

Pattern Recognition: Homework 9

Due date: 2023.4.25

In this homework we will focus on multilayer perceptron. You need to first draw the computing graph, deduce the gradient of each parameter recursively, then fill in the blank of the file `train.py`

Problem 1 (20 pt)

Consider a three layer MLP with ReLU activation $g(x) = \max(x, 0)$. For the ℓ_{th} layer, denote the weight matrix as $\mathbf{W}_\ell \in \mathbb{R}^{d_\ell \times d_{\ell+1}}$, bias term as $\mathbf{b}_\ell \in \mathbb{R}^{d_\ell}$, where d_ℓ is the dimensionality of the ℓ_{th} layer. For the forward pass of ℓ_{th} layer, the model first takes the last layer's output \mathbf{a}_ℓ ($\mathbf{a}_0 = \mathbf{x}_0$ is the input), then send it through a linear function as

$$\mathbf{z}_{\ell+1} = \mathbf{W}_\ell \mathbf{a}_\ell + \mathbf{b}_\ell. \quad \ell = 0, 1, 2$$

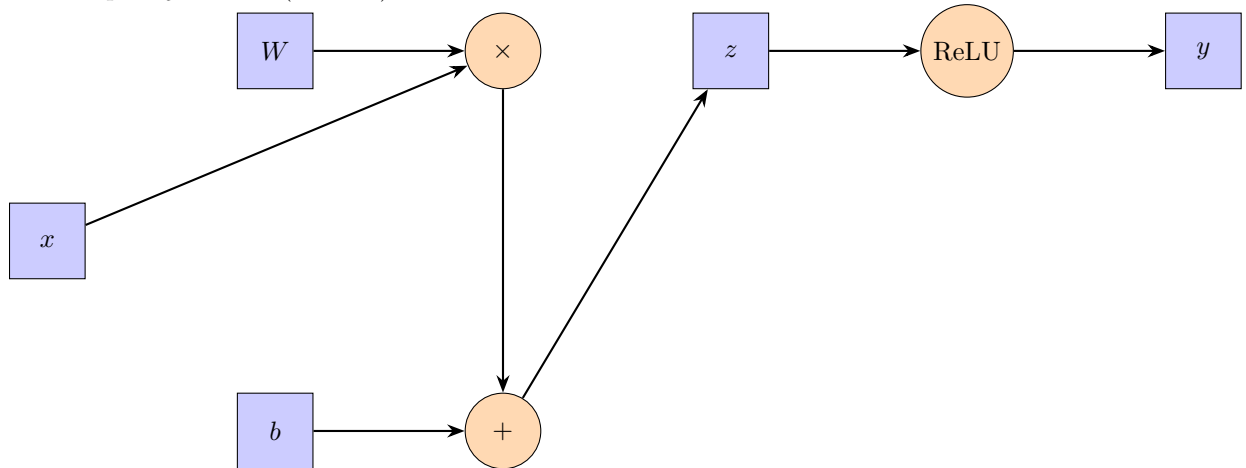
Then apply a non-linear activation g element-wise as

$$\mathbf{a}_{\ell+1} = g(\mathbf{z}_{\ell+1}). \quad \ell = 0, 1$$

Specially, for the last layer, we do not need to take activation $g(\mathbf{z}_3)$ and directly output \mathbf{z}_3 . In one line, the function parametrized by $\boldsymbol{\theta} = (\mathbf{W}_0, \mathbf{b}_0, \mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2)$ is

$$f_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbf{W}_2 g(\mathbf{W}_1 g(\mathbf{W}_0 \mathbf{x} + \mathbf{b}_0) + \mathbf{b}_1) + \mathbf{b}_2 \quad (1)$$

Give a data point (\mathbf{x}, y) and loss function as $J(\boldsymbol{\theta}, \mathbf{x}) = \frac{1}{2} \|f(\mathbf{x}) - y\|^2$, draw a computing graph of the loss. The graph should contain every parameter $\mathbf{w}_\ell, \mathbf{b}_\ell$, data \mathbf{x} and intermediate output \mathbf{a}_ℓ as node. Draw an arrow $a \rightarrow b$ if computing b depends on a . For example, here is the simplest computational graph of forward pass $y = \text{ReLU}(Wx + b)$



Problem 2 (20 pt)

Based on the computing graph above, recursively derive the gradient $\frac{\partial J}{\partial \mathbf{n}}$ for each node \mathbf{n} in the graph. Note that given an edge $a \rightarrow b$, you can directly use gradient $\frac{\partial J}{\partial b}$ as known in your expression for $\frac{\partial J}{\partial a}$ and do not need to expand it. **hint:** draw it from downstream loss J to upstream \mathbf{x} .

Problem 3 (60 pt)

Why draw the graph? Because we will use back propagation to compute the gradient for each learnable parameter. You may notice that it is very complicated if we directly compute the gradient $\frac{\partial J}{\partial \mathbf{W}_0}$ from equation (1). But we can recursively compute it back from J to \mathbf{W}_0 with the relationship you derived from problem 2. Fill in the blank in `train.py` to complete the computation of each parameter's gradient, and run the toy regression code to see how neural network fit a function. You should achieve low training and test loss.