

Deep Learning 1

Lecturer: Changshui Zhang

zcs@mail.tsinghua.edu.cn

Hong Zhao

vzhao@tsinghua.edu.cn

Student: Jingxuan Yang

yangjx20@mails.tsinghua.edu.cn

Backward Propagation Algorithm

1. We have a k -class classification problem. The input \mathbf{x} is a d -dimensional vector, and the corresponding label \mathbf{y} is a one-hot k -dimensional vector with one element being 1 and others being 0. We define the following neural network $\hat{\mathbf{y}} = f(\mathbf{x}; \mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2)$:

$$\mathbf{h}_1 = \mathbf{W}_1^\top \mathbf{x} + \mathbf{b}_1, \quad \mathbf{h}_1 \in \mathbb{R}^s \quad (1)$$

$$\mathbf{a}_1 = \text{ReLU}(\mathbf{h}_1), \quad \mathbf{a}_1 \in \mathbb{R}^s \quad (2)$$

$$\mathbf{h}_2 = \text{Concat}(\mathbf{a}_1, \mathbf{x}), \quad \mathbf{h}_2 \in \mathbb{R}^{s+d} \quad (3)$$

$$\mathbf{m} \sim \text{Bernoulli}(p), \quad \mathbf{m} \in \mathbb{R}^{s+d} \quad (4)$$

$$\mathbf{a}_2 = \mathbf{h}_2 \odot \mathbf{m}, \quad \mathbf{a}_2 \in \mathbb{R}^{s+d} \quad (5)$$

$$\mathbf{h}_3 = \mathbf{W}_2^\top \mathbf{a}_2 + \mathbf{b}_2, \quad \mathbf{h}_3 \in \mathbb{R}^k \quad (6)$$

$$\hat{\mathbf{y}} = \text{Softmax}(\mathbf{h}_3), \quad \hat{\mathbf{y}} \in \mathbb{R}^k \quad (7)$$

where $\text{ReLU}(\cdot)$ is element-wise ReLU activation function, $\text{Concat}(\mathbf{a}, \mathbf{b})$ means concatenating vector \mathbf{b} after vector \mathbf{a} to make a vector with larger dimensionality, \mathbf{m} is current dropout mask, \odot means element-wise multiplication (Specifically, equation(5) randomly zeroes some of the elements of the input tensor \mathbf{h}_2 with probability p using \mathbf{m} sampled from a Bernoulli distribution. For example, if $\mathbf{h}_2 = [a, b, c, d]^\top$, and $\mathbf{m} = [e, f, g, h]^\top$, then $\mathbf{a}_2 = [ae, bf, cg, dh]^\top$ according to equation(5)), and $\text{Softmax}(\cdot)$ is softmax activation function. The dimensionalities of parameters and internal variables could be inferred from the neural network definition.

We use the following cross-entropy loss:

$$L = -\mathbf{y}^\top \log \hat{\mathbf{y}}. \quad (8)$$

1.1. Please use backward propagation algorithm to find the following derivatives:

$$\frac{\partial L}{\partial \hat{\mathbf{y}}}, \quad \frac{\partial L}{\partial \mathbf{h}_3}, \quad \frac{\partial L}{\partial \mathbf{W}_2}, \quad \frac{\partial L}{\partial \mathbf{b}_2}, \quad \frac{\partial L}{\partial \mathbf{a}_2}, \quad \frac{\partial L}{\partial \mathbf{h}_2}, \quad \frac{\partial L}{\partial \mathbf{a}_1}, \quad \frac{\partial L}{\partial \mathbf{h}_1}, \quad \frac{\partial L}{\partial \mathbf{W}_1}, \quad \frac{\partial L}{\partial \mathbf{b}_1}, \quad \frac{\partial L}{\partial \mathbf{x}}. \quad (9)$$

Please express your results in matrices and vectors instead of sum of individual elements.

解: 原题目中 $\mathbf{w}_1, \mathbf{w}_2$ 实为矩阵, 故改用 $\mathbf{W}_1, \mathbf{W}_2$ 符号代替之. 为避免符号下标混淆, 本题中一律以右上角小括号内指标表示向量, 矩阵和张量的元素, 如 $\mathbf{a}^{(1)}$ 表示向量 \mathbf{a} 的第 1 个元素, $\mathbf{A}^{(2,3)}$ 表示矩阵 \mathbf{A} 的第 2 行第 3 列的元素, $\mathbf{A}^{(1,2,3)}$ 表示张量 \mathbf{A} 的第 1 行第 2 列第 3 竖的元素.

交叉熵损失函数 L 对 $\hat{\mathbf{y}}$ 的偏导数为

$$\frac{\partial L}{\partial \hat{\mathbf{y}}} = \frac{\partial(-\mathbf{y}^\top \log \hat{\mathbf{y}})}{\partial \hat{\mathbf{y}}} = -\frac{\partial \log \hat{\mathbf{y}}}{\partial \hat{\mathbf{y}}} \mathbf{y} = -\text{diag}(\mathbf{1} \oslash \hat{\mathbf{y}}) \mathbf{y} = -\mathbf{y} \oslash \hat{\mathbf{y}} \quad (10)$$

其中, 符号 \oslash 表示逐项除法 (element-wise division).

函数 $\hat{\mathbf{y}}$ 对 \mathbf{h}_3 的偏导数为

$$\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{h}_3} = \frac{\partial}{\partial \mathbf{h}_3} \left(\frac{\exp(\mathbf{h}_3)}{\mathbf{1}^\top \exp(\mathbf{h}_3)} \right) \quad (11)$$

其中, $\forall i, j = 1, 2, \dots, k$, 当 $i \neq j$ 时, 有

$$\begin{aligned} \frac{\partial \hat{\mathbf{y}}^{(i)}}{\partial \mathbf{h}_3^{(j)}} &= \frac{\partial}{\partial \mathbf{h}_3^{(j)}} \left(\frac{\exp(\mathbf{h}_3^{(i)})}{\mathbf{1}^\top \exp(\mathbf{h}_3)} \right) \\ &= \frac{0 - \exp(\mathbf{h}_3^{(i)}) \exp(\mathbf{h}_3^{(j)})}{[\mathbf{1}^\top \exp(\mathbf{h}_3)]^2} \\ &= -\frac{\exp(\mathbf{h}_3^{(i)})}{\mathbf{1}^\top \exp(\mathbf{h}_3)} \frac{\exp(\mathbf{h}_3^{(j)})}{\mathbf{1}^\top \exp(\mathbf{h}_3)} \\ &= -\hat{\mathbf{y}}^{(i)} \hat{\mathbf{y}}^{(j)} \end{aligned} \quad (12)$$

当 $i = j$ 时, 有

$$\begin{aligned} \frac{\partial \hat{\mathbf{y}}^{(j)}}{\partial \mathbf{h}_3^{(j)}} &= \frac{\partial}{\partial \mathbf{h}_3^{(j)}} \left(\frac{\exp(\mathbf{h}_3^{(j)})}{\mathbf{1}^\top \exp(\mathbf{h}_3)} \right) \\ &= \frac{\exp(\mathbf{h}_3^{(j)}) \mathbf{1}^\top \exp(\mathbf{h}_3) - \exp(\mathbf{h}_3^{(j)}) \exp(\mathbf{h}_3^{(j)})}{[\mathbf{1}^\top \exp(\mathbf{h}_3)]^2} \\ &= -\frac{\exp(\mathbf{h}_3^{(j)})}{\mathbf{1}^\top \exp(\mathbf{h}_3)} \frac{\exp(\mathbf{h}_3^{(j)})}{\mathbf{1}^\top \exp(\mathbf{h}_3)} + \frac{\exp(\mathbf{h}_3^{(j)})}{\mathbf{1}^\top \exp(\mathbf{h}_3)} \\ &= -\hat{\mathbf{y}}^{(j)} \hat{\mathbf{y}}^{(j)} + \hat{\mathbf{y}}^{(j)} \end{aligned} \quad (13)$$

综上, $\hat{\mathbf{y}}$ 对 \mathbf{h}_3 的偏导数为

$$\begin{aligned} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{h}_3} &= \frac{\partial}{\partial \mathbf{h}_3} \left(\frac{\exp(\mathbf{h}_3)}{\mathbf{1}^\top \exp(\mathbf{h}_3)} \right) \\ &= -\hat{\mathbf{y}} \hat{\mathbf{y}}^\top + \text{diag}(\hat{\mathbf{y}}) \end{aligned} \quad (14)$$

注意到 \mathbf{y} 是 one-hot 向量, 则交叉熵损失函数 L 对 \mathbf{h}_3 的偏导数为

$$\begin{aligned}
 \frac{\partial L}{\partial \mathbf{h}_3} &= \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{h}_3} \frac{\partial L}{\partial \hat{\mathbf{y}}} \\
 &= [-\hat{\mathbf{y}}\hat{\mathbf{y}}^\top + \text{diag}(\hat{\mathbf{y}})](-\mathbf{y} \oslash \hat{\mathbf{y}}) \\
 &= \hat{\mathbf{y}}\hat{\mathbf{y}}^\top(\mathbf{y} \oslash \hat{\mathbf{y}}) - \text{diag}(\hat{\mathbf{y}})(\mathbf{y} \oslash \hat{\mathbf{y}}) \\
 &= \hat{\mathbf{y}} \sum_{i=1}^k \mathbf{y}^{(i)} - \mathbf{y} \\
 &= \hat{\mathbf{y}} - \mathbf{y}
 \end{aligned} \tag{15}$$

函数 \mathbf{h}_3 对 \mathbf{W}_2 的偏导数为 3 阶张量 $\mathbf{W} \in \mathbb{R}^{(s+d) \times k \times k}$, 其元素为

$$\mathbf{W}^{(m,n,i)} = \frac{\partial \mathbf{h}_3^{(i)}}{\partial \mathbf{W}_2^{(m,n)}} = \frac{\partial [(\mathbf{W}_2^\top)^{(i,\cdot)} \mathbf{a}_2 + \mathbf{b}_2^{(i)}]}{\partial \mathbf{W}_2^{(m,n)}} = \delta_{in} \mathbf{a}_2^{(m)}, \quad \forall i, n = 1, 2, \dots, k, m = 1, 2, \dots, s+d \tag{16}$$

其中, δ_{in} 为 Kronecker delta 符号.

故交叉熵损失函数 L 对 $\mathbf{W}_2^{(m,n)}$ 的偏导数为

$$\begin{aligned}
 \frac{\partial L}{\partial \mathbf{W}_2^{(m,n)}} &= \sum_{i=1}^k \frac{\partial \mathbf{h}_3^{(i)}}{\partial \mathbf{W}_2^{(m,n)}} \frac{\partial L}{\partial \mathbf{h}_3^{(i)}} \\
 &= \sum_{i=1}^k \delta_{in} \mathbf{a}_2^{(m)} [\hat{\mathbf{y}}^{(i)} - \mathbf{y}^{(i)}] \\
 &= \mathbf{a}_2^{(m)} [\hat{\mathbf{y}}^{(n)} - \mathbf{y}^{(n)}]
 \end{aligned} \tag{17}$$

所以, 交叉熵损失函数 L 对 \mathbf{W}_2 的偏导数为

$$\frac{\partial L}{\partial \mathbf{W}_2} = \mathbf{a}_2(\hat{\mathbf{y}} - \mathbf{y})^\top \tag{18}$$

交叉熵损失函数 L 对 \mathbf{b}_2 的偏导数为

$$\frac{\partial L}{\partial \mathbf{b}_2} = \frac{\partial \mathbf{h}_3}{\partial \mathbf{b}_2} \frac{\partial L}{\partial \mathbf{h}_3} = \frac{\partial (\mathbf{W}_2^\top \mathbf{a}_2 + \mathbf{b}_2)}{\partial \mathbf{b}_2} \frac{\partial L}{\partial \mathbf{h}_3} = \mathbf{I}(\hat{\mathbf{y}} - \mathbf{y}) = \hat{\mathbf{y}} - \mathbf{y} \tag{19}$$

交叉熵损失函数 L 对 \mathbf{a}_2 的偏导数为

$$\frac{\partial L}{\partial \mathbf{a}_2} = \frac{\partial \mathbf{h}_3}{\partial \mathbf{a}_2} \frac{\partial L}{\partial \mathbf{h}_3} = \mathbf{W}_2(\hat{\mathbf{y}} - \mathbf{y}) \tag{20}$$

交叉熵损失函数 L 对 \mathbf{h}_2 的偏导数为

$$\frac{\partial L}{\partial \mathbf{h}_2} = \frac{\partial \mathbf{a}_2}{\partial \mathbf{h}_2} \frac{\partial L}{\partial \mathbf{a}_2} = \frac{\partial (\mathbf{h}_2 \odot \mathbf{m})}{\partial \mathbf{h}_2} \frac{\partial L}{\partial \mathbf{a}_2} = \text{diag}(\mathbf{m}) \mathbf{W}_2(\hat{\mathbf{y}} - \mathbf{y}) = \mathbf{m} \odot [\mathbf{W}_2(\hat{\mathbf{y}} - \mathbf{y})] \tag{21}$$

交叉熵损失函数 L 对 \mathbf{a}_1 的偏导数为

$$\frac{\partial L}{\partial \mathbf{a}_1} = \frac{\partial \mathbf{h}_2}{\partial \mathbf{a}_1} \frac{\partial L}{\partial \mathbf{h}_2} = \frac{\partial \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{x} \end{bmatrix}}{\partial \mathbf{a}_1} \frac{\partial L}{\partial \mathbf{h}_2} = [\mathbf{I}_{s \times s} \quad \mathbf{0}_{s \times d}] \text{diag}(\mathbf{m}) \mathbf{W}_2 (\hat{\mathbf{y}} - \mathbf{y}) \quad (22)$$

交叉熵损失函数 L 对 \mathbf{h}_1 的偏导数为

$$\frac{\partial L}{\partial \mathbf{h}_1} = \frac{\partial \mathbf{a}_1}{\partial \mathbf{h}_1} \frac{\partial L}{\partial \mathbf{a}_1} = \frac{\partial \text{ReLU}(\mathbf{h}_1)}{\partial \mathbf{h}_1} \frac{\partial L}{\partial \mathbf{a}_1} = \text{diag} \left(\frac{1 + \text{sgn}(\mathbf{h}_1)}{2} \right) \frac{\partial L}{\partial \mathbf{a}_1} \quad (23)$$

函数 \mathbf{h}_1 对 \mathbf{W}_1 的偏导数为

$$\frac{\partial \mathbf{h}_1^{(i)}}{\partial \mathbf{W}_1^{(m,n)}} = \frac{\partial \left[(\mathbf{W}_1^\top)^{(i,:)} \mathbf{x} + \mathbf{b}_1^{(i)} \right]}{\partial \mathbf{W}_1^{(m,n)}} = \delta_{in} \mathbf{x}^{(m)}, \quad \forall i, n = 1, 2, \dots, s, m = 1, 2, \dots, d \quad (24)$$

故交叉熵损失函数 L 对 $\mathbf{W}_1^{(m,n)}$ 的偏导数为

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{W}_1^{(m,n)}} &= \sum_{i=1}^s \frac{\partial \mathbf{h}_1^{(i)}}{\partial \mathbf{W}_1^{(m,n)}} \frac{\partial L}{\partial \mathbf{h}_1^{(i)}} \\ &= \sum_{i=1}^s \frac{\partial L}{\partial \mathbf{h}_1^{(i)}} \delta_{in} \mathbf{x}^{(m)} \\ &= \mathbf{x}^{(m)} \frac{\partial L}{\partial \mathbf{h}_1^{(n)}} \end{aligned} \quad (25)$$

所以, 交叉熵损失函数 L 对 \mathbf{W}_1 的偏导数为

$$\frac{\partial L}{\partial \mathbf{W}_1} = \mathbf{x} \left(\frac{\partial L}{\partial \mathbf{h}_1} \right)^\top \quad (26)$$

交叉熵损失函数 L 对 \mathbf{b}_1 的偏导数为

$$\frac{\partial L}{\partial \mathbf{b}_1} = \frac{\partial \mathbf{h}_1}{\partial \mathbf{b}_1} \frac{\partial L}{\partial \mathbf{h}_1} = \frac{\partial (\mathbf{W}_1^\top \mathbf{x} + \mathbf{b}_1)}{\partial \mathbf{b}_1} \frac{\partial L}{\partial \mathbf{h}_1} = \mathbf{I} \frac{\partial L}{\partial \mathbf{h}_1} = \frac{\partial L}{\partial \mathbf{h}_1} \quad (27)$$

交叉熵损失函数 L 对 \mathbf{x} 的偏导数为

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{x}} &= \frac{\partial \mathbf{h}_1}{\partial \mathbf{x}} \frac{\partial L}{\partial \mathbf{h}_1} + \frac{\partial \mathbf{h}_2}{\partial \mathbf{x}} \frac{\partial L}{\partial \mathbf{h}_2} \\ &= \frac{\partial (\mathbf{W}_1^\top \mathbf{x} + \mathbf{b}_1)}{\partial \mathbf{x}} \frac{\partial L}{\partial \mathbf{h}_1} + \frac{\partial \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{x} \end{bmatrix}}{\partial \mathbf{x}} \frac{\partial L}{\partial \mathbf{h}_2} \\ &= \mathbf{W}_1 \frac{\partial L}{\partial \mathbf{h}_1} + [\mathbf{0}_{d \times s} \quad \mathbf{I}_{d \times d}] \frac{\partial L}{\partial \mathbf{h}_2} \end{aligned} \quad (28)$$

1.2. Derive the back-propagation updates for convolutional neural networks: In this situation, the input \mathbf{x} is an image with size $C_{\text{in}} \times H \times W$, where the three dimensions represent channel, height and width respectively. We define the following simple convolutional neural network:

$$\mathbf{u}_1 = \text{Conv2d}(C_{\text{in}}, C_{\text{out}}, k)(\mathbf{x}) \quad (29)$$

$$\mathbf{h}_1 = \text{MaxPool2d}(N)(\mathbf{u}_1) \quad (30)$$

$$\mathbf{a}_1 = \text{ReLU}(\mathbf{h}_1) \quad (31)$$

$$\mathbf{u}_2 = \text{Flatten}(\mathbf{a}_1) \quad (32)$$

$$\mathbf{h}_2 = \mathbf{W}_2^\top \mathbf{u}_2 + \mathbf{b}_2 \quad (33)$$

$$\hat{\mathbf{y}} = \text{Softmax}(\mathbf{h}_2) \quad (34)$$

where **Conv2d** (There are learnable parameters $\mathbf{W}_1, \mathbf{b}_1$ in this function. See this [link](#) for detailed definitions of the function arguments and this [link](#) for visualization) and **MaxPool2d** are convolutional function and pooling function defined in a PyTorch style, **Flatten** means reshaping the input vector into a one-dimensional vector, and the final loss function is the same with Eq. (8). What is the derivatives for the network parameters (i.e. convolutional and linear weights and biases)?

Hint: Read the following reference material: [Notes on Convolutional Neural Networks](#).

解: 首先确定各个变量的维数, 输入图片 $\mathbf{x} \in \mathbb{R}^{C_{\text{in}} \times H \times W}$, Conv2d 函数输出 $\mathbf{u}_1 \in \mathbb{R}^{C_{\text{out}} \times H_{\text{out}} \times W_{\text{out}}}$, 其中

$$H_{\text{out}} = H - k + 1, \quad W_{\text{out}} = W - k + 1 \quad (35)$$

MaxPool2d 函数输出 $\mathbf{h}_1 \in \mathbb{R}^{C_{\text{out}} \times H_{\text{mp}} \times W_{\text{mp}}}$, 其中

$$H_{\text{mp}} = \left\lfloor \frac{H_{\text{out}}}{N} \right\rfloor, \quad W_{\text{mp}} = \left\lfloor \frac{W_{\text{out}}}{N} \right\rfloor \quad (36)$$

ReLU 函数输出 $\mathbf{a}_1 \in \mathbb{R}^{C_{\text{out}} \times H_{\text{mp}} \times W_{\text{mp}}}$, Flatten 函数输出 $\mathbf{u}_2 \in \mathbb{R}^d$, 其中

$$d = C_{\text{out}} + H_{\text{mp}} + W_{\text{mp}} \quad (37)$$

设 $\mathbf{W}_2 \in \mathbb{R}^{d \times s}$, 则 $\mathbf{b}_2 \in \mathbb{R}^s$, $\mathbf{h}_2 \in \mathbb{R}^s$, $\hat{\mathbf{y}} \in \mathbb{R}^s$.

由 1.1 题可知, 交叉熵损失函数 L 对 \mathbf{W}_2 的偏导数为

$$\frac{\partial L}{\partial \mathbf{W}_2} = \mathbf{a}_2(\hat{\mathbf{y}} - \mathbf{y})^\top \quad (38)$$

交叉熵损失函数 L 对 \mathbf{b}_2 的偏导数为

$$\frac{\partial L}{\partial \mathbf{b}_2} = \hat{\mathbf{y}} - \mathbf{y} \quad (39)$$

交叉熵损失函数 L 对 \mathbf{u}_2 的偏导数为

$$\frac{\partial L}{\partial \mathbf{u}_2} = \mathbf{W}_2(\hat{\mathbf{y}} - \mathbf{y}) \quad (40)$$

函数 \mathbf{u}_2 对 \mathbf{a}_1 的偏导数为 4 阶张量 $\mathbf{A} \in \mathbb{R}^{C_{\text{out}} \times H_{\text{mp}} \times W_{\text{mp}} \times d}$, 其元素为

$$\mathbf{A}^{(i,j,k,m)} = \frac{\partial \mathbf{u}_2^{(m)}}{\partial \mathbf{a}_1^{(i,j,k)}} = \frac{\partial [\text{Flatten}(\mathbf{a}_1)]^{(m)}}{\partial \mathbf{a}_1^{(i,j,k)}} = \delta_{mn(i,j,k)} \quad (41)$$

其中 $n(\cdot, \cdot, \cdot)$ 是三元函数,

$$n(i, j, k) \triangleq (i-1)H_{\text{mp}}W_{\text{mp}} + (j-1)W_{\text{mp}} + k \quad (42)$$

交叉熵损失函数 L 对 \mathbf{a}_1 的偏导数为

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{a}_1^{(i,j,k)}} &= \sum_{m=1}^d \frac{\partial \mathbf{u}_2^{(m)}}{\partial \mathbf{a}_1^{(i,j,k)}} \frac{\partial L}{\partial \mathbf{u}_2^{(m)}} \\ &= \sum_{m=1}^d \delta_{mn(i,j,k)} \mathbf{W}_2^{(m,:)} (\hat{\mathbf{y}} - \mathbf{y}) \\ &= \mathbf{W}_2^{(n(i,j,k),:)} (\hat{\mathbf{y}} - \mathbf{y}) \end{aligned} \quad (43)$$

函数 \mathbf{a}_1 对 \mathbf{h}_1 的偏导数为 6 阶张量 $\mathbf{H} \in \mathbb{R}^{C_{\text{out}} \times H_{\text{mp}} \times W_{\text{mp}} \times C_{\text{out}} \times H_{\text{mp}} \times W_{\text{mp}}}$, 其元素为

$$\mathbf{H}^{(r,s,t,i,j,k)} = \frac{\partial \mathbf{a}_1^{(i,j,k)}}{\partial \mathbf{h}_1^{(r,s,t)}} = \frac{\partial \text{ReLU}(\mathbf{h}_1^{(i,j,k)})}{\partial \mathbf{h}_1^{(r,s,t)}} = \delta_{ir} \delta_{js} \delta_{kt} \frac{1 + \text{sgn}(\mathbf{h}_1^{(r,s,t)})}{2} \quad (44)$$

交叉熵损失函数 L 对 \mathbf{h}_1 的偏导数为

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{h}_1^{(r,s,t)}} &= \sum_{i=1}^{C_{\text{out}}} \sum_{j=1}^{H_{\text{mp}}} \sum_{k=1}^{W_{\text{mp}}} \frac{\partial \mathbf{a}_1^{(i,j,k)}}{\partial \mathbf{h}_1^{(r,s,t)}} \frac{\partial L}{\partial \mathbf{a}_1^{(i,j,k)}} \\ &= \sum_{i=1}^{C_{\text{out}}} \sum_{j=1}^{H_{\text{mp}}} \sum_{k=1}^{W_{\text{mp}}} \delta_{ir} \delta_{js} \delta_{kt} \frac{1 + \text{sgn}(\mathbf{h}_1^{(r,s,t)})}{2} \mathbf{W}_2^{(n(i,j,k),:)} (\hat{\mathbf{y}} - \mathbf{y}) \\ &= \frac{1 + \text{sgn}(\mathbf{h}_1^{(r,s,t)})}{2} \mathbf{W}_2^{(n(r,s,t),:)} (\hat{\mathbf{y}} - \mathbf{y}) \end{aligned} \quad (45)$$

函数 \mathbf{h}_1 对 \mathbf{u}_1 的偏导数为 6 阶张量 $\mathbf{U} \in \mathbb{R}^{C_{\text{out}} \times H_{\text{out}} \times W_{\text{out}} \times C_{\text{out}} \times H_{\text{mp}} \times W_{\text{mp}}}$, 其元素为

$$\begin{aligned} \mathbf{U}^{(p,q,l,r,s,t)} &= \frac{\partial \mathbf{h}_1^{(r,s,t)}}{\partial \mathbf{u}_1^{(p,q,l)}} \\ &= \frac{\partial \text{MaxPool2d}(\mathbf{h}_1^{(r,s,t)})}{\partial \mathbf{u}_1^{(p,q,l)}} \\ &= \begin{cases} \delta_{\mathbf{h}_1^{(r,s,t)} \mathbf{u}_1^{(p,q,l)}}, & \text{if } p-1 < \frac{r}{N} \leq p, \quad q-1 < \frac{s}{N} \leq q, \quad l-1 < \frac{t}{N} \leq l \\ 0, & \text{o.w.} \end{cases} \end{aligned} \quad (46)$$

交叉熵损失函数 L 对 \mathbf{u}_1 的偏导数为

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{u}_1^{(p,q,l)}} &= \sum_{r=1}^{C_{\text{out}}} \sum_{s=1}^{H_{\text{mp}}} \sum_{t=1}^{W_{\text{mp}}} \frac{\partial \mathbf{h}_1^{(r,s,t)}}{\partial \mathbf{u}_1^{(p,q,l)}} \frac{\partial L}{\partial \mathbf{h}_1^{(r,s,t)}} \\ &= \sum_{r=1}^{C_{\text{out}}} \sum_{s=1}^{H_{\text{mp}}} \sum_{t=1}^{W_{\text{mp}}} \frac{\partial \mathbf{h}_1^{(r,s,t)}}{\partial \mathbf{u}_1^{(p,q,l)}} \frac{1 + \text{sgn}(\mathbf{h}_1^{(r,s,t)})}{2} \mathbf{W}_2^{(n(r,s,t),:)} (\hat{\mathbf{y}} - \mathbf{y}) \end{aligned} \quad (47)$$

由函数 Conv2d 定义可知

$$\mathbf{u}_1^{(j,:,:) } = \mathbf{b}_1^{(j,:,:) } + \sum_{k=1}^{C_{\text{in}}} \mathbf{W}_1^{(j,k,:,:) } \star \mathbf{x}^{(k,:,:) } \quad (48)$$

其中, \star 为二维互相关 (cross-correlation) 运算符号,

$$\mathbf{b}_1 \in \mathbb{R}^{C_{\text{out}} \times H_{\text{out}} \times W_{\text{out}}}, \quad \mathbf{W}_1 \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in}} \times k \times k} \quad (49)$$

交叉熵损失函数 L 对 \mathbf{b}_1 的偏导数为

$$\frac{\partial L}{\partial \mathbf{b}_1} = \frac{\partial L}{\partial \mathbf{u}_1} \quad (50)$$

函数 \mathbf{u}_1 对 \mathbf{W}_1 的偏导数为 7 阶张量 $\mathbf{W} \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in}} \times k \times k \times C_{\text{out}} \times H_{\text{out}} \times W_{\text{out}}}$, 难以具体表出, 根据参考文献 [2] 可知可以使用 MATLAB 内置函数直接计算交叉熵损失函数 L 对 \mathbf{W}_1 的偏导数.

交叉熵损失函数 L 对 \mathbf{W}_1 的偏导数为

$$\frac{\partial L}{\partial \mathbf{W}_1^{(m,n,:,:)}} = \text{rot180} \left\{ \text{conv2} \left[\mathbf{x}^{(n,:,:)}, \text{rot180} \left(\frac{\partial L}{\partial \mathbf{u}_1^{(m,:,:)}} \right), \text{'valid'} \right] \right\} \quad (51)$$

其中, rot180 和 conv2 均为 MATLAB 中的函数.

Vanishing and Exploding Gradients

2. In this section we will investigate the vanishing/exploding gradient problem and see why ReLU and ResNet can mitigate this problem. The vanishing/exploding gradient problem often occurs during training very deep neural networks and could make the training process failed.

Suppose the input \mathbf{x} and output $\hat{\mathbf{y}}$ are both d -dimensional vectors, and the gradient of loss w.r.t. $\hat{\mathbf{y}}$: $\frac{\partial L}{\partial \hat{\mathbf{y}}} \in \mathbb{R}^d$ is known. Suppose the deep feed-forward neural network has l layers (i.e., $l > 100$), and all internal nodes are also d -dimensional. All weights are initialized with zero-mean Gaussian distribution and all biases are initialized to zero. We now analyze the gradients in the first optimization iteration.

Effect of ReLU. The neural network is defined by:

$$\begin{aligned} \mathbf{a}_1 &= \text{Sigmoid}(\mathbf{W}_1^\top \mathbf{x} + \mathbf{b}_1) \\ \mathbf{a}_2 &= \text{Sigmoid}(\mathbf{W}_2^\top \mathbf{a}_1 + \mathbf{b}_2) \\ &\vdots \\ \hat{\mathbf{y}} = \mathbf{a}_l &= \text{Sigmoid}(\mathbf{W}_l^\top \mathbf{a}_{l-1} + \mathbf{b}_l) \end{aligned} \quad (52)$$

2.1. Find $\frac{\partial L}{\partial \mathbf{W}_1}$.

解: 对 $z \in \mathbb{R}$, Sigmoid 函数定义为

$$\text{Sigmoid}(z) = \frac{1}{1 + \exp(-z)} \quad (53)$$

其对 z 的偏导数为

$$\begin{aligned} \frac{\partial \text{Sigmoid}(z)}{\partial z} &= \frac{-[-\exp(-z)]}{[1 + \exp(z)]^2} \\ &= \frac{1}{1 + \exp(-z)} \frac{\exp(-z)}{1 + \exp(-z)} \\ &= \text{Sigmoid}(z)[1 - \text{Sigmoid}(z)] \end{aligned} \quad (54)$$

对 $\mathbf{z} \in \mathbb{R}^d$, Sigmoid 函数可用逐项除法 \odot 符号表示为

$$\text{Sigmoid}(\mathbf{z}) = \mathbf{1} \odot [1 + \exp(-\mathbf{z})] \quad (55)$$

则由式 (54) 可知 $\text{Sigmoid}(\mathbf{z})$ 对 \mathbf{z} 的偏导数为

$$\frac{\partial \text{Sigmoid}(\mathbf{z})}{\partial \mathbf{z}} = \text{diag} \left\{ \text{Sigmoid}(\mathbf{z}) \odot [\mathbf{1} - \text{Sigmoid}(\mathbf{z})] \right\} \quad (56)$$

记 $\mathbf{a}_0 \triangleq \mathbf{x}$, 令

$$\mathbf{h}_i \triangleq \mathbf{W}_i^\top \mathbf{a}_{i-1} + \mathbf{b}_i, \quad \forall i = 1, 2, \dots, l \quad (57)$$

则

$$\frac{\partial \mathbf{h}_i}{\partial \mathbf{a}_{i-1}} = \mathbf{W}_i, \quad \forall i = 1, 2, \dots, l \quad (58)$$

又由式 (56) 可知

$$\frac{\partial \mathbf{a}_i}{\partial \mathbf{h}_i} = \frac{\partial \text{Sigmoid}(\mathbf{h}_i)}{\partial \mathbf{h}_i} = \text{diag}[\mathbf{a}_i \odot (\mathbf{1} - \mathbf{a}_i)], \quad \forall i = 1, 2, \dots, l \quad (59)$$

则由链式法则及式 (18) 可得

$$\begin{aligned}
\frac{\partial L}{\partial \mathbf{W}_1} &= \frac{\partial \mathbf{a}_1}{\partial \mathbf{W}_1} \frac{\partial \mathbf{a}_2}{\partial \mathbf{a}_1} \cdots \frac{\partial \mathbf{a}_{l-1}}{\partial \mathbf{a}_{l-2}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}_{l-1}} \frac{\partial L}{\partial \hat{\mathbf{y}}} \\
&= \left(\frac{\partial \mathbf{h}_1}{\partial \mathbf{W}_1} \frac{\partial \mathbf{a}_1}{\partial \mathbf{h}_1} \right) \left(\frac{\partial \mathbf{h}_2}{\partial \mathbf{a}_1} \frac{\partial \mathbf{a}_2}{\partial \mathbf{h}_2} \right) \cdots \left(\frac{\partial \mathbf{h}_{l-1}}{\partial \mathbf{a}_{l-2}} \frac{\partial \mathbf{a}_{l-1}}{\partial \mathbf{h}_{l-1}} \right) \left(\frac{\partial \mathbf{h}_l}{\partial \mathbf{a}_{l-1}} \frac{\partial \mathbf{a}_l}{\partial \mathbf{h}_l} \right) \frac{\partial L}{\partial \hat{\mathbf{y}}} \\
&= \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}_1} \frac{\partial \mathbf{a}_1}{\partial \mathbf{h}_1} \left(\prod_{i=2}^l \frac{\partial \mathbf{h}_i}{\partial \mathbf{a}_{i-1}} \frac{\partial \mathbf{a}_i}{\partial \mathbf{h}_i} \right) \frac{\partial L}{\partial \hat{\mathbf{y}}} \\
&= \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}_1} \text{diag} [\mathbf{a}_1 \odot (\mathbf{1} - \mathbf{a}_1)] \left(\prod_{i=2}^l \mathbf{W}_i \text{diag} [\mathbf{a}_i \odot (\mathbf{1} - \mathbf{a}_i)] \right) \frac{\partial L}{\partial \hat{\mathbf{y}}} \\
&= \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}_1} \left(\prod_{i=2}^l \mathbf{W}_i \right) \left(\prod_{j=1}^l \text{diag} [\mathbf{a}_j \odot (\mathbf{1} - \mathbf{a}_j)] \right) \frac{\partial L}{\partial \hat{\mathbf{y}}} \\
&= \mathbf{x} \left[\left(\prod_{i=2}^l \mathbf{W}_i \right) \left(\prod_{j=1}^l \text{diag} [\mathbf{a}_j \odot (\mathbf{1} - \mathbf{a}_j)] \right) \frac{\partial L}{\partial \hat{\mathbf{y}}} \right]^\top
\end{aligned} \tag{60}$$

2.2. Suppose $z \in \mathbb{R}$ is a scalar, find the maximum of sigmoid function's derivative: $g(z) = \text{Sigmoid}'(z)$.

解: 由式 (54) 可知

$$g(z) = \text{Sigmoid}'(z) = \text{Sigmoid}(z)[1 - \text{Sigmoid}(z)] \tag{61}$$

根据基本不等式可得

$$g(z) = \text{Sigmoid}(z)[1 - \text{Sigmoid}(z)] \leq \frac{[\text{Sigmoid}(z) + 1 - \text{Sigmoid}(z)]^2}{4} = \frac{1}{4} \tag{62}$$

等号成立当且仅当

$$\text{Sigmoid}(z) = 1 - \text{Sigmoid}(z) \tag{63}$$

即

$$\text{Sigmoid}(z) = \frac{1}{2} \tag{64}$$

即

$$z = 0 \tag{65}$$

故当 $z = 0$ 时, Sigmoid 函数的梯度 $g(z)$ 取到最大值

$$\max_{z \in \mathbb{R}} g(z) = g(0) = \frac{1}{4} \tag{66}$$

2.3. Since many sigmoid functions' derivatives are multiplied together in $\frac{\partial L}{\partial \mathbf{W}_1}$, the values in $\frac{\partial L}{\partial \mathbf{W}_1}$ will tend to vanish to zero. Is there any possibility that we can initialize weight matrices carefully to recover the values into a good range for float32 (e.g., -0.01 to 1000)? (*Hint: Since the sigmoid function saturates in both sides,*

in order to keep the gradient scale “reasonable”, one needs to make sure the values before sigmoid function lie in non-saturation region.)

解: 由上题可知 Sigmoid(z) 函数的最大梯度位于 $z = 0$ 处, 且 Sigmoid(0) = 1/2, 即要求

$$\begin{aligned} \mathbf{W}_1^\top \mathbf{x} + \mathbf{b}_1 &= \mathbf{0} \\ \mathbf{W}_i^\top \frac{\mathbf{1}}{2} + \mathbf{b}_i &= \mathbf{0}, \quad \forall i = 2, 3, \dots, l \end{aligned} \quad (67)$$

初始化时偏置 $\mathbf{b}_i = \mathbf{0}$, $\forall i = 1, 2, \dots, l$, 则初始化权重矩阵需满足

$$\begin{aligned} \mathbf{W}_1^\top \mathbf{x} &= \mathbf{0} \\ \mathbf{W}_i^\top \mathbf{1} &= \mathbf{0}, \quad \forall i = 2, 3, \dots, l \end{aligned} \quad (68)$$

此时获得的 Sigmoid 函数梯度连乘是最大的, 但由于

$$\max_{z \in \mathbb{R}} g(z) = \frac{1}{4} \quad (69)$$

又 $l > 100$, $\forall z_j \in \mathbb{R}$, $j = 1, 2, \dots, l$ 有

$$\prod_{j=1}^l g(z_j) < \prod_{j=1}^l \max_{z_j \in \mathbb{R}} g(z_j) = 2^{-2l} < 2^{-200} \ll 2^{-126} \quad (70)$$

故当层数很多时, 在梯度都取最大值这种最好的情况下也几乎不可能把梯度值恢复到对 float32 精度友好的数值范围 (e.g., -0.01 to 1000) 内, 甚至会远远超出 float32 精度的表示范围 ($2^{-126} \sim 2^{127}$).

2.4. If our computer supports efficient arbitrary precision floating point calculation, do we need to care about vanishing gradient problem anymore?

解: 需要考虑, 因为即便计算精确, 梯度 $\frac{\partial L}{\partial \mathbf{W}_1}$ 过小也会使得权重矩阵的更新过程

$$\mathbf{W}_1 \leftarrow \mathbf{W}_1 - \eta \frac{\partial L}{\partial \mathbf{W}_1} \quad (71)$$

极为缓慢, 最终导致整个神经网络长时间无法收敛.

2.5. Can we alleviate the vanishing problem by replacing the activation function from sigmoid to ReLU? Why?

解: 可以缓解一部分. 此时

$$\frac{\partial \mathbf{a}_i}{\partial \mathbf{h}_i} = \frac{\partial \text{ReLU}(\mathbf{h}_i)}{\partial \mathbf{h}_i} = \text{diag} \left(\frac{\mathbf{1} + \text{sgn}(\mathbf{a}_j)}{2} \right), \quad \forall i = 1, 2, \dots, l \quad (72)$$

则

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{W}_1} &= \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}_1} \frac{\partial \mathbf{a}_1}{\partial \mathbf{h}_1} \left(\prod_{i=2}^l \frac{\partial \mathbf{h}_i}{\partial \mathbf{a}_{i-1}} \frac{\partial \mathbf{a}_i}{\partial \mathbf{h}_i} \right) \frac{\partial L}{\partial \hat{\mathbf{y}}} \\ &= \mathbf{x} \left[\left(\prod_{i=2}^l \mathbf{W}_i \right) \left(\prod_{j=1}^l \text{diag} \left(\frac{\mathbf{1} + \text{sgn}(\mathbf{a}_j)}{2} \right) \right) \frac{\partial L}{\partial \hat{\mathbf{y}}} \right]^\top \end{aligned} \quad (73)$$

所以, ReLU 函数梯度连乘的部分

$$\prod_{j=1}^l \text{diag} \left(\frac{\mathbf{1} + \text{sgn}(\mathbf{a}_j)}{2} \right)$$

不会消失, 但若 $\det(\mathbf{W}_i)$ 过小, 则权重矩阵连乘的部分

$$\prod_{i=2}^l \mathbf{W}_i$$

会消失, 这也会导致梯度的消失. 因此, 把 Sigmoid 函数换成 ReLU 函数可以缓解一部分梯度消失的问题, 但不能完全避免这个问题.

Effect of ResNet. The neural network now has skip connections:

$$\begin{aligned} \mathbf{a}_1 &= \text{Sigmoid}(\mathbf{W}_1^\top \mathbf{x} + \mathbf{b}_1) + \mathbf{x} \\ \mathbf{a}_2 &= \text{Sigmoid}(\mathbf{W}_2^\top \mathbf{a}_1 + \mathbf{b}_2) + \mathbf{a}_1 \\ &\vdots \\ \hat{\mathbf{y}} = \mathbf{a}_l &= \text{Sigmoid}(\mathbf{W}_l^\top \mathbf{a}_{l-1} + \mathbf{b}_l) + \mathbf{a}_{l-1} \end{aligned} \quad (74)$$

2.7. Find $\frac{\partial L}{\partial \mathbf{W}_1}$ and explain why the gradients can not easily vanish.

解: 与 2.1 题类似, 由链式法则及式 (18) 可得

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{W}_1} &= \frac{\partial \mathbf{a}_1}{\partial \mathbf{W}_1} \frac{\partial \mathbf{a}_2}{\partial \mathbf{a}_1} \cdots \frac{\partial \mathbf{a}_{l-1}}{\partial \mathbf{a}_{l-2}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{a}_{l-1}} \frac{\partial L}{\partial \hat{\mathbf{y}}} \\ &= \left(\frac{\partial \mathbf{h}_1}{\partial \mathbf{W}_1} \frac{\partial \mathbf{a}_1}{\partial \mathbf{h}_1} \right) \left(\frac{\partial \mathbf{h}_2}{\partial \mathbf{a}_1} \frac{\partial \mathbf{a}_2}{\partial \mathbf{h}_2} + \mathbf{I} \right) \cdots \left(\frac{\partial \mathbf{h}_{l-1}}{\partial \mathbf{a}_{l-2}} \frac{\partial \mathbf{a}_{l-1}}{\partial \mathbf{h}_{l-1}} + \mathbf{I} \right) \left(\frac{\partial \mathbf{h}_l}{\partial \mathbf{a}_{l-1}} \frac{\partial \mathbf{a}_l}{\partial \mathbf{h}_l} + \mathbf{I} \right) \frac{\partial L}{\partial \hat{\mathbf{y}}} \\ &= \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}_1} \frac{\partial \mathbf{a}_1}{\partial \mathbf{h}_1} \left[\prod_{i=2}^l \left(\frac{\partial \mathbf{h}_i}{\partial \mathbf{a}_{i-1}} \frac{\partial \mathbf{a}_i}{\partial \mathbf{h}_i} + \mathbf{I} \right) \right] \frac{\partial L}{\partial \hat{\mathbf{y}}} \\ &= \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}_1} \text{diag} [\mathbf{a}_1 \odot (\mathbf{1} - \mathbf{a}_1)] \left[\prod_{i=2}^l \left(\mathbf{W}_i \text{diag} [\mathbf{a}_i \odot (\mathbf{1} - \mathbf{a}_i)] + \mathbf{I} \right) \right] \frac{\partial L}{\partial \hat{\mathbf{y}}} \\ &= \mathbf{x} \left\{ \text{diag} [\mathbf{a}_1 \odot (\mathbf{1} - \mathbf{a}_1)] \left[\prod_{i=2}^l \left(\mathbf{W}_i \text{diag} [\mathbf{a}_i \odot (\mathbf{1} - \mathbf{a}_i)] + \mathbf{I} \right) \right] \frac{\partial L}{\partial \hat{\mathbf{y}}} \right\}^\top \end{aligned} \quad (75)$$

所以, 连乘部分

$$\prod_{i=2}^l \left(\mathbf{W}_i \text{diag} [\mathbf{a}_i \odot (\mathbf{1} - \mathbf{a}_i)] + \mathbf{I} \right)$$

中的每一项都加上了一个单位矩阵, 则此部分轻易不会消失, 因此 ResNet 可以有效防止梯度消失的问题.

Programming: Image Classification with CIFAR-10

Please install PyTorch and run the [CIFAR-10 tutorial](#). If you want to use TensorFlow or other deep learning frameworks, please find corresponding CIFAR-10 tutorial for that framework and run it.

If you successfully run the tutorial, write “Success” in the report and otherwise please write “Failed”. For this assignment report, only this single word is required, no details/accuracies/codes are needed. The assignment intends to get you familiar with the deep learning library. If you have already trained some deep classifiers using any deep learning library before, you can safely ignore the tutorial and directly write “Success” in the report.

解: Success.

参考文献

- [1] Dimitri P. Bertsekas. Nonlinear Programming: 3rd Edition[M]. Athena Scientific, 2016. ISBN: 978-1-886529-05-2, <http://www.athenasc.com/nonlinbook.html>.
- [2] Bouvrie, Jake. Notes on Convolutional Neural Networks. CogPrints.org, 2006. http://cogprints.org/5869/1/cnn_tutorial.pdf.