

# CS324 Assignment 1

Li Xuanran  
12312110

October 26, 2025

## 1 Part I: The perceptron (20 pts)

### 1.1 Task 1

Simply use `numpy`. Setting `np.random.seed(1)` is for reproducing the same random numbers, and the parameter can be any number. Labels distinguish the different 2D Gaussian distributions.

---

```
1 def generate(mean=2, cov=1):
2     np.random.seed(11)
3     mean1 = [mean, mean]
4     cov1 = [[cov, cov], [cov, cov]]
5     mean2 = [-mean, -mean]
6     cov2 = [[cov, -cov], [-cov, cov]]
7     data1 = np.random.multivariate_normal(mean1, cov1, 100)
8     data2 = np.random.multivariate_normal(mean2, cov2, 100)
9     labels1 = np.zeros(100)
10    labels2 = np.ones(100)
11    train_data = np.vstack((data1[:80], data2[:80]))
12    train_labels = np.hstack((labels1[:80], labels2[:80]))
13    test_data = np.vstack((data1[80:], data2[80:]))
14    test_labels = np.hstack((labels1[80:], labels2[80:]))
```

---

Listing 1: Generate a dataset

### 1.2 Task 2

We compute the sign function `y`, which determines the mask operation.

---

```
1 def forward(self, input_vec):
2     label_0 = np.ones([input_vec.shape[0], 1])
3     matrix = np.hstack((input_vec, label_0))
4     y = np.dot(self.weights, matrix.T)
5     y = np.array([1 if i >= 0 else -1 for i in y])
6     return y
```

---

Listing 2: forward function

According to **perceptron\_tutorial.pdf**, we should only keep the sample  $i$  where  $\text{pred}[i] * \text{labels}[i] < 0$  to calculate the gradient, So we define `mask` matrix, and `True` means we pick up the sample.

---

```

1 def train(self, training_inputs, labels):
2     for _ in range(self.max_epochs):
3         p = self.forward(training_inputs)
4         mask = p * labels < 0
5         X, y = training_inputs[mask], labels[mask]
6         grad = -np.sum(X * y.reshape(y.shape[0], 1)) / training_inputs.
            shape[0]
7         self.weights = self.weights - self.learning_rate * grad

```

---

Listing 3: train function

### 1.3 Task 3

We train the model with  $\text{epoch} = 20$  and  $\alpha = 0.01$ , setting  $\text{mean} = 2$  and  $\text{cov} = 1$  by default. For accuracy, we can directly invoke `accuracy_score` function in `sklearn.metrics`. **Listing 4** only shows the main part.

---

```

1 if __name__ == '__main__':
2     model = Perceptron(n_inputs=2, max_epochs=20, learning_rate=0.01)
3     X_train, X_test, y_train, y_test = generate(2, 1)
4     model.train(X_train, y_train)
5     p = model.forward(X_test)
6     print(f'mean={2}, cov={1}, Accuracy: {accuracy_score(p, y_test)}')

```

---

Listing 4: Training and accuracy

### 1.4 Task 4

We took  $\mu = \{20, 2, 0.2, 0.02\}$  and  $\sigma = \{10, 1, 0.1, 0.01\}$ .

---

```

1 means, covs = [20, 2, 0.2, 0.02], [10, 1, 0.1, 0.01]
2 for mean in means:
3     for cov in covs:
4         model = Perceptron(n_inputs=2, max_epochs=20, learning_rate=0.01)
5         X_train, X_test, y_train, y_test = generate(mean, cov)
6         model.train(X_train, y_train)
7         p = model.forward(X_test)
8         print(f'mean={mean}, cov={cov}, Accuracy: {accuracy_score(p,
            y_test)}')

```

---

Listing 5: Experiment with different sets of points

The result shows in **Table 1**.

We found if the means are too close (which means  $\mu$  is too small), or their variances is too high (which means  $\sigma$  is too big), the accuracy score will be much lower.

$\mu \backslash \sigma$	10	1	0.1	0.01
20	1.0	1.0	1.0	1.0
2	0.875	1.0	1.0	1.0
0.2	0.7	0.35	0.5	0.5
0.02	0.275	0.65	0.45	0.5

Table 1: Accuracy of different means and variances

## 1.5 Theoretical Analysis

The perceptron is a binary classifier with:

1. Input:  $\mathbf{x} \in \mathbb{R}^n$
2. Weighed vector:  $\mathbf{w} \in \mathbb{R}^n$
3. Bias:  $b$
4. Output:

$$\hat{y}_i = \begin{cases} 1, & \text{if } \mathbf{w}^T \mathbf{x} + b \geq 0, \\ -1, & \text{otherwise.} \end{cases}$$

For each training example, we have forward and backward process:

1. Forward:

$$\hat{y}_i = \text{sign}(\mathbf{w}^T \mathbf{x} + b).$$

2. Backward: Update  $\mathbf{w}$  and  $b$  by

$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} + \eta(y_i - \hat{y}_i)\mathbf{x}_i, \\ b &\leftarrow b + \eta(y_i - \hat{y}_i). \end{aligned}$$

## 2 Part II: The mutli-layer perceptron (60 pts)

### 2.1 Task 1 & Theoretical Analysis

1. Hidden Layers (Layer 1 to  $N - 1$ ): **ReLU**

$$x^{(l)} = \max(0, \tilde{x}^{(l)}), \quad \frac{\partial x^{(l)}}{\partial \tilde{x}^{(l)}} = \mathbf{1}_{\tilde{x}^{(l)} > 0}.$$

2. Output Layer (Layer  $N$ ): **Softmax**

$$x^{(N)} = \frac{\exp(\tilde{x}^{(N)})}{\sum_{i=1}^{d_N} \exp(\tilde{x}^{(N)})_i}, \quad \frac{\partial x_i^{(N)}}{\partial \tilde{x}_j^{(N)}} = x_i^{(N)}(\delta_{ij} - x_j^{(N)}).$$

3. Loss function: **Cross-Entropy**

$$L(x^{(N)}, t) = - \sum_i t_i \log x_i^{(N)}, \quad \frac{\partial L}{\partial \tilde{x}^{(N)}} = x^{(N)} - t.$$

4. **Backpropagation:** For the last linear layer  $L$ , the gradient of the loss with respect to the output  $\hat{y}$  is:

$$\frac{\partial L}{\partial \hat{y}} = \frac{\partial L}{\partial x^{(L)}}$$

Then, the gradient w.r.t. weights is:

$$\frac{\partial L}{\partial W^{(l)}} = \frac{\partial L}{\partial \tilde{x}^{(l)}} (x^{(l-1)})^T.$$

The gradient w.r.t. bias is:

$$\frac{\partial L}{\partial b^{(l)}} = \sum_i \frac{\partial L}{\partial \tilde{x}_i^{(l)}}.$$

Finally, we have the gradient w.r.t. the previous layer:

$$\frac{\partial L}{\partial x^{(l-1)}} = (W^{(l)})^T \frac{\partial L}{\partial \tilde{x}^{(l)}}.$$

For hidden layers with ReLU:

$$\frac{\partial L}{\partial \tilde{x}^{(l)}} = \frac{\partial L}{\partial x^{(l)}} \odot \mathbf{1}_{\tilde{x}^{(l)} > 0}.$$

---

```

1 class ReLU(object):
2     def forward(self, x):
3         self.x = x
4         return np.maximum(0, x)
5
6     def backward(self, dout):
7         dx = dout * (self.x > 0).astype(float)
8         return dx

```

---

Listing 6: ReLU function

---

```

1 class SoftMax(object):
2     def forward(self, x):
3         shift_x = x - np.max(x, axis=1, keepdims=True)
4         exp_x = np.exp(shift_x)
5         return exp_x / np.sum(exp_x, axis=1, keepdims=True)
6
7     def backward(self, dout):
8         return dout

```

---

Listing 7: Softmax function

---

```

1 class CrossEntropy(object):
2     def forward(self, x, y):
3         pred_clipped = np.clip(x, 1e-12, 1.-1e-12)
4         loss = -np.sum(y * np.log(pred_clipped), axis=1)
5         return np.mean(loss)
6
7     def backward(self, x, y):
8         return x-y

```

---

Listing 8: Cross-Entropy function

## 2.2 Task 2

We use `make_moons()` to generate the training and testing data.

---

```

1 dnn_hidden_units = [int(x) for x in dnn_hidden_units.split(',')]
2 X, y = make_moons(n_samples=1000, shuffle=True)
3 n_values = np.max(y) + 1
4 y = np.eye(n_values)[y]
5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
6     shuffle=True)
7
8 model = MLP(n_inputs=2, n_hidden=dnn_hidden_units, n_classes=2)
9 loss = CrossEntropy()

```

---

Listing 9: Load data and initialize the MLP model

Line 3 and line 4 is to transform the 0/1 classification to one-hot encoding possibility.

---

```

1     for step in range(max_steps):
2         y = model.forward(X_train)
3         l = loss.forward(y, y_train)
4         dout = loss.backward(y, y_train)
5         model.backward(dout)
6         model.head.params['weight'] = model.head.params['weight'] -
7             learning_rate * model.head.grads['weight']
8         model.head.params['bias'] = model.head.params['bias'] -
9             learning_rate * model.head.grads['bias']
10        for linear in model.linears:
11            linear.params['weight'] = linear.params['weight'] -
12                learning_rate * linear.grads['weight']
13            linear.params['bias'] = linear.params['bias'] - learning_rate
14            * linear.grads['bias']
15
16        if step % eval_freq == 0 or step == max_steps - 1:
17            y = model.forward(X_test)
18            print(f"Step: {step}, Loss: {l}, Accuracy: {accuracy(y, y_test
19                )}")

```

---

Listing 10: Implement the training loop

## 2.3 Task 3

Simply run `mlp.ipynb` in `Part_2` folder. Here we set `MAX_EPOCHS_DEFAULT = 1500`. The visualized result is as follows.

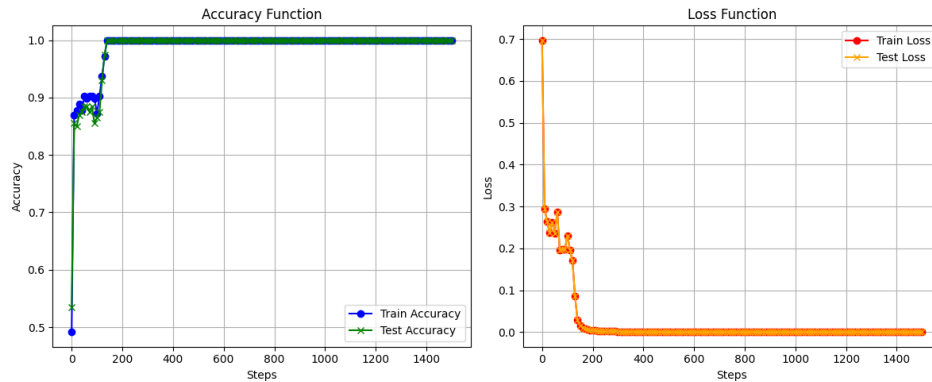


Figure 1: Accuracy function and loss function

From **Figure 1**, we know that after about **180 steps**, the accuracy and loss both converge.

## 3 Part III: Stochastic gradient descent (15 pts)

### 3.1 Task 1

Notion: There are two more parsers added in `train_mlp_numpy.py` in `Part_3`, so the way to run `train_mlp_numpy.py` becomes

---

```
1 python train_mlp_numpy.py --dnn_hidden_units 20 --learning_rate 0.01 --
  max_steps 1500 --eval_freq 10 --sgd True --batch_size 1
```

---

Listing 11: How to run Part 3

where the default `sgd` is `True`, the default `batch_size` is 1.

When using SGD, the training loop becomes

---

```
1 batches = shuffle_and_batch(X_train, y_train, batch_size=batch_size)
2 for X, y in batches:
3     X = X.reshape((1, -1))
4     y_hat = model.forward(X)
5     l = loss.forward(y_hat, y)
6     dout = loss.backward(y_hat, y)
7     model.backward(dout)
8     model.head.params['weight'] = model.head.params['weight'] -
  learning_rate * model.head.grads['weight']
9     model.head.params['bias'] = model.head.params['bias'] - learning_rate
  * model.head.grads['bias']
10     for linear in model.linears:
11         linear.params['weight'] = linear.params['weight'] - learning_rate
  * linear.grads['weight']
```

---

```

12         linear.params['bias'] = linear.params['bias'] - learning_rate *
           linear.grads['bias']

```

Listing 12: sgd

### 3.2 Task 2

We took `batch_size = {1, 10, 50, 200, 500}`. Simply run `train_mlp.ipynb` in `Part_3` folder, and the visualized result is as follows:

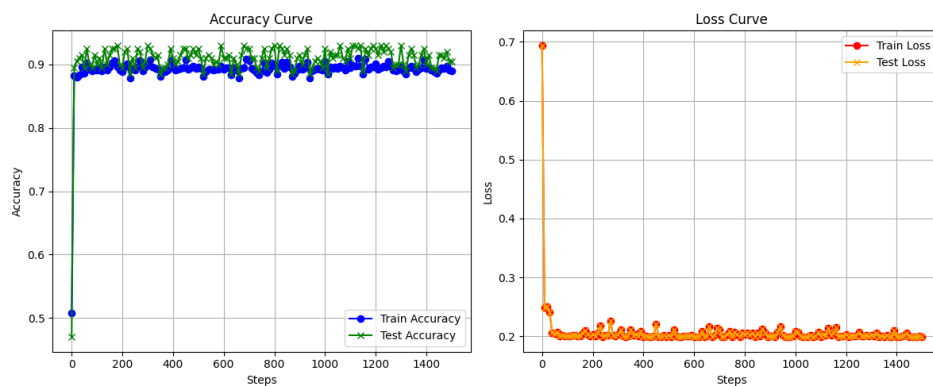


Figure 2: Accuracy and loss function when batch size is 1

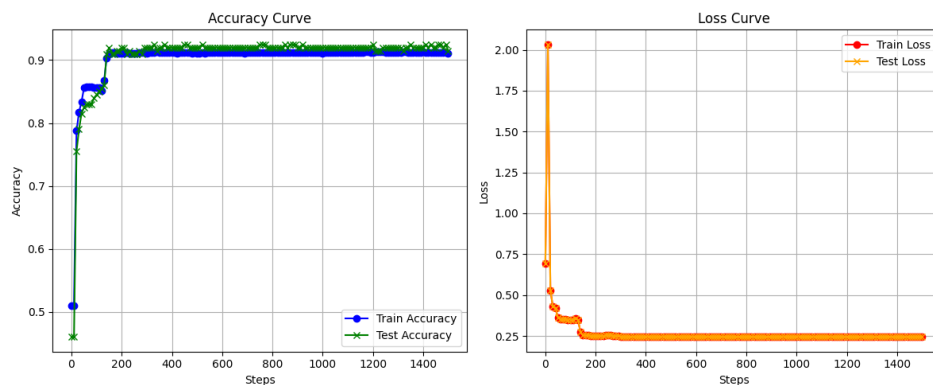


Figure 3: Accuracy and loss function when batch size is 10

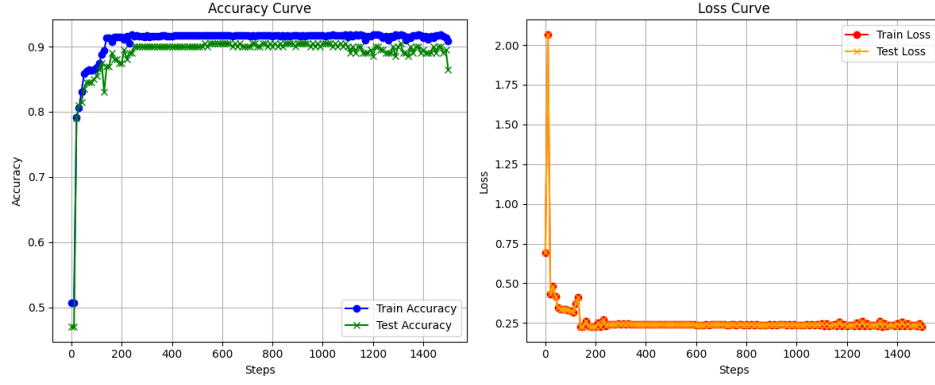


Figure 4: Accuracy and loss function when batch size is 50

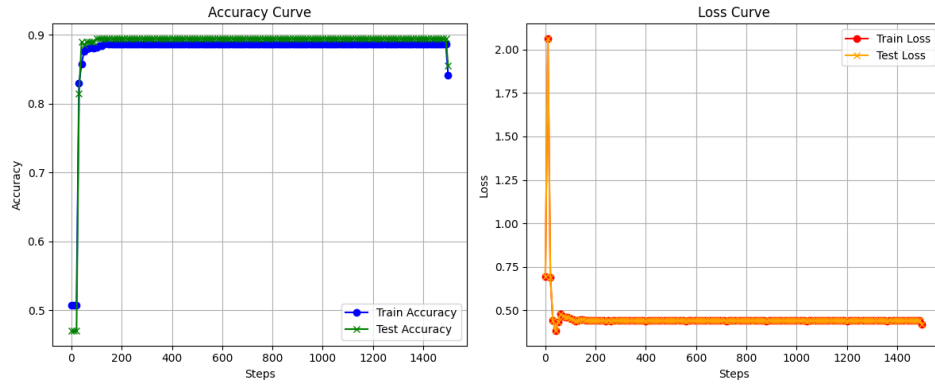


Figure 5: Accuracy and loss function when batch size is 200

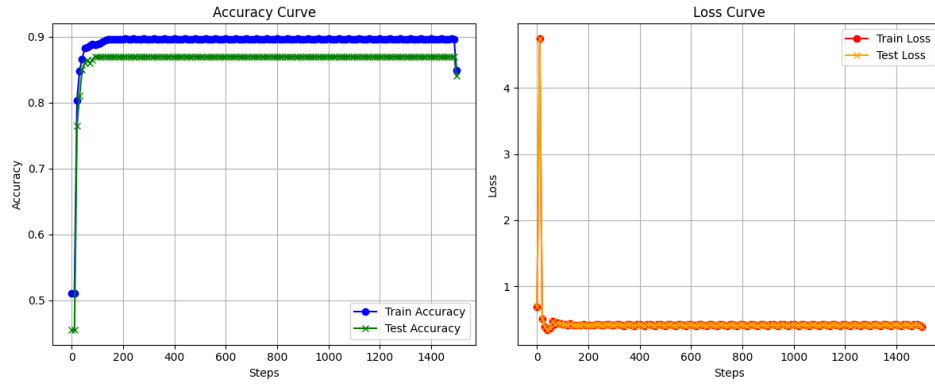


Figure 6: Accuracy and loss function when batch size is 500

We found that when the batch size gets larger, the oscillation of accuracy and loss will decrease. As a result, setting batch\_size= 200 is relatively the best choice.