

Homework Assignment #3

Instructor: Chixiao Chen

Name: Chunyu Wang, FudanID: 20210860017

- This HW counts 15% of your final score, please treat it carefully.
- Please submit the electronic copy via mail: faet_english@126.com before 06/11/2020 11:59pm.
- It is encouraged to use L^AT_EX to edit it, the source code of the assignment is available via: <https://www.overleaf.com/read/mrhqrdztstdzs>
- You can also open it by Office Word, and save it as a .doc file for easy editing. Also, you can print it out, complete it and scan it by your cellphone.
- Problem 2 needs python and numpy. If you do not have a local python environment, please use an online version <https://colab.research.google.com/>.
- You can answer the assignment either in Chinese or English

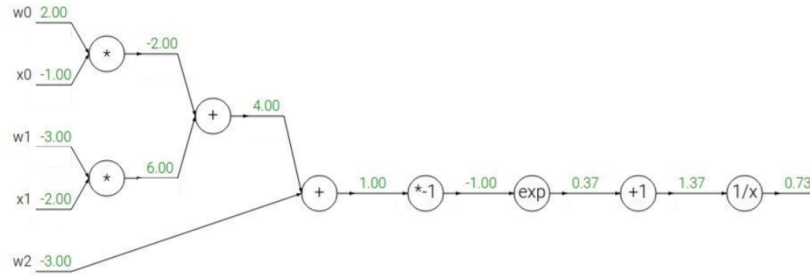
Problem 1: Gradient Computing

(30 points)

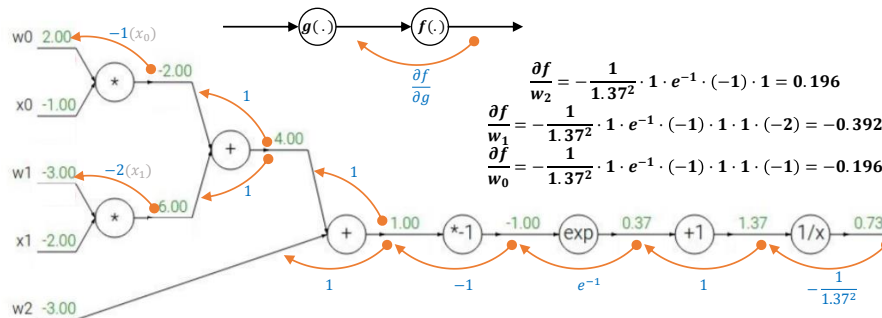
Assuming that one loss function in a classifier has the following output expression:

$$f(x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}},$$

and the current state is shown below:



Please compute all the weight gradients $\frac{\partial f}{\partial w_i}, i = 0, 1, 2$.



Problem 2: Training a two-layer neural network using Numpy

(70 points)

Assuming you have a tiny dataset which has 8 inputs, 4 classes and 500 samples. Please design a two-layer neural network as the classifier. Both forward (inference) and backward (training) propagation are required. The first 400 samples are for training, and the last 100 samples are for test. The dataset is available via: <https://cihlab.github.io/course/dataset.txt>. The activation function is ReLU in the case.

The following table is an example interpretation of the dataset file. (The first two lines of the file is illustrated.)

x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	Class Label
0.4812	0.7790	0.8904	0.7361	0.9552	0.2119	0.7992	0.2409	4
0.4472	0.5985	0.7859	0.5035	0.6912	0.4038	0.0787	0.2301	1

Please submit your code and a brief report with the loss function definition, the final accuracy results, the neuron number in the hidden layers, etc. Also include your strategy for batch size and learning rate. (Hint: It is encouraged to use python and numpy (<https://www.numpy.org/>). You can refer to the slides 34 in the lecture 7 notes. The problem does not encourage you to use Tensorflow/caffe/pytorch, but if you have no idea about numpy, you can also using these frameworks.)

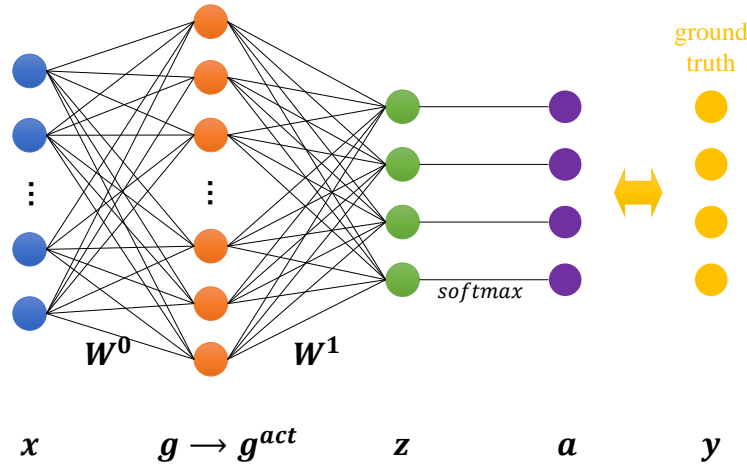


Figure 1: Network Structure

Network Structure A two-layer neural network as shown in Fig.1 is implemented using **pure numpy**. I will give the definition of this network and the specific parameter settings. It has 32 neurons in the hidden layer (g). The entire network can be written in the following form:

$$a = F(x) = (\sigma(xW^{(0)} + b^{(0)})W^{(1)} + b^{(1)})$$

where, $x \in \mathbb{R}^{1 \times 8}$ is the input, $a \in \mathbb{R}^{1 \times 4}$ is the output of F . y is the ground truth in the form of onehot. $W^{(0)} \in \mathbb{R}^{8 \times H}$, $W^{(1)} \in \mathbb{R}^{H \times 4}$, $b^{(0)} \in \mathbb{R}^{1 \times H}$, $b^{(1)} \in \mathbb{R}^{1 \times 4}$ are weights and bias, H indicates the number of neurons in hidden layer ($H = 32$ in our experiment). $\sigma(\cdot)$ is sigmoid function or ReLU as activation function.

For ease of representation, we have the following definition:

$$\begin{aligned} g &= xW^{(0)} + b^{(0)} \\ g^{act} &= \sigma(g) \\ z &= g^{act}W^{(1)} + b^{(1)} \\ a &= S(z) = softmax(z) \end{aligned}$$

Loss Function we use cross entropy loss as the loss function:

$$L = - \sum_i y_i \log a_i$$

Gradient Calculation

$$\begin{aligned}
\frac{\partial L}{\partial \mathbf{z}_i} &= \mathbf{a}_i - \mathbf{y}_i \\
\rightarrow \frac{\partial L}{\partial \mathbf{W}_{ij}^{(1)}} &= \frac{\partial L}{\partial \mathbf{z}_j} \frac{\partial \mathbf{z}_j}{\partial \mathbf{W}_{ij}^{(1)}} = \frac{\partial L}{\partial \mathbf{z}_j} \mathbf{g}_i^{act} \\
\rightarrow \frac{\partial L}{\partial \mathbf{g}_i^{act}} &= \sum_j \frac{\partial L}{\partial \mathbf{z}_j} \frac{\partial \mathbf{z}_j}{\partial \mathbf{g}_i^{act}} = \sum_j \frac{\partial L}{\partial \mathbf{z}_j} \mathbf{W}_{ij}^{(1)} \\
\rightarrow \frac{\partial L}{\partial \mathbf{g}_i} &= \frac{\partial L}{\partial \mathbf{g}_i^{act}} \frac{\partial \mathbf{g}_i^{act}}{\partial \mathbf{g}_i} = \frac{\partial L}{\partial \mathbf{g}_i^{act}} \sigma_i^{-1} \\
\rightarrow \frac{\partial L}{\partial \mathbf{W}_{ij}^{(0)}} &= \frac{\partial L}{\partial \mathbf{g}_j} \frac{\partial \mathbf{g}_j}{\partial \mathbf{W}_{ij}^{(0)}} = \frac{\partial L}{\partial \mathbf{g}_j} \mathbf{x}_i \\
\rightarrow \frac{\partial L}{\partial \mathbf{b}_i^{(1)}} &= \frac{\partial L}{\partial \mathbf{z}_i}, \quad \frac{\partial L}{\partial \mathbf{b}_i^{(0)}} = \frac{\partial L}{\partial \mathbf{g}_i}
\end{aligned}$$

Parameters Update

We use gradients of mini batch with `batch-size=16` to update parameter θ ($\mathbf{W}^{(0)}, \mathbf{W}^{(1)}, \mathbf{b}^{(0)}, \mathbf{b}^{(1)}$):

$$\theta_i := \theta_i - \alpha \frac{1}{m} \frac{\partial L}{\partial \theta_i}$$

where m is `batch-size=16`, α is learning rate. We used an initial learning rate of 1 ($\alpha = 1$), and the learning rate is reduced to its half for every 100 epoch.

$$\alpha := \frac{1}{2} \alpha$$

Experiments & Results

The original dataset was randomly shuffled and then divided into 70% training set and 30% test set. During the training process, the average loss of each epoch on the test set and training set and the accuracy rate on the test set are shown in the figure below.

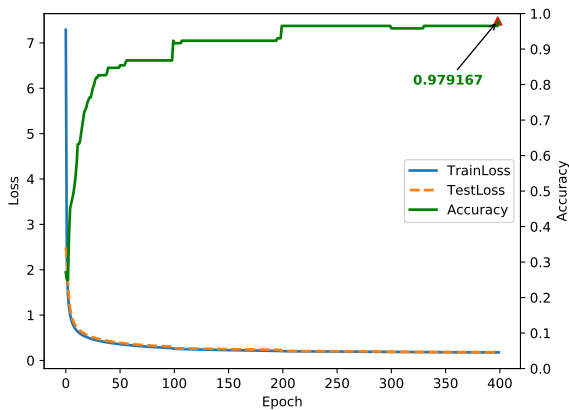


Figure 2: $\sigma = \text{sigmoid}$

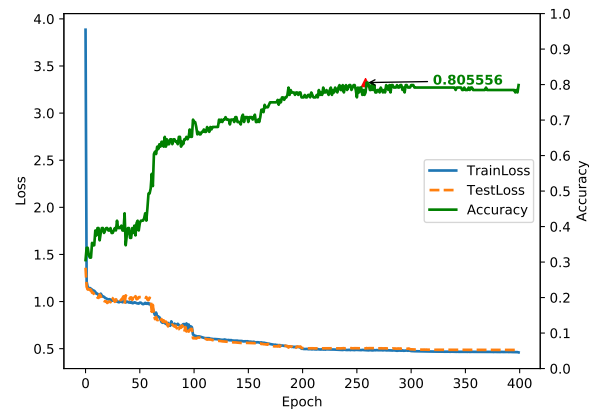


Figure 3: $\sigma = \text{ReLU}$

σ	hidden-size	batch-size	init-lr	accuracy
sigmoid	32	16	1.0	0.979167
ReLU	32	16	1e-1	0.805556

Sigmoid activation function performs better than ReLU under the same conditions, and the final accuracy reaches **0.979167**. Sigmoid performs better on shallow networks such as this task, while ReLU is widely used in deeper networks with the ability of sparsity and handling vanishing gradient problem.

Source Code <https://github.com/yiwangchunyu/AI-RISC/tree/hw3/HW3>

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from matplotlib.ticker import MultipleLocator
4
5 def ReLU(x):
6     return np.maximum(x,0)
7
8 #The derivative of ReLU
9 def ReLU_d(x):
10     return np.where(x > 0, 1, 0)
11
12 def sigmoid(x):
13     return 1/(1+np.exp(-x))
14
15 #The derivative of sigmoid
16 def sigmoid_d(x):
17     return sigmoid(x)*(1-sigmoid(x))
18
19 def softmax(z):
20     t = np.exp(z)
21     a = np.exp(z) / np.sum(t, axis=1).reshape(-1,1)
22     return a
23
24 class DataLoader():
25     def __init__(self,filename='dataset.txt',batch_size=16, shuffle=True, train=True,
26                 train_ratio=0.7):
27         self.filename=filename
28         self.batch_size=batch_size
29         self.inputs,self.labels=self.load_data()
30         if shuffle:
31             self.shuffle()
32         if train:
33             self.inputs, self.labels = self.inputs[:int(self.inputs.shape[0]*train_ratio)], \
34                                     self.labels[:int(self.labels.shape[0]*train_ratio)]
35         else:
36             self.inputs, self.labels = self.inputs[int(self.inputs.shape[0] * train_ratio):], \
37                                     self.labels[int(self.labels.shape[0] * train_ratio):]
38         self.length = self.inputs.shape[0] // self.batch_size
39
40     def shuffle(self):
41         shuffle_ix = np.random.permutation(np.arange(len(self.labels)))
42         self.inputs = self.inputs[shuffle_ix]
43         self.labels = self.labels[shuffle_ix]
44
45     def load_data(self):
46         inputs = []
47         labels = []
48         with open(self.filename, 'r') as f:
49             for line in f:
50                 row = line.split()
51                 labels.append(int(row[-1]))
52                 inputs.append(list(map(float, row[:-1])))
53         return np.array(inputs), np.array(labels)-1
54
55     def __iter__(self):
56         self.id = 0
57         return self
58
59     def __next__(self):
60         if self.id<self.length:
61             inputs = self.inputs[self.id*self.batch_size:(self.id+1)*self.batch_size]

```

```

62         labels = self.labels[self.id * self.batch_size:(self.id + 1) * self.batch_size]
63         self.id+=1
64         return inputs,labels
65     else:
66         raise StopIteration
67
68
69 class Net():
70     def __init__(self, input_size=8, hidden_size=32, bias=True, num_class=4, lr = None,
71         act='sigmoid'):
72         self.hidden_size = hidden_size
73         self.input_size=input_size
74         self.num_class=num_class
75         self.bias=bias
76         self.lr=lr
77         # self.W0=np.zeros((self.input_size, self.hidden_size))
78         # self.W1 = np.zeros((self.hidden_size,self.num_class))
79         self.W1 = np.random.rand(self.hidden_size,self.num_class)
80         self.W0 = np.random.rand(self.input_size, self.hidden_size)
81
82         # self.b0=np.zeros((1,self.hidden_size))
83         # self.b1=np.zeros((1,self.num_class))
84         self.b0 = np.random.rand(1, self.hidden_size)
85         self.b1 = np.random.rand(1, self.num_class)
86         if act=='sigmoid':
87             self.activate=sigmoid
88             self.activate_d = sigmoid_d
89             if not self.lr:
90                 self.lr=1.0
91         elif act=='relu':
92             self.activate = ReLU
93             self.activate_d = ReLU_d
94             if not self.lr:
95                 self.lr=1e-1
96         else:
97             exit(-1)
98
99     # compute cross entropy loss
100     def loss(self,a,y,reduction='mean'):
101         self.y = np.eye(self.num_class)[y] # onehot
102         loss = -np.sum(self.y*np.log(a),axis=1)
103         if reduction=='mean':
104             loss = np.sum(loss) / self.y.shape[0]
105         return loss
106
107     def forword(self,x): # x: batch_size x 8
108         self.x=x
109         self.g = self.x.dot(self.W0)+self.b0
110         self.g_act = self.activate(self.g)
111         self.z = self.g_act.dot(self.W1) + self.b1
112         self.a = softmax(self.z)
113         return self.a
114
115     # calculate gradients
116     def backword(self,y):
117         self.y = np.eye(self.num_class)[y] # onehot
118         self.grad_z=self.a-self.y
119         self.grad_W1=self.g_act.T.dot(self.grad_z)/self.x.shape[0]
120         self.grad_g_act=self.grad_z.dot(self.W1.T)
121         self.grad_g=self.activate_d(self.g)*self.grad_g_act
122         self.grad_W0=self.x.T.dot(self.grad_g)/self.x.shape[0]
123         if self.bias:
124             self.grad_b1=self.grad_z.copy()

```

```

124         self.grad_b0 = self.grad_g.copy()
125
126     # update params (gradient descent)
127     def step(self, lr_shrink=1):
128         lr=lr_shrink*self.lr
129         self.W0=self.W0-lr*self.grad_W0
130         self.W1 = self.W1 - lr * self.grad_W1
131         if self.bias:
132             self.b0=self.b0-lr*self.grad_b0
133             self.b1 = self.b1 - lr * self.grad_b1
134
135     def __call__(self, x):
136         return self.forward(x)
137
138     def train():
139         train_Loader=DataLoader()
140         test_Loader = DataLoader(train=False)
141         net=Net()
142         train_losses,test_losses=[],[]
143         pos,best_acc,accs=0,0,[]
144         lr_shrink = 1
145         for epoch in range(n_epoch):
146             train_loss,test_loss=0,0
147             if (epoch+1)%100==0:
148                 lr_shrink*=0.5
149             for i,(inputs,labels) in enumerate(train_Loader):
150                 outputs=net(inputs)
151                 loss=net.loss(outputs,labels)
152                 train_loss+=loss
153                 net.backward(labels)
154                 net.step(lr_shrink)
155             # test
156             correct=0
157             for i,(inputs,labels) in enumerate(test_Loader):
158                 outputs=net(inputs)
159                 loss=net.loss(outputs,labels)
160                 test_loss+=loss
161                 preds = np.argmax(outputs,axis=1)
162                 correct += (preds==labels).sum()
163
164             # logging...
165             train_loss/=train_Loader.length
166             test_loss/=test_Loader.length
167             train_losses.append(train_loss)
168             test_losses.append(test_loss)
169             acc=correct/(test_Loader.length*test_Loader.batch_size)
170             if acc>best_acc:
171                 best_acc=acc
172                 pos=len(accs)-1
173             accs.append(acc)
174             print("epoch:%d, train_loss:%f, test_loss:%f, acc=%f (%d,%d), best_acc:%f" % (
175                 epoch, train_loss, test_loss,
176                 acc, correct, test_Loader.length*test_Loader.batch_size,best_acc))
177
178         print('best accuracy:', best_acc)
179
180     #plot
181     fig = plt.figure()
182     ax1 = fig.add_subplot(111)
183     plot11=ax1.plot(np.arange(0,len(train_losses)),train_losses
184                     ,linewidth = '2',label='TrainLoss')
185     plot12=ax1.plot(np.arange(0, len(test_losses)), test_losses
186                     ,linewidth = '2',linestyle='--', label='TestLoss')

```

```
187     ax1.set_xlabel('Epoch')
188     ax1.set_ylabel('Loss')
189
190     ax2 = ax1.twinx()
191     plot2=ax2.plot(np.arange(0, len(accs)), accs
192                   ,color='g',linewidth = '2',linestyle='-', label='Accuracy')
193     ax2.set_ylabel('Accuracy')
194     ax2.set_ylim(0,1)
195     y_major_locator = MultipleLocator(0.1)
196     ax2.yaxis.set_major_locator(y_major_locator)
197     ax2.annotate('%f'%(best_acc), (pos,best_acc)
198                 ,xytext=(n_epoch*0.8,0.8),weight='heavy',color='g',
199                 arrowprops=dict(arrowstyle='->'))
200     ax2.scatter(pos,best_acc,color='r',marker='^')
201     lines=plot11+plot12+plot2
202     ax1.legend(lines, [l.get_label() for l in lines],loc='center right')
203     plt.savefig('loss.pdf', dpi=300)
204     plt.show()
205
206 n_epoch=400
207 if __name__=="__main__":
208     train()
```