

向量函数的导数

2021年2月7日

1 一元函数的导数

导数对于函数而言的，具体而言，它是一个极限：

定义1.1 (导数) 对于函数 $f: D \rightarrow \mathbb{R}$ ，设 x_0 是 D 的一个内点，若极限

$$\lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h} \quad (1.1)$$

存在，则称 f 在 x_0 处可导，将这个极限记为 $f'(x_0)$ ，并称 $f'(x_0)$ 是 f 在 x_0 处的导数。

对于一元函数，有一个有用的事实：

定理1.1 (复合函数的导数) 对于函数 $f: A \rightarrow B, g: C \rightarrow \mathbb{R}$ ，设 $x \in A$ ，记 $y = f(x)$ ，那么只要 f 在 x 可导并且 g 在 y 可导，就一定可以说明复合函数 $g \circ f$ 在 x 可导，并且我们有

$$(g \circ f)'(x) = g'(y)f'(x) \quad (1.2)$$

对于复合函数的写法，我们约定 $g \circ f: x \mapsto g(f(x))$ 。

证明： 设 x' 在 x 某邻域内：

$$f(x') = f(x) + f'(x)(x' - x) + o(x' - x) \quad (1.3)$$

由此可知当 x' 足够接近 x 时，存在常数 C ，使得

$$|f(x') - f(x)| \leq C|x' - x| \quad (1.4)$$

于是有

$$\begin{aligned} g(f(x')) &= g(f(x)) + g'(f(x))(f(x') - f(x)) + o(f(x') - f(x)) \\ &= g(f(x)) + g'(f(x))f'(x)(x' - x) + g'(f(x))o(x' - x) + o(x' - x) \\ &= g(f(x)) + g'(f(x))f'(x)(x' - x) + o(x' - x) \end{aligned} \quad (1.5)$$

观察式 1.5 我们发现： $g'(f(x))f'(x)$ 就是 $g \circ f$ 在 x 处的导数。 \square

2 向量函数及其导数

首先介绍向量：

定义2.1 (向量) 若干个数字写成一行叫做**行向量**，若干个数字写成一列叫做**列向量**，行向量和列向量统称**向量**。

例2.1 我们称

$$\begin{bmatrix} 109 \\ 24 \end{bmatrix}$$

是一个列向量，并且说它是 2 维的。称

$$\begin{bmatrix} 10 & 201 & 34 & 75 \end{bmatrix}$$

是一个行向量，并且说它是 4 维的。它有几个元素就说它是几维的。我们暂不涉及无穷维的向量。

那么输出值（也可能有多个）依赖于多个输入值（受其影响）的函数就可以说是一个向量函数：

定义2.2（向量函数） 以向量作为输入，以向量作为输出的函数叫做向量函数。

例2.2 令

$$f: \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \mapsto \begin{bmatrix} x_1^2 + x_2^2 \\ x_1^2 - x_2^2 \\ x_1 x_2 \end{bmatrix}$$

则这里的 f 就是一个堂堂正正的向量函数。

由于向量函数有多个输出，以及多个输入，所以要完全刻画它的导数性质，非得要用到一个矩阵不可，具体地，设 $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$ ，对于 \mathbb{R}^m 中的一点 \mathbf{x} ，记 $\mathbf{y} = f(\mathbf{x})$ ，我们这么样描述 f 的导数：

$$\mathbf{J}^T(f) = \begin{bmatrix} \frac{\partial \mathbf{y}}{\partial x_1} \\ \frac{\partial \mathbf{y}}{\partial x_2} \\ \vdots \\ \frac{\partial \mathbf{y}}{\partial x_m} \end{bmatrix} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \cdots & \frac{\partial y_n}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_n}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_m} & \frac{\partial y_2}{\partial x_m} & \cdots & \frac{\partial y_n}{\partial x_m} \end{bmatrix} \quad (2.1)$$

如果还有一个函数 $g: \mathbb{R}^n \rightarrow \mathbb{R}^q$ ，并且我们还知道

$$\mathbf{J}^T(g) = \begin{bmatrix} \frac{\partial z_1}{\partial y_1} & \frac{\partial z_2}{\partial y_1} & \cdots & \frac{\partial z_q}{\partial y_1} \\ \frac{\partial z_1}{\partial y_2} & \frac{\partial z_2}{\partial y_2} & \cdots & \frac{\partial z_q}{\partial y_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_1}{\partial y_n} & \frac{\partial z_2}{\partial y_n} & \cdots & \frac{\partial z_q}{\partial y_n} \end{bmatrix} \quad (2.2)$$

前面我们演示了：对于一元函数，求导运算的链式法则成立，现在用矩阵来表述多元函数的导数，我们发现链式法则对于多元函数的求导运算也成立，具体地：设 f 将 \mathbb{R}^m 中的 \mathbf{x} 映射为 \mathbb{R}^n 中的 \mathbf{y} ，而后 g 又将 \mathbf{y} 映到 \mathbb{R}^q 中的 \mathbf{z} ，类似于：

$$\mathbf{x} := \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \xrightarrow{f} \mathbf{y} := \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \xrightarrow{g} \mathbf{z} := \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_q \end{bmatrix} \quad (2.3)$$

那么，我们有

$$\mathbf{J}_x^T(g \circ f) = \mathbf{J}_x^T(f) \mathbf{J}_y^T(g) \quad (2.4)$$

也就是说

$$\begin{bmatrix} \frac{\partial z_1}{\partial x_1} & \frac{\partial z_2}{\partial x_1} & \cdots & \frac{\partial z_q}{\partial x_1} \\ \frac{\partial z_1}{\partial x_2} & \frac{\partial z_2}{\partial x_2} & \cdots & \frac{\partial z_q}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_1}{\partial x_m} & \frac{\partial z_2}{\partial x_m} & \cdots & \frac{\partial z_q}{\partial x_m} \end{bmatrix} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \cdots & \frac{\partial y_n}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_n}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_m} & \frac{\partial y_2}{\partial x_m} & \cdots & \frac{\partial y_n}{\partial x_m} \end{bmatrix} \begin{bmatrix} \frac{\partial z_1}{\partial y_1} & \frac{\partial z_2}{\partial y_1} & \cdots & \frac{\partial z_q}{\partial y_1} \\ \frac{\partial z_1}{\partial y_2} & \frac{\partial z_2}{\partial y_2} & \cdots & \frac{\partial z_q}{\partial y_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_1}{\partial y_n} & \frac{\partial z_2}{\partial y_n} & \cdots & \frac{\partial z_q}{\partial y_n} \end{bmatrix} \quad (2.5)$$

我们可以通过回忆矩阵乘法的运算规则和多元函数的求导法则来理解这个式子。

3 梯度下降

梯度下降是指：在给定的学习率 α 下，按照梯度的指引，一步步地调整函数的输入，使得其输出减少的一种函数最小化的方法。所谓的梯度其实就是导数，除了笔划数量不同并无差别。那么既然要「按照梯度的指引」肯定首先要知道梯度，而经过我们刚才对向量函数的导数知识的学习，相信这不在话下。

3.1 一个简单的例子

考虑函数

$$f(x_1, x_2) = (x_1 x_2 - 3)^2 + 1 \quad (3.1)$$

找到一点 (x_1, x_2) 使得 $f(x_1, x_2)$ 尽可能地小。

我们可将它看做一个 \mathbb{R}^2 到 \mathbb{R} 的映射，首先来计算梯度：

$$\mathbf{J}(f) = \begin{bmatrix} 2x_2(x_1 x_2 - 3) & 2x_1(x_1 x_2 - 3) \end{bmatrix} \quad (3.2)$$

然后我们首先**随机**选一点（随机就是随便的意思），比如说，我选 $\mathbf{x}^{(1)} = (1, 1)$ ，先代入看看此时的函数值是多少：

$$f(\mathbf{x}^{(1)}) = (1 - 3)^2 + 1 = 5 \quad (3.3)$$

设定一个学习率 $\alpha = 1 \times 10^{-2}$ ，然后计算 f 在 $\mathbf{x}^{(1)}$ 的梯度：

$$\mathbf{J}_{\mathbf{x}^{(1)}}(f) = \begin{bmatrix} -4 & -4 \end{bmatrix} \quad (3.4)$$

然后更新它：

$$\begin{aligned} \mathbf{x}^{(2)} &= \mathbf{x}^{(1)} - \alpha \mathbf{J}_{\mathbf{x}^{(1)}}(f) \\ &= \begin{bmatrix} 1 & 1 \end{bmatrix} - 1 \times 10^{-2} \begin{bmatrix} -4 & -4 \end{bmatrix} \\ &= \begin{bmatrix} 1.04 & 1.04 \end{bmatrix} \end{aligned} \quad (3.5)$$

再看此时的函数值：

$$f(\mathbf{x}^{(2)}) \approx 4.68025856 \quad (3.6)$$

再做一次更新：

$$\begin{aligned} \mathbf{x}^{(3)} &= \mathbf{x}^{(2)} - \alpha \mathbf{J}_{\mathbf{x}^{(2)}}(f) \\ &= \begin{bmatrix} 1.04 & 1.04 \end{bmatrix} - 1 \times 10^{-2} \begin{bmatrix} -1.995136 & -1.995136 \end{bmatrix} \\ &= \begin{bmatrix} 1.05995136 & 1.05995136 \end{bmatrix} \end{aligned} \quad (3.7)$$

再看此时的函数值：

$$f(\mathbf{x}^{(3)}) = 4.5212639385 \quad (3.8)$$

这个过程的重复就是梯度下降。有时候学习率选得太大，或者初始点选得不合适，可能会导致效果不好。

3.2 神经网络

神经网络就可以看做是多个向量函数的复合，可以认为，神经网络的基本构成 (building-block) 是层，与之对应的数学对象也就是向量函数，为此，现在我们就来定义几个「层」，第一层添加哑节点并且对输入做线性组合：

$$f_1 : M[4, 3] \times \mathbb{R}^2 \longrightarrow \mathbb{R}^4$$

$$\left(\boldsymbol{\theta}^{(1)}, \mathbf{x}^{(1)} \right) \mapsto \mathbf{x}^{(2)} := \boldsymbol{\theta}^{(1)} \begin{bmatrix} 1 \\ \mathbf{x}^{(1)} \end{bmatrix} \quad (3.9)$$

其中 $M[4, 3]$ 表示全体 4 行 3 列矩阵构成的集合，第三层对第二层的输出做按元素 Sigmoid 处理，并且添加一个偏置：

$$f_2 : \mathbb{R}^4 \longrightarrow \mathbb{R}^5$$

$$\mathbf{x}^{(2)} \mapsto \mathbf{x}^{(3)} := \begin{bmatrix} 1 \\ \text{Sigmoid}(\mathbf{x}^{(2)}) \end{bmatrix} \quad (3.10)$$

第三层对第二层的输出做线性组合并且进行 Sigmoid 处理：

$$f_3 : M[1, 5] \times \mathbb{R}^5 \longrightarrow \mathbb{R}$$

$$\left(\boldsymbol{\theta}^{(2)}, \mathbf{x}^{(3)} \right) \mapsto x^{(4)} := \text{Sigmoid} \left(\left(\boldsymbol{\theta}^{(2)} \right)^\top \mathbf{x}^{(3)} \right) \quad (3.11)$$

如果将这个神经网络的神经元一个个画出来，我们会看到图 3-1：

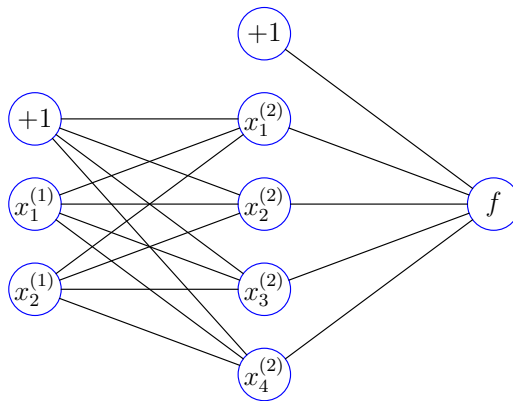


图 3-1: 神经网络示意图

如果仅仅是做预测，即假设正确的参数 $\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}$ 是已知的，那么对于一个数据点 $\mathbf{x}^{(1)}$ ，仅仅计算到 f_3 的输出值，是足够的，然而如果我们不知道参数 $\boldsymbol{\theta}$ ，就要再加一层用来计算误差，当然这仅仅是用来计算误差和做梯度下降，做预测的时候用不到这一层：

$$f_4 : \mathbb{R} \longrightarrow \{1, 0\}$$

$$x^{(4)} \mapsto \frac{1}{2} (x^{(4)} - t(\mathbf{x}^{(1)}))^2 \quad (3.12)$$

其中 $t(\cdot)$ 对应真实标签。

为了更加清楚地搞清楚对于这种多层的神经网络，梯度下降（也就是后馈传播）应该怎么做，我们可以按照数值与数值之间的影响关系画出一张图（图 3-2）：

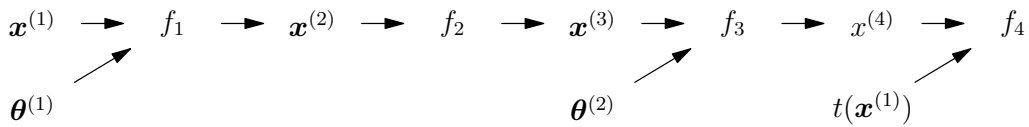


图 3-2: 计算图

从图 3-2中, 我们可以看到: $\mathbf{x}^{(1)}$ 和 $\boldsymbol{\theta}^{(1)}$ 共同透过 f_1 影响着 $\mathbf{x}^{(2)}$, 而 $\mathbf{x}^{(2)}$ 透过 f_2 影响 $\mathbf{x}^{(3)}$, 由于 f_2 其实是按元素 Sigmoid, 所以没有参数, 然后 $\mathbf{x}^{(3)}$ 和 $\boldsymbol{\theta}^{(2)}$ 共同影响 $\mathbf{x}^{(4)}$, 也就是神经网络最终输出的概率, 而 $t(\mathbf{x}^{(1)})$ (也就是真实标签) 和 $\mathbf{x}^{(4)}$ (也就是输出概率值) 共同影响着 f_4 的计算结果 (也就是损失值). 我们希望通过链式法则计算出 f_4 的输出值分别关于 $\boldsymbol{\theta}^{(2)}$ 与 $\boldsymbol{\theta}^{(1)}$ 的导数 (也就是梯度), 然后根据这个梯度和学习率一起, 去微调之, 从而使得 f_4 的输出值降低.

3.3 梯度的计算

我们的目的是找到 $\boldsymbol{\theta}^{(1)}$ 和 $\boldsymbol{\theta}^{(2)}$ 使得对于 $\mathbf{x}^{(1)}$ 和 $t(\mathbf{x}^{(1)})$ 的取值分别如表 3-1所示的时候

表 3-1: 自变量和标签的取值

$\mathbf{x}^{(1)}$	$t(\mathbf{x}^{(1)})$
(0, 0)	0
(0, 1)	1
(1, 0)	1
(1, 1)	0

f_4 的输出都比较小. 为此, 我们这样做: 固定 $\mathbf{x}^{(1)}$ 和 $t(\mathbf{x}^{(1)})$, 然后 f_4 关于各 $\boldsymbol{\theta}$ 的梯度, 然后根据这些梯度去微调 $\boldsymbol{\theta}$ 使得 f_4 下降 (梯度下降), 然后让 $\mathbf{x}^{(1)}$ 和 $t(\mathbf{x}^{(1)})$ 都取到表格的下一行 (如果已是最后一行就返回第一行), 重复这些步骤直到 f_4 的输出相当小.

观察图3-2, 我们发现: 顺着箭头方向的计算都是比较容易进行, 稍微有点麻烦的其实是求导, 如果只是求 f 对于 \mathbf{x} 的导数还好说, 稍微有点麻烦的是对于那些 $\boldsymbol{\theta}$ 的导数, $\boldsymbol{\theta}$ 是矩阵, 而我们只知道在「以向量作为输入」且「以向量作为输出」的情况下如何求导, 那么如果是「以矩阵作为」输入的时候该怎么办呢? 我们主要是将矩阵也视作一个折叠起来的向量, 要将矩阵写成一个向量, 我们首先取矩阵的第一行, 然后取第二行接在第一行的后头, 然后第三行接在第二行的后头, 这样一来, 就建立了有限维矩阵和向量之间的一一对应关系. 譬如说, 对于矩阵

$$\begin{bmatrix} 14 & 32 & 66 \\ 72 & 86 & 97 \end{bmatrix} \quad (3.13)$$

我们就将它写成一个行向量

$$[14 \quad 32 \quad 66 \quad 72 \quad 86 \quad 97] \quad (3.14)$$

有了这种对应关系, 就可以求一个函数的向量输出关于矩阵输入的导数, 例如说, 对于函数

$$\mathbf{g} : M[2, 3] \rightarrow \mathbb{R}^2$$

$$\mathbf{x} := \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \end{bmatrix} \mapsto \begin{bmatrix} g_1(\mathbf{x}) \\ g_2(\mathbf{x}) \end{bmatrix} \quad (3.15)$$

其中 $x_0^{(3)} = 1$. 然后计算雅克比矩阵, 首先计算关于 $\mathbf{x}^{(3)}$ 的

$$\mathbf{J}_{\mathbf{x}^{(3)}}(f_3) = \begin{bmatrix} \theta_{1,0}^{(2)}s(1-s) & \theta_{1,1}^{(2)}s(1-s) & \theta_{1,2}^{(2)}s(1-s) & \theta_{1,3}^{(2)}s(1-s) & \theta_{1,4}^{(2)}s(1-s) \end{bmatrix} \quad (3.24)$$

然后计算关于 $\boldsymbol{\theta}^{(2)}$ 的

$$\mathbf{J}_{\boldsymbol{\theta}^{(2)}}(f_3) = \begin{bmatrix} s(1-s) & x_1^{(3)}s(1-s) & x_1^{(3)}s(1-s) & x_1^{(3)}s(1-s) & x_1^{(3)}s(1-s) \end{bmatrix} \quad (3.25)$$

这里的 s 等于 $\text{Sigmoid}\left(\sum_{j=0}^4 \theta_{1,j}^{(2)}x_j^{(3)}\right)$. 可以看到梯度 (雅克比矩阵) 的计算主要是在前向传播时进行. 接下来我们主要关心 f_4 关于 $x^{(4)}$ 的梯度, 首先有

$$f_4(x^{(4)}) = \frac{1}{2} (x^{(4)} - t(\mathbf{x}^{(1)}))^2 \quad (3.26)$$

那么就可以写出导数:

$$\frac{df_4}{dx^{(4)}} = x^{(4)} - t(\mathbf{x}^{(1)}) \quad (3.27)$$

导数的符号计算暂时告一段落.

3.4 编程实现

我们并不打算给出哪怕是稍微完整一点的神经网络的实现, 我们主要是想展现损失函数随着梯度下降的进行而下降的现象, 这样做对于理解梯度下降有帮助.

首先是初始化一个神经网络, 主要是初始化它的参数:

```

1 import numpy as np
2 from scipy.linalg import block_diag
3
4 class NeuralNet:
5
6     def __init__(self):
7
8         self.x1 = np.zeros(shape=(2,1), dtype=np.float)
9         self.x2 = np.zeros(shape=(4,1), dtype=np.float)
10        self.x3 = np.zeros(shape=(5,1), dtype=np.float)
11        self.x4 = np.zeros(shape=(1,1), dtype=np.float)
12        self.x5 = np.zeros(shape=(1,1), dtype=np.float)
13
14        self.theta1 = np.zeros(shape=(4,3), dtype=np.float)
15        self.theta2 = np.zeros(shape=(1,5), dtype=np.float)
16        self.tx1 = np.zeros(shape=(1,1), dtype=np.float)
17
18        self.jacobian_f1_x1 = np.zeros(shape=(4,2), dtype=np.float)
19        self.jacobian_f1_theta1 = np.zeros(shape=(4,12), dtype=np.float)
20        self.jacobian_f2_x2 = np.zeros(shape=(5,4), dtype=np.float)
21        self.jacobian_f3_x3 = np.zeros(shape=(1,5), dtype=np.float)
22        self.jacobian_f3_theta2 = np.zeros(shape=(1,5), dtype=np.float)
23        self.jacobian_f4_x4 = np.zeros(shape=(1,1), dtype=np.float)

```

符号与我们上面的数学推导所用到的基本一致. 然后定义前馈过程, 主要是计算雅克比矩阵和误差:

```

1 def feed_forward(self, x: np.ndarray, y: np.ndarray) -> np.ndarray:
2     if self.x1.shape == x.shape and self.tx1.shape == y.shape:
3         self.x1 = x

```

```

4         self.tx1 = y
5
6         self.x2 = self.x1
7         self.x2 = np.concatenate(
8             ([[1]], self.x2),
9             axis=0
10        )
11        self.x2 = self.theta1 @ self.x2
12        self.jacobian_f1_x1 = self.theta1[:, [1,2]].copy()
13
14        temp = np.concatenate([[1]], self.x1, axis=0).T
15        self.jacobian_f1_theta1 = block_diag(temp, temp, temp, temp)
16
17        self.x3 = self.x2
18        self.x3 = 1 / (1 + np.exp(0-self.x3))
19        self.x3 = np.concatenate(
20            ([[1]], self.x3, ),
21            axis=0
22        )
23
24        self.jacobian_f2_x2 = (self.x3 * (1 - self.x3)).flatten()
25        self.jacobian_f2_x2 = np.diag(self.jacobian_f2_x2)[: , 1:]
26
27        temp = self.theta2 @ self.x3
28        self.x4 = 1 / (1 + np.exp(0-temp))
29
30        self.jacobian_f3_x3 = self.theta2 * self.x4 * (1-self.x4)
31        self.jacobian_f3_theta2 = self.x3.T * self.x4 * (1-self.x4)
32
33        self.x5 = (1/2)*np.power(self.x4-self.tx1, 2)
34        self.jacobian_f4_x4 = self.x4 - self.tx1
35
36        return self.x5
37    else:
38        raise ValueError()

```

算好了雅克比矩阵就可以进行误差逆传播了：

```

1 def back_propagate(self) -> np.ndarray:
2     theta2_gradient = self.jacobian_f4_x4 @ self.jacobian_f3_theta2
3     theta1_gradient = self.jacobian_f4_x4 @ \
4     self.jacobian_f3_x3 @ \
5     self.jacobian_f2_x2 @ \
6     self.jacobian_f1_theta1
7     theta1_gradient = theta1_gradient.reshape(4,3)
8
9     return (theta1_gradient, theta2_gradient,)

```

误差逆传播给出的其实是梯度，所以我们还要定义一个过程去更新参数：

```

1 def update_parameter(self, lr: float = 0.02) -> None:
2     theta1_gradient, theta2_gradient = self.back_propagate()
3     self.theta1 = self.theta1 - lr * theta1_gradient

```



```
4     self.theta2 = self.theta2 - lr * theta2_gradient
```

当然了，在真正开始训练之前，还有随机设置它的参数，通过一个参数设置方法：

```
1 def set_parameter(self, theta1: np.ndarray, theta2: np.ndarray):
2     self.theta1 = theta1
3     self.theta2 = theta2
```

现在可以调用这些函数了：

```
1 net = NeuralNet()
2 theta1 = np.random.rand(4,3,)
3 theta2 = np.random.rand(1,5,)
4 net.set_parameter(theta1, theta2)
5
6 for i in range(100000):
7     x1 = np.array([[0], [0]])
8     y = np.array([[0]])
9     e1 = net.feed_forward(x1, y)[0,0]
10    net.back_propagate()
11    net.update_parameter()
12
13    x1 = np.array([[0], [1]])
14    y = np.array([[1]])
15    e2 = net.feed_forward(x1, y)[0,0]
16    net.back_propagate()
17    net.update_parameter()
18
19    x1 = np.array([[1], [0]])
20    y = np.array([[1]])
21    e3 = net.feed_forward(x1, y)[0,0]
22    net.back_propagate()
23    net.update_parameter()
24
25    x1 = np.array([[1], [1]])
26    y = np.array([[0]])
27    e4 = net.feed_forward(x1, y)[0,0]
28    net.back_propagate()
29    net.update_parameter()
30
31    if (i+1) % 10000 == 0:
32        error = e1 + e2 + e3 + e4
33        print(error)
```

如果不收敛，建议多试几组随机参数。

4 总结

在本篇文章中，我们主要讲了向量函数的链式求导法则及其应用，尤其是，通过以神经网络为例，讲解了向量函数的求导在误差逆传播算法中的体现，我们衷心地希望您通过读懂这篇文章使得自己无论是在对微积分知识的理解上还是在神经网络梯度下降算法的理解上都能够有看得见的提升！