

一份不太专业的多层感知机介绍

Day Z

目录

摘要	2
1 模型结构	2
1.1 样本与数据	2
1.2 预测与学习	3
1.3 数学模型	4
1.3.1 神经网络	4
1.3.2 损失函数	6
1.3.3 激活函数	7
1.4 通用近似定理	9
2 优化方法	10
2.1 梯度下降法	10
2.2 反向传播公式	12
3 训练过程	15
参考文献	17
代码	18
神经网络	18
激活函数	24
MNIST 测试	25

摘要

多层感知机可用于解决分类与回归问题，它是两个向量集之间的一个映射。根据通用近似定理，多层感知机隐藏层的输出的加权和可被视为某个多元函数的值。经过参数优化，该多元函数能在一定程度上拟合诸多连续函数。

小批量梯度下降法是一种机器学习领域常见的参数优化方法，它可被用于训练多层感知机。梯度下降的难点在于计算损失函数对各层参数的偏导数，而基于链式法则实现的反向传播算法则较好地解决了这一问题。

关键词 多层感知机; 神经网络; 通用近似定理; 梯度下降; 反向传播

1 模型结构

1.1 样本与数据

样本是具体或抽象的事物，是研究对象。分类（Classification）任务是确定给定样本所属的类别，是一种定性的任务，如鸢尾花品种分类。One-hot 向量是仅有一个分量为 1，其余分量为 0 的向量。在样本被划分为 $n(n \geq 2)$ 个类别的多分类任务中，不同的 n 维 One-hot 向量表示不同的类别。从样本的角度进一步看，样本的 One-hot 向量的不同分量实际上对应不同的类别。一个分量可被视为一个真值或概率值，它表明样本是否属于该分量对应的类别。当一个分量为 1 时，样本属于该分量对应的类别，否则不属于。回归（Regression）任务是确定给定样本在特定情形下对应的具体数值或数值组，是一种定量的任务，如波士顿房价预测。

特征（Feature）是样本的属性，其值是对样本本身的描述。特征向量是样本全体共有的一组给定特征的值按给定顺序组成的向量，它表征了样本。标签（Label）是样本全体被赋予的具有意义的属性或属性组，其值由样本的特征向量决定。标签向量是标签值的数学抽象。在分类任务中，样本标签值是样本所属的类别，标签向量是 One-hot 向量。在回归任务中，标签值是若干有序实数，标签向量是分量为这些实数的实向量。若标签值为一个实数，则标签向量为 1 维实向量。

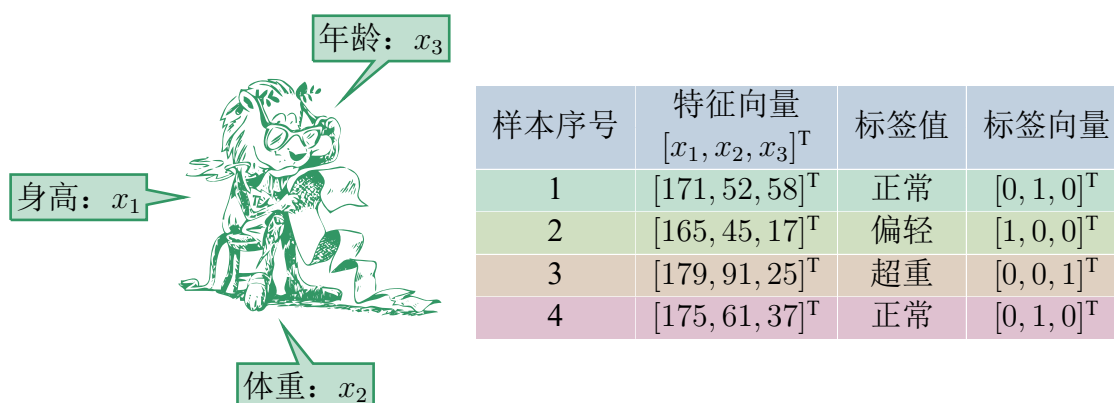


图 1 分类任务中的样本数据示例¹

¹图中 TeX Lion 为 Duane Bibby 原创，后文不再说明。

特征空间是所有可能的特征向量的集合，标签空间是所有可能的标签向量的集合。存在一个真实的映射，它是一个从特征空间到标签空间的映射，记为 f_t 。由此可知，每个特征向量都有唯一一个与之对应的标签向量。样本空间是所有可能的特征向量及其对应的标签向量的组合的集合。 f_t 未必是已知的，于是对给定的特征向量，其对应的标签向量是存在的，但可能是未知的。

1.2 预测与学习

多层感知机（Multi-layer Perceptron, MLP）是一种机器学习（Machine Learning）模型，它常用于分类任务与回归任务中。通俗而言，多层感知机的任务是根据给定样本的特征向量预测其对应的标签向量。在针对一个样本的预测任务中，多层感知机的输入是该样本的特征向量，输出是该样本的标签向量预测值。

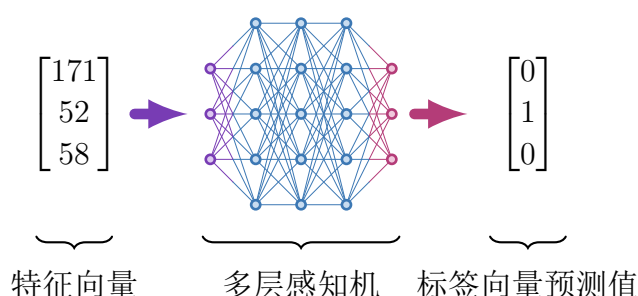


图2 多层感知机的输入与输出示例

在简单任务中，样本特征向量和标签向量之间的关系 f_t 是容易被人为确定并设计为固定算法的，而复杂任务中则不然。复杂任务要解决的是诸如手写数字识别等不易发现 f_t 的问题，这些问题是传统的算法逻辑基本固定的确定性程序难以应付的。而多层感知机能较好地解决这些问题，它能从数据中学习，从而构建出接近或吻合 f_t 的关系。

多层感知机模型内部有大量参数，它们参与了处理特征向量从而生成标签向量预测值的过程。学习的本质是利用大量已知特征向量及其对应标签向量的样本数据来调整模型内部参数，这一过程称为训练（Training），这些用于训练的已知样本的集合称为训练集（Training Set）。训练集可以认为是样本空间的一个很小的子集。在训练的某一过程中，多层感知机接收某个训练集子集的特征向量组及其对应的标签向量组，同时计算出标签向量预测值组。之后，它采用特定的参数优化方法，根据标签向量组和标签向量预测值组调整内部参数，从而提高预测的准确率。

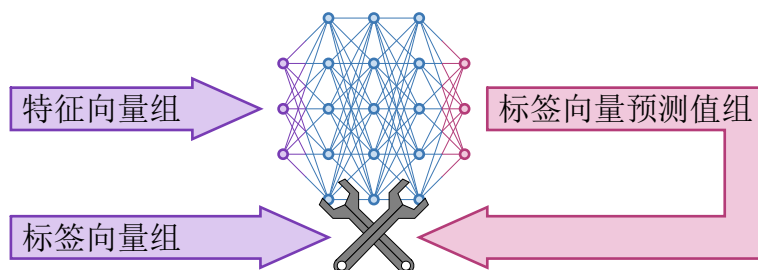


图3 多层感知机的参数调整

1.3 数学模型

1.3.1 神经网络

多层感知机是一种人工神经网络（Artificial Neural Network）模型，其设计思想最初源于生物学意义上的神经网络^{[1][2]}。从任务和计算角度看，多层感知机是从特征空间到另一个向量集的映射，记为 \hat{f} 。

假设在一个任务中，所有样本独立同分布（Independent and Identically Distributed），样本的特征向量是 $a(a \geq 1)$ 维的，标签向量是 $q(q \geq 1)$ 维的。现针对任务，分别给出 $M(M \geq 2)$ 个 $n^{(m)} \times n^{(m-1)}$ 矩阵 $\mathbf{W}^{(m)}$ 、 $n^{(m)} \times 1$ 列向量 $\mathbf{b}^{(m)}$ 和函数 $\varphi^{(m)}$ （其中， $m = 1, 2, \dots, M$ ， $n^{(0)} = a$ ， $n^{(M)} = q$ ， $n^{(m)} \geq 1$ ）。方便起见，记 $\mathbf{W}^{(m)}$ 第 $i(1 \leq i \leq n^{(m)})$ 行的转置为 $n^{(m-1)} \times 1$ 列向量 $\mathbf{w}_i^{(m)}$ ，而 $\mathbf{w}_i^{(m)}$ 的第 $j(1 \leq j \leq n^{(m-1)})$ 个分量为 $w_{i,j}^{(m)}$ 。又记 $\mathbf{b}^{(m)}$ 的第 $i(1 \leq i \leq n^{(m)})$ 个分量为 $b_i^{(m)}$ 。现给定 $s(s \geq 1)$ 个样本，并声明如下几个变量。

符号	格式	意义
\mathbf{x}	$n^{(0)} \times 1$ 列向量	特征向量
x_j	标量	特征的值， \mathbf{x} 的第 $j(0 \leq j \leq n^{(0)})$ 个分量
\mathbf{y}	$n^{(M)} \times 1$ 列向量	标签向量
y_j	标量	\mathbf{y} 的第 $j(0 \leq j \leq n^{(M)})$ 个分量
$\hat{\mathbf{y}}$	$n^{(M)} \times 1$ 列向量	标签向量预测值
\hat{y}_j	标量	$\hat{\mathbf{y}}$ 的第 $j(0 \leq j \leq n^{(M)})$ 个分量
\mathbf{X}	$s \times n^{(0)}$ 矩阵	特征矩阵，多层感知机的输入，每一行为一个样本特征向量的转置
\mathbf{x}_i	$n^{(0)} \times 1$ 列向量	特征向量， \mathbf{X} 第 $i(1 \leq i \leq s)$ 行的转置
$x_{i,j}$	标量	特征的值， \mathbf{x}_i 的第 $j(0 \leq j \leq n^{(0)})$ 个分量
\mathbf{Y}	$s \times n^{(M)}$ 矩阵	标签矩阵，每一行为一个样本标签向量的转置
\mathbf{y}_i	$n^{(M)} \times 1$ 列向量	标签向量， \mathbf{Y} 第 $i(1 \leq i \leq s)$ 行的转置
$y_{i,j}$	标量	\mathbf{y}_i 的第 $j(0 \leq j \leq n^{(M)})$ 个分量
$\hat{\mathbf{Y}}$	$s \times n^{(M)}$ 矩阵	多层感知机的输出，每一行为一个样本标签向量预测值的转置
$\hat{\mathbf{y}}_i$	$n^{(M)} \times 1$ 列向量	标签向量预测值， $\hat{\mathbf{Y}}$ 第 $i(1 \leq i \leq s)$ 行的转置
$\hat{y}_{i,j}$	标量	$\hat{\mathbf{y}}_i$ 的第 $j(0 \leq j \leq n^{(M)})$ 个分量
$\mathbf{X}^{(m)}$	$s \times n^{(m)}$ 矩阵	
$x_{i,j}^{(m)}$	标量	$\mathbf{X}^{(m)}$ 第 $i(1 \leq i \leq s)$ 行、第 $j(1 \leq j \leq n^{(m)})$ 列元素
$\mathbf{x}^{(m)}$	$n^{(m)} \times 1$ 列向量	$s = 1$ 时的 $\mathbf{X}^{(m)\top}$
$x_j^{(m)}$	标量	$\mathbf{x}^{(m)}$ 的第 $j(1 \leq j \leq n^{(m)})$ 个分量
$\mathbf{Y}^{(m)}$	$s \times n^{(m)}$ 矩阵	
$y_{i,j}^{(m)}$	标量	$\mathbf{Y}^{(m)}$ 第 $i(1 \leq i \leq s)$ 行、第 $j(1 \leq j \leq n^{(m)})$ 列元素
$\mathbf{y}^{(m)}$	$n^{(m)} \times 1$ 列向量	$s = 1$ 时的 $\mathbf{Y}^{(m)\top}$
$y_j^{(m)}$	标量	$\mathbf{y}^{(m)}$ 的第 $j(1 \leq j \leq n^{(m)})$ 个分量

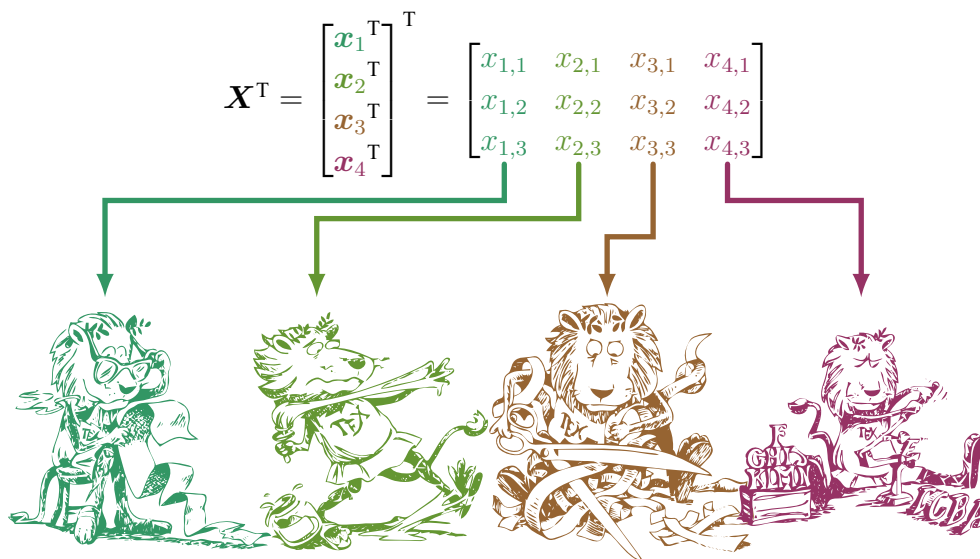


图 4 4×3 特征矩阵 \mathbf{X} 的转置示例

现给出 M 个每一列均为 $\mathbf{b}^{(m)}$ 的 $n^{(m)} \times s$ 矩阵 $\mathbf{B}^{(m)} = [\mathbf{b}^{(m)}, \dots, \mathbf{b}^{(m)}]$, 那么多层感知机的输出 $\hat{\mathbf{Y}} = \hat{f}(\mathbf{X})$ 可计算如下

$$\begin{cases} \mathbf{X}^{(m)\top} = \mathbf{W}^{(m)}\mathbf{Y}^{(m-1)\top} + \mathbf{B}^{(m)} \\ \mathbf{Y}^{(m)} = \varphi^{(m)}(\mathbf{X}^{(m)}) \\ \mathbf{Y}^{(0)} = \mathbf{X} \\ \mathbf{Y}^{(M)} = \hat{\mathbf{Y}} \end{cases} \quad (1-1)$$

其中, $\mathbf{W}^{(m)}$ 被称为权重 (Weight) 矩阵, $\mathbf{b}^{(m)}$ 被称为偏置 (Bias) 列向量, $\varphi^{(m)}$ 被称为激活函数 (Activation Function), 且约定 $(\varphi^{(m)}(\mathbf{X}^{(m)}))_{i,j} = \varphi^{(m)}(x_{i,j}^{(m)})$ 。 m 是多层感知机的层号, 第 0 层是输入层, 第 M 层是输出层, 其余层是隐藏层。

当 $s = 1$ 时, 由(1-1)式, 有

$$y_j^{(m)} = \varphi^{(m)}(\mathbf{w}_j^{(m)\top} \mathbf{y}^{(m-1)} + b_j^{(m)}) \quad (1-2)$$

称(1-2)式为多层感知机第 m 层的第 j 个神经元 (Neuron)。

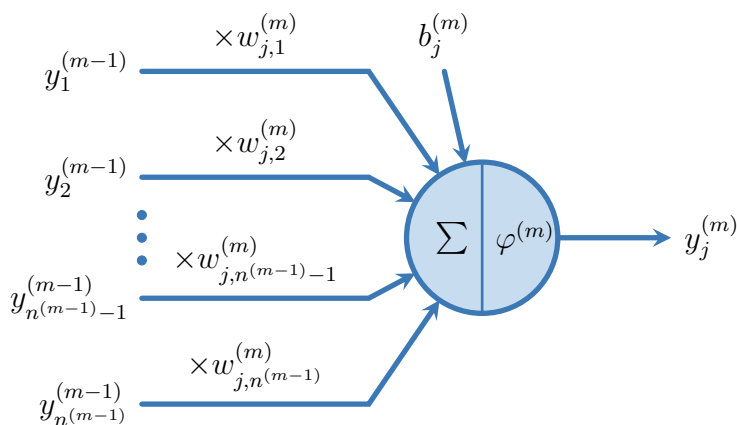


图 5 多层感知机第 m 层的第 j 个神经元结构

神经元是多层感知机的基本单位。由于多层感知机的神经元分层排列，每个神经元只接收前一层神经元的输出作为输入，各层神经元之间不存在反馈，故多层感知机属于前馈神经网络（Feedforward Neural Network）。由于每个神经元均接收前一层所有神经元的输出，故这种神经网络结构又是全连接的（Fully Connected）。图 6 所示的神经元连接线上的数据为多层感知机第 m 层的第 $n^{(m)}$ 个神经元上待求和的加权输入。其中 $y_1^{(m-1)}$ 、 $y_2^{(m-1)}$ 、 $y_{n^{(m-1)}-1}^{(m-1)}$ 和 $y_{n^{(m-1)}}^{(m-1)}$ 分别为多层感知机第 $m-1$ 层的第 1、2、 $n^{(m-1)}-1$ 和 $n^{(m-1)}$ 个神经元的输出。

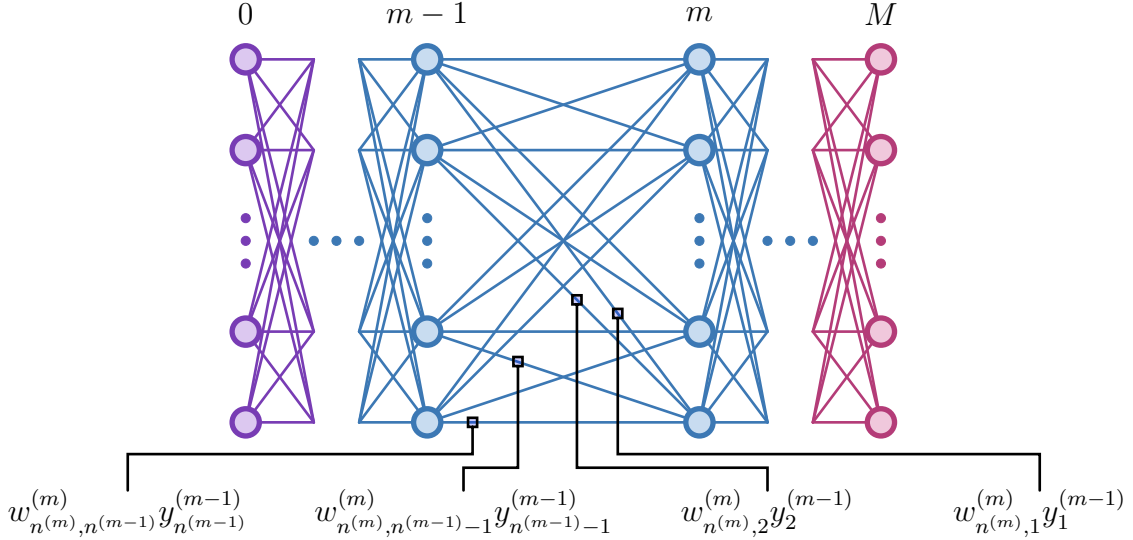


图 6 多层感知机中第 m 层的第 $n^{(m)}$ 个神经元的加权输入

1.3.2 损失函数

损失函数（Loss Function）用于衡量多层感知机的输出与实际的样本标签向量之间的差异程度，其值为非负实数。一般而言，损失函数是一个关于标签向量预测值组 $\hat{\mathbf{Y}}$ 和标签向量组 \mathbf{Y} 的函数，记为 \mathcal{L} 。令 $L \in \mathbb{R}^*$ 为损失函数值，有

$$L = \mathcal{L}(\hat{\mathbf{Y}}, \mathbf{Y}) \quad (2-1)$$

由于多层感知机是映射 \hat{f} ，且 $\hat{\mathbf{Y}} = \hat{f}(\mathbf{X})$ ，于是有

$$L = \mathcal{L}(\hat{f}(\mathbf{X}), \mathbf{Y}) \quad (2-2)$$

由此可知，损失函数是一个从样本空间到 \mathbb{R}^* 的函数。对一个较好的多层感知机模型，它在整个特征空间上的输出与对应的标签向量之间的差异应该尽可能小，换言之，损失函数在整个样本空间上的值应该尽可能小。从损失函数的角度看，训练就是调整参数 $\mathbf{W}^{(m)}$ 和 $\mathbf{b}^{(m)}$ ，从而找到一个使得 \mathcal{L} 在整个样本空间上的函数值尽可能小的映射 \hat{f} 。

一般地，训练只能基于样本规模十分有限的训练集。因此，可能出现损失函数在训练集上的值很小，而在训练集以外的部分样本空间上的值较大的现象，这种现象一般为过拟合（Overfitting）。而损失函数在整个样本空间上的值均较大的现象一般为欠拟合（Underfitting）。

1.3.2.1 交叉熵损失函数

在分类任务中，常用交叉熵（Cross Entropy）作为损失函数。对前面给定的 \mathbf{y} 和 $\hat{\mathbf{y}}$ ， $\hat{\mathbf{y}}$ 一般是一个属于空间 $(0, 1)^{n^{(M)}}$ 的向量，其第 j ($1 \leq j \leq n^{(M)}$) 个分量可被视为样本属于第 j 个类

别的概率，于是交叉熵损失为

$$l_{CE} = \sum_{j=1}^{n^{(M)}} y_j \ln \hat{y}_j$$

对于(1-1)式，有交叉熵损失

$$\mathcal{L}_{CE} = \frac{1}{s} \sum_{i=1}^s \sum_{j=1}^{n^{(M)}} y_{i,j} \ln \hat{y}_{i,j} \quad (2-3-1)$$

\mathcal{L}_{CE} 对元素 $\hat{y}_{p,q}$ 的偏导数为

$$\frac{\partial}{\partial \hat{y}_{p,q}} \mathcal{L}_{CE} = -\frac{1}{s} \frac{y_{p,q}}{\hat{y}_{p,q}} \quad (2-3-2)$$

1.3.2.2 均方差损失函数

均方误差（Mean Squared Error, MSE），简称均方差。在回归任务中，常用均方差作为损失函数。对同属实向量的 \mathbf{y} 和 $\hat{\mathbf{y}}$ ，均方差损失为

$$l_{MSE} = \frac{1}{n^{(M)}} \sum_{j=1}^{n^{(M)}} (\hat{y}_j - y_j)^2$$

对于(1-1)式，有均方差损失

$$\mathcal{L}_{MSE} = \frac{1}{n^{(M)}s} \sum_{i=1}^s \sum_{j=1}^{n^{(M)}} (\hat{y}_{i,j} - y_{i,j})^2 \quad (2-4-1)$$

\mathcal{L}_{MSE} 对元素 $\hat{y}_{p,q}$ 的偏导数为

$$\frac{\partial}{\partial \hat{y}_{p,q}} \mathcal{L}_{MSE} = \frac{2}{n^{(M)}s} (\hat{y}_{p,q} - y_{p,q}) \quad (2-4-2)$$

1.3.3 激活函数

激活函数一般指非线性激活函数，其主要作用是多层感知机引入非线性能力或控制神经元输出范围。

在(1-1)式中，若所有的激活函数 $\varphi^{(m)}(x) = x$ ，则

$$\begin{aligned} \hat{\mathbf{Y}}^T &= \mathbf{W}^{(M)} (\dots \mathbf{W}^{(2)} (\mathbf{W}^{(1)} \mathbf{X}^T + \mathbf{B}^{(1)}) + \mathbf{B}^{(2)} \dots) + \mathbf{B}^{(M)} \\ &= \mathbf{U} \mathbf{X}^T + \mathbf{V} \end{aligned}$$

其中

$$\mathbf{U} = \mathbf{W}^{(M)} \dots \mathbf{W}^{(2)} \mathbf{W}^{(1)}$$

$$\mathbf{V} = \mathbf{W}^{(M)} \dots \mathbf{W}^{(2)} \mathbf{B}^{(1)} + \mathbf{W}^{(M)} \dots \mathbf{W}^{(3)} \mathbf{B}^{(2)} + \dots + \mathbf{W}^{(M)} \mathbf{B}^{(M-1)} + \mathbf{B}^{(M)}$$

可见，当多层感知机的所有神经元均不被施以非线性激活函数时，该多层感知机实际上是一个仿射变换（Affine Transformation）。无论如何调整参数 $\mathbf{W}^{(m)}$ 和 $\mathbf{b}^{(m)}$ ，它始终与不属于仿射变换的非线性变换（Nonlinear Transformation）存在本质区别。换言之，它无法更好地拟合非线性变换，故需要依靠非线性激活函数引入非线性能力。

在训练过程中，随着参数 $\mathbf{W}^{(m)}$ 和 $\mathbf{b}^{(m)}$ 的不断调整，神经元的输出可能超出计算机数据类型表示范围，即可能发生数据溢出现象。此时应当适当“挤压”输出数据，将其映射到某个有界区间上。而在分类任务中，输出层神经元的输出通常也被控制在区间 $(0, 1)$ 上，并被视为样本属于某一类别的概率。故此，多层感知机需要依靠非线性激活函数控制神经元输出范围。

非线性激活函数种类较多，常见的有以 Logistic 和 \tanh 等为代表的 Sigmoid 型和以 ReLU、Leaky ReLU 和 Parametric ReLU 等为代表的 ReLU 型。激活函数有各自的优缺点，对训练速度和模型精度等方面的影响各不相同。例如 ReLU 型计算简单，但容易造成数据溢出，而 Sigmoid 型反之。下面给出几个较为常见的激活函数及其导数的计算式。

1.3.3.1 Logistic 函数

Logistic 函数是一个经典的激活函数，其函数图像形如“S”，故常被称为 Sigmoid 函数。以下二式分别为 Logistic 及其导数的计算方法。

$$\varphi_{lg} = \frac{1}{1 + e^{-x}} \quad (3-1-1)$$

$$\varphi'_{lg} = \varphi_{lg}(1 - \varphi_{lg}) \quad (3-1-2)$$

1.3.3.2 tanh 函数

\tanh 即双曲正切函数（Hyperbolic Tangent Function），它与 Logistic 同属 Sigmoid 型函数，但其输出以 0 为均值。以下二式分别为 \tanh 及其导数的计算方法。

$$\varphi_{th} = \frac{1 - e^{-2x}}{1 + e^{-2x}} \quad (3-2-1)$$

$$\varphi'_{th} = 1 - \varphi_{th}^2 \quad (3-2-2)$$

1.3.3.3 ReLU 函数

ReLU 即修正线性单元（Rectified Linear Unit）。以下二式分别为 ReLU 及其导数的计算方法。

$$\varphi_r = \max(0, x) \quad (3-3-1)$$

$$\varphi'_r = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases} \quad (3-3-2)$$

1.3.3.3 Leaky ReLU 函数

Leaky ReLU 在 ReLU 的基础上保留了负方向的信息。令 $\alpha > 0$ ，则以下二式分别为 Leaky ReLU 及其导数的计算方法。

$$\varphi_{kr} = \begin{cases} x, & x \geq 0 \\ \alpha x, & x < 0 \end{cases} \quad (3-4-1)$$

$$\varphi'_{kr} = \begin{cases} 1, & x \geq 0 \\ \alpha, & x < 0 \end{cases} \quad (3-4-2)$$

一般地，取 α 为 0.01 左右。需要指明的是，Leaky ReLU 和 ReLU 在 $x = 0$ 时均不可导，因而(3-4-2)和(3-3-2)在 $x = 0$ 处的值均是人为规定的。

1.3.3.2 Softmax 函数

Softmax 函数是一个比较特殊的激活函数，常用在分类任务中，一般为多层感知机输出层的激活函数。在(1-1)式中，对 $i(1 \leq i \leq s)$ 、 $p(1 \leq p \leq n^{(M)})$ 以及 $q(1 \leq q \leq n^{(M)}, q \neq p)$ ，有 Softmax 及其偏导数的计算方法如下

$$\hat{y}_{i,p} = \frac{e^{x_{i,p}^{(M)}}}{\sum_{j=1}^{n^{(M)}} e^{x_{i,j}^{(M)}}} \quad (3-5-1)$$

$$\frac{\partial}{\partial x_{i,p}^{(M)}} \hat{y}_{i,p} = \hat{y}_{i,p}(1 - \hat{y}_{i,p}) \quad (3-5-2)$$

$$\frac{\partial}{\partial x_{i,q}^{(M)}} \hat{y}_{i,p} = -\hat{y}_{i,p}\hat{y}_{i,q} \quad (3-5-3)$$

1.4 通用近似定理

通用近似定理 (Universal Approximation Theorem) [3][4][5] 是多层感知机的理论依据，它表明多层感知机对连续函数具有拟合能力。下面是通用近似定理的一种表述。

令 $\varphi(\bullet)$ 为一个连续、有界、单调递增 (区别于严格单调递增) 的非常数函数， \mathbf{I}_D 为 $D(D \geq 1)$ 维单位超立方体 $[0, 1]^D$ ， $C(\mathbf{I}_D)$ 为定义在 \mathbf{I}_D 上的连续函数的集合。对给定的任意实数 $\varepsilon > 0$ 和任意函数 $g \in C(\mathbf{I}_D)$ ，存在一个自然数 $N \in \mathbb{N}^+$ ，一组实数 $v_i \in \mathbb{R}$ 、 $b_i \in \mathbb{R}$ 和实向量 $\mathbf{w}_i \in \mathbb{R}^m$ (其中 $i = 1, 2, \dots, N$)，以至于对任意 $\mathbf{x} \in \mathbf{I}_D$ ，可以定义

$$F(\mathbf{x}) = \sum_{i=1}^N v_i \varphi(\mathbf{w}_i^T \mathbf{x} + b_i)$$

使得

$$|F(\mathbf{x}) - g(\mathbf{x})| < \varepsilon$$

换言之，形如 $F(\mathbf{x})$ 的函数在 $C(\mathbf{I}_D)$ 中是稠密的。

回顾(1-2)的神经元模型，令多层感知机第 $m+1(1 \leq m \leq M-1)$ 层的第 $k(0 \leq k \leq n^{(m+1)})$ 个神经元的激活函数的输入为 $x_k^{(m+1)}$ ，则

$$\begin{aligned} x_k^{(m+1)} &= \mathbf{w}_k^{(m+1)T} \mathbf{y}^{(m)} + b_k^{(m+1)} \\ &= \sum_{j=1}^{n^{(m)}} w_{k,j}^{(m+1)} y_j^{(m)} + b_k^{(m+1)} \\ &= \sum_{j=1}^{n^{(m)}} w_{k,j}^{(m+1)} \varphi^{(m)}(\mathbf{w}_j^{(m)T} \mathbf{y}^{(m-1)} + b_j^{(m)}) + b_k^{(m+1)} \end{aligned}$$

假设 $\mathbf{y}^{(m-1)} \in [-C, C]^{n^{(m-1)}} (C \in \mathbb{R}^+)$ ，令

$$\mathbf{y}^* = \frac{1}{2C} \mathbf{y}^{(m-1)} + [\frac{1}{2}, \dots, \frac{1}{2}]^T$$

则 $\mathbf{y}^* \in [0, 1]^{n^{(m-1)}}$ ，于是有

$$x_k^{(m+1)} = \sum_{j=1}^{n^{(m)}} w_{k,j}^{(m+1)} \varphi_j^{(m)} (2C \mathbf{w}_j^{(m)T} \mathbf{y}^* + (b_j^{(m)} - \mathbf{w}_j^{(m)T} [C, \dots, C]^T)) + b_k^{(m+1)} \quad (4^*)$$

现令 $g(\bullet)$ 是定义在 $[-C, C]^{n^{(m-1)}}$ 上的连续函数，则

$$g^*(\bullet) = g(2C(\bullet) - [C, \dots, C]^T)$$

是定义在 $[0, 1]^{n^{(m-1)}}$ 上的连续函数。忽略偏置 $b_k^{(m+1)}$ ，易知(4*)式符合通用近似定理，即对连续函数 $g^*(\bullet)$ ，只要参数 $n^{(m)}$ 、 $w_{k,j}^{(m+1)}$ 、 $\mathbf{w}_j^{(m)T}$ 和 $b_j^{(m)}$ 恰当，(4*)式就能在某种程度上拟合连续函数 $g^*(\bullet)$ 。由前文已知，参数 $w_{k,j}^{(m+1)}$ 、 $\mathbf{w}_j^{(m)T}$ 和 $b_j^{(m)}$ 可通过训练不断调整，而多层感知机第 m 个隐藏层的神经元个数 $n^{(m)}$ 则不能调整。 $n^{(m)}$ 一般是一个固定的常数，它在多层感知机构建之初便确定下来，而其具体值的确定是一项经验性工作。

在激活函数中，ReLU 型函数并不是有界函数，这在直觉上不符合通用近似定理。但事实上，一个输出层为线性激活层 ($\varphi^{(M)}(x) = x$)，隐藏层为 ReLU 激活层的二层感知机已经具有极强的拟合能力。一种观点认为，这是因为仅由 ReLU 便可以构造出一个新的连续、有界且单调递增的非常数函数。令 $p < q (p \in \mathbb{R}, q \in \mathbb{R})$ ，则该新函数如下

$$\begin{aligned} \varphi_{p,q} &= \varphi_r(x - p) - \varphi_r(x - q) \\ &= \max(0, x - p) - \max(0, x - q) \\ &= \begin{cases} 0, & x < p \\ x, & p \leq x < q \\ q - p, & x \geq q \end{cases} \end{aligned}$$

2 优化方法

2.1 梯度下降法

回顾(2-2)式，现要求在 \mathbf{X} 为训练集特征向量组和 \mathbf{Y} 为对应的训练集标签向量组的情况下，找出使得 \mathcal{L} 值尽可能小的权重矩阵参数 $\mathbf{W}^{(m)}$ 和偏置列向量参数 $\mathbf{b}^{(m)}$ 。由此，视 \mathcal{L} 为一个关于权重参数组 $w_{1,1}^{(1)}, w_{1,2}^{(1)}, \dots, w_{n^{(M)}, n^{(M-1)}}^{(M)}$ 和偏置参数组 $b_1^{(1)}, b_2^{(1)}, \dots, b_{n^{(M)}}^{(M)}$ 的多元函数

$$L = \mathcal{L}(w_{1,1}^{(1)}, w_{1,2}^{(1)}, \dots, w_{n^{(M)}, n^{(M-1)}}^{(M)}, b_1^{(1)}, b_2^{(1)}, \dots, b_{n^{(M)}}^{(M)}) = \mathcal{L}(\mathbf{p})$$

对多元函数 \mathcal{L} 而言，其极值可能为最值，而其所有偏导数均为 0 的点可能为极值点。于是存在一种方法为，计算出使得 \mathcal{L} 的偏导数均为 0 的权重组和偏置组，即解算方程组

$$\frac{\partial}{\partial \mathbf{p}} \mathcal{L} = \mathbf{0}$$

但事实上，这种方法在计算上几乎不可行，故引入梯度下降法 (Gradient Descent Method)。

根据函数在其梯度的反方向上的方向导数最小的原理可知，函数在该方向上的值下降最快。对 \mathcal{L} 而言，只要顺着其梯度反方向前进，就能在局部上找到使得 \mathcal{L} 尽可能小的权重组和

偏置组。在一次梯度下降的迭代过程中，设当前权重组和偏置组为点 \mathbf{p} ，则从当前位置向梯度反方向前进一小段距离，并更新 \mathbf{p} 到下一个位置，即

$$\mathbf{p} \leftarrow \mathbf{p} - \eta \frac{\partial}{\partial \mathbf{p}} \mathcal{L} \quad (5-1)$$

其中， $\eta (0 < \eta < 1)$ 称为学习率，它决定了前进的幅度。重复多轮执行(5-1)， \mathcal{L} 在点 \mathbf{p} 上的值在总体上将呈现为一个下降的趋势。针对(1-1)式，(5-1)式也可以具体拆分为

$$\begin{cases} \mathbf{W}^{(m)} \leftarrow \mathbf{W}^{(m)} - \eta \frac{\partial}{\partial \mathbf{W}^{(m)}} \mathcal{L} \\ \mathbf{b}^{(m)} \leftarrow \mathbf{b}^{(m)} - \eta \frac{\partial}{\partial \mathbf{b}^{(m)}} \mathcal{L} \end{cases} \quad (5-2)$$

梯度下降法并不能保证得到全局最优解。图 7 所示为两个梯度下降算法和学习率均相同，但参数初始值不同的梯度下降过程。以圆形标记为起点的绿色路径找到了全局最优解，而以十字标记为起点的红色路径则陷入了局部最优。解决局部最优的方案较多，可以从自适应学习率、动量梯度下降、参数初始化等方面入手。

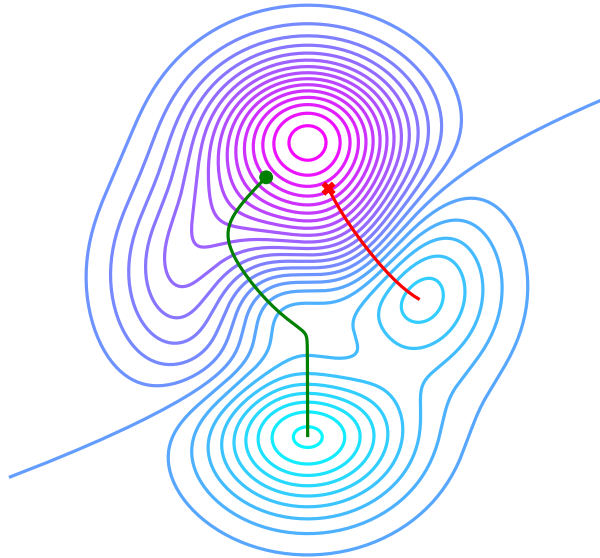


图 7 全局最优与局部最优示例

在上文中已交代， \mathbf{X} 和 \mathbf{Y} 分别是整个训练集的特征向量组和标签向量组，这种每次迭代过程都基于整个训练集来更新参数的梯度下降称为批量梯度下降（Batch Gradient Descent, BGD）。对计算机而言，整个训练集的样本数据是相对庞大的，因而批量梯度下降较为突出的缺点之一是计算开销大。不同于批量梯度下降，随机梯度下降（Stochastic Gradient Descent, SGD）的每次迭代过程只随机采取一个样本进行参数更新。换言之，此时的 \mathbf{X} 和 \mathbf{Y} 均为行向量，而非具有多行的矩阵。然而，单个样本不足以代表整个训练集，随机梯度下降具有一定盲目性，容易陷入局部最优。综上，可以考虑在批量梯度下降和随机梯度下降之间进行折中处理，引入小批量梯度下降（Mini-batch Gradient Descent, MBGD），它的每次迭代过程基于训练集中的小部分样本来更新参数。此时，矩阵 \mathbf{X} 和 \mathbf{Y} 分别为训练集小部分样本的特征向量组和标签向量组。

2.2 反向传播公式

如(5-2)式，梯度下降法的难点在于计算 \mathcal{L} 对各层参数 $\mathbf{W}^{(m)}$ 和 $\mathbf{b}^{(m)}$ 的偏导数，故而引入误差反向传播（Backward Propagation of Errors）算法，也称为反向传播（Backpropagation, BP）算法^{[6][7]}。反向传播算法基于链式法则实现，它利用被称为误差项的 $\partial\mathcal{L}/\partial\mathbf{X}^{(m)}$ 由后至前逐层计算出损失函数 \mathcal{L} 对权重参数和偏置参数的偏导数，这一过程即为反向传播过程^[8]。

现针对前文(1-1)式及(2-1)式，给出以下几个公式及其证明。

2.2.0.1 损失函数对神经元输出的偏导数

$$\frac{\partial\mathcal{L}}{\partial\mathbf{Y}^{(m)}} = \frac{\partial\mathcal{L}}{\partial\mathbf{X}^{(m+1)}} \mathbf{W}^{(m+1)} \quad (6-1)$$

其中， $1 \leq m \leq M-1$ 。现给出如下证明。

首先，计算 $x_{i,v}^{(m+1)}$ 对 $y_{i,j}^{(m)}$ 的偏导数如下

$$\begin{aligned} \frac{\partial x_{i,v}^{(m+1)}}{\partial y_{i,j}^{(m)}} &= \frac{\partial}{\partial y_{i,j}^{(m)}} \left(\sum_{k=1}^{n^{(m)}} w_{v,k}^{(m+1)} y_{i,k}^{(m)} + b_v^{(m+1)} \right) \\ &= w_{v,j}^{(m+1)} \end{aligned}$$

由上式，再依据链式法则，有

$$\begin{aligned} \left(\frac{\partial\mathcal{L}}{\partial\mathbf{Y}^{(m)}} \right)_{i,j} &= \frac{\partial\mathcal{L}}{\partial y_{i,j}^{(m)}} \\ &= \sum_{u=1}^s \sum_{v=1}^{n^{(m+1)}} \frac{\partial x_{u,v}^{(m+1)}}{\partial y_{i,j}^{(m)}} \frac{\partial\mathcal{L}}{\partial x_{u,v}^{(m+1)}} \\ &= \sum_{u \neq i} \sum_{v=1}^{n^{(m+1)}} \frac{\partial x_{u,v}^{(m+1)}}{\partial y_{i,j}^{(m)}} \frac{\partial\mathcal{L}}{\partial x_{u,v}^{(m+1)}} + \sum_{v=1}^{n^{(m+1)}} \frac{\partial x_{i,v}^{(m+1)}}{\partial y_{i,j}^{(m)}} \frac{\partial\mathcal{L}}{\partial x_{i,v}^{(m+1)}} \\ &= 0 + \sum_{v=1}^{n^{(m+1)}} w_{v,j}^{(m+1)} \frac{\partial\mathcal{L}}{\partial x_{i,v}^{(m+1)}} \\ &= \sum_{v=1}^{n^{(m+1)}} \left(\frac{\partial\mathcal{L}}{\partial\mathbf{X}^{(m+1)}} \right)_{i,v} (\mathbf{W}^{(m+1)})_{v,j} \end{aligned}$$

由此，(6-1)式成立。

2.2.0.2 损失函数对权重的偏导数

$$\frac{\partial\mathcal{L}}{\partial\mathbf{W}^{(m)}} = \left(\frac{\partial\mathcal{L}}{\partial\mathbf{X}^{(m)}} \right)^T \mathbf{Y}^{(m-1)} \quad (6-2)$$

现给出如下证明。

首先，计算 $x_{v,i}^{(m)}$ 对 $w_{i,j}^{(m)}$ 的偏导数如下

$$\begin{aligned} \frac{\partial x_{v,i}^{(m)}}{\partial w_{i,j}^{(m)}} &= \frac{\partial}{\partial w_{i,j}^{(m)}} \left(\sum_k^{n^{(m-1)}} w_{i,k}^{(m)} y_{v,k}^{(m-1)} + b_i^{(m)} \right) \\ &= y_{v,j}^{(m-1)} \end{aligned}$$

由上式，再依据链式法则，有

$$\begin{aligned}
\left(\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{(m)}} \right)_{i,j} &= \frac{\partial \mathcal{L}}{\partial w_{i,j}^{(m)}} \\
&= \sum_{u=1}^{n^{(m)}} \sum_{v=1}^s \frac{\partial \mathcal{L}}{\partial x_{v,u}^{(m)}} \frac{\partial x_{v,u}^{(m)}}{\partial w_{i,j}^{(m)}} \\
&= \sum_{u \neq i} \sum_{v=1}^s \frac{\partial \mathcal{L}}{\partial x_{v,u}^{(m)}} \frac{\partial x_{v,u}^{(m)}}{\partial w_{i,j}^{(m)}} + \sum_{v=1}^s \frac{\partial \mathcal{L}}{\partial x_{v,i}^{(m)}} \frac{\partial x_{v,i}^{(m)}}{\partial w_{i,j}^{(m)}} \\
&= 0 + \sum_{v=1}^s \frac{\partial \mathcal{L}}{\partial x_{v,i}^{(m)}} y_{v,j}^{(m-1)} \\
&= \sum_{v=1}^s \left(\left(\frac{\partial \mathcal{L}}{\partial \mathbf{X}^{(m)}} \right)^T \right)_{i,v} (\mathbf{Y}^{(m-1)})_{v,j}
\end{aligned}$$

由此，(6-2)式成立。

2.2.0.3 损失函数对偏置的偏导数

$$\left(\frac{\partial \mathcal{L}}{\partial \mathbf{b}^{(m)}} \right)_i = \sum_{v=1}^s \left(\frac{\partial \mathcal{L}}{\partial \mathbf{X}^{(m)}} \right)_{v,i} \quad (6-3)$$

换言之， $\partial \mathcal{L} / \partial \mathbf{b}^{(m)}$ 的第 i 个分量是 $\partial \mathcal{L} / \partial \mathbf{X}^{(m)}$ 的第 i 列之和。现给出如下证明。

首先，计算 $x_{v,i}^{(m)}$ 对 $b_i^{(m)}$ 的偏导数如下

$$\begin{aligned}
\frac{\partial x_{v,i}^{(m)}}{\partial b_i^{(m)}} &= \frac{\partial}{\partial b_i^{(m)}} \left(\sum_{k=1}^{n^{(m-1)}} w_{i,k}^{(m)} y_{v,k}^{(m-1)} + b_i^{(m)} \right) \\
&= 1
\end{aligned}$$

由上式，再依据链式法则，有

$$\begin{aligned}
\left(\frac{\partial \mathcal{L}}{\partial \mathbf{b}^{(m)}} \right)_i &= \frac{\partial \mathcal{L}}{\partial b_i^{(m)}} \\
&= \sum_{u=1}^{n^{(m)}} \sum_{v=1}^s \frac{\partial \mathcal{L}}{\partial x_{v,u}^{(m)}} \frac{\partial x_{v,u}^{(m)}}{\partial b_i^{(m)}} \\
&= \sum_{u \neq i} \sum_{v=1}^s \frac{\partial \mathcal{L}}{\partial x_{v,u}^{(m)}} \frac{\partial x_{v,u}^{(m)}}{\partial b_i^{(m)}} + \sum_{v=1}^s \frac{\partial \mathcal{L}}{\partial x_{v,i}^{(m)}} \frac{\partial x_{v,i}^{(m)}}{\partial b_i^{(m)}} \\
&= \sum_{v=1}^s \frac{\partial \mathcal{L}}{\partial x_{v,i}^{(m)}} \\
&= \sum_{v=1}^s \left(\frac{\partial \mathcal{L}}{\partial \mathbf{X}^{(m)}} \right)_{v,i}
\end{aligned}$$

2.2.0.4 隐藏层误差项

$$\frac{\partial \mathcal{L}}{\partial \mathbf{X}^{(m)}} = \varphi^{(m)'}(\mathbf{X}^{(m)}) \odot \left(\frac{\partial \mathcal{L}}{\partial \mathbf{X}^{(m+1)}} \cdot \mathbf{W}^{(m+1)} \right) \quad (6-4)$$

其中, $1 \leq m \leq M-1$, \odot 是哈达玛积 (Hadamard Product), $\varphi^{(m)'}$ 是 $\varphi^{(m)}$ 的导数且约定 $(\varphi^{(m)' }(\mathbf{X}^{(m)}))_{i,j} = \varphi^{(m)'}(x_{i,j}^{(m)})$ 。同时, 需要注意, 第 m 层任意一个神经元的输出 $y_{u,v}^{(m)}$ 仅与 $x_{u,v}^{(m)}$ 相关, 即 $y_{u,v}^{(m)} = \varphi^{(m)}(x_{u,v}^{(m)})$ 仅有唯一变量 $x_{u,v}^{(m)}$ 。现给出如下证明。

首先, 由链式法则, 有

$$\begin{aligned} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{X}^{(m)}} \right)_{i,j} &= \frac{\partial \mathcal{L}}{\partial x_{i,j}^{(m)}} \\ &= \sum_{u=1}^s \sum_{v=1}^{n^{(m)}} \frac{\partial y_{u,v}^{(m)}}{\partial x_{i,j}^{(m)}} \frac{\partial \mathcal{L}}{\partial y_{u,v}^{(m)}} \\ &= \sum_{u \neq i} \sum_{v=1}^{n^{(m)}} \frac{\partial y_{u,v}^{(m)}}{\partial x_{i,j}^{(m)}} \frac{\partial \mathcal{L}}{\partial y_{u,v}^{(m)}} + \sum_{v=1}^{n^{(m)}} \frac{\partial y_{i,v}^{(m)}}{\partial x_{i,j}^{(m)}} \frac{\partial \mathcal{L}}{\partial y_{i,v}^{(m)}} \\ &= \sum_{v=1}^{n^{(m)}} \frac{\partial y_{i,v}^{(m)}}{\partial x_{i,j}^{(m)}} \frac{\partial \mathcal{L}}{\partial y_{i,v}^{(m)}} \end{aligned} \quad (6^*)$$

已知当 $v \neq j$ 时, $y_{i,v}^{(m)}$ 与 $x_{i,j}^{(m)}$ 不相关, 于是由(6*), 有

$$\begin{aligned} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{X}^{(m)}} \right)_{i,j} &= \sum_{v \neq j} \frac{\partial y_{i,v}^{(m)}}{\partial x_{i,j}^{(m)}} \frac{\partial \mathcal{L}}{\partial y_{i,v}^{(m)}} + \frac{\partial y_{i,j}^{(m)}}{\partial x_{i,j}^{(m)}} \frac{\partial \mathcal{L}}{\partial y_{i,j}^{(m)}} \\ &= 0 + \varphi^{(m)'}(x_{i,j}^{(m)}) \frac{\partial \mathcal{L}}{\partial y_{i,j}^{(m)}} \\ &= \varphi^{(m)'}(\mathbf{X}^{(m)})_{i,j} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{Y}^{(m)}} \right)_{i,j} \end{aligned}$$

再结合(6-1)式, 可知(6-4)式成立。

2.2.0.5 基于交叉熵损失的 Softmax 输出层误差项

$$\frac{\partial \mathcal{L}}{\partial \mathbf{X}^{(M)}} = \frac{1}{s} (\hat{\mathbf{Y}} - \mathbf{Y}) \quad (6-5)$$

(6-5)式适用于损失函数为交叉熵、输出层激活函数为 Softmax 的用于解决分类问题的多层感知机。现给出如下证明。

首先, 由(6*)式, 有

$$\begin{aligned} \left(\frac{\partial \mathcal{L}}{\partial \mathbf{X}^{(M)}} \right)_{i,j} &= \sum_{v=1}^{n^{(M)}} \frac{\partial y_{i,v}^{(M)}}{\partial x_{i,j}^{(M)}} \frac{\partial \mathcal{L}}{\partial y_{i,v}^{(M)}} \\ &= \sum_{v \neq j} \frac{\partial y_{i,v}^{(M)}}{\partial x_{i,j}^{(M)}} \frac{\partial \mathcal{L}}{\partial y_{i,v}^{(M)}} + \frac{\partial y_{i,j}^{(M)}}{\partial x_{i,j}^{(M)}} \frac{\partial \mathcal{L}}{\partial y_{i,j}^{(M)}} \\ &= \sum_{v \neq j} \frac{\partial \hat{y}_{i,v}}{\partial x_{i,j}^{(M)}} \frac{\partial \mathcal{L}}{\partial \hat{y}_{i,v}} + \frac{\partial \hat{y}_{i,j}}{\partial x_{i,j}^{(M)}} \frac{\partial \mathcal{L}}{\partial \hat{y}_{i,j}} \end{aligned}$$

同时, 由于 \mathbf{Y} 各行均为 One-hot 向量, 于是有

$$\sum_{v=1}^{n^{(M)}} y_{i,v} = 1$$

再由(2-3-2)式、(3-5-2)式和(3-5-3)式，有

$$\begin{aligned}
\left(\frac{\partial \mathcal{L}}{\partial \mathbf{X}^{(M)}}\right)_{i,j} &= \sum_{v \neq j} (-\hat{y}_{i,v} \hat{y}_{i,j}) \left(-\frac{1}{s} \frac{y_{i,v}}{\hat{y}_{i,v}}\right) + \hat{y}_{i,j} (1 - \hat{y}_{i,j}) \left(-\frac{1}{s} \frac{y_{i,j}}{\hat{y}_{i,j}}\right) \\
&= \sum_{v \neq j} \frac{1}{s} \hat{y}_{i,j} y_{i,v} + \frac{1}{s} \hat{y}_{i,j} y_{i,j} - \frac{1}{s} y_{i,j} \\
&= \sum_{v=1}^{n^{(M)}} \frac{1}{s} \hat{y}_{i,j} y_{i,v} - \frac{1}{s} y_{i,j} \\
&= \frac{1}{s} \hat{y}_{i,j} \left(\sum_{v=1}^{n^{(M)}} y_{i,v}\right) - \frac{1}{s} y_{i,j} \\
&= \frac{1}{s} \hat{y}_{i,j} - \frac{1}{s} y_{i,j} \\
&= \frac{1}{s} ((\hat{\mathbf{Y}})_{i,j} - (\mathbf{Y})_{i,j})
\end{aligned}$$

由此，(6-5)式成立。

2.2.0.6 基于均方差损失的线性激活输出层误差项

$$\frac{\partial \mathcal{L}}{\partial \mathbf{X}^{(M)}} = \frac{2}{n^{(M)} s} (\hat{\mathbf{Y}} - \mathbf{Y}) \quad (6-6)$$

(6-6)式适用于损失函数为均方差、输出层激活函数为线性函数 $\varphi^{(M)}(x) = x$ 的用于解决分类问题的多层感知机。现给出如下证明。

由于 $\hat{\mathbf{Y}} = \varphi^{(M)}(\mathbf{X}^{(M)}) = \mathbf{X}^{(M)}$ ，于是依据(2-4-2)式，有

$$\begin{aligned}
\left(\frac{\partial \mathcal{L}}{\partial \mathbf{X}^{(M)}}\right)_{i,j} &= \left(\frac{\partial \mathcal{L}}{\partial \hat{\mathbf{Y}}}\right)_{i,j} \\
&= \frac{\partial \mathcal{L}}{\partial \hat{y}_{i,j}} \\
&= \frac{2}{n^{(M)} s} (\hat{y}_{i,j} - y_{i,j}) \\
&= \frac{2}{n^{(M)} s} ((\hat{\mathbf{Y}})_{i,j} - (\mathbf{Y})_{i,j})
\end{aligned}$$

由此，(6-6)成立。

3 训练过程

一般地，机器学习需要经过数据采集、数据分析、数据处理、模型选择、模型训练和模型评估等流程，现简化为数据准备、模型准备和模型训练。

数据准备的主要工作是特征工程（Feature Engineering）。简而言之，特征工程是将原始数据转化为能够更好地表征实体并且适用于特定机器学习模型的特征的过程。例如，在对 28×28 像素灰度图上的手写数字的识别任务中，每张图片由 28×28 矩阵表征，该矩阵的元素为像素灰度值，是 $[0, 255]$ 上的整数。由于多层感知机的输入是 1 维的，故需要将 28×28 矩阵转化为含有 784 个分量的向量。此外，为避免出现收敛困难、数据溢出等问题，还需要对数据进行归一化处理。如前例，所有像素灰度值应被映射到 $[0, 1]$ 上。

模型准备的主要工作是超参数（Hyperparameter）选择。一般而言，超参数是用于定义模型结构和优化方法的一类参数，如多层感知机层数、各层神经元个数、各层激活函数以及梯度下降学习率等。超参数很难通过优化方法自动调整，通常需要人为设置。综上，超参数选择是一项经验性极强的工作。

一般地，在基于小批量梯度下降的训练过程中，对模型进行一次参数调整的一个过程称为一次迭代（Iteration），用于一次迭代的一小批样本数据称为一个批量（Batch），整个训练集全部数据均已被用于参数调整的一个完整过程称为一个时期（Epoch）。合理设置批量大小（Batch Size）、迭代次数和时期次数也是训练的关键环节。现给出一种训练方案的主要过程，见如下算法。

算法: 基于误差反向传播的小批量梯度下降训练主过程

输入: $\mathcal{D} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^{|\mathcal{D}|}$,

$n_{ep}(n_{ep} \geq 1)$,

$s(s \leq |\mathcal{D}|)$

```

1 随机或按特定方法初始化参数组  $\mathcal{P} = \{\mathbf{W}^{(m)}, \mathbf{b}^{(m)}\}_{m=1}^M$ ;
2  for  $e \leftarrow 1$  to  $n_{ep}$  do                                     //  $n_{ep}$ : epoch number
3       $n_{it} \leftarrow \lceil |\mathcal{D}|/s \rceil$ ;                          //  $n_{it}$ : iteration number;  $s$ : batch size
4      对  $\mathcal{D}$  中的所有有序对  $(\mathbf{x}_i, \mathbf{y}_i)$  进行随机排序;
5      按  $s$  将  $\mathcal{D}$  顺序划分为  $n_{it}$  批数据  $\{\mathbf{X}_i, \mathbf{Y}_i\}_{i=1}^{n_{it}}$ ;
6      for  $i \leftarrow 1$  to  $n_{it}$  do
7           $\mathbf{Y}^{(0)} \leftarrow \mathbf{X}_i$ ;
8          for  $m \leftarrow 1$  to  $M$  do
9              依据(1-1)式计算  $\mathbf{X}^{(m)}$  和  $\mathbf{Y}^{(m)}$ ;
10         end
11          $\mathbf{Y} \leftarrow \mathbf{Y}_i$ ;
12         依据(6-5)式或(6-6)式计算输出层误差项  $\partial\mathcal{L}/\partial\mathbf{X}^{(M)}$ ;
13          $\mathbf{W} \leftarrow \mathbf{W}^{(M)}$ ;
14         依据(6-2)式计算  $\partial\mathcal{L}/\partial\mathbf{W}^{(M)}$ ;
15         依据(6-3)式计算  $\partial\mathcal{L}/\partial\mathbf{b}^{(M)}$ ;
16         依据(5-2)式更新  $\mathbf{W}^{(M)}$  和  $\mathbf{b}^{(M)}$ ;
17         for  $m \leftarrow (M-1)$  to 1 do
18             依据(6-4)式和更新前的第  $m+1$  层权重矩阵  $\mathbf{W}$  计算误差项  $\partial\mathcal{L}/\partial\mathbf{X}^{(m)}$ ;
19              $\mathbf{W} \leftarrow \mathbf{W}^{(m)}$ ;
20             依据(6-2)式计算  $\partial\mathcal{L}/\partial\mathbf{W}^{(m)}$ ;
21             依据(6-3)式计算  $\partial\mathcal{L}/\partial\mathbf{b}^{(m)}$ ;
22             依据(5-2)式更新  $\mathbf{W}^{(m)}$  和  $\mathbf{b}^{(m)}$ ;
23         end
24     end
25 end
输出:  $\mathcal{P} = \{\mathbf{W}^{(m)}, \mathbf{b}^{(m)}\}_{m=1}^M$ 

```

参考文献

- [1] Warren S McCulloch, Walter Pitts. A Logical Calculus of the Ideas Immanent in Nervous Activity[J]. Bulletin of Mathematical Biophysics, 1943, 5(1): 115-133.
- [2] F Rosenblatt. The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain[J]. Psychological Review, 1958, 65(6): 386-408.
- [3] G Cybenko. Approximation by Superpositions of a Sigmoidal Function[J]. Mathematics of Control, Signals, and Systems, 1989, 2(1): 303-314.
- [4] Kurt Hornik. Multilayer Feedforward Networks are Universal Approximators[J]. Neural Networks, 1989, 2(5): 359-366.
- [5] Kurt Hornik. Approximation Capabilities of Multilayer Feedforward Networks[J]. Neural Networks, 1991, 4(2): 251-257.
- [6] Paul J Werbos. Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Sciences[D]. Cambridge, MA: Harvard University, 1974.
- [7] David E Rumelhart, Geoffrey E Hinton, Ronald J Williams. Learning Representations by Back-propagating Errors[J]. Nature, 1986, 323(9): 533-536.
- [8] 邱锡鹏. 神经网络与深度学习 [M]. 北京: 机械工业出版社, 2020: 92-95.

代码

神经网络

```
1  #-*- coding: UTF-8 -*-
2  """
3  File Name: neural_network.py
4  """
5  # Python 3.8.16
6
7  import numpy as np  # 1.18.5
8
9
10 class FullyConnectedNN(object):
11     """Fully connected neural network"""
12
13     @property
14     def num_layers(self):
15         """Number of layers"""
16         return self.__num_lyr
17
18     def __init__(self, layer_dims, activ_funcs):
19         # Number of layers
20         self.__num_lyr = len(layer_dims)
21         # Input matrices of hidden and output layers
22         self.lyr_in = dict()
23         # Output matrices of all layers
24         self.lyr_out = {0: np.array([])}
25         # Weight matrices of hidden and output layers
26         self.weight = dict()
27         # Bias column vectors of hidden and output layers
28         self.bias = dict()
29         # Activation functions of hidden and output layers
30         self.__ac_func = dict()
31         # Initialization
32         for i in range(1, self.__num_lyr, 1):
33             self.lyr_in[i] = np.array([])
34             self.lyr_out[i] = np.array([])
35             self.weight[i] = np.random.uniform(
36                 -1, 1, (layer_dims[i], layer_dims[i - 1]))
37         )
38         self.bias[i] = np.random.uniform(-1, 1, (layer_dims[i], 1))
39         self.__ac_func[i] = activ_funcs[i - 1]
```

```

40         # Error matrix of the backward propagation process
41         self.bp_err = np.array([])
42         # Weight matrix stored for the backward propagation process before updating
43         self.bp_wgt = np.array([])
44
45     def grad_loss_out(self):
46         """The partial derivatives of the loss function
47         with respect to elements of an output matrix"""
48         return self.bp_err @ self.bp_wgt
49
50     def grad_loss_weight(self, layer_id):
51         """The partial derivatives of the loss function
52         with respect to elements of a weight matrix"""
53         return self.bp_err.T @ self.lyr_out[layer_id - 1]
54
55     def grad_loss_bias(self):
56         """The partial derivatives of the loss function
57         with respect to elements of a bias column vector"""
58         err0_sum = self.bp_err.sum(0)
59         return err0_sum.reshape((err0_sum.shape[0], 1))
60
61     def forward(self, features):
62         """Forward propagation"""
63         self.lyr_out[0] = features
64         for i in range(1, self.__num_lyr, 1):
65             self.lyr_in[i] = (
66                 self.weight[i] @ self.lyr_out[i - 1].T + self.bias[i]
67             ).T
68             self.lyr_out[i] = self.__ac_func[i](self.lyr_in[i])
69
70
71     class BackpropagationNN(FullyConnectedNN):
72         """Back propagation (backward propagation of errors) neural network"""
73
74     def set_learning_rate(self, rate):
75         """Initialize the learning rate."""
76         self.__lr = rate if 0 < rate < 1 else 0.01
77
78     def cal_out_err(self):
79         """Calculate the error matrix of the output layer."""
80         pass
81
82     def __init__(self, layer_dims, activ_funcs, hid_grad_activ_funcs):

```

```

83     super().__init__(layer_dims, activ_funcs)
84     # Label matrix
85     self.label = None
86     # Learning rate
87     self.__lr = 0.01
88     # Derivatives of all hidden layer activation functions
89     self.__hid_grd_acf = dict()
90     # Output layer id, a key of lyr_out
91     self.max_lid = self.num_layers - 1
92     for i in range(1, self.max_lid, 1):
93         self.__hid_grd_acf[i] = hid_grad_activ_funcs[i - 1]
94
95     def __adjust_weight(self, layer_id):
96         self.weight[layer_id] -= self.__lr * self.grad_loss_weight(layer_id)
97
98     def __adjust_bias(self, layer_id):
99         self.bias[layer_id] -= self.__lr * self.grad_loss_bias()
100
101     def __cal_hid_err(self, layer_id):
102         """Calculate an error matrix of a hidden layer."""
103         self.bp_err = self.__hid_grd_acf[layer_id](
104             self.lyr_in[layer_id]) * self.grad_loss_out()
105         self.bp_wgt = self.weight[layer_id]
106
107     def __backward(self):
108         self.bp_err = self.cal_out_err()
109         self.bp_wgt = self.weight[self.max_lid]
110         self.__adjust_weight(self.max_lid)
111         self.__adjust_bias(self.max_lid)
112         for i in range(self.max_lid - 1, 0, -1):
113             self.__cal_hid_err(i)
114             self.__adjust_weight(i)
115             self.__adjust_bias(i)
116
117     def train_one_batch(self, features, labels):
118         """Train one batch."""
119         self.forward(features)
120         self.label = labels
121         self.__backward()
122
123     def train_one_epoch(self, features, labels, batch_size):
124         """Train one epoch."""
125         num_sp = labels.shape[0]

```

```

126         iteration = (num_sp + batch_size - 1) // batch_size
127         index = np.arange(0, num_sp, 1)
128         # Disrupt the entire training set order.
129         # And this is actually not mandatory to execute.
130         np.random.shuffle(index)
131         ite_idx = 0
132         for i in range(iteration):
133             sub_idx = index[ite_idx:ite_idx + batch_size:1]
134             sub_features = features[sub_idx]
135             sub_labels = labels[sub_idx]
136             self.train_one_batch(sub_features, sub_labels)
137             ite_idx += batch_size
138
139
140 class SoftmaxClassifier(BackpropagationNN):
141     """Classifier based on cross-entropy & softmax"""
142
143     @staticmethod
144     def softmax(np2d_arr):
145         """Softmax function"""
146         row_max = np.amax(np2d_arr, axis=1)
147         exp_matrix = np.exp(np2d_arr - row_max.reshape((row_max.shape[0], 1)))
148         row_sum = np.sum(exp_matrix, axis=1)
149         return exp_matrix / row_sum.reshape((row_sum.shape[0], 1))
150
151     @staticmethod
152     def ce_loss(outputs, labels):
153         """Cross-entropy loss function"""
154         return -np.sum(labels * np.log(
155             np.where(outputs == 0, 0.001, outputs))) / labels.shape[0]
156
157     def cal_out_err(self):
158         """Calculate the error matrix of the output layer."""
159         return (self.lyr_out[self.max_lid] - self.label) / self.label.shape[0]
160
161     def __init__(self, layer_dims, hid_activ_funcs, hid_grad_activ_funcs):
162         """
163         Parameters
164         --
165         layer_dims: tuple or list
166             Numbers of neurons in each layer
167         hid_activ_funcs: tuple or list
168             Activation functions of all hidden layers

```

```

169     hid_grad_activ_funcs: tuple or list
170         Derivatives of all hidden layer activation functions
171     """
172     activ_funcs = [func for func in hid_activ_funcs]
173     activ_funcs.append(self.softmax)
174     super().__init__(layer_dims, activ_funcs, hid_grad_activ_funcs)
175     self.loss = []
176     self.accuracy = []
177
178     def train_mul_epoch(self, features, labels, batch_size, epoch_size,
179                        test_features=None, test_labels=None):
180         """Train multiple epochs."""
181         test_x = features if test_features is None else test_features
182         test_y = labels if test_labels is None else test_labels
183         for i in range(epoch_size):
184             self.train_one_epoch(features, labels, batch_size)
185             #
186             # ** **
187             _features = test_x
188             _labels = test_y
189             self.forward(_features)
190             loss = self.ce_loss(self.lyr_out[self.max_lid], _labels)
191             outputs = np.argmax(self.lyr_out[self.max_lid], axis=1)
192             _labels = np.argmax(_labels, axis=1)
193             accu = np.sum(outputs == _labels) / _labels.shape[0]
194             print('{: <6}'.format(i), end='')
195             print('{: <15.6f}'.format(loss), end='')
196             print('{:.6f}'.format(accu))
197             self.loss.append(loss)
198             self.accuracy.append(accu)
199             # ** **
200             #
201
202     def predict(self, features):
203         """
204         Returns
205         --
206         k_id: numpy.ndarray
207             The ordinal numbers that distinguish their corresponding sample labels
208         """
209         self.forward(features)
210         return np.argmax(self.lyr_out[self.max_lid], axis=1)

```

```

212     def test_accuracy(self, features, labels):
213         """
214         Returns
215         --
216         num: integer
217             The total number of samples whose prediction results matched their labels
218         """
219         self.forward(features)
220         outputs = np.argmax(self.lyr_out[self.max_lid], axis=1)
221         labels = np.argmax(labels, axis=1)
222         return np.sum(outputs == labels)
223
224
225     class MSERegressor(BackpropagationNN):
226         """Regressor based on mean squared error"""
227
228         @staticmethod
229         def ms_loss(outputs, labels):
230             """Mean squared loss function"""
231             delta = outputs - labels
232             return np.sum(delta * delta) / (outputs.shape[0] * outputs.shape[1])
233
234         def cal_out_err(self):
235             """Calculate the error matrix of the output layer."""
236             outs = self.lyr_out[self.max_lid]
237             return 2 * (outs - self.label) / (outs.shape[0] * outs.shape[1])
238
239         def __init__(self, layer_dims, hid_activ_funcs, hid_grad_activ_funcs):
240             """
241             Parameters
242             --
243             layer_dims: tuple or list
244                 Numbers of neurons in each layer
245             hid_activ_funcs: tuple or list
246                 Activation functions of all hidden layers
247             hid_grad_activ_funcs: tuple or list
248                 Derivatives of all hidden layer activation functions
249             """
250             activ_funcs = [func for func in hid_activ_funcs]
251             activ_funcs.append(lambda arr: arr)
252             super().__init__(layer_dims, activ_funcs, hid_grad_activ_funcs)
253             self.loss = []
254

```



```

255     def train_mul_epoch(self, features, labels, batch_size, epoch_size,
256                          test_features=None, test_labels=None):
257         """Train multiple epochs."""
258         test_x = features if test_features is None else test_features
259         test_y = labels if test_labels is None else test_labels
260         for i in range(epoch_size):
261             self.train_one_epoch(features, labels, batch_size)
262             #
263             # ** **
264             _features = test_x
265             _labels = test_y
266             self.forward(_features)
267             loss = self.ms_loss(self.lyr_out[self.max_lid], _labels)
268             print('{: <6}'.format(i), end='')
269             print('{: <15.6f}'.format(loss))
270             self.loss.append(loss)
271             # ** **
272             #
273
274     def predict(self, features):
275         """
276         Returns
277         --
278         k_id: numpy.ndarray
279         """
280         self.forward(features)
281         return self.lyr_out[self.max_lid]

```

激活函数

```

1  # -*- coding: UTF-8 -*-
2  """
3  File Name: activation_function.py
4  """
5  # Python 3.8.16
6
7  import numpy as np  # 1.18.5
8
9
10 def sigmoid(np2d_arr):
11     """Sigmoid function"""
12     return 1 / (1 + np.exp(-np2d_arr))

```

```

13
14
15 def d_sigmoid(np2d_arr):
16     """The partial derivative of the sigmoid function"""
17     sgm = sigmoid(np2d_arr)
18     return sgm * (1 - sgm)
19
20
21 def relu(np2d_arr):
22     """Rectified linear unit function"""
23     return np.maximum(np2d_arr, 0)
24
25
26 def d_relu(np2d_arr):
27     """The partial derivative of the rectified linear unit function"""
28     return np.where(np2d_arr > 0, 1, 0)
29
30
31 def leaky_relu(np2d_arr):
32     """Leaky rectified linear unit function"""
33     return np.minimum(np2d_arr, 0) * 0.01 + np.maximum(np2d_arr, 0)
34
35
36 def d_leaky_relu(np2d_arr):
37     """The partial derivative of the leaky rectified linear unit function"""
38     mtx = np.where(np2d_arr < 0, np2d_arr, 1)
39     return np.where(mtx > 0, mtx, 0.01)
40
41
42 def tan_h(np2d_arr):
43     """Hyperbolic tangent function"""
44     exp = np.exp(-2 * np2d_arr)
45     return (1 - exp) / (1 + exp)
46
47
48 def d_tan_h(np2d_arr):
49     """The partial derivative of the hyperbolic tangent function"""
50     th = tan_h(np2d_arr)
51     return 1 - th * th

```

MNIST 测试

```

1  # -*- coding: UTF-8 -*-
2  """
3  File Name: test_mnist.py
4  """
5  # Python 3.8.16
6
7  import numpy as np  # 1.18.5
8  from tensorflow.keras.datasets import mnist  # tensorflow-gpu 2.3.0
9  import activation_function as af
10 import neural_network as nn
11
12
13 def get_one_hot(one_dim_digit):
14     """Represent a digit with a One-hot."""
15     one_hot = np.zeros(shape=(10,))
16     one_hot[one_dim_digit[0]] = 1
17     return one_hot
18
19
20 # DATA
21
22 (train_features, train_labels), (test_features, test_labels) = mnist.load_data()
23 feature_dim = train_features.shape[1] * train_features.shape[2]
24 # train data
25 train_num = train_features.shape[0]
26 train_features = train_features.reshape((train_num, feature_dim)) / 255
27 train_labels = np.apply_along_axis(
28     get_one_hot, axis=1, arr=train_labels.reshape((train_num, 1))
29 )
30 # test data
31 test_num = test_features.shape[0]
32 test_features = test_features.reshape((test_num, feature_dim)) / 255
33 test_labels = np.apply_along_axis(
34     get_one_hot, axis=1, arr=test_labels.reshape((test_num, 1))
35 )
36
37 # MODEL
38
39 # There are 81 neurons in the hidden layer.
40 classifier = nn.SoftmaxClassifier((feature_dim, 81, 10),
41                                   (af.sigmoid,),
42                                   (af.d_sigmoid,))
43

```

```
44 # TRAINING
45
46 classifier.set_learning_rate(0.55)
47 # batch_size = 9, epoch_size = 11
48 classifier.train_mul_epoch(train_features, train_labels, 9, 11,
49                             test_features, test_labels)
50 classifier.set_learning_rate(0.1)
51 # batch_size = 27, epoch_size = 9
52 classifier.train_mul_epoch(train_features, train_labels, 27, 9,
53                             test_features, test_labels)
54 classifier.set_learning_rate(0.001)
55 # batch_size = 31, epoch_size = 13
56 classifier.train_mul_epoch(train_features, train_labels, 31, 13,
57                             test_features, test_labels)
58
59 # TEST
60
61 print('\naccuracy = ', end='')
62 print(classifier.test_accuracy(test_features, test_labels) / test_num)
```

整理不易，转载请注明作者及出处。

作者：Deng Z

联系：3030299946

出处：binisbull@GitHub