

CSED490F Lab: Autograd

Team 4: Yunkyu Lee (20210733), Hyeonu Cho (20230740)

September 14, 2025

1 Automatic Differentiation

Reverse-mode automatic differentiation performs back-to-front accumulation of local gradients based on the chain rule. Consider a function $f : \mathbb{R}^n \mapsto \mathbb{R}^m$, and its $m \times n$ Jacobian \mathbf{J}_f . For $\mathbf{y} = f(\mathbf{x})$, suppose we know $\nabla_{\mathbf{y}}\mathcal{L}$, the gradient of some scalar value \mathcal{L} w.r.t. \mathbf{y} . Then, we can obtain $\nabla_{\mathbf{x}}\mathcal{L}$ as a vector-Jacobian product (VJP):

$$\nabla_{\mathbf{x}}\mathcal{L} = (\nabla_{\mathbf{y}}\mathcal{L})^\top \mathbf{J}_f$$

As $\mathbf{J}_{f_1 \circ f_2} = \mathbf{J}_{f_1} \circ \mathbf{J}_{f_2}$, we can chain the above VJP in order to obtain gradients of composed functions. Therefore, knowing the local gradients of function outputs w.r.t. their inputs is sufficient for reverse-mode automatic differentiation. We describe the local gradients of some operators below.

1.1 Basic Operations

$$\text{Add}(x, y) = x + y, \quad \frac{\partial \text{Add}}{\partial x} = 1, \quad \frac{\partial \text{Add}}{\partial y} = 1 \quad (1)$$

$$\text{Mul}(x, y) = x \times y, \quad \frac{\partial \text{Mul}}{\partial x} = y, \quad \frac{\partial \text{Mul}}{\partial y} = x \quad (2)$$

$$\text{Pow}(x, y) = x^y, \quad \frac{\partial \text{Pow}}{\partial x} = yx^{y-1}, \quad \frac{\partial \text{Pow}}{\partial y} = x^y \log x \quad (3)$$

$$\text{Log}(x) = \log x, \quad \frac{\partial \text{Log}}{\partial x} = x^{-1}, \quad (4)$$

$$\text{ReLU}(x) = \max(x, 0), \quad \frac{\partial \text{ReLU}}{\partial x} = \begin{cases} 1, & x > 0 \\ 0, & x \leq 0 \end{cases}, \quad (5)$$

$$\text{Sum}(x, \text{axis}) = \sum_{i \in \text{axis}} x_i, \quad \frac{\partial \text{Sum}}{\partial x} = 1, \quad (6)$$

$$\text{MatMul}(x, y) = xy, \quad \frac{\partial \text{MatMul}}{\partial x} = y, \quad \frac{\partial \text{MatMul}}{\partial y} = x \quad (7)$$

1.2 Classification

Following [1], we use a shifted softmax for numerical stability:

$$\text{Softmax}(x_j) = \frac{\exp(x_j - s)}{\sum_i \exp(x_i - s)}$$

The negative log-likelihood loss (with log-probability input)

$$\text{NLLLoss}(\log \hat{y}, y) = \sum_i -y_i \log \hat{y}_i$$

The cross-entropy loss (with logit inputs) can be seen as a combination of **Softmax** and **NLLoss**.

$$\text{CrossEntropyLoss}(x, y) = \text{NLLoss}(\log(\text{Softmax}(x)), y) = \sum_i -y_i \log \text{Softmax}(x_i)$$

1.3 Matrix Multiplication

Consider $y = xw$ for $w \in \mathbb{R}^{D \times M}$, $x \in \mathbb{R}^{N \times D}$, and $y \in \mathbb{R}^{N \times M}$. Then, each element of y is computed as

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}.$$

In matrix form,

$$\begin{bmatrix} y_{1,1} & \cdots & y_{1,M} \\ \vdots & \ddots & \vdots \\ y_{N,1} & \cdots & y_{N,M} \end{bmatrix} = \begin{bmatrix} x_{1,1} & \cdots & x_{1,D} \\ \vdots & \ddots & \vdots \\ x_{N,1} & \cdots & x_{N,D} \end{bmatrix} \begin{bmatrix} w_{1,1} & \cdots & w_{1,M} \\ \vdots & \ddots & \vdots \\ w_{D,1} & \cdots & w_{D,M} \end{bmatrix}$$

Suppose $\frac{\partial \mathcal{L}}{\partial y}$ is given as

$$\begin{bmatrix} \frac{\partial \mathcal{L}}{\partial y_{1,1}} & \cdots & \frac{\partial \mathcal{L}}{\partial y_{1,M}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \mathcal{L}}{\partial y_{N,1}} & \cdots & \frac{\partial \mathcal{L}}{\partial y_{N,M}} \end{bmatrix}$$

Considering any element of x , $x_{i,j}$, it only contributes to all y_n : as

$$y_{i:} = [\sum_d x_{i,d} w_{d,1} \quad \cdots \quad \sum_d x_{i,d} w_{d,M}]$$

This leads to $\frac{\partial \mathcal{L}}{\partial x}$ as

$$\frac{\partial \mathcal{L}}{\partial x_{n,d}} = \sum_m \frac{\partial \mathcal{L}}{\partial y_{n,m}} \frac{\partial y_{n,m}}{\partial x_{n,d}} = \sum_m \frac{\partial \mathcal{L}}{\partial y_{n,m}} w_{d,m} = \frac{\partial \mathcal{L}}{\partial y_{n:}} (w_{d:})^\top$$

Consequently, we can get $\frac{\partial \mathcal{L}}{\partial x} = \frac{\partial \mathcal{L}}{\partial y} w^\top$. Similarly, we can get $\frac{\partial \mathcal{L}}{\partial w} = x^\top \frac{\partial \mathcal{L}}{\partial y}$.

1.4 Cross Entropy Loss

2 MNIST Classification

Utilizing the operators described in section 1, we perform testing on the MNIST dataset. We use a 3-Layer MLP as shown below.

The model was trained for 10 epochs with a batch size of 100 and a learning rate of 0.1. The result is shown in figure below.

Consequently, the model achieved an accuracy of acc, confirming the correctness of its implementation.

References

- [1] Pierre Blanchard, Desmond J Higham, and Nicholas J Higham. “Accurately computing the log-sum-exp and softmax functions”. In: *IMA Journal of Numerical Analysis* 41.4 (Aug. 2020), pp. 2311–2330. ISSN: 0272-4979. DOI: 10.1093/imanum/draa038. eprint: <https://academic.oup.com/imanum/article-pdf/41/4/2311/40758053/draa038.pdf>. URL: <https://doi.org/10.1093/imanum/draa038>.