

# CMSC5743 2021F Homework 1

**Due:** Oct. 14, 2021

All solutions should be submitted to the blackboard in the format of **PDF/MS Word**.

## Q1 (12%)

- (a) (4%) We provide a very simple neural network as shown in Figure 1, please calculate the result in the blank neuron.
- (b) (4%) If we choose to prune one weight, which weight do you choose to achieve the best result? What's your evaluation metric?
- (c) (4%) If you have a chance to prune any weights, what's your pruning plan to make a better tradeoff between accuracy and the number of weights?

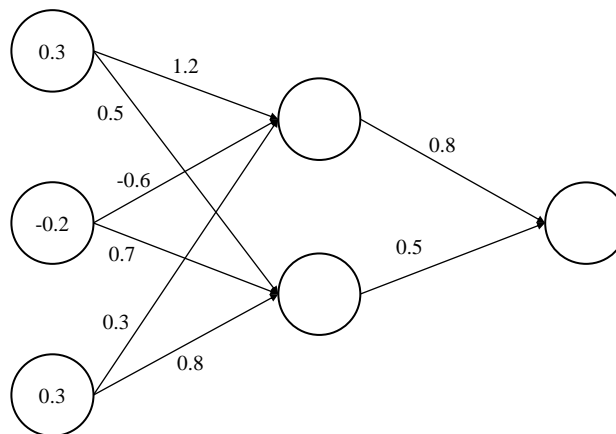


Figure 1: A simple 2-layer neural network

## Q2 (13%)

- (a) (5%) Some people may argue that why we do not simply train a smaller neural network instead of pruning a neural network with huge amount of parameters. Have you thought about this problem? What are the advantages of network pruning over training a smaller network? Please list two points and provide as much support as possible.
- (b) (4%) If someone want to apply structured pruning of fixed proportion for each layer, is it necessary?
- (c) (4%) If someone want to apply unstructured pruning of unfixed proportion for each layer, is it necessary?

## Q3 (12%) **Regularization.** In this question, you are going to solve a toy problem. Consider a function $J(x) = (x - 2)^2, x \in \mathbb{R}$ .

- (a) (3%) Find the global minimum of function  $J(x) + 6x^2$ . Justify your answer.

- (b) (3%) Find the global minimum of function  $J(x) + 6|x|$ . Justify your answer.  
 (c) (3%) Consider the following optimization problem.

$$\min_{x \in \mathbb{R}} J(\alpha; x) = (x - 2)^2 + \alpha|x|.$$

How should we determine  $\alpha$  so that the minimizer is at  $x = 0$ ?

- (d) (3%) How do you get inspired from the above questions about  $\ell_1$ ,  $\ell_2$  and sparsity? Please explain *briefly*.

**Q4** (13%)  $\ell_0$ -norm. Consider the  $p$ -norm (or  $\ell_p$ -norm) of a vector  $\mathbf{x} = [x_1, \dots, x_n]^\top$

$$\|\mathbf{x}\|_p = \sqrt[p]{|x_1|^p + |x_2|^p + \dots + |x_n|^p}, \quad (1)$$

where integer  $p \geq 1$  and the dimension  $n$  is fixed.

- (a) (4%) If we have a vector  $\mathbf{x}$  whose  $\ell_2$ -norm  $\|\mathbf{x}\|_2 \leq 1$ , will its  $\ell_1$ -norm  $\|\mathbf{x}\|_1$  be bounded? Justify your answer.  
 (b) (4%) Generalize (1) by letting  $p$  be any positive number. Show that the following limit exists, and give the result.

$$\lim_{p \rightarrow 0^+} |x_1|^p + |x_2|^p + \dots + |x_n|^p.$$

- (c) (5%) A vector norm function  $f(\mathbf{x})$  must satisfy *absolute homogeneity*, that is, for any scalar  $\alpha$  and vector  $\mathbf{x}$ , we must have  $f(\alpha\mathbf{x}) = |\alpha|f(\mathbf{x})$ . Can the result of (b) be a proper vector norm? Justify your answer.

**Q5** (12%)

- (a) (3%) A concrete formulation is shown as follows.

$$\min_{\beta_1, \beta_2, \beta_3} \left\| \begin{bmatrix} 1 \\ 2 \end{bmatrix} - \begin{bmatrix} 3 & 4 & 5 \\ 6 & 7 & 8 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} \right\|_2^2 + 3(|\beta_1| + |\beta_2| + |\beta_3|) + 8(\beta_1^2 + \beta_2^2 + \beta_3^2).$$

Let  $\beta'_1 = 3\beta_1$ ,  $\beta'_2 = 3\beta_2$ ,  $\beta'_3 = 3\beta_3$ , please transfer the above formulation as an equivalent formulation with respect to  $\beta'_1$ ,  $\beta'_2$  and  $\beta'_3$ .

- (b) (6%) Considering a more general formulation as follows.

$$\min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda_1\|\boldsymbol{\beta}\|_1 + \lambda_2\|\boldsymbol{\beta}\|_2^2.$$

Try to explain or prove how this formulation can be converted into an equivalent LASSO problem with  $\lambda_1$  and  $\lambda_2$  positive numbers.

- (c) (3%) Compare the formulation in (b) and typical LASSO formulation by discussing the advantages and disadvantages.

**Q6** (13%)

(a) (3%) Consider the typical LASSO formulation

$$\min_{\beta_1, \beta_2, \beta_3} \left\| \begin{bmatrix} 1 \\ 2 \end{bmatrix} - \begin{bmatrix} 3 & 4 & 5 \\ 6 & 7 & 8 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} \right\|_2^2 + 3(|\beta_1| + |\beta_2| + |\beta_3|).$$

Let  $\beta'_1 = 2\beta_1$ ,  $\beta'_2 = 3\beta_2$  and  $\beta'_3 = 4\beta_3$ , please transfer the above formulation as an equivalent formulation with respect to  $\beta'_1$ ,  $\beta'_2$  and  $\beta'_3$ .

(b) (7%) Consider a formulation as follows.

$$\min_{\beta} \left( \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \sum_{i=1}^p w_i |\beta_i| \right),$$

where  $\beta = [\beta_1, \beta_2, \dots, \beta_p]^\top$  and  $w_i > 0$ . We expect to apply same algorithms for solving LASSO problems to handle above object function. Considering that, try to convert the above formulation into the standard LASSO problem (i.e.,  $\min_{\beta} (\|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda_1 \|\beta\|_1)$ ) under a general assumption that  $w_k \neq w_j$ , if  $k \neq j$ .

(c) (3%) Compared with the typical LASSO formulation, what are advantages for the formulation in (b)?

**Q7** (10%) Convolution is the most important operation in CNN. As shown in Figure 2, the input activation tensor is  $\mathcal{X} \in \mathbb{R}^{H \times W \times C}$ . Weight tensor is  $\mathcal{W} \in \mathbb{R}^{R \times S \times C \times K}$ . The output activation tensor is  $\mathcal{Y} \in \mathbb{R}^{P \times Q \times K}$ . Here we set  $H = W = 5$ ,  $C = 8$ ,  $R = S = 3$ ,  $K = 6$  and  $P = Q = 3$ . Besides, the stride number is 1 and the padding number is 0.

(a) (2%) Write down direct convolution by C++ language style.

(b) (4%) The loop unrolling is one of loop optimization techniques to make full use of the hardware on-chip storage resources. Write down the loop unrolling at input channel level and output channel level, respectively, by C++ language style.

(c) (4%) Sketch the corresponding computing hardware architectures for the two loop unrolling strategies in (b), respectively.

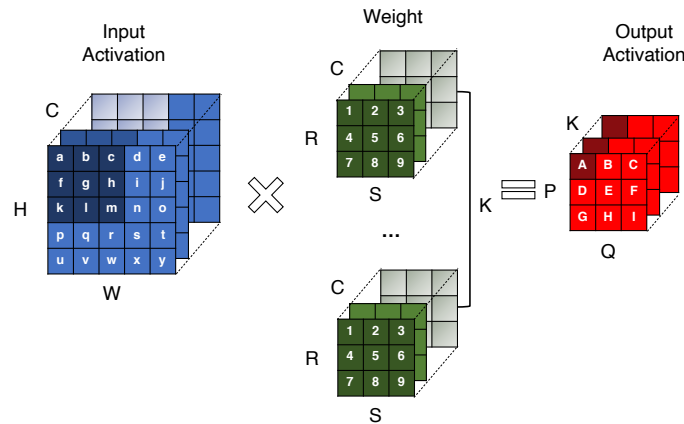


Figure 2: Convolution.

**Q8** (15%) Convolution can be equivalently represented as matrix matrix multiplication. Here we consider a special case:  $\mathbf{Y} = \mathbf{X} \cdot \mathbf{W} + \mathbf{V}$ , where  $\mathbf{Y} \in \mathbb{R}^{N \times K}$  and  $\mathbf{X} \in \mathbb{R}^{N \times M}$  are known input and output matrices.  $\mathbf{V} \in \mathbb{R}^{N \times K}$  is an unknown model error matrix.  $\mathbf{W} \in \mathbb{R}^{M \times K}$  is an unknown model coefficient matrix. In particular, indexes of nonzero elements of each row in  $\mathbf{W}$  is identical to achieve structure sparsity. Let

$$\mathbf{Y} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \\ 7 & 8 \end{bmatrix}, \mathbf{X} = \begin{bmatrix} 9 & 13 & 17 \\ 10 & 14 & 18 \\ 11 & 15 & 19 \\ 12 & 16 & 20 \end{bmatrix}, \mathbf{W} = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix}.$$

(a) (7%) Write down the coordinate descent method to handle the formulation as follows.

$$\min_{\mathbf{W}} \left\| \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \\ 7 & 8 \end{bmatrix} - \begin{bmatrix} 9 & 13 & 17 \\ 10 & 14 & 18 \\ 11 & 15 & 19 \\ 12 & 16 & 20 \end{bmatrix} \cdot \mathbf{W} \right\|_2^2 + \sum_{i=1}^3 \lambda_i \|\mathbf{w}_{i,\cdot}\|_2,$$

where  $\mathbf{w}_{i,\cdot}$  denotes the  $i$ -th row in  $\mathbf{W}$ .  $\lambda_1 = 1$ ,  $\lambda_2 = 100$  and  $\lambda_3 = 1$ . The initial matrix  $\mathbf{W}^{(0)} = \mathbf{O}$ . The stopping criterion is set to 2 iterations. Please show the final numerical result.

(b) (8%) Obtaining this structure sparse model coefficient matrix can be formulated as

$$\min_{\mathbf{W}} \left\| \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \\ 7 & 8 \end{bmatrix} - \begin{bmatrix} 9 & 13 & 17 \\ 10 & 14 & 18 \\ 11 & 15 & 19 \\ 12 & 16 & 20 \end{bmatrix} \cdot \mathbf{W} \right\|_2^2$$

$$\text{s.t. } \sum_{i=1}^3 \mathcal{I}[\|\mathbf{w}_{i,\cdot}\| > 0] \leq 2,$$

where  $\mathcal{I}[\cdot]$  denotes the indicator function and  $\|\cdot\|$  is an any vector norm. In fact,  $\sum_{i=1}^3 \mathcal{I}[\|\mathbf{w}_{i,\cdot}\| > 0]$  denotes the number of nonzero rows in the matrix  $\mathbf{W}$ . In the constraint, 2 is given to determine the number of nonzero rows in the matrix  $\mathbf{W}$ . Orthogonal matching pursuit, as a heuristics method, is widely used in one-dimension sparse vector reconstruction. Please extend the typical orthogonal matching pursuit to handle the formulation and show the final numerical result.