begin{cases} 0 & \mathbf{z}_{1i} < 0 \\ 1 & \mathbf{z}_{1i} > 0 \end{cases} \text{ for each } i \in \{1, \dots, K\} \\
\overline{\mathbf{z}_2} &= \frac{\partial \mathbf{h}_2^\top}{\partial \mathbf{z}_2} \overline{\mathbf{h}_2} \\
&= \overline{\mathbf{h}_2} \circ \sigma'(\mathbf{z}_2) \\
&= \overline{\mathbf{h}_2} \circ \sigma(\mathbf{z}_2) \circ (1 - \sigma(\mathbf{z}_2)) \\
\overline{\mathbf{x}} &= \frac{\partial \mathbf{z}_1^\top}{\partial \mathbf{x}} \overline{\mathbf{z}_1} + \frac{\partial \mathbf{z}_2^\top}{\partial \mathbf{x}} \overline{\mathbf{z}_2} + \frac{\partial \mathbf{y}^\top}{\partial \mathbf{x}} \overline{\mathbf{y}} \\
&= [\mathbf{W}^{(1)}]^\top \overline{\mathbf{z}_1} + [\mathbf{W}^{(2)}]^\top \overline{\mathbf{z}_2} + [\mathbf{W}^{(4)}]^\top \overline{\mathbf{y}}
\end{aligned}$$

## 2.2 Gradient Norm Computation

Many deep learning algorithms require you to compute the  $L^2$  norm of the gradient of a loss function with respect to the model parameters for every example in a minibatch. Unfortunately, most differentiation functionality provided by most software frameworks (Tensorflow, PyTorch) does not support computing gradients for individual samples in a minibatch. Instead, they only give one gradient per minibatch that aggregates individual gradients for you. A naive way to get the per-example gradient norm is to use a batch size of 1 and repeat the back-propagation  $N$  times, where  $N$  is the minibatch size. After that, you can compute the  $L^2$  norm of each gradient vector. As you can imagine, this approach is very inefficient. It can not exploit the parallelism of minibatch operations provided by the framework.

In this question, we will investigate a more efficient way to compute the per-example gradient norm and reason about its complexity compared to the naive method. For simplicity, let us consider the following two-layer neural network.

$$\begin{aligned}\mathbf{z} &= \mathbf{W}^{(1)}\mathbf{x} \\ \mathbf{h} &= \text{ReLU}(\mathbf{z}) \\ \mathbf{y} &= \mathbf{W}^{(2)}\mathbf{h},\end{aligned}$$

where  $\mathbf{W}^{(1)} = \begin{pmatrix} 1 & 2 & 1 \\ -2 & 1 & 0 \\ 1 & -2 & -1 \end{pmatrix}$  and  $\mathbf{W}^{(2)} = \begin{pmatrix} -2 & 4 & 1 \\ 1 & -2 & -3 \\ -3 & 4 & 6 \end{pmatrix}$ .

### 2.2.1 Naive Computation [1pt]

Let us assume the input  $x = (1 \ 3 \ 1)^\top$  and the error vector  $\bar{\mathbf{y}} = \frac{\partial \mathcal{J}}{\partial \mathbf{y}}^\top = (1 \ 1 \ 1)^\top$ . In this question, write down the Jacobian matrix (numerical value)  $\frac{\partial \mathcal{J}}{\partial \mathbf{W}^{(1)}}$  and  $\frac{\partial \mathcal{J}}{\partial \mathbf{W}^{(2)}}$  using back-propagation. Then, compute the square of Frobenius Norm of the two Jacobian matrices,  $\|A\|_F^2$ . The square of Frobenius norm of a matrix  $A$  is defined as follows:

$$\|A\|_F^2 = \sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 = \text{trace}(A^\top A)$$

*Hints: Be careful about the transpose. Show all your work for partial marks.*

**Answer**

$$\begin{aligned}\mathbf{x} &= (1 \ 3 \ 1)^\top, \mathbf{z} = (8 \ 1 \ -6)^\top, \mathbf{h} = (8 \ 1 \ 0)^\top \\ \bar{\mathbf{y}} &= (1 \ 1 \ 1)^\top, \bar{\mathbf{h}} = (-4 \ 6 \ 4)^\top, \bar{\mathbf{z}} = (-4 \ 6 \ 0)^\top \\ \frac{\partial \mathcal{J}}{\partial \mathbf{W}^{(1)}} &= \bar{\mathbf{z}}\mathbf{x}^\top = (-4 \ 6 \ 0)^\top (1 \ 3 \ 1) = \begin{pmatrix} -4 & -12 & -4 \\ 6 & 18 & 6 \\ 0 & 0 & 0 \end{pmatrix} \\ \frac{\partial \mathcal{J}}{\partial \mathbf{W}^{(2)}} &= \bar{\mathbf{y}}\mathbf{h}^\top = (1 \ 1 \ 1)^\top (8 \ 1 \ 0) = \begin{pmatrix} 8 & 1 & 0 \\ 8 & 1 & 0 \\ 8 & 1 & 0 \end{pmatrix}\end{aligned}$$



$$\left\| \frac{\partial \mathcal{J}}{\partial \mathbf{W}^{(1)}} \right\|_F^2 = 572$$

$$\left\| \frac{\partial \mathcal{J}}{\partial \mathbf{W}^{(2)}} \right\|_F^2 = 195$$

### 2.2.2 Efficient Computation [0.5pt]

Notice that weight Jacobian can be expressed as the outer product of the error vector and activation  $\frac{\partial \mathcal{J}}{\partial \mathbf{W}^{(1)}} = (\bar{\mathbf{z}} \mathbf{x}^\top)^\top$  and  $\frac{\partial \mathcal{J}}{\partial \mathbf{W}^{(2)}} = (\bar{\mathbf{y}} \mathbf{h}^\top)^\top$ . We can compute the Jacobian norm more efficiently using the following trick:

$$\begin{aligned} \left\| \frac{\partial \mathcal{J}}{\partial \mathbf{W}^{(1)}} \right\|_F^2 &= \text{trace} \left( \frac{\partial \mathcal{J}}{\partial \mathbf{W}^{(1)}}^\top \frac{\partial \mathcal{J}}{\partial \mathbf{W}^{(1)}} \right) && \text{(Definition)} \\ &= \text{trace} \left( \bar{\mathbf{z}} \mathbf{x}^\top \mathbf{x} \bar{\mathbf{z}}^\top \right) \\ &= \text{trace} \left( \mathbf{x}^\top \mathbf{x} \bar{\mathbf{z}}^\top \bar{\mathbf{z}} \right) && \text{(Cyclic Property of Trace)} \\ &= \left( \mathbf{x}^\top \mathbf{x} \right) \left( \bar{\mathbf{z}}^\top \bar{\mathbf{z}} \right) && \text{(Scalar Multiplication)} \\ &= \|\mathbf{x}\|_2^2 \|\bar{\mathbf{z}}\|_2^2 \end{aligned}$$

Compute the **square** of the Frobenius Norm of the two Jacobian matrices by plugging the value into the above trick.

*Hints: Verify the solution is the same as naive computation. Show all your work for partial marks.*

**Answer**

$$\left\| \frac{\partial \mathcal{J}}{\partial \mathbf{W}^{(1)}} \right\|_F^2 = 572$$

$$\left\| \frac{\partial \mathcal{J}}{\partial \mathbf{W}^{(2)}} \right\|_F^2 = 195$$

### 2.2.3 Complexity Analysis [1.5pt]

Now, let us consider a general neural network with  $K - 1$  hidden layers ( $K$  weight matrices). All input units, output units, and hidden units have a dimension of  $D$ . Assume we have  $N$  input vectors. How many scalar multiplications  $T$  (integer) do we need to compute the per-example gradient norm using naive and efficient computation, respectively? And, what is the memory cost  $M$  (big  $\mathcal{O}$  notation)?

For simplicity, you can ignore the activation function and loss function computation. You can assume the network does not have a bias term. You can also assume there are no in-place operations. Please fill up the table below.

	T (Naive)	T (Efficient)	M (Naive)	M (Efficient)
Forward Pass				
Backward Pass				
Gradient Norm Computation				

*Hints: The forward pass computes all the activations and needs memory to store model parameters and activations. The backward pass computes all the error vectors. Moreover, you also need to*

compute the parameter's gradient in naive computation. During the Gradient Norm Computation, the naive method needs to square the gradient before aggregation. In contrast, the efficient method relies on the trick. Thinking about the following questions may be helpful. 1) Do we need to store all activations in the forward pass? 2) Do we need to store all error vectors in the backward pass? 3) Why standard backward pass is twice more expensive than the forward pass? Don't forget to consider  $K$  and  $N$  in your answer.

### Answer

	T (Naive)	T (Efficient)	M (Naive)	M (Efficient)
Forward Pass	$NKD^2$	$NKD^2$	$\mathcal{O}(KD^2 + NKD)$	$\mathcal{O}(KD^2 + NKD)$
Backward Pass	$2NKD^2$	$NKD^2$	$\mathcal{O}(NKD^2)$	$\mathcal{O}(KD^2 + NKD)$
Gradient Norm Computation	$NKD^2$	$NK(2D + 1)$	$\mathcal{O}(NKD^2)$	$\mathcal{O}(NKD)$

1) **Forward Pass:** Both the naive and efficient methods need to compute and store all activations. For each training input, we will compute  $K$  activations. Each activation is computed by a matrix-vector product of a  $D \times D$  matrix and a  $D$ -dimension vector, resulting in  $D^2$  scalar multiplications. Thus, for a  $K$ -layer network and  $N$  training inputs, we will have  $NKD^2$  scalar multiplications. As for memory, we need  $KD^2$  (as there are  $K$   $D \times D$  weight matrices) for parameters and  $NKD$  for  $NK$  activations each of size  $D$ .

2) **Backward Pass:** Both the naive and efficient methods need to compute all activation vectors. Moreover, the naive method needs to construct the gradient vector explicitly using the outer product of activation and error vector. This extra computation is the cause of the discrepancy between the number of scalar multiplications ( $2NKD^2$  vs.  $NKD^2$ ) between the backward passes for the naive and efficient methods. Computing the error vector is very similar to what we do in the forward pass. Each error vector is computed by a matrix-vector product of a  $D \times D$  matrix and a  $D$ -dimension vector, resulting in  $D^2$  scalar multiplications. Thus, for a  $K$ -layer network and  $N$  training inputs, we will have  $NKD^2$  scalar multiplications. As for the memory costs, both methods need to store the activations and error vectors, which takes  $\mathcal{O}(NKD)$  memory. Both methods also need to store the model parameters, which, as mentioned earlier, is  $\mathcal{O}(KD^2)$ . This results in a memory cost of  $\mathcal{O}(NKD + KD^2)$  for the efficient method. However, the naive method needs to store all the parameters' gradients, which takes  $\mathcal{O}(NKD^2)$ , which dominates all the other terms so far. This means  $M$  for the backward pass in the naive method is  $\mathcal{O}(NKD^2)$ .

3) **Gradient Norm Computation:** To calculate the per-example gradient norm using the naive method, we would need to square all the gradients. This would take  $D^2$  scalar multiplications for each of the Jacobian matrices (there are  $NK$  of them), resulting in a total of  $NKD^2$ . Furthermore, since there are  $NK$   $D \times D$  matrices to deal with, the memory cost for this step in the naive method is  $\mathcal{O}(NKD^2)$ . In contrast, the efficient method determines the Frobenius norm of each weight Jacobian by multiplying the norms of two  $D$ -dimensional vectors. Calculating the norm of a  $D$ -dimensional vector takes  $D$  scalar multiplications; doing this twice and multiplying the norms together would therefore take  $2D + 1$  multiplications. This results in  $NK(2D + 1)$  scalar multiplications in total. Finally, the gradient norm computation in the efficient method deals with  $NK$  vectors, each of which is  $D$ -dimensional, resulting in a memory cost of  $\mathcal{O}(NKD)$ .

4) **Alternative acceptable answers:** In forward pass, the above answer assumes we compute all the activations including the final output, which is unnecessary for this question as we provide the final error vector. Therefore,  $N(K - 1)D^2$  is also correct. In backward pass, the above answer assumes we compute all error vectors including the error vector for the input, which is unnecessary

for this question. Therefore,  $N(2K - 1)D^2$  and  $N(K - 1)D^2$  are also correct.

### 2.3 Inner product of Jacobian: JVP and VJP [0pt]

A more general case of computing the gradient norm is to compute the inner product of the Jacobian matrices computed using two different examples. Let  $f_1, f_2$  and  $y_1, y_2$  be the final outputs and layer outputs of two different examples respectively. The inner product  $\Theta$  of Jacobian matrices of layer parameterized by  $\theta$  is defined as:

$$\Theta_{\theta}(f_1, f_2) := \frac{\partial f_1}{\partial \theta} \frac{\partial f_2}{\partial \theta}^{\top} = \frac{\partial f_1}{\partial y_1} \frac{\partial y_1}{\partial \theta} \frac{\partial y_2}{\partial \theta}^{\top} \frac{\partial f_2}{\partial y_2}^{\top} = \underbrace{\frac{\partial f_1}{\partial y_1}}_{\mathbf{O} \times \mathbf{Y}} \underbrace{\frac{\partial y_1}{\partial \theta}}_{\mathbf{Y} \times \mathbf{P}} \underbrace{\frac{\partial y_2}{\partial \theta}^{\top}}_{\mathbf{P} \times \mathbf{Y}} \underbrace{\frac{\partial f_2}{\partial y_2}^{\top}}_{\mathbf{Y} \times \mathbf{O}},$$

Where  $\mathbf{O}, \mathbf{Y}, \mathbf{P}$  represent the dimension of the final output, layer output, model parameter respectively. How to formulate the above computation using Jacobian Vector Product (JVP) and Vector Jacobian Product (VJP)? What are the computation cost using the following three ways of contracting the above equation?

- (a) Outside-in:  $M_1 M_2 M_3 M_4 = ((M_1 M_2)(M_3 M_4))$
- (b) Left-to-right and right-to-left:  $M_1 M_2 M_3 M_4 = (((M_1 M_2)(M_3)M_4) = (M_1(M_2(M_3 M_4))))$
- (c) Inside-out-left and inside-out-right:  $M_1 M_2 M_3 M_4 = ((M_1(M_2 M_3))M_4) = (M_1((M_2 M_3)M_4))$

## 3 Linear Regression

The reading on linear regression located at <https://uoft-csc413.github.io/2022/assets/readings/L01a.pdf> may be useful for this question.

Given  $n$  pairs of input data with  $d$  features and scalar label  $(\mathbf{x}_i, t_i) \in \mathbb{R}^d \times \mathbb{R}$ , we wish to find a linear model  $f(\mathbf{x}) = \hat{\mathbf{w}}^{\top} \mathbf{x}$  with  $\hat{\mathbf{w}} \in \mathbb{R}^d$  that minimizes the squared error of prediction on the training samples defined below. This is known as an empirical risk minimizer. For concise notation, denote the data matrix  $X \in \mathbb{R}^{n \times d}$  and the corresponding label vector  $\mathbf{t} \in \mathbb{R}^n$ . The training objective is to minimize the following loss:

$$\min_{\hat{\mathbf{w}}} \frac{1}{n} \sum_{i=1}^n (\hat{\mathbf{w}}^{\top} \mathbf{x}_i - t_i)^2 = \min_{\hat{\mathbf{w}}} \frac{1}{n} \|X \hat{\mathbf{w}} - \mathbf{t}\|_2^2.$$

We assume  $X$  is full rank:  $X^{\top} X$  is invertible when  $n > d$ , and  $X X^{\top}$  is invertible otherwise. Note that when  $d > n$ , the problem is *underdetermined*, i.e. there are less training samples than parameters to be learned. This is analogous to learning an *overparameterized* model, which is common when training of deep neural networks.

### 3.1 Deriving the Gradient [0pt]

Write down the gradient of the loss w.r.t. the learned parameter vector  $\hat{\mathbf{w}}$ .

### 3.2 Underparameterized Model

#### 3.2.1 [0.5pt]

First consider the underparameterized  $d < n$  case. Show that the solution obtained by gradient descent is  $\hat{\mathbf{w}} = (X^{\top} X)^{-1} X^{\top} \mathbf{t}$ , assuming training converges. Show your work.

**Answer**

$$\begin{aligned}
\frac{2}{n} X^\top (X \hat{\mathbf{w}} - \mathbf{t}) &= 0; \\
\Leftrightarrow X^\top X \hat{\mathbf{w}} &= X^\top \mathbf{t}; \\
\Leftrightarrow \hat{\mathbf{w}} &= (X^\top X)^{-1} X^\top \mathbf{t},
\end{aligned}$$

where we used the invertibility of  $X^\top X$  in the last step. This solution is unique.

**3.2.2 [0.5pt]**

Now consider the case of noisy linear regression. The training labels  $t_i = \mathbf{w}^* \top \mathbf{x}_i + \epsilon_i$  are generated by a ground truth linear target function, where the noise term,  $\epsilon_i$ , is generated independently with zero mean and variance  $\sigma^2$ . The final training error can be derived as a function of  $n$  and  $\epsilon$ , as:

$$Error = \frac{1}{n} \|(X(X^\top X)^{-1} X^\top - I) \epsilon\|_2^2,$$

Show this is true by substituting your answer from 3.2.1 into  $\frac{1}{n} \|X \hat{\mathbf{w}} - \mathbf{t}\|_2^2$ . Also, find the expectation of the above training error in terms of  $n, d$  and  $\sigma$ .

*Hints: you might find the cyclic property<sup>5</sup> of trace useful.*

**Answer**

Starting with the loss function, and substituting for  $\mathbf{t}$ , in vector form, with  $X \mathbf{w}^* + \epsilon$ , we have:

$$\begin{aligned}
&\frac{1}{n} \|X \hat{\mathbf{w}} - \mathbf{t}\|_2^2; \\
&\Leftrightarrow \frac{1}{n} \|X(X^\top X)^{-1} X^\top \mathbf{t} - \mathbf{t}\|_2^2; \\
&\Leftrightarrow \frac{1}{n} \|(X(X^\top X)^{-1} X^\top - I) \mathbf{t}\|_2^2; \\
&\Leftrightarrow \frac{1}{n} \|(X(X^\top X)^{-1} X^\top - I)(X \mathbf{w}^* + \epsilon)\|_2^2; \\
&\Leftrightarrow \frac{1}{n} \|X(X^\top X)^{-1} X^\top X \mathbf{w}^* - X \mathbf{w}^* + X(X^\top X)^{-1} X^\top \epsilon - \epsilon\|_2^2; \\
&\Leftrightarrow \frac{1}{n} \|X(X^\top X)^{-1} X^\top \epsilon - \epsilon\|_2^2,
\end{aligned}$$

from which the first result directly follows. For the expectation, expanding the last expression, we have:

$$\begin{aligned}
&\frac{1}{n} \|X(X^\top X)^{-1} X^\top \epsilon - \epsilon\|_2^2; \\
&\Leftrightarrow \frac{1}{n} (X(X^\top X)^{-1} X^\top \epsilon - \epsilon)^\top (X(X^\top X)^{-1} X^\top \epsilon - \epsilon); \\
&\Leftrightarrow \frac{1}{n} (\epsilon \epsilon^\top - \epsilon^\top X(X^\top X)^{-1} X^\top \epsilon),
\end{aligned}$$

<sup>5</sup>[https://en.wikipedia.org/wiki/Trace\\_\(linear\\_algebra\)#Cyclic\\_property](https://en.wikipedia.org/wiki/Trace_(linear_algebra)#Cyclic_property)

and applying the trace operator (to a scalar) and expectation operator, we have:

$$\begin{aligned}
 \mathbb{E}[\text{Error}] &= \frac{1}{n} \mathbb{E}[\text{Tr}(\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T - \boldsymbol{\epsilon}^T X (X^T X)^{-1} X^T \boldsymbol{\epsilon})] \\
 &= \frac{1}{n} \text{Tr}(\mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T]) - \frac{1}{n} \text{Tr}(\mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T] X (X^T X)^{-1} X^T) \\
 &= \frac{1}{n} (n\sigma^2) - \frac{\sigma^2}{n} \text{Tr}(X (X^T X)^{-1} X^T) \\
 &= \sigma^2 - \frac{\sigma^2}{n} \text{Tr}(I_d) \\
 &= \sigma^2 \left(1 - \frac{d}{n}\right),
 \end{aligned}$$

noting that  $\mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T] = \sigma^2 I_n$  since noise terms are independent.

### 3.3 Overparameterized Model

#### 3.3.1 [0pt]

Now consider the overparameterized  $d > n$  case. We first illustrate that there exist multiple empirical risk minimizers. For simplicity we let  $n = 1$  and  $d = 2$ . Choose  $\mathbf{x}_1 = [1; 1]$  and  $t_1 = 3$ , i.e. the one data point and all possible  $\hat{\mathbf{w}}$  lie on a 2D plane. Show that there exists infinitely many  $\hat{\mathbf{w}}$  satisfying  $\hat{\mathbf{w}}^T \mathbf{x}_1 = y_1$  on a real line. Write down the equation of the line.

**Answer**

$$\hat{\mathbf{w}}^T \mathbf{x}_1 = t_1 \Rightarrow \hat{w}_2 = -\hat{w}_1 + 3$$

Every  $\hat{\mathbf{w}}$  on the line specified above achieves zero training error.

#### 3.3.2 [1pt]

Now, let's generalize the previous 2D case to the general  $d > n$ . Show that gradient descent from zero initialization i.e.  $\hat{\mathbf{w}}(0) = 0$  finds a unique minimizer if it converges. Show that the solution by gradient descent is  $\hat{\mathbf{w}} = X^T (X X^T)^{-1} \mathbf{t}$ . Show your work.

*Hints: You can assume that the gradient is spanned by the rows of  $X$  and write  $\hat{\mathbf{w}} = X^T \mathbf{a}$  for some  $\mathbf{a} \in \mathbb{R}^n$ .*

**Answer**

Since the gradient is always spanned by the rows of  $X$ , starting from zero initialization, every iterate and thus the final solution of gradient descent can also be written as a linear combination of rows of  $X$ . We can thus write  $\hat{\mathbf{w}} = X^T \mathbf{a}$  for some  $\mathbf{a} \in \mathbb{R}^n$ .

From the stationary condition, setting the gradient to zero we have:

$$\begin{aligned}
 X^T (X \hat{\mathbf{w}} - \mathbf{t}) &= 0; \\
 \Leftrightarrow X \hat{\mathbf{w}} - \mathbf{t} &= 0,
 \end{aligned}$$

since  $X$  is full rank by assumption. Plugging this into the stationary condition yields:

$$\begin{aligned} X\hat{\mathbf{w}} - \mathbf{t} &= XX^\top \mathbf{a} - \mathbf{t} = 0; \\ \Leftrightarrow \mathbf{a} &= (XX^\top)^{-1}\mathbf{t}; \\ \Leftrightarrow \hat{\mathbf{w}} &= X^\top (XX^\top)^{-1}\mathbf{t}, \end{aligned}$$

in which we utilized the invertibility of  $XX^\top$  when  $d > n$ . Due to this and iterative update property, the solution is unique.

### 3.3.3 [0pt]

Repeat part 3.2.2 for the overparameterized case.

### 3.3.4 [0.5pt]

Visualize and compare underparameterized with overparameterized polynomial regression: [https://colab.research.google.com/github/uoft-csc413/2022/blob/master/assets/assignments/LS\\_polynomial\\_regression.ipynb](https://colab.research.google.com/github/uoft-csc413/2022/blob/master/assets/assignments/LS_polynomial_regression.ipynb) Include your code snippets for the `fit_poly` function in the write-up. Does overparameterization (higher degree polynomial) always lead to overfitting, i.e. larger test error?

#### Answer

Overparameterization does not always lead to worse test error. Instead the test error peaks when degree of freedom of the model is roughly the same as the number of training samples. This surprising peak, known as *double descent* <https://openai.com/blog/deep-double-descent>, requires us to rethink the conventional view of bias-variance tradeoff in machine learning models.

### 3.3.5 [0pt]

Give  $n_1, n_2$  with  $n_1 \leq n_2$ , and fixed dimension  $d$  for which  $L_2 \geq L_1$ , i.e. the loss with  $n_2$  data points is greater than loss with  $n_1$  data points. Explain the underlying phenomenon. Be sure to also include the error values  $L_1$  and  $L_2$  or provide visualization in your solution.

*Hints: use your code to experiment with relevant parameters, then vary to find region and report one such setting.*

