

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

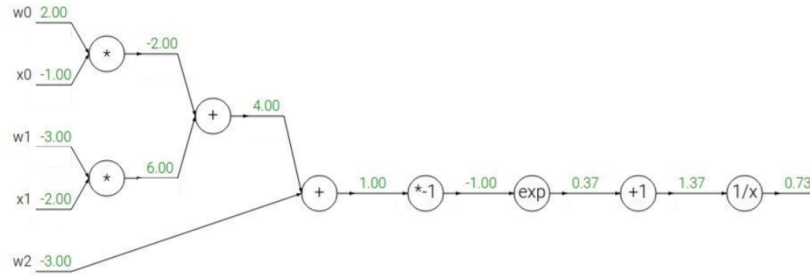
Problem 1: Gradient Computing

(30 points)

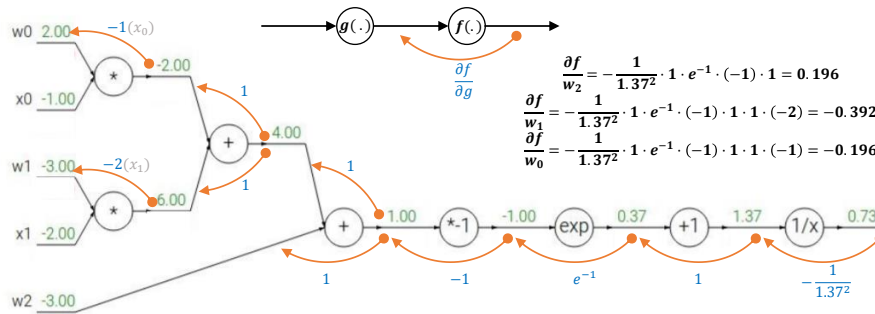
Assuming that one loss function in a classifier has the following output expression:

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

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/cafe/pytorch, but if you have no idea about numpy, you can also using these frameworks.)

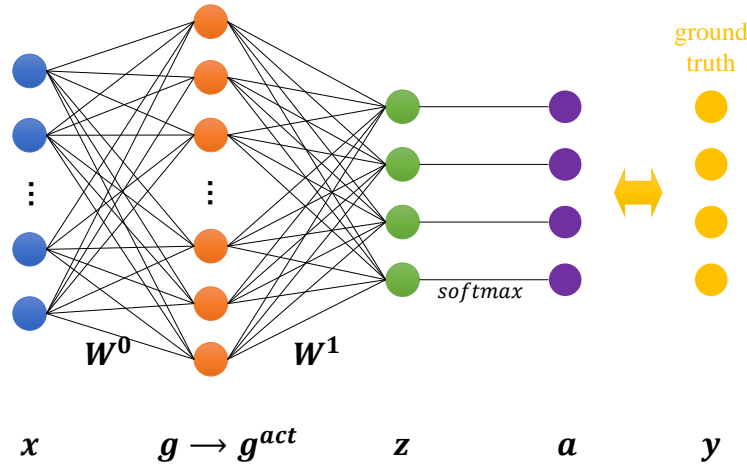


Figure 1: Network Structure

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:

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

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.

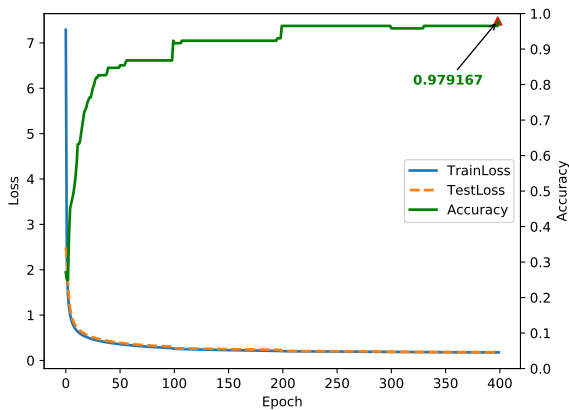


Figure 2: $\sigma = \text{sigmoid}$

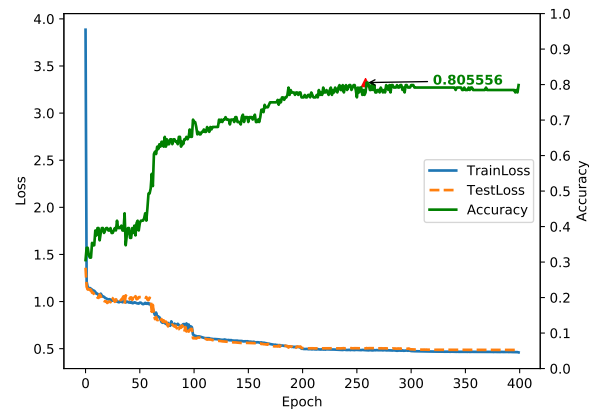


Figure 3: $\sigma = \text{ReLU}$

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 reach **0.979167**. Sigmoid performs better on shallow networks such as this task, while ReLU is widely used in deep networks with the ability of sparsity and handling vanishing gradient problem.

Source Code

```

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

```

```

62     self.input_size=input_size
63     self.num_class=num_class
64     self.bias=bias
65     self.lr=lr
66     # self.W0=np.zeros((self.input_size, self.hidden_size))
67     # self.W1 = np.zeros((self.hidden_size,self.num_class))
68     self.W1 = np.random.rand(self.hidden_size,self.num_class)
69     self.W0 = np.random.rand(self.input_size, self.hidden_size)
70
71     # self.b0=np.zeros((1,self.hidden_size))
72     # self.b1=np.zeros((1,self.num_class))
73     self.b0 = np.random.rand(1, self.hidden_size)
74     self.b1 = np.random.rand(1, self.num_class)
75     self.activate=sigmoid
76
77 # compute cross entropy loss
78 def loss(self,a,y,reduction='mean'):
79     self.y = np.eye(self.num_class)[y] # onehot
80     loss = -np.sum(self.y*np.log(a),axis=1)
81     if reduction=='mean':
82         loss = np.sum(loss) / self.y.shape[0]
83     return loss
84
85 def forward(self,x): # x: batch_size x 8
86     self.x=x
87     self.g = self.x.dot(self.W0)+self.b0
88     self.g_act = self.activate(self.g)
89     self.z = self.g_act.dot(self.W1) + self.b1
90     self.a = softmax(self.z)
91     return self.a
92
93 # calculate gradients
94 def backward(self,y):
95     self.y = np.eye(self.num_class)[y] # onehot
96     self.grad_z=self.a-self.y
97     self.grad_W1=self.g_act.T.dot(self.grad_z)/self.x.shape[0]
98     self.grad_g_act=self.grad_z.dot(self.W1.T)
99     self.grad_g=self.g_act*(1-self.g_act)*self.grad_g_act
100    self.grad_W0=self.x.T.dot(self.grad_g)/self.x.shape[0]
101    if self.bias:
102        self.grad_b1=self.grad_z.copy()
103        self.grad_b0 = self.grad_g.copy()
104
105 # update params (gradient descent)
106 def step(self, lr_shrink=1):
107     lr=lr_shrink*self.lr
108     self.W0=self.W0-lr*self.grad_W0
109     self.W1 = self.W1 - lr * self.grad_W1
110     if self.bias:
111         self.b0=self.b0-lr*self.grad_b0
112         self.b1 = self.b1 - lr * self.grad_b1
113
114 def __call__(self, x):
115     return self.forward(x)
116
117 def train():
118     train_loader=DataLoader()
119     test_loader = DataLoader(train=False)
120     net=Net()
121     train_losses,test_losses=[],[]
122     pos,best_acc,accs=0,0,[]
123     lr_shrink = 1
124     for epoch in range(n_epoch):

```

```

125     train_loss,test_loss=0,0
126     if (epoch+1)%100==0:
127         lr_shrink*=0.5
128     for i,(inputs,labels) in enumerate(train_Loader):
129         outputs=net(inputs)
130         loss=net.loss(outputs,labels)
131         train_loss+=loss
132         net.backward(labels)
133         net.step(lr_shrink)
134     # test
135     correct=0
136     for i,(inputs,labels) in enumerate(test_Loader):
137         outputs=net(inputs)
138         loss=net.loss(outputs,labels)
139         test_loss+=loss
140         preds = np.argmax(outputs,axis=1)
141         correct += (preds==labels).sum()
142
143     # logging...
144     train_loss/=train_Loader.length
145     test_loss/=test_Loader.length
146     train_losses.append(train_loss)
147     test_losses.append(test_loss)
148     acc=correct/(test_Loader.length*test_Loader.batch_size)
149     if acc>best_acc:
150         best_acc=acc
151         pos=len(accs)-1
152     accs.append(acc)
153     print("epoch:%d, train_loss:%f, test_loss:%f, acc=%f (%d,%d), best_acc:%f" % (
154         epoch, train_loss, test_loss,
155         acc, correct, test_Loader.length*test_Loader.batch_size,best_acc))
156
157     print('best accuracy:', best_acc)
158
159     #plot
160     fig = plt.figure()
161     ax1 = fig.add_subplot(111)
162     plot11=ax1.plot(np.arange(0,len(train_losses)),train_losses
163         ,linewidth = '2',label='TrainLoss')
164     plot12=ax1.plot(np.arange(0, len(test_losses)), test_losses
165         ,linewidth = '2',linestyle='--', label='TestLoss')
166     ax1.set_xlabel('Epoch')
167     ax1.set_ylabel('Loss')
168
169     ax2 = ax1.twinx()
170     plot2=ax2.plot(np.arange(0, len(accs)), accs
171         ,color='g',linewidth = '2',linestyle='-', label='Accuracy')
172     ax2.set_ylabel('Accuracy')
173     ax2.set_ylim(0,1)
174     y_major_locator = MultipleLocator(0.1)
175     ax2.yaxis.set_major_locator(y_major_locator)
176     ax2.annotate('%f'%(best_acc),(pos,best_acc)
177         ,xytext=(n_epoch*0.8,0.8),weight='heavy',color='g',
178         arrowprops=dict(arrowstyle='->'))
179     ax2.scatter(pos,best_acc,color='r',marker='^')
180     lines=plot11+plot12+plot2
181     ax1.legend(lines, [l.get_label() for l in lines],loc='center right')
182     plt.savefig('loss.pdf', dpi=300)
183     plt.show()
184
185     n_epoch=400
186     if __name__=="__main__":
187         train()

```