

计算两层神经网络参数迭代的 JACOBI 矩阵：以 MNIST 数据集为例

ZEYU XIE¹, ANGXIU NI^{2,3}

1. 变量定义

1.1. 输入输出.

- (a) 学习率 $\alpha = 0.1$
- (b) 一个 batch 的大小 $k = 60000$
- (c) 输入层的输入 $X \in \mathbb{R}^{k \times 784}$

$$(1) \quad X = \begin{bmatrix} x^1 \\ x^2 \\ \vdots \\ x^k \end{bmatrix} = \begin{bmatrix} x_1^1 & x_2^1 & \cdots & x_{784}^1 \\ x_1^2 & x_2^2 & \cdots & x_{784}^2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^k & x_2^k & \cdots & x_{784}^k \end{bmatrix}$$

- (d) 隐藏层的输入 $A_1 \in \mathbb{R}^{k \times 50}$

$$(2) \quad A_1 = \begin{bmatrix} a_1^1 & a_2^1 & \cdots & a_{50}^1 \\ a_1^2 & a_2^2 & \cdots & a_{50}^2 \\ \vdots & \vdots & \ddots & \vdots \\ a_1^k & a_2^k & \cdots & a_{50}^k \end{bmatrix}$$

- (e) 隐藏层的输出 $Z_1 \in \mathbb{R}^{k \times 50}$

$$(3) \quad Z_1 = \begin{bmatrix} z_1^1 & z_2^1 & \cdots & z_{50}^1 \\ z_1^2 & z_2^2 & \cdots & z_{50}^2 \\ \vdots & \vdots & \ddots & \vdots \\ z_1^k & z_2^k & \cdots & z_{50}^k \end{bmatrix}$$

¹ DEPARTMENT OF MATHEMATICS, TSINGHUA UNIVERSITY, BEIJING, CHINA.

² DEPARTMENT OF MATHEMATICS, UNIVERSITY OF CALIFORNIA, IRVINE, USA

³ YAU MATHEMATICAL SCIENCES CENTER, TSINGHUA UNIVERSITY, BEIJING, CHINA.

E-mail address: niangxiu@gmail.com.

Date: 2024 年 4 月 5 日.

(f) 输出层的输入 $A_2 \in \mathbb{R}^{k \times 10}$ ($Z_1 = A_2$, 是同一个矩阵)

$$(4) \quad A_2 = \begin{bmatrix} a'_1{}^1 & a'_2{}^1 & \cdots & a'_{10}{}^1 \\ a'_1{}^2 & a'_2{}^2 & \cdots & a'_{10}{}^2 \\ \vdots & \vdots & \ddots & \vdots \\ a'_1{}^k & a'_2{}^k & \cdots & a'_{10}{}^k \end{bmatrix}$$

(g) 输出层的输出 $Y \in \mathbb{R}^{k \times 10}$

$$(5) \quad Y = \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^k \end{bmatrix} = \begin{bmatrix} y_1^1 & y_2^1 & \cdots & y_{10}^1 \\ y_1^2 & y_2^2 & \cdots & y_{10}^2 \\ \vdots & \vdots & \ddots & \vdots \\ y_1^k & y_2^k & \cdots & y_{10}^k \end{bmatrix}$$

(h) 正确答案 $T \in \mathbb{R}^{k \times 10}$

$$(6) \quad T = \begin{bmatrix} t^1 \\ t^2 \\ \vdots \\ t^k \end{bmatrix} = \begin{bmatrix} t_1^1 & t_2^1 & \cdots & t_{10}^1 \\ t_1^2 & t_2^2 & \cdots & t_{10}^2 \\ \vdots & \vdots & \ddots & \vdots \\ t_1^k & t_2^k & \cdots & t_{10}^k \end{bmatrix}$$

1.2. 神经网络参数.

(a) 输入层到隐藏层的权重 $W_1 \in \mathbb{R}^{784 \times 50}$

$$(7) \quad W_1 = \begin{bmatrix} w_1^1 & w_2^1 & \cdots & w_{50}^1 \\ w_1^2 & w_2^2 & \cdots & w_{50}^2 \\ \vdots & \vdots & \ddots & \vdots \\ w_1^{784} & w_2^{784} & \cdots & w_{50}^{784} \end{bmatrix}$$

(b) 输入层到隐藏层的偏置 $b_1 \in \mathbb{R}^{50}$

$$(8) \quad b_1 = [b_1^1 \quad b_2^1 \quad \cdots \quad b_{50}^1]$$

(c) 隐藏层到输出层的权重 $W_2 \in \mathbb{R}^{50 \times 10}$

$$(9) \quad W_2 = \begin{bmatrix} w_1^1 & w_2^1 & \cdots & w_{10}^1 \\ w_1^2 & w_2^2 & \cdots & w_{10}^2 \\ \vdots & \vdots & \ddots & \vdots \\ w_1^{50} & w_2^{50} & \cdots & w_{10}^{50} \end{bmatrix}$$

(d) 隐藏层到输出层的偏置 $b_2 \in \mathbb{R}^{10}$

$$(10) \quad b_2 = [b_2^1 \quad b_2^2 \quad \cdots \quad b_{10}^1]$$

1.3. 神经网络参数的处理.

(a) 第 t 步迭代的参数向量 $\theta^t \in \mathbb{R}^{39760}$ ¹

$$(11) \quad \theta^t = (\theta_1^t, \theta_2^t, \dots, \theta_{39760}^t)$$

(b) 每一步迭代的参数向量的 Jacobi 矩阵 $Df(\theta^t) \in \mathbb{R}^{39760 \times 39760}$

$$(12) \quad Df(\theta^t) = \begin{bmatrix} \frac{\partial \theta_1^{t+1}}{\partial \theta_1^t} & \frac{\partial \theta_1^{t+1}}{\partial \theta_2^t} & \cdots & \frac{\partial \theta_1^{t+1}}{\partial \theta_{39760}^t} \\ \frac{\partial \theta_2^{t+1}}{\partial \theta_1^t} & \frac{\partial \theta_2^{t+1}}{\partial \theta_2^t} & \cdots & \frac{\partial \theta_2^{t+1}}{\partial \theta_{39760}^t} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \theta_{39760}^{t+1}}{\partial \theta_1^t} & \frac{\partial \theta_{39760}^{t+1}}{\partial \theta_2^t} & \cdots & \frac{\partial \theta_{39760}^{t+1}}{\partial \theta_{39760}^t} \end{bmatrix}$$

2. 神经网络的结构

2.1. **输入层.** 输入层的维度为 784, 即 28×28 的图片的 784 个像素。

2.2. **隐藏层.** 隐藏层的维度为 50。

先进行线性变换

$$(13) \quad A_1 = XW_1 + b_1$$

再经 *sigmoid* 激活函数处理

¹其中 $39760 = 784 \times 50 + 50 + 50 \times 10 + 10$

(14)

$$Z_1 = \text{sigmoid}(A_1) = \text{sigmoid}\left(\begin{bmatrix} a_1^1 & a_2^1 & \cdots & a_{50}^1 \\ a_1^2 & a_2^2 & \cdots & a_{50}^2 \\ \vdots & \vdots & \ddots & \vdots \\ a_1^k & a_2^k & \cdots & a_{50}^k \end{bmatrix}\right) = \begin{bmatrix} \frac{1}{1+\exp(a_1^1)} & \frac{1}{1+\exp(a_2^1)} & \cdots & \frac{1}{1+\exp(a_{50}^1)} \\ \frac{1}{1+\exp(a_1^2)} & \frac{1}{1+\exp(a_2^2)} & \cdots & \frac{1}{1+\exp(a_{50}^2)} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{1}{1+\exp(a_1^k)} & \frac{1}{1+\exp(a_2^k)} & \cdots & \frac{1}{1+\exp(a_{50}^k)} \end{bmatrix}$$

2.3. **输出层.** 输出层的维度为 10

先进行线性变换

(15)

$$A_2 = Z_1 W_2 + b_2$$

再经过 *softmax* 层

(16)

$$Y = \text{softmax}(A_2) = \text{softmax}\left(\begin{bmatrix} a'^1 \\ a'^2 \\ \vdots \\ a'^k \end{bmatrix}\right) = \begin{bmatrix} \frac{\exp(a'^1_1)}{\sum_{l=1}^{10} \exp(a'^1_l)} & \frac{\exp(a'^1_2)}{\sum_{l=1}^{10} \exp(a'^1_l)} & \cdots & \frac{\exp(a'^1_{10})}{\sum_{l=1}^{10} \exp(a'^1_l)} \\ \frac{\exp(a'^2_1)}{\sum_{l=1}^{10} \exp(a'^2_l)} & \frac{\exp(a'^2_2)}{\sum_{l=1}^{10} \exp(a'^2_l)} & \cdots & \frac{\exp(a'^2_{10})}{\sum_{l=1}^{10} \exp(a'^2_l)} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\exp(a'^k_1)}{\sum_{l=1}^{10} \exp(a'^k_l)} & \frac{\exp(a'^k_2)}{\sum_{l=1}^{10} \exp(a'^k_l)} & \cdots & \frac{\exp(a'^k_{10})}{\sum_{l=1}^{10} \exp(a'^k_l)} \end{bmatrix}$$

2.4. **损失函数.** 我们用交叉熵作为损失函数，即

(17)

$$L = -\frac{1}{k} \sum_{i=1}^k \sum_{j=1}^{10} y_j^i \log y_j^i$$

3. 梯度和 JACOBI 矩阵

第 t 步迭代的参数为

(18)

$$\theta^t = (W_1(t), b_1(t), W_2(t), b_2(t)) = (\theta_1^t, \theta_2^t, \cdots, \theta_{39760}^t)$$

此时对应的梯度为

(19)

$$\text{grad}(t) = \left(\frac{\partial L}{\partial \theta_1^t}, \frac{\partial L}{\partial \theta_2^t}, \cdots, \frac{\partial L}{\partial \theta_{39760}^t}\right) = \left(\frac{\partial L}{\partial W_1^t}, \frac{\partial L}{\partial b_1^t}, \frac{\partial L}{\partial W_2^t}, \frac{\partial L}{\partial b_2^t}\right)$$

神经网络迭代过程即

(20)

$$\theta^{t+1} = \theta^t - \alpha \cdot \text{grad}(t)$$

也即

$$\begin{aligned}
 W_1^{t+1} &= W_1^t - \alpha \cdot \frac{\partial L}{\partial W_1^t} \\
 b_1^{t+1} &= b_1^t - \alpha \cdot \frac{\partial L}{\partial b_1^t} \\
 W_2^{t+1} &= W_2^t - \alpha \cdot \frac{\partial L}{\partial W_2^t} \\
 b_2^{t+1} &= b_2^t - \alpha \cdot \frac{\partial L}{\partial b_2^t}
 \end{aligned}
 \tag{21}$$

因此迭代的 Jacobi 矩阵为

$$J(t) = \begin{bmatrix} \frac{\partial \theta_1^{t+1}}{\partial \theta_1^t} & \frac{\partial \theta_1^{t+1}}{\partial \theta_2^t} & \cdots & \frac{\partial \theta_1^{t+1}}{\partial \theta_{39760}^t} \\ \frac{\partial \theta_2^{t+1}}{\partial \theta_1^t} & \frac{\partial \theta_2^{t+1}}{\partial \theta_2^t} & \cdots & \frac{\partial \theta_2^{t+1}}{\partial \theta_{39760}^t} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \theta_{39760}^{t+1}}{\partial \theta_1^t} & \frac{\partial \theta_{39760}^{t+1}}{\partial \theta_2^t} & \cdots & \frac{\partial \theta_{39760}^{t+1}}{\partial \theta_{39760}^t} \end{bmatrix}
 \tag{22}$$

也即

$$J(t) = \begin{bmatrix} \frac{\partial W_1^{t+1}}{\partial W_1^t} & \frac{\partial W_1^{t+1}}{\partial b_1^t} & \frac{\partial W_1^{t+1}}{\partial W_2^t} & \frac{\partial W_1^{t+1}}{\partial b_2^t} \\ \frac{\partial b_1^{t+1}}{\partial W_1^t} & \frac{\partial b_1^{t+1}}{\partial b_1^t} & \frac{\partial b_1^{t+1}}{\partial W_2^t} & \frac{\partial b_1^{t+1}}{\partial b_2^t} \\ \frac{\partial W_2^{t+1}}{\partial W_1^t} & \frac{\partial W_2^{t+1}}{\partial b_1^t} & \frac{\partial W_2^{t+1}}{\partial W_2^t} & \frac{\partial W_2^{t+1}}{\partial b_2^t} \\ \frac{\partial b_2^{t+1}}{\partial W_1^t} & \frac{\partial b_2^{t+1}}{\partial b_1^t} & \frac{\partial b_2^{t+1}}{\partial W_2^t} & \frac{\partial b_2^{t+1}}{\partial b_2^t} \end{bmatrix}
 \tag{23}$$

4. JACOBI 矩阵的计算

4.1. 基本性质.

4.1.1. *sigmoid* 函数的 *Jacobi* 矩阵. 假设 $A \in \mathbb{R}^{n \times m}$ 是同一个矩阵, 则 $\frac{\partial(\text{sigmoid}(A))}{\partial A}$ 的计算公式为

$$\frac{\partial(\text{sigmoid}(A))}{\partial A} = \text{diag}\left\{\frac{d(\text{sigmoid}(x_1^1))}{dx_1^1}, \frac{d(\text{sigmoid}(x_2^1))}{dx_2^1}, \dots, \frac{d(\text{sigmoid}(x_m^n))}{dx_m^n}\right\}
 \tag{24}$$

4.1.2. *softmax* 函数的 *Jacobi* 矩阵. 假设 A 是同一个矩阵, 则 $\frac{\partial(\text{softmax}(A))}{\partial A}$ 的计算公式为

$$\frac{\partial(\text{softmax}(A))}{\partial A} = \text{softmax}(A) \cdot (I - \text{softmax}(A))
 \tag{25}$$

$$\begin{aligned}
(26) \quad \frac{\partial W_1(t+1)}{\partial W_1(t)} &= \frac{\partial(W_1(t) - \alpha \cdot \frac{\partial L}{\partial W_1(t)})}{\partial W_1(t)} \\
&= I - \alpha \frac{\partial^2 L}{\partial W_1 \partial W_1} \\
&= I - \alpha \frac{\partial}{\partial W_1} \left(\frac{\partial L}{\partial W_1} \right) \\
&= I - \alpha \frac{\partial}{\partial W_1} \left(\frac{\partial L}{\partial Y} \frac{\partial Y}{\partial A_2} \frac{\partial A_2}{\partial W_1} \right) \\
&= I - \alpha \frac{\partial}{\partial W_1} \left(\frac{\partial L}{\partial Y} \frac{\partial Y}{\partial A_2} \right) W_2^T Z_1^T \\
&= I - \alpha \left(\frac{\partial}{\partial W_1} \left(\frac{\partial L}{\partial Y} \right) \frac{\partial Y}{\partial A_2} + \frac{\partial L}{\partial Y} \frac{\partial}{\partial W_1} \left(\frac{\partial Y}{\partial A_2} \right) \right) W_2^T Z_1^T
\end{aligned}$$

$$(27) \quad \frac{\partial W_1'}{\partial b_1} = -\alpha \left(softmax'(A_2) W_2^T sigmoid'(A_1) \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 1 \end{bmatrix} \right) X$$

$$(28) \quad \frac{\partial W_1'}{\partial W_2} = -\alpha \left(softmax'(A_2) W_2^T sigmoid'(A_1) \begin{bmatrix} x_1^1 & x_2^1 & \cdots & x_{784}^1 \\ x_1^2 & x_2^2 & \cdots & x_{784}^2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^k & x_2^k & \cdots & x_{784}^k \end{bmatrix} \right)$$

$$\frac{\partial W_1'}{\partial b_2} = -\alpha \left(softmax'(A_2) W_2^T sigmoid'(A_1) \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 1 \end{bmatrix} \right)$$

$$\begin{aligned}
(29) \quad \frac{\partial b_1'}{\partial W_1} &= -\alpha \left(softmax'(A_2) W_2^T sigmoid'(A_1) \begin{bmatrix} x_1^1 & x_2^1 & \cdots & x_{784}^1 \\ x_1^2 & x_2^2 & \cdots & x_{784}^2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^k & x_2^k & \cdots & x_{784}^k \end{bmatrix} \right) \\
\frac{\partial b_1'}{\partial b_1} &= -\alpha \left(softmax'(A_2) W_2^T sigmoid'(A_1) \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 1 \end{bmatrix} \right)
\end{aligned}$$

$$\begin{aligned}
 (30) \quad \frac{\partial b_1'}{\partial W_2} &= -\alpha \left(\text{softmax}'(A_2) W_2^T \text{sigmoid}'(A_1) \begin{bmatrix} x_1^1 & x_2^1 & \cdots & x_{784}^1 \\ x_1^2 & x_2^2 & \cdots & x_{784}^2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^k & x_2^k & \cdots & x_{784}^k \end{bmatrix} \right) \\
 \frac{\partial b_1'}{\partial b_2} &= -\alpha \left(\text{softmax}'(A_2) W_2^T \text{sigmoid}'(A_1) \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 1 \end{bmatrix} \right)
 \end{aligned}$$

$$\begin{aligned}
 (31) \quad \frac{\partial W_2'}{\partial W_1} &= -\alpha \left(\text{softmax}'(A_2) \begin{bmatrix} a_1^1 & a_2^1 & \cdots & a_{50}^1 \\ a_1^2 & a_2^2 & \cdots & a_{50}^2 \\ \vdots & \vdots & \ddots & \vdots \\ a_1^k & a_2^k & \cdots & a_{50}^k \end{bmatrix} \right) \\
 \frac{\partial W_2'}{\partial b_1} &= -\alpha \left(\text{softmax}'(A_2) \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 1 \end{bmatrix} \right)
 \end{aligned}$$

$$\begin{aligned}
 (32) \quad \frac{\partial W_2'}{\partial W_2} &= I - \alpha \frac{\partial^2 Y}{\partial W_2 \partial W_2} \\
 &= I - \alpha \frac{\partial}{\partial W_2} (\text{softmax}'(A_2) A_1) \\
 &= I - \alpha \frac{\partial}{\partial W_2} (\text{softmax}'(A_2)) A_1 \\
 &= I - \alpha \left(\frac{\partial \text{softmax}'(A_2)}{\partial A_2} \frac{\partial A_2}{\partial W_2} \right) A_1 \\
 &= I - \alpha \left(\text{softmax}'(A_2) \frac{\partial A_2}{\partial W_2} \right) A_1 \\
 &= I - \alpha \left(\text{softmax}'(A_2) \begin{bmatrix} a_1^1 & a_2^1 & \cdots & a_{50}^1 \\ a_1^2 & a_2^2 & \cdots & a_{50}^2 \\ \vdots & \vdots & \ddots & \vdots \\ a_1^k & a_2^k & \cdots & a_{50}^k \end{bmatrix} \right) A_1
 \end{aligned}$$

$$(33) \quad \frac{\partial W_2'}{\partial b_2} = -\alpha \left(softmax'(A_2) \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 1 \end{bmatrix} \right)$$

$$(34) \quad \frac{\partial b_2'}{\partial W_1} = -\alpha \left(softmax'(A_2) \begin{bmatrix} a_1^1 & a_2^1 & \cdots & a_{50}^1 \\ a_1^2 & a_2^2 & \cdots & a_{50}^2 \\ \vdots & \vdots & \ddots & \vdots \\ a_1^k & a_2^k & \cdots & a_{50}^k \end{bmatrix} \right)$$

$$\frac{\partial b_2'}{\partial b_1} = -\alpha \left(softmax'(A_2) \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 1 \end{bmatrix} \right)$$

$$(35) \quad \frac{\partial b_2'}{\partial W_2} = -\alpha \left(softmax'(A_2) \begin{bmatrix} a_1^1 & a_2^1 & \cdots & a_{50}^1 \\ a_1^2 & a_2^2 & \cdots & a_{50}^2 \\ \vdots & \vdots & \ddots & \vdots \\ a_1^k & a_2^k & \cdots & a_{50}^k \end{bmatrix} \right)$$

$$\frac{\partial b_2'}{\partial b_2} = -\alpha \left(softmax'(A_2) \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 1 \end{bmatrix} \right)$$

将以上 $4 \times 4 = 16$ 个 Jacobi 矩阵计算公式代入 $Df(W_1, b_1, W_2, b_2)$ 的定义, 即得到 $Df(W_1, b_1, W_2, b_2)$ 的表达式