

# 作业

## Homework3

庄镇华 502022370071

A Neural Networks Homework Assignment



南京大學  
NANJING UNIVERSITY

2023 年 5 月 15 日

2023 年 5 月 15 日

注意：对于 4、5、6 题，请在 pdf 文件中贴代码，并配以必要的注释。

### ✔ 题目一

什么是单层感知机的局限性，造成其局限性的原因是什么，该如何解决？

**解答：**单层感知机的局限指其无法解决线性不可分问题，比如异或问题。对于非线性问题，需要增加神经元，以拟合输入-输出关系。

解决方法是增加网络层数，即增加隐藏层，这样单层感知机就变成了多层感知机，不仅可以解决异或问题，而且具有很好的非线性分类效果。

### ✔ 题目二

相比单层感知机，多层感知机有什么优点，又带来什么问题？

**解答：**多层感知机的优点有：

1. **更好地解决非线性问题。**多层感知机可以通过训练权值参数来学习非线性函数，这对于解决很多实际问题是很有用的。

2. **具有较高的准确率。**多层感知机可以通过训练来获得较高的准确率，尤其是当数据集较大时。

3. **具有较好的泛化能力。**多层感知机可以很好地适应新的数据，对于未来的数据具有较好的泛化能力。

4. **训练方式灵活。**多层感知机可以使用不同的优化算法和损失函数来训练，这使得它可以应用于各种不同的场景。

随着层数的增多带来的问题有：

1. **过拟合**，可通过 Dropout 解决。

2. **参数难以调试**，可通过 Adagrad、Adam、Adadelta 等自适应的梯度下降方法降低调试参数的负担。

3. **梯度弥散。**使用 Sigmoid 在反向传播中梯度值会逐渐减少，经过多层的传递后会呈指数级的剧烈减少，因此梯度值在传递到前面几层时就变得非常小了，神经网络参数的更新将会非常缓慢。

### ✔ 题目三

比较绝对误差与平方误差的优劣，并介绍其他的任意两种误差函数。

**解答：**绝对误差函数曲线呈 V 字型，在 0 处不可导，计算机求解导数比较困难，而且该函数大部分情况下梯度都是相等的，这意味着即使对于小的损失值，其梯度也是大的，不利于模型的收敛和学习。但绝对误差对离群点不那么敏感，更有包容性，因为该函数计算的是误差的绝对值，无论是误差大于 1 还是小于 1，没有平方项的作用，惩罚力度都是一样的。

平方误差处处可导，计算机求解梯度较为容易，且其梯度也根据误差值而动态变化，能较快准确达到收敛。但是从离群点角度来看，平方误差很容易受到异常点的影响，对异常点会赋予较大的权重，如果异常点不属于考虑范围，是由于某种错误导致的，则此函数指导方向将出现偏差。

2023 年 5 月 15 日

其他的误差函数：

### Huber 损失

$$L_{\delta}(y, f(x)) = \begin{cases} \frac{1}{2}(y - f(x))^2 & \text{for } |y - f(x)| \leq \delta, \\ \delta|y - f(x)| - \frac{1}{2}\delta^2 & \text{otherwise.} \end{cases}$$

优点：对异常值更加鲁棒；在最优点附近由于调整为 MSE，梯度更新会随着误差减小而减小，有利于收敛。缺点：引入额外的超参，需要调试；临界点处不可导。

### Log-Cosh 损失

$$L(y, y^p) = \sum_{i=1}^n \log(\cosh(y_i^p - y_i))$$

优点：具有 huber 损失具备的所有优点；二阶处处可微，许多机器学习算法比如 XGBoost 算法采用牛顿法逼近最优点，而牛顿法要求损失函数二阶可微。缺点：误差很大情况下，一阶梯度和 Hessian 会变成定值，导致 XGBoost 出现缺少分裂点的情况。

## ✓ 题目四

描述梯度下降的优缺点和局部最小值的定义，使用梯度下降方法优化 Himmelblau 函数：

$$f(x, y) = (x^2 + y - 11)^2 + (x + y^2 - 7)^2$$

可使用 pytorch 框架。（Himmelblau 函数可视化代码如下）

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from mpl_toolkits.mplot3d import Axes3D
4
5 def him(x):
6     return (x[0]**2+x[1]-11)**2+(x[0]+x[1]**2-7)**2
7
8 x = np.arange(-6,6,0.1)
9 y = np.arange(-6,6,0.1)
10 X,Y = np.meshgrid(x,y)
11 Z = him([X,Y])
12
13 fig = plt.figure()
14 ax = fig.gca(projection='3d')
15 ax.plot_surface(X,Y,Z,cmap='rainbow')
16 ax.set_xlabel('x[0]')
17 ax.set_ylabel('x[1]')
18 ax.set_zlabel('f')
19 fig.show()
```

### 解答：梯度下降的优势

1. 相比大规模数值矩阵，梯度下降算法的迭代求解效率更高。
2. 对于无法计算全域唯一优解的情况，梯度下降仍然能够有效进行最小值点求解。

### 梯度下降的缺点

1. 局部最小陷阱：指的是该点左右两端取值都大于该点，但是该点不是全域最小值点。

2023 年 5 月 15 日

2. 鞍点陷阱：鞍点是指那些不是极值点但梯度为 0 的点
3. 当网络层数特别深时，会存在梯度爆炸和梯度消失现象。

**局部极小值**：如果存在一个  $\epsilon > 0$ ，使的所有满足  $|x - x^*| < \epsilon$  的  $x$  都有  $f(x^*) \leq f(x)$  我们就把点  $x^*$  对应的函数值  $f(x^*)$  称为一个函数  $f$  的局部最小值。

```

1 import torch
2 lr = 1e-3 # 学习率
3 T = 20000 # 迭代次数
4 if __name__ == '__main__':
5     # x代表坐标值(x,y)
6     x = torch.tensor([0., 0.], requires_grad=True)
7     # 定义Adam优化器，学习速率是1e-3
8     optimizer = torch.optim.Adam([x], lr=lr)
9     for step in range(T):
10         # 输入坐标，得到预测值
11         pred = him(x)
12         # 梯度清零
13         optimizer.zero_grad()
14         # 梯度回传，获取坐标的梯度信息
15         pred.backward()
16         # 沿梯度方向更新梯度，优化坐标值
17         optimizer.step()
18
19         if step % (T // 10) == 0:
20             print('step {}: x = {}, f(x) = {}'.format(step, x.tolist(), pred.item()))
21
22 # 输出情况
23 '''
24 step 0:      x = [0.0009999999310821295, 0.0009999999310821295], f(x) = 170.0
25 step 2000:   x = [2.3331809043884277, 1.9540693759918213], f(x) =
26             13.730910301208496
27 step 4000:   x = [2.9820079803466797, 2.0270984172821045], f(x) =
28             0.014858869835734367
29 step 6000:   x = [2.999983549118042, 2.0000221729278564], f(x) = 1.1074007488787174
30             e-08
31 step 8000:   x = [2.9999938011169434, 2.0000083446502686], f(x) =
32             1.5572823031106964e-09
33 step 10000:  x = [2.999997854232788, 2.000002861022949], f(x) = 1.8189894035458565e
34             -10
35 step 12000:  x = [2.9999992847442627, 2.0000009536743164], f(x) =
36             1.6370904631912708e-11
37 step 14000:  x = [2.999999761581421, 2.000000238418579], f(x) = 1.8189894035458565e
38             -12
39 step 16000:  x = [3.0, 2.0], f(x) = 0.0
40 step 18000:  x = [3.0, 2.0], f(x) = 0.0

```

函数图像可视化如下图所示：

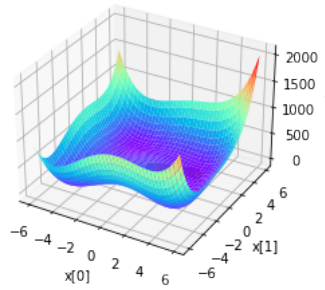


图 1: Himmelblau 函数可视化

### ✓ 题目五

代码设计单层感知机，使其满足  $y = \sigma(Wx + b)$ 。(其函数定义如以下代码，只需完成其实现)

```
1 class Layer():
2     def __init__(self, input_dim, output_dim, bias=True):
3         None
4     def __call__(self, x, train=False):
5         None
6     def forward(self, x, train):
7         None
8     # 反向传播函数
9     def backward(self, error, eta):
10        None
```

解答：单层感知机实现代码如下，具体实现细节见注释。

```
1 import numpy as np
2
3 # 激活函数
4 def relu(x):
5     return np.maximum(x, 0)
6
7 # 激活函数导数
8 def relu_derivative(x):
9     grad = np.array(x, copy=True)
10    grad[x > 0] = 1.
11    grad[x <= 0] = 0.
12    return grad
13
14 class Layer():
15     def __init__(self, input_dim, output_dim, bias=True):
16         super(Layer, self).__init__()
17         # 输入维度
18         self.input_dim = input_dim
19         # 输出维度
20         self.output_dim = output_dim
21         # 随机初始化权重
22         self.w = np.random.randn(self.input_dim, self.output_dim)
```

2023 年 5 月 15 日

```

23     self.bias = bias
24     # 初始化激活函数
25     self.activation = relu
26     self.activation_derivative = relu_derivative
27     # 随机初始化偏置
28     if self.bias:
29         self.b = np.random.randn(self.output_dim)
30
31     def __call__(self, x, train=False):
32         # __call__ 函数使得类对象具有类似函数的功能
33         # 直接内部调用forward函数即可
34         z = self.forward(x, train)
35         return z
36
37     def forward(self, x, train):
38         # 计算wx + b
39         y = np.dot(x, self.w) + self.b
40         # 计算\sigma(wx + b)
41         z = self.activation(y)
42         # 如果是训练模式，则需要计算梯度
43         if train:
44             self.pre_output = x
45             self.grad_z = self.activation_derivative(y)
46         return z
47
48     # 反向传播函数
49     def backward(self, error, eta):
50         '''
51         反向传播过程 error->out->net->w,b
52         '''
53         # d_out_net代表out对net的求导结果
54         d_out_net = self.grad_z
55         # d_net_w代表net对w的求导结果
56         d_net_w = self.pre_output
57
58
59         # d_error_net代表error对net的求导结果，这里应用了链式法则
60         d_error_net = np.multiply(d_out_net, error)
61
62         # d_error_w代表error对w的求导结果，这里应用了链式法则
63         d_error_w = np.dot(d_net_w.T, d_error_net)
64
65         self.dw = d_error_w
66
67         # 计算偏置的权重
68         if self.bias:
69             d_error_b = d_error_net.sum(axis=0)
70             self.db = d_error_b
71
72         # 给下一层的error对out的求导的结果为上一层的加权和
73         self.d_error_out_ = np.dot(d_error_net, self.w.T)
74
75         # eta为学习率，进行梯度更新
76         self.w -= eta * self.dw

```

2023 年 5 月 15 日

```

77     self.b -= eta * self.db
78
79     return self.d_error_out_

```

测试代码如下所示，经测试可以收敛

```

1  import matplotlib.pyplot as plt
2  from sklearn import datasets
3  from sklearn.model_selection import train_test_split
4
5  x, y = datasets.make_moons(n_samples=1000, noise=0.2, random_state=100)
6  x, y = x.reshape(1000, 2), y.reshape(1000, 1)
7  x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.3,
8  random_state=42)
9  print(x.shape, y.shape)
10
11 np.random.seed(2023)
12 layer = Layer(2, 1, True)
13
14 for i in range(800):
15     y_pred = layer.forward(x_train, True)
16     error = y_pred - y_train
17     layer.backward(error, 0.0001)
18     if i % 80 == 0:
19         mse = np.mean(np.square(error))
20         print(mse)

```

## ✓ 题目六

在题目五的基础上，使用 Layer 类完成多层感知机的构造，其成员函数和 Layer 一致，并用于对  $\cos(2\pi x)$  函数的拟合。

解答：多层感知机实现代码如下，具体实现细节见注释。

```

1  def swish(x):
2      beta = 0.2
3      return x * (1 / (1 + np.exp(-beta * x)))
4
5  def swish_derivative(x):
6      beta = 0.2
7      sigma = (1 / (1 + np.exp(-beta * x)))
8      fx = x * sigma
9      return beta * fx + sigma * (1 - beta * fx)
10
11 # 多层感知机
12 class MultiLayer():
13     def __init__(self, layer_dim_list, lr, bias=True):
14         super(MultiLayer, self).__init__()
15         # 单层感知机列表
16         self.layer_list = []
17         # 学习率
18         self.lr = lr
19         for i in range(len(layer_dim_list)-1):
20             input_dim = layer_dim_list[i]

```

2023 年 5 月 15 日

```

21         output_dim = layer_dim_list[i+1]
22         self.layer_list.append(Layer(input_dim,output_dim,bias))
23         # 多层感知机层数
24         self.layer_num = len(self.layer_list)
25
26     def __call__(self,x,train=False):
27         # __call__函数使得类对象具有类似函数的功能
28         # 直接内部调用forward函数即可
29         z = self.forward(x,train)
30         return z
31
32     def forward(self,x,train):
33         x = x.reshape(x.shape[0], -1)
34         out = x
35         for layer in self.layer_list:
36             out = layer.forward(x,True)
37             x = out
38         return out
39
40     # 反向传播函数
41     def backward(self,error,eta):
42         # 从后往前进行梯度更新
43         for idx in range(self.layer_num-1,-1,-1):
44             error = self.layer_list[idx].backward(error,eta)
45         return None

```

经过实验验证，当激活函数为 swish 函数时，对  $\cos(2\pi x)$  函数的拟合效果较好。测试代码如下所示，经测试可以收敛

```

1 x = np.linspace(-0.5, 0.6, 1000, endpoint=True)
2 y = np.cos(2*np.pi*x)
3 x_train, y_train = x.reshape(x.shape[0], 1), y.reshape(y.shape[0], 1)
4
5 np.random.seed(2028)
6 layer = MultiLayer([1,8,1], 0.0004, True)
7 print(layer.layer_num)
8
9 for i in range(80000):
10     y_pred = layer.forward(x_train, True)
11     error = y_pred - y_train
12     layer.backward(error, 0.0004)
13     if i % 800 == 0:
14         mse = np.mean(np.square(error))
15         print(mse)

```

拟合效果如下图所示



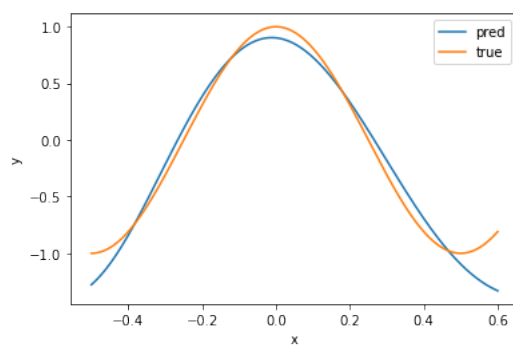


图 2: 余弦函数拟合效果