

hw4

Wenhua Bao 2512664  
Xinrui Chen 2503398

August 2, 2020

**Task 1: Neural Networks**

1a)

1.

Get the data: firstly load the data. The output of the original dataset is 1 dimension which includes number 0 9. We change it to 10 dimension, and describe the number with the index. The index of value 1 is the number.

Code:

```

1 mnist_train_in = np.loadtxt("./dataSets/mnist_small_train_in.txt", delimiter=',')
2 mnist_train_out = np.loadtxt("./dataSets/mnist_small_train_out.txt")
3 mnist_test_in = np.loadtxt("./dataSets/mnist_small_test_in.txt", delimiter=',')
4 mnist_test_out = np.loadtxt("./dataSets/mnist_small_test_out.txt")
5 train_y = np.zeros((mnist_train_out.shape[0], 10))
6 test_y = np.zeros((mnist_test_out.shape[0], 10))
7 for k in range(train_y.shape[0]):
8     g = int(mnist_train_out[k])
9     train_y[k][g] = 1
10 for k in range(test_y.shape[0]):
11     g = int(mnist_test_out[k])
12     test_y[k][g] = 1

```

2.

We built a network, in which input has 784 nodes, hidden layer has 30 nodes and output has 10 nodes.

Code:

```

1 class neuralNetwork:
2     def __init__(self, num_input, num_output, hidden_layer):
3         self.input = num_input
4         self.output = num_output
5         self.hidden = hidden_layer
6         self.biases = [np.random.randn(hidden_layer, 1), np.random.randn(num_output, 1)]
7         self.weights = [np.random.randn(num_input, hidden_layer),
8                         np.random.randn(hidden_layer, num_output)]

```

3.

Forward propagation: We choose sigmoid activation for hidden layer and output layer and  $L(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$  as the loss function. Here Nonlinear log-likelihood can also be used, because it is a classification problem.

Code:

```

1 def fw(self, x):
2     a0 = np.dot(x, self.weights[0]).T + self.biases[0]
3     z0 = sigmoid(a0)
4     a1 = np.dot(z0, self.weights[1]).T + self.biases[1]
5     z1 = sigmoid(a1)
6     return z1

```

```

1 def sigmoid(x):
2     return 1.0 / (1.0 + np.exp(-x))
3 def d_sigmoid(x):
4     return sigmoid(x) * (1 - sigmoid(x))

```

4.

Backward propagation: update the weight and bias.

Code:

```

1  def bw(self, x, y):
2      a0 = np.dot(x.T, self.weights[0])+self.biases[0].T
3      z0 = sigmoid(a0)
4      a1 = np.dot(z0, self.weights[1])+self.biases[1].T
5      z1 = sigmoid(a1)
6      # backward pass
7      db1 = (z1 - y) * d_sigmoid(a1)
8      db0 = np.dot(db1, self.weights[1].T) * d_sigmoid(a0)
9      x = np.mat(x).T
10     z0 = np.mat(z0).T
11     dw0 = np.dot(x, db0)
12     dw1 = np.dot(z0, db1)
13     db = [db0, db1]
14     dw = [dw0, dw1]
15     return db, dw

```

```

1  def update_parameter(self, mini_batch, learn_rate):
2      b_new= [np.zeros(b.shape) for b in self.biases]
3      w_new = [np.zeros(w.shape) for w in self.weights]
4      number = len(mini_batch[0])
5      for k in range(number):
6          print('k', k)
7          x = mini_batch[0][k]
8          y = mini_batch[1][k]
9          db, dw = self.bw(x, y)
10
11         b_new = [b + db.T for b, db in zip(b_new, db)]
12         w_new = [w + dw for w, dw in zip(w_new, dw)]
13     self.weights = [w - learn_rate * nw
14                     for w, nw in zip(self.weights, w_new)]
15     self.biases = [b - learn_rate * nb
16                   for b, nb in zip(self.biases, b_new)]

```

5. Train the data: we use the Mini-Batch Gradient Descent to train the data, with the mini-batch size 10, 30, 100. Code:

```

1  def random_mini_batch(self, x, y, size):
2      N = self.input
3      N = x.shape[0]
4      # shuffle
5      x_y = list(zip(x, y))
6      np.random.shuffle(x_y)
7      x, y = zip(*x_y)
8      x = np.array(x)
9      y = np.array(y)
10     mini_batch_num = N//size
11     print('mini_batch_num:', mini_batch_num)
12     mini_batch_list = []
13     for i in range(mini_batch_num):
14         mini_batch_x = x[i * size:(i+1) * size, :]
15         mini_batch_y = y[i * size:(i+1) * size]
16         mini_batch = (mini_batch_x, mini_batch_y) # mini_batch is a tuple
17         mini_batch_list.append(mini_batch)
18     return mini_batch_list

```

```

1  def train(self, x, y, iter_num, mini_batch_size, learn_rate, test_x, test_y):

```

```
2     for j in range(iter_num):
3         mini_batches = self.random_mini_batch(x, y, mini_batch_size)
4         for mini_batch in mini_batches:
5             self.update_parameter(mini_batch, learn_rate)
```

6.

set the network and train with the learning rate 1, 2 and 3. But unfortunately, there are still several errors when we run the code. So we don't have a evaluation for this network.

---

1b)

---

Difference:

What multi-layer neural networks do is feature mapping to values. Features are picked manually. What deep learning does is extract features from the signal and get results/value from the features. The features are selected by the network itself instead of people. In traditional neural networks, back propagation is used, which is an iterative algorithm to train the entire network, randomly setting the initial value, calculating the output of the current network, and then changing the parameters of the previous layers according to the difference between the current output and the label until convergence. In deep learning, we don't use this method, because there are too many layers in deep learning and then the gradient descent method makes less sense.

Which limitations of classical NN does deep learning overcome?

1. when the dataset is very big, it needs to extract features of the raw data as input manually. Irrelevant variables must be ignored, while useful information is retained. But what can be considered as the usefull information is hard to decide. Then using the deep learning is better.
2. Classical NN cannot process time series data (e.g. audio) because classical neural networks do not contain time parameters.
3. In the area of image recognition ,natural language processing and semantic analysis translation, deep learning is better.

# 1 Task2

## 1.0.1 a

1. SVM is the linear classifier with the largest spacing defined on the feature space . 2. Advantages: The classifier only depends on a few data points. Both the original SVM formulation (primal) as well as the derived dual formulation are quadratic programming problems (quadratic cost, linear constraints), which have unique solutions that can be computed efficiently

## 1.0.2 b

$$\begin{aligned} \arg \min_{w, b} \frac{1}{2} \|w\|^2 \\ \text{s.t. } y_i (w^\top x_i + b) - 1 \geq 0 \quad \forall i \end{aligned}$$

## 1.0.3 c

In reality the modes have noise to make models not linear separable, slack variables make the model can tolerant noise and become linear separable.

$$\begin{aligned} \arg \min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \zeta_i \\ \text{s.t. } y_i (w^\top x_i + b) \geq 1 - \zeta_i \quad (i = 1, 2, \dots, N) \\ \zeta_i \geq 0 \end{aligned}$$

## 1.0.4 d

$$\begin{aligned} L(w, b, \alpha) &= \frac{1}{2} \|w\|^2 - \alpha_i \sum_{i=1}^N (y_i (w^\top x_i + b) - 1) \\ &= \frac{1}{2} \|w\|^2 - \sum_{i=1}^N \alpha_i y_i w^\top x_i - \sum_{i=1}^N \alpha_i y_i b + \sum_{i=1}^N \alpha_i \\ \frac{\partial L}{\partial w} &= w - \sum \alpha_i y_i x_i = 0 \\ w &= \sum \alpha_i y_i x_i \\ \frac{\partial L}{\partial b} &= - \sum \alpha_i y_i = 0 \\ \sum \alpha_i y_i &= 0 \\ L(w, b, \alpha) &= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_j^\top x_i) - \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_j^\top x_i) + \sum_{i=1}^N a_i \\ &= \sum_{i=1}^N a_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_j^\top x_i) \end{aligned}$$

result

$$\begin{aligned} \min \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_i^\top x_j) \\ \text{s.t. } \alpha_i \geq 0 \\ \sum_{i=1}^N \alpha_i y_i = 0 \end{aligned}$$

### 1.0.5 e

we can skip the derivation

### 1.0.6 f

1. Replace every occurrence of a scalar product between features with a kernel function

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$$

2. A kernel  $K(\mathbf{x}_i, \mathbf{x}_j)$  that allows us to compute the scalar product without making the mapping explicit

### 1.0.7 g

## 2 Task3

### 2.1 a

```
def _exponential_cov(u1, u2, theta):
    return theta[0] * np.exp(-5*theta[1]*np.subtract.outer(u1, u2).sum(1))

def conditional(u_new, x, theta):
    A = _exponential_cov(u_new, x, theta)
    B = _exponential_cov(x, x, theta)
    C = _exponential_cov(u_new, u_new, theta)
    m0 = np.linalg.inv(B).dot(A.T).T.dot(C)
    sigma = C-A.dot(np.linalg.inv(B).dot(A.T))
    return np.squeeze(m0), np.squeeze(sigma)

theta = [1, 1]
delta_0 = _exponential_cov(0, 0, theta)
u = [1]
y = [np.random.normal(scale=delta_0)]
delta_1 = _exponential_cov(u, x, theta)
def predict(u, data, kernel, theta, sigma, y):
    k = [kernel(u, y, theta) for y in data]
    sim = np.linalg.inv(sigma)
    z_pre = np.dot(k, sim).dot(y)
    sigma_new = kernel(u, x, theta)*np.dot(k, sim).dot(k)
    return y_pre, sigma_new
x_pre = np.arange(0.5*np.pi, 0.6*np.pi)
y_pre = np.sin(x_pre)*np.sin(x_pre)
predictions = [predict(i, x, _exponential_cov, theta, delta_1, y) for i in x_pre]
y_pre, sigma = np.transpose(predictions)
```

