

基于 tensor 实现的神经网络

2333105 Ning Zhiheng

1 理论基础

1.1 矩阵求导

众所周知，对于实值标量函数 $y = f(\mathbf{X})$ 其中 $y \in \mathbb{R}$, $\mathbf{X} \in \mathbb{R}^{p \times q}$, $f: \mathbb{R}^{p \times q} \rightarrow \mathbb{R}$ 。我们有下式成立：

$$dy = \text{tr}(\frac{\partial f}{\partial \mathbf{X}^T} d\mathbf{X}) \quad (1.1)$$

由于没有既成的矩阵复合求导公式，并不能简单遵循实数求导的链式法则认为 $\frac{\partial y}{\partial \mathbf{X}} = \frac{\partial f}{\partial \mathbf{G}(\mathbf{X})} \frac{\partial \mathbf{G}(\mathbf{X})}{\partial \mathbf{X}}$ ，其中 $f: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$, $\mathbf{G}: \mathbb{R}^{p \times q} \rightarrow \mathbb{R}^{m \times n}$ 是成立的。

下面我们从微分的本质来推导矩阵复合函数的求导法则，给定 $y = f(\mathbf{M})$, $\mathbf{M} = \mathbf{A}\mathbf{X}\mathbf{B}$ ，求 $\frac{\partial y}{\partial \mathbf{X}}$

$$dy = \text{tr}(\frac{\partial f}{\partial \mathbf{M}^T} d\mathbf{M}) = \text{tr}(\frac{\partial f}{\partial \mathbf{M}^T} \mathbf{A} d\mathbf{X} \mathbf{B}) = \text{tr}((\mathbf{A}^T \frac{\partial f}{\partial \mathbf{M}} \mathbf{B}^T)^T d\mathbf{X}) \quad (1.2)$$

即 $\frac{\partial y}{\partial \mathbf{X}} = \mathbf{A}^T \frac{\partial f}{\partial \mathbf{M}} \mathbf{B}^T$ 。

(为了节省篇幅，本文后面给定的抽象函数都认为是实值标量函数)，

1.2 监督学习

先回忆一下监督学习的范式，定义样本空间 $\mathcal{X} \in \mathcal{X}$ ，标记空间 $\mathcal{Y} \in \mathcal{Y}$ 以及包含了从 \mathcal{X} 到 \mathcal{Y} 的假设空间 $\mathcal{H} = \mathcal{X} \times \mathcal{Y}$ ，同时考虑损失函数 $\mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ ， P 为样本空间中 \mathcal{H} 的测度，我们要做的就是最小化损失函数，即

$$\min_{w,b} E_p(L(\hat{\mathbf{Y}}, h_{w,b}(\mathbf{X}))) \quad (1.3)$$

其中, $\hat{\mathbf{Y}}$ 为真实值, $h \in \mathcal{H}$ 。

给定样本集 $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, 其中 $\mathbf{x}_i \in \mathbb{R}^d, y_i \in \mathbb{R}, |\mathcal{D}| = n$ 。并且 \mathcal{D} 中每个样本集都是独立从样本空间采样的, 我们有 $\lim_{n \rightarrow +\infty} \frac{1}{n} L(\hat{\mathbf{Y}}, h_{\mathbf{w}, \mathbf{b}}(\mathbf{X})) = E_p(L(\hat{\mathbf{Y}}, h_{\mathbf{w}, \mathbf{b}}(\mathbf{X})))$, 这是根据 Hoeffding 不等式给出的, 其中 $\frac{1}{n} L(\hat{\mathbf{Y}}, h_{\mathbf{w}, \mathbf{b}}(\mathbf{X}))$ 是经验损失函数。

接着, 问题便转换成优化经验损失函数, 即:

$$\mathbf{w}, \mathbf{b} = \arg \min_{\mathbf{w}, \mathbf{b}} L(\hat{\mathbf{Y}}, h_{\mathbf{w}, \mathbf{b}}(\mathbf{X})) \quad (1.4)$$

下面使用梯度下降法来求解这个问题 (完整思想不再赘述, 请读者自行查阅)。梯度下降法的核心是 **沿着梯度反方向的方向, 函数局部下降最快**。即可以通过迭代的方式, 不断优化参数逼近最小值, 表达式如下:

$$\begin{aligned} \mathbf{w}^{(k+1)} &= \mathbf{w}^{(k)} - \alpha \nabla_{\mathbf{w}} L \\ \mathbf{b}^{(k+1)} &= \mathbf{b}^{(k)} - \alpha \nabla_{\mathbf{b}} L \end{aligned} \quad (1.5)$$

其中 α 是步长, $\mathbf{w}^{(k)}$ 、 $\mathbf{b}^{(k)}$ 是第 k 次迭代的参数。

按照上述流程, 假如我们一直迭代, 便可以得到函数取**全局最小值**时的参数 \mathbf{w}, \mathbf{b} 。细心的读者可能会发现, 既然用这么简单的方法就可以求解, 那为什么深度学习 (神经网络) 中还有如此多调参的 “艺术”?

这需要解释三个核心的点:

- 1) 沿着梯度反方向函数值确实会减少, 但如果 α 太大, 可能就适得其反, 如图 1.1, 所以正确的调参模式应该为:

$$\begin{aligned} \mathbf{w}^{(k+1)} &= \arg \min_{\mathbf{w}} L(\mathbf{w}^{(k)} - \alpha \nabla_{\mathbf{w}} L) \\ \mathbf{b}^{(k+1)} &= \arg \min_{\mathbf{b}} L(\mathbf{b}^{(k)} - \alpha \nabla_{\mathbf{b}} L) \end{aligned} \quad (1.6)$$

由于 \mathbf{w}, \mathbf{b} 是变值, 每一次迭代取得最小值的 α 便不是唯一的, 可以使用精确步长搜索方法确定 α 。为了简单实践, 本文采用简单批量随机梯度下降法, 并设置 $\alpha = 0.01$, 这可以近似保证函数的值呈现下降趋势。

- 2) 简单梯度下降法是一个贪心的算法, 因此不能保证函数在**全局**上下降得最快, 若求解一个可解问题所需的时间很多, 那么算法便失去了意义, 这时候我们说这个问题是 “无解” 的。

- 3) 通常来说我们希望损失函数为凸函数，因为梯度下降法一旦收敛便是该函数的全局最小值点；而对于非凸函数来说，由于存在多个局部最小值点，一旦算法跳进局部陷阱，便很难跳出来。

所以为了更好更快地使函数收敛到最小值，神经网络超参数调优也是一个值得研究的领域。

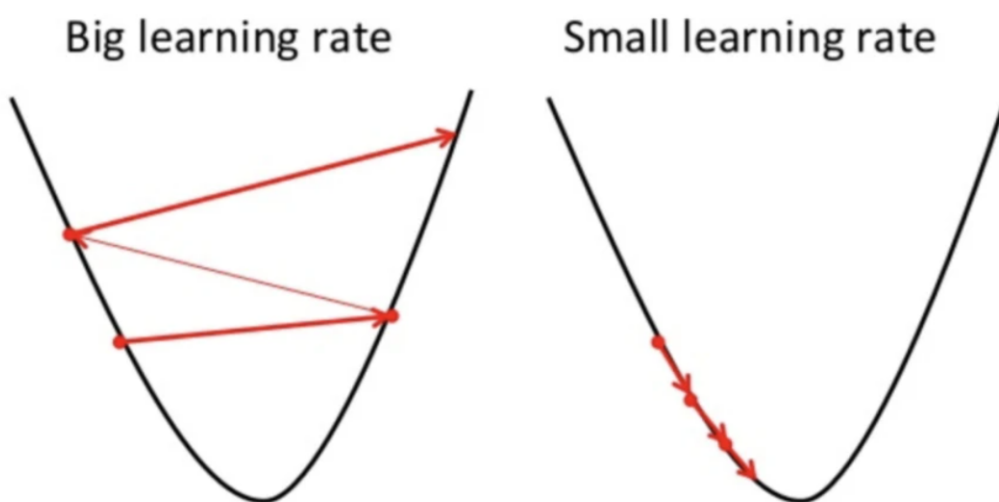


图 1.1: 步长选择策略

2 网络结构

2.1 全连接层

首先给出全连接层的数学模型，设输入该层的数据为 $\mathbf{X} \in \mathbb{R}^{n \times d}$ ，权重矩阵为 $\mathbf{W} \in \mathbb{R}^{d \times p}$ ，偏置向量为 $\mathbf{b} \in \mathbb{R}^{1 \times p}$ ，输出为 $\mathbf{Y} \in \mathbb{R}^{n \times p}$ 。则：

$$\mathbf{Y}_{n \times p} = \mathbf{X}_{n \times d} \mathbf{W}_{d \times p} + \mathbf{1}_{n \times 1} \mathbf{b}_{1 \times p} \quad (2.1)$$

这实际完成了特征从 d 维空间到 p 维空间的映射。

2.1.1 前向传播

因此，我们可以写出前向传播对应的函数：

```
def forward(self, x):
    self.x = x
    sample_dim = x.shape[0]
    self.tensor_one = torch.ones([sample_dim, 1], dtype=torch.
                                  float64)

    x = torch.mm(x, self.w) + torch.mm(self.tensor_one, self.b)
    return x
```

2.1.2 反向传播

根据 1.1，容易给出 L 对 (W, b, X) 的偏导：

$$\begin{aligned}\nabla_W L &= \mathbf{X}^T \frac{\partial L}{\partial \mathbf{Y}} \\ \nabla_b L &= \mathbf{1}_{n \times 1}^T \frac{\partial L}{\partial \mathbf{Y}} \\ \nabla_X L &= \frac{\partial L}{\partial \mathbf{Y}} \mathbf{W}^T\end{aligned}\quad (2.2)$$

因此我们可以写出反向传播对应的函数：

```
def backward(self, top_gard):
    self.dw = torch.mm(self.x.transpose(0, 1), top_gard)
    self.db = torch.mm(self.tensor_one.transpose(0, 1), top_gard)

    bottom_gard = torch.mm(top_gard, self.w.transpose(0, 1))
    return bottom_gard
```

2.2 Relu 层

给出 Relu 表达式：

$$\text{Relu}(\mathbf{X}) = \max(\mathbf{X}, 0) = \begin{cases} \mathbf{X} & \mathbf{X} \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (2.3)$$

给出 Relu 的图像，如下图所示：

可以看到这是一个非线性的函数，可以作为全连接层的激活函数。

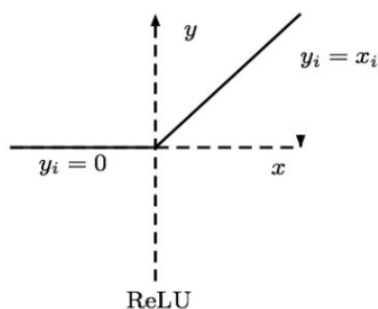


图 2.1: Relu 函数示意图

2.2.1 前向传播

因此，我们可以写出前向传播对应的函数：

```
def forward(self, x):
    self.x = x
    x = torch.maximum(x, torch.tensor([0.0]))
    return x
```

2.2.2 反向传播

根据式(2.3)，求出 $\text{Relu}(x)$ 的导数：

$$\text{Relu}'(x) = \begin{cases} 1 & x \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (2.4)$$

接着写出反向传播对应的公式：

$$\nabla_{\mathbf{X}} L = \nabla_{\mathbf{Y}} L \circ \text{Relu}'(\mathbf{X}) = \begin{cases} \nabla_{\mathbf{Y}} L & \mathbf{X} \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (2.5)$$

其中 \circ 为哈达玛积。

因此我们可以写出反向传播对应的函数：

```
def backward(self, top_grad):
    bottom_grad = top_grad
    bottom_grad[self.x < 0] = 0
    return bottom_grad
```

2.3 softmax 层

给出多分类情况下大家常用的 softmax 表达式:

$$\text{softmax}(\mathbf{x}) = \frac{\exp(\mathbf{x})}{\sum_{i=1}^d \exp(\mathbf{x})} = \frac{\exp(\mathbf{x})}{\mathbf{1}^T \exp(\mathbf{x})} \quad (2.6)$$

其中 $\mathbf{x} \in \mathbb{R}^d$, 可以看到本质上就是一个将向量中每个元素映射到区间 $[0, 1]$ 的函数, 即 $f: \mathbb{R}^n \rightarrow \mathbb{R}_{[0,1]}$ 。若 $\mathbf{X} \in \mathbb{R}^{n \times d}$ 为矩阵, 则公式变为:

$$\begin{aligned} \text{softmax}(\mathbf{X}) &= (\text{softmax}(\mathbf{x}_1), \dots, \text{softmax}(\mathbf{x}_n))^T \\ &= \text{diag}^{-1}\{\exp(\mathbf{X})\mathbf{1}\} \exp(\mathbf{X}) \end{aligned} \quad (2.7)$$

其中 $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$ 。

读者看到这里可能会感觉疑惑, $\text{softmax}(\mathbf{X})$ 与 \mathbf{X} 都是矩阵的形式, 而我们从未提及 $\frac{\partial \text{softmax}(\mathbf{X})}{\partial \mathbf{X}}$ 这样矩阵对矩阵的求导形式, 该怎么办?

事实上, 在神经网络中, softmax 与交叉熵函数是成双成对出现的。为了便于后续说明, 在此先给出交叉熵函数的数学表达式:

$$L(\hat{\mathbf{y}}, \mathbf{y}) = \sum_{i=1}^d -\hat{y}_i \log(y_i) = -\hat{\mathbf{Y}}^T \log(\mathbf{Y}) \quad (2.8)$$

其中 $\mathbf{y} = (y_1, \dots, y_d)^T, \hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_d)^T, \hat{\mathbf{Y}} = (\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_n)^T, \hat{y}_i = \begin{cases} 1 & i = k \\ 0 & i \neq k \end{cases}$,

表明样本属于第 k 类。

考察式(2.8), 容易发现其为一个标量。于是尝试把式(2.7)与式(2.8)结合起来以实现标量 L 对矩阵 \mathbf{X} 的求导。至此, 给出重新定义的 softmax 层的公式:

$$\begin{aligned} L(\hat{\mathbf{Y}}, h(\mathbf{X})) &= \sum_{i=1}^n -\hat{\mathbf{y}}_i \log(\text{softmax}(\mathbf{x}_i)) = \text{tr}(-\hat{\mathbf{Y}} \log \text{softmax}(\mathbf{X})) \\ &= -\text{tr}(\hat{\mathbf{Y}} \log(\text{diag}^{-1}\{\exp(\mathbf{X})\mathbf{1}\})) \end{aligned} \quad (2.9)$$

2.3.1 前向传播公式

因此, 我们可以写出前向传播对应的函数 (这包含了损失函数):

