

**Due Date : February 16th, 2019**

Instructions

- For all questions, show your work!
- Use a document preparation system such as LaTeX.
- Submit your answers electronically via the course studium page, and via Gradescope.

**Question 1.** Using the following definition of the derivative and the definition of the Heaviside step function :

$$\frac{d}{dx}f(x) = \lim_{\epsilon \rightarrow 0} \frac{f(x + \epsilon) - f(x)}{\epsilon} \quad H(x) = \begin{cases} 1 & \text{if } x > 0 \\ \frac{1}{2} & \text{if } x = 0 \\ 0 & \text{if } x < 0 \end{cases}$$

1. Show that the derivative of the rectified linear unit  $g(x) = \max\{0, x\}$ , **wherever it exists**, is equal to the Heaviside step function.
2. Give two alternative definitions of  $g(x)$  using  $H(x)$ .
3. Show that  $H(x)$  can be well approximated by the sigmoid function  $\sigma(x) = \frac{1}{1+e^{-kx}}$  asymptotically (i.e for large  $k$ ), where  $k$  is a parameter.
- \*4. Although the Heaviside step function is not differentiable, we can define its **distributional derivative**. For a function  $F$ , consider the functional  $F[\phi] = \int_{\mathbb{R}} H(x)\phi(x)dx$ , where  $\phi$  is a smooth function (infinitely differentiable) with compact support ( $\phi(x) = 0$  whenever  $|x| \geq A$ , for some  $A > 0$ ).

Show that whenever  $F$  is differentiable,  $F'[\phi] = -\int_{\mathbb{R}} F(x)\phi'(x)dx$ . Using this formula as a definition in the case of non-differentiable functions, show that  $H'[\phi] = \phi(0)$ . ( $\delta[\phi] \doteq \phi(0)$  is known as the Dirac delta function.)

**Answer 1.** Write your answer here.

**Question 2.** Let  $\mathbf{x}$  be an  $n$ -dimensional vector. Recall the softmax function :  $S : \mathbf{x} \in \mathbb{R}^n \mapsto S(\mathbf{x}) \in \mathbb{R}^n$  such that  $S(\mathbf{x})_i = \frac{e^{\mathbf{x}_i}}{\sum_j e^{\mathbf{x}_j}}$ ; the diagonal function :  $\text{diag}(\mathbf{x})_{ij} = \mathbf{x}_i$  if  $i = j$  and  $\text{diag}(\mathbf{x})_{ij} = 0$  if  $i \neq j$ ; and the Kronecker delta function :  $\delta_{ij} = 1$  if  $i = j$  and  $\delta_{ij} = 0$  if  $i \neq j$ .

1. Show that the derivative of the softmax function is  $\frac{dS(\mathbf{x})_i}{d\mathbf{x}_j} = S(\mathbf{x})_i (\delta_{ij} - S(\mathbf{x})_j)$ .
2. Express the Jacobian matrix  $\frac{\partial S(\mathbf{x})}{\partial \mathbf{x}}$  using matrix-vector notation. Use  $\text{diag}(\cdot)$ .
3. Compute the Jacobian of the sigmoid function  $\sigma(\mathbf{x}) = 1/(1 + e^{-\mathbf{x}})$ .
4. Let  $\mathbf{y}$  and  $\mathbf{x}$  be  $n$ -dimensional vectors related by  $\mathbf{y} = f(\mathbf{x})$ ,  $L$  be an unspecified differentiable loss function. According to the chain rule of calculus,  $\nabla_{\mathbf{x}} L = (\frac{\partial \mathbf{y}}{\partial \mathbf{x}})^{\top} \nabla_{\mathbf{y}} L$ , which takes up  $\mathcal{O}(n^2)$  computational time in general. Show that if  $f(\mathbf{x}) = \sigma(\mathbf{x})$  or  $f(\mathbf{x}) = S(\mathbf{x})$ , the above matrix-vector multiplication can be simplified to a  $\mathcal{O}(n)$  operation.

**Answer 2.** Write your answer here.

**Question 3.** Recall the definition of the softmax function :  $S(\mathbf{x})_i = e^{\mathbf{x}_i} / \sum_j e^{\mathbf{x}_j}$ .

1. Show that softmax is translation-invariant, that is :  $S(\mathbf{x} + c) = S(\mathbf{x})$ , where  $c$  is a scalar constant.

2. Show that softmax is not invariant under scalar multiplication. Let  $S_c(\mathbf{x}) = S(c\mathbf{x})$  where  $c \geq 0$ . What are the effects of taking  $c$  to be 0 and arbitrarily large?
3. Let  $\mathbf{x}$  be a 2-dimensional vector. One can represent a 2-class categorical probability using softmax  $S(\mathbf{x})$ . Show that  $S(\mathbf{x})$  can be reparameterized using sigmoid function, i.e.  $S(\mathbf{x}) = [\sigma(z), 1 - \sigma(z)]^\top$  where  $z$  is a scalar function of  $\mathbf{x}$ .
4. Let  $\mathbf{x}$  be a  $K$ -dimensional vector ( $K \geq 2$ ). Show that  $S(\mathbf{x})$  can be represented using  $K - 1$  parameters, i.e.  $S(\mathbf{x}) = S([0, y_1, y_2, \dots, y_{K-1}]^\top)$  where  $y_i$  is a scalar function of  $\mathbf{x}$  for  $i \in \{1, \dots, K - 1\}$ .

**Answer 3.** Write your answer here.

**Question 4.** Consider a 2-layer neural network  $y : \mathbb{R}^D \rightarrow \mathbb{R}^K$  of the form :

$$y(x, \Theta, \sigma)_k = \sum_{j=1}^M \omega_{kj}^{(2)} \sigma \left( \sum_{i=1}^D \omega_{ji}^{(1)} x_i + \omega_{j0}^{(1)} \right) + \omega_{k0}^{(2)}$$

for  $1 \leq k \leq K$ , with parameters  $\Theta = (\omega^{(1)}, \omega^{(2)})$  and logistic sigmoid activation function  $\sigma$ . Show that there exists an equivalent network of the same form, with parameters  $\Theta' = (\tilde{\omega}^{(1)}, \tilde{\omega}^{(2)})$  and tanh activation function, such that  $y(x, \Theta', \tanh) = y(x, \Theta, \sigma)$  for all  $x \in \mathbb{R}^D$ , and express  $\Theta'$  as a function of  $\Theta$ .

**Answer 4.** Write your answer here.

**Question 5.** Given  $N \in \mathbb{Z}^+$ , we want to show that for any  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and any sample set  $\mathcal{S} \subset \mathbb{R}^n$  of size  $N$ , there is a set of parameters for a two-layer network such that the output  $y(\mathbf{x})$  matches  $f(\mathbf{x})$  for all  $\mathbf{x} \in \mathcal{S}$ . That is, we want to interpolate  $f$  with  $y$  on any finite set of samples  $\mathcal{S}$ .

1. Write the generic form of the function  $y : \mathbb{R}^n \rightarrow \mathbb{R}^m$  defined by a 2-layer network with  $N - 1$  hidden units, with linear output and activation function  $\phi$ , in `te(rmsofitsweightsandbiases(W(1), b(1))` and `(W(2), b(2))`.
2. In what follows, we will restrict  $\mathbf{W}^{(1)}$  to be  $\mathbf{W}^{(1)} = [\mathbf{w}, \dots, \mathbf{w}]^T$  for some  $\mathbf{w} \in \mathbb{R}^n$  (so the rows of  $\mathbf{W}^{(1)}$  are all the same). Show that the interpolation problem on the sample set  $\mathcal{S} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\} \subset \mathbb{R}^n$  can be reduced to solving a matrix equation :  $\mathbf{M}\tilde{\mathbf{W}}^{(2)} = \mathbf{F}$ , where  $\tilde{\mathbf{W}}^{(2)}$  and  $\mathbf{F}$  are both  $N \times m$ , given by

$$\tilde{\mathbf{W}}^{(2)} = [\mathbf{W}^{(2)}, \mathbf{b}^{(2)}]^\top \quad \mathbf{F} = [f(\mathbf{x}^{(1)}), \dots, f(\mathbf{x}^{(N)})]^\top$$

Express the  $N \times N$  matrix  $\mathbf{M}$  in terms of  $\mathbf{w}$ ,  $\mathbf{b}^{(1)}$ ,  $\phi$  and  $\mathbf{x}^{(i)}$ .

- \*3. **Proof with Relu activation.** Assume  $\mathbf{x}^{(i)}$  are all distinct. Choose  $\mathbf{w}$  such that  $\mathbf{w}^\top \mathbf{x}^{(i)}$  are also all distinct (Try to prove the existence of such a  $\mathbf{w}$ , although this is not required for the assignment - See Assignment 0). Set  $\mathbf{b}_j^{(1)} = -\mathbf{w}^\top \mathbf{x}^{(j)} + \epsilon$ , where  $\epsilon > 0$ . Find a value of  $\epsilon$  such that  $\mathbf{M}$  is triangular with non-zero diagonal elements. Conclude. (Hint : assume an ordering of  $\mathbf{w}^\top \mathbf{x}^{(i)}$ .)
- \*4. **Proof with sigmoid-like activations.** Assume  $\phi$  is continuous, bounded,  $\phi(-\infty) = 0$  and  $\phi(0) > 0$ . Decompose  $\mathbf{w}$  as  $\mathbf{w} = \lambda \mathbf{u}$ . Set  $\mathbf{b}_j^{(1)} = -\lambda \mathbf{u}^\top \mathbf{x}^{(j)}$ . Fixing  $\mathbf{u}$ , show that  $\lim_{\lambda \rightarrow +\infty} \mathbf{M}$  is triangular with non-zero diagonal elements. Conclude. (Note that doing so preserves the distinctness of  $\mathbf{w}^\top \mathbf{x}^{(i)}$ .)

**Answer 5.** Write your answer here.

**Question 6.** Compute the *full*, *valid*, and *same* convolution (with kernel flipping) for the following 1D matrices :  $[1, 2, 3, 4] * [1, 0, 2]$

**Answer 6.** Write your answer here.

**Question 7.** Consider a convolutional neural network. Assume the input is a colorful image of size  $256 \times 256$  in the RGB representation. The first layer convolves 64  $8 \times 8$  kernels with the input, using a stride of 2 and no padding. The second layer downsamples the output of the first layer with a  $5 \times 5$  non-overlapping max pooling. The third layer convolves 128  $4 \times 4$  kernels with a stride of 1 and a zero-padding of size 1 on each border.

1. What is the dimensionality (scalar) of the output of the last layer ?
2. Not including the biases, how many parameters are needed for the last layer ?

**Answer 7.** Write your answer here.

**Question 8.** Assume we are given data of size  $3 \times 64 \times 64$ . In what follows, provide the correct configuration of a convolutional neural network layer that satisfies the specified assumption. Answer with the window size of kernel ( $k$ ), stride ( $s$ ), padding ( $p$ ), and dilation ( $d$ , with convention  $d = 0$  for no dilation). Use square windows only (e.g. same  $k$  for both width and height).

1. The output shape of the first layer is  $(64, 32, 32)$ .
  - (a) Assume  $k = 8$  without dilation.
  - (b) Assume  $d = 6$ , and  $s = 2$ .
2. The output shape of the second layer is  $(64, 8, 8)$ . Assume  $p = 0$  and  $d = 0$ .
  - (a) Specify  $k$  and  $s$  for pooling with non-overlapping window.
  - (b) What is output shape if  $k = 8$  and  $s = 4$  instead ?
3. The output shape of the last layer is  $(128, 4, 4)$ .
  - (a) Assume we are not using padding or dilation.
  - (b) Assume  $d = 1$ ,  $p = 2$ .
  - (c) Assume  $p = 1$ ,  $d = 0$ .

**Answer 8.** Write your answer here.