

Chapter 9

A Tutorial on Neural Networks

9.1 Learning Objectives

This chapter gives a comprehensive tutorial to demonstrate how we implement some key practical considerations (discussed in chapter 8) in real projects through python programming. Examples in chapter 8 have already shown the implementation of weight regularization and dropout in python. The tutorial in this chapter will focus on mini-batch gradient descent and optimization methods. After completing this chapter, one should be able to

- Develop a general neural network from the scratch;
- Implement mini-batch training;
- Implement momentum and Adam optimizations from the scratch;
- Be aware of the effects of hyperparameters on neural networks;

9.2 Architecture of Neural Network

We will build a generic 3-layer neural network, as shown in Fig.1, and apply it to two different tasks throughout this chapter. The size of each layer is defined by the corresponding element in the list *layers_dims*. For example, *layers_dims* =[12288, 25,12,6] specifies 12288 features of input, 25 units in the first hidden layer and 12 units in the second hidden layer, 6 units in the output layer. All hidden layers use ReLU activation functions while the output layer uses softmax for classification. Note that binary classification can be implemented by assigning *layers_dims*[3]=2 with one-hot-coded labels.

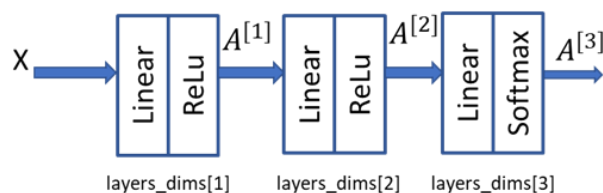


Fig.1 A generic 3-layer neural network

We will train the neural network using mini-batch strategy to feed data, and using three different optimization options: basic gradient descent, momentum and Adam. The network model in python is given by

```
model(X, Y, layers_dims, optimizer, learning_rate = 0.0007, mini_batch_size = 64, beta
    = 0.9, beta1 = 0.9, beta2 = 0.999, epsilon = 1e-8, num_epochs = 10000, print_cost
    = True)
```

The values for arguments are default. The meaning of each argument will be explained in the subsequent texts.

The design flowchart of the model is shown in Fig.2, giving an overall picture of the model. The red fonts indicate functions defined separately. This flowchart can help a reader understand the whole picture of mini-batch concept and its implementation. It is helpful to refer to this flowchart when reading through the subsequent texts.

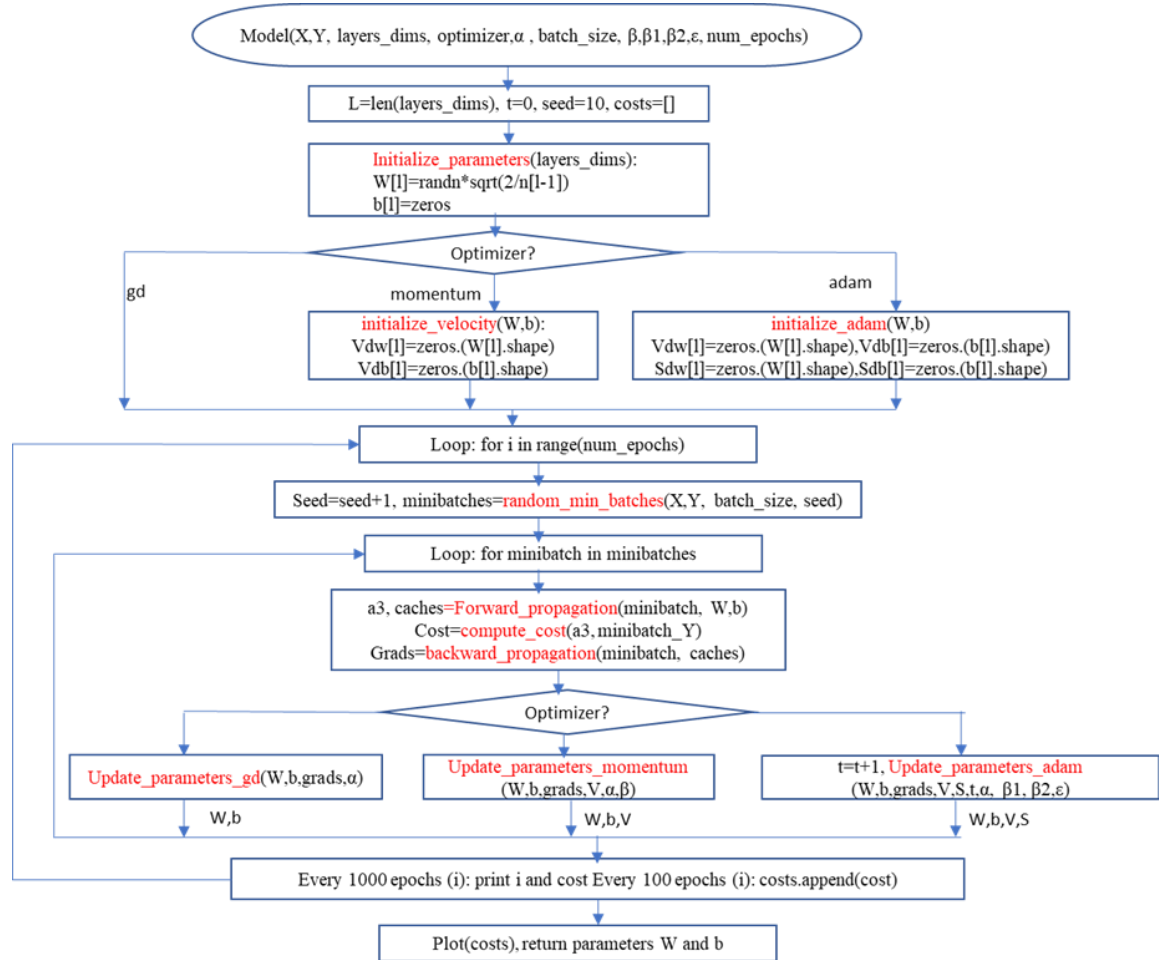


Fig. 2 Design flowchart of the neural network model

The model has the following arguments:

X -- input data, of shape (number of features, number of examples)
Y -- one-hot-label, of shape (K, number of examples), where K is the number of classes.
layers_dims -- python list, containing the size of each layer
learning_rate α -- the learning rate, scalar.
mini_batch_size -- the size of a mini batch
beta -- Momentum hyperparameter
beta1 -- Exponential decay hyperparameter for the past gradients estimates
beta2 -- Exponential decay hyperparameter for the past squared gradients estimates
epsilon -- hyperparameter preventing division by zero in Adam updates
num_epochs -- number of epochs
print_cost -- True to print the cost every 100 epochs

Returns: parameters -- python dictionary containing your updated parameters

The parameters W and b are updated in two nested loops. The outer loop is designed for different epochs while the inner loop works for different batches within an epoch. Each epoch is obtained by randomly shuffling the original dataset, and then partitioned into mini-batches. Each update of parameters is based on one mini-batch. For example, if the total number of epochs is 10,000, and 5 mini-batches in each epoch, then the cost and the parameters W and b will be updated 50,000 times. Please note that we can specify to save and plot some costs (not all) which includes enough information for monitoring the behavior of training. When the mini-batch size is equal to one, the model implements stochastic gradient descent. When the mini-batch size is m , the model implements batch gradient descent. There are two levels of initializations: 1) basically initialize W and b ; and 2) extra initialize V for momentum and (V, S) for Adam.

9.3 Details of Key Functions

In this section, we will describe the design details of some key functions which are used for neural network model construction. Forward propagation, cost function, and backward propagation have been detailed and standardized in previous chapters. Here we focus on mini-batch process, momentum and Adam optimizations.

9.3.1 Mini-batch gradient descent

```
def random_mini_batches(X, Y, mini_batch_size = 64, seed = 0):
```

The function `random_mini_batches ()` divides the entire training set into equally sized portions (i.e. batches). To make each epoch different, the training set needs to be randomly shuffled before being divided:

- 1) **Shuffle:** Create a shuffled version of the training set (X, Y) as shown in Fig.3. Each column of X and Y represents a training example. Note that the random shuffling is done synchronously between X and Y . Such that after the shuffling the i th column of X is the example corresponding to the i th label in Y . The shuffling step ensures that examples will be split randomly into different mini-batches.

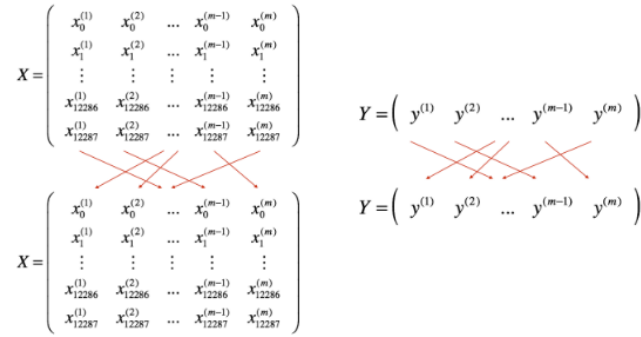


Fig.3 Shuffle training set

- 2) **Partition:** Partition the shuffled (X, Y) into mini-batches of size mini_batch_size (e.g. 64). Note that the number of training examples is not always divisible by mini_batch_size. The last mini batch might be smaller, but you don't need to worry about this. When the final mini-batch is smaller than the full mini_batch_size, it will look like this:

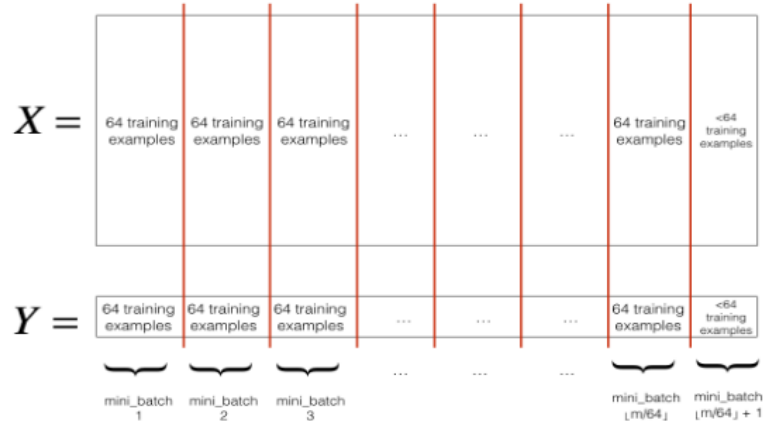


Fig.4 mini-batches

The function `random_mini_batches()` returns a Python list `mini_batches` whose elements are mini-batches.

9.3.2 Optimization with momentum

Because mini-batch gradient descent makes a parameter update after seeing just a subset of examples, the direction of the update has some variance, and so the path taken by mini-batch gradient descent will "zigzag" toward convergence. Using momentum can reduce this zigzagging, as shown in Fig.5. Momentum takes into account the past gradients to smooth out the update. We will store the 'direction' of the previous gradients in the variable `v`. Formally, this will be the exponentially weighted average of the gradient on previous steps. From the previous chapter, we know that optimization with momentum can be described by (17) in the previous chapter (chapter 8.)

Initialization: $v_{dw} = 0, v_{db} = 0$

On iteration t:

Compute dW, db on the current mini-batch

$$v_{dw} = \beta v_{dw} + (1 - \beta) dW$$

$$v_{db} = \beta v_{db} + (1 - \beta) db$$

$$W := W - \alpha v_{dw}$$

$$b := b - \alpha v_{db}$$

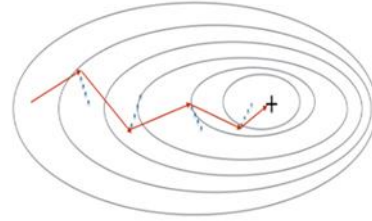


Fig.5 The red arrows show the direction taken by one step of mini-batch gradient descent with momentum. The blue points show the direction of the gradient (with respect to the current mini-batch) on each step.

To implement this momentum, we need to define the following two functions:

```
def initialize_velocity(parameters):
```

```
def update_parameters_with_momentum(parameters, grads, v, beta, learning_rate):
```

where beta (i.e. β) is the momentum and learning_rate (i.e. α) is the learning rate. All parameters (i.e. W, b) should be stored in the parameters dictionary.

Remarks:

- The velocity is initialized with zeros. So the algorithm will take a few iterations to "build up" velocity and start to take bigger steps.
- If $\beta=0$, then this just becomes standard gradient descent without momentum.
- The larger the momentum β is, the smoother the update because the more we take the past gradients into account. But if β is too big, it could also smooth out the updates too much.
- Common values for β range from 0.8 to 0.999. If you don't feel inclined to tune this, $\beta=0.9$ is often a reasonable default.
- Tuning the optimal β for your model might need trying several values to see what works best in term of reducing the value of the cost function J .
- Momentum takes past gradients into account to smooth out the steps of gradient descent. It can be applied with batch gradient descent, mini-batch gradient descent or stochastic gradient descent.
- You usually need to tune a momentum hyperparameter β and a learning rate α .

9.3.3 Optimization with Adam

Adam is one of the most effective optimization algorithms for training neural networks. It combines ideas from RMSProp and Momentum. The algorithm is given by (19) in chapter 8 as follows:

Initialization: $v_{dw} = 0, s_{dw} = 0, v_{db} = 0, s_{db} = 0$

On batch iteration t :

Compute dW, db on the current mini-batch

$$v_{dw} = \beta_1 v_{dw} + (1 - \beta_1) dw,$$

$$v_{db} = \beta_1 v_{db} + (1 - \beta_1) db$$

$$s_{dw} = \beta_2 s_{dw} + (1 - \beta_2)(dw * dw),$$

$$s_{db} = \beta_2 s_{db} + (1 - \beta_2)(db * db)$$

$$v_{dw} := \frac{v_{dw}}{1 - \beta_1^t}$$

$$v_{db} := \frac{v_{db}}{1 - \beta_1^t}$$

$$s_{dw} := \frac{s_{dw}}{1 - \beta_2^t}$$

$$s_{db} := \frac{s_{db}}{1 - \beta_2^t}$$

$$W := W - \alpha \frac{v_{dw}}{\sqrt{s_{dw} + \epsilon}}, \quad b := b - \alpha \frac{v_{db}}{\sqrt{s_{db} + \epsilon}}$$

The following two functions are designed to perform the parameter update for Adam algorithm.

```
def initialize_adam(parameters) :
```

```
def update_parameters_with_adam(parameters, grads, v, s, t, learning_rate = 0.01,
```

```
    beta1 = 0.9, beta2 = 0.999, epsilon = 1e-8):
```

The algorithm performs the following computations for parameters of all layers:

1. It calculates an exponentially weighted average of past gradients, and stores it in variables v (before bias correction) and $v_corrected$ (with bias correction).
2. It calculates an exponentially weighted average of the squares of the past gradients, and stores it in variables s (before bias correction) and $s_corrected$ (with bias correction).
3. It updates parameters in a direction based on combining information from "1" and "2".
4. It returns parameters, v and s .

where:

- t counts the number of steps taken of Adam, i.e. how many times the parameters have been updated so far. It is not the index of epochs, nor the index of mini-batches. It is equal to the product of epoch number and mini-batch number.
- L is the number of layers
- β_1 and β_2 are hyperparameters that control the two exponentially weighted averages.
- α is the learning rate
- ϵ is a very small number to avoid dividing by zero

9.4 Python Programming

9.4.1 Python files

This section attaches Python codes which demonstrate how the ideas are exactly implemented in the neural network. The first file is `opt_utils.py` for some basic functions. The second file `ch8.py` includes all other required functions, top-level neural network model, and applications to two classification tasks.

`\machine learing\NN_m_class\opt_utils.py`

```
1.  # -*- coding: utf-8 -*-
2.  """
3.  Created on Tue Jun 11 15:57:44 2019
4.  opt_utils.py
5.  @author: weido
6.  """
7.
8.  import numpy as np
9.  import matplotlib.pyplot as plt
10. import h5py
11. import scipy.io
12. import sklearn
13. import sklearn.datasets
14.
15. def sigmoid(x):
16.     """
17.     Compute the sigmoid of x
18.     Arguments:
19.     x -- A scalar or numpy array of any size.
20.     Return:
21.     s -- sigmoid(x)
22.     """
23.     s = 1/(1+np.exp(-x))
24.     return s
25. def softmax(x):
26.     """
27.     Compute the sigmoid of x
28.     Arguments:
29.     x -- A scalar or numpy array of any size.
30.     Return:
31.     s -- softmax(x)
32.     """
33.     t=np.exp(x)
34.     s = t/np.sum(t, axis=0)
35.
36.     return s
37. def relu(x):
38.     """
39.     Compute the relu of x
40.     Arguments:
41.     x -- A scalar or numpy array of any size.
42.     Return:
43.     s -- relu(x)
44.     """
45.     s = np.maximum(0,x)
46.
```

```

47.     return s
48.
49. def load_params_and_grads(seed=1):
50.     np.random.seed(seed)
51.     W1 = np.random.randn(2,3)
52.     b1 = np.random.randn(2,1)
53.     W2 = np.random.randn(3,3)
54.     b2 = np.random.randn(3,1)
55.
56.     dW1 = np.random.randn(2,3)
57.     db1 = np.random.randn(2,1)
58.     dW2 = np.random.randn(3,3)
59.     db2 = np.random.randn(3,1)
60.
61.     return W1, b1, W2, b2, dW1, db1, dW2, db2
62.
63.
64. def initialize_parameters(layer_dims):
65.     """
66.     Arguments:
67.         layer_dims -
68.         - python array (list) containing the dimensions of each layer in our network
69.
70.     Returns:
71.         parameters -
72.         - python dictionary containing your parameters "W1", "b1", ..., "WL", "bL":
73.             W1 -- weight matrix of shape (layer_dims[l], layer_dims[l-
74.             1])
75.             b1 -- bias vector of shape (layer_dims[l], 1)
76.             Wl -- weight matrix of shape (layer_dims[l-
77.             1], layer_dims[l])
78.             bl -- bias vector of shape (1, layer_dims[l])
79.
80.     Tips:
81.         - For example: the layer_dims for the "Planar Data classification model" would have been [2,2,1].
82.         This means W1's shape was (2,2), b1 was (1,2), W2 was (2,1) and b2 was (1,1).
83.         Now you have to generalize it!
84.         - In the for loop, use parameters['W' + str(l)] to access Wl, where l is the
85.         iterative integer.
86.     """
87.     np.random.seed(3)
88.     parameters = {}
89.     L = len(layer_dims) # number of layers in the network
90.
91.     for l in range(1, L):
92.         parameters['W' + str(l)] = np.random.randn(layer_dims[l], layer_dims[l-
93.         1])* np.sqrt(2 / layer_dims[l-1])
94.         parameters['b' + str(l)] = np.zeros((layer_dims[l], 1))
95.
96.         assert(parameters['W' + str(l)].shape == layer_dims[l], layer_dims[l-
97.         1])
98.         assert(parameters['W' + str(l)].shape == layer_dims[l], 1)
99.
100.    return parameters
101.
102.
103. def compute_cost(a3, Y):
104.     """

```



```

99.     Implement the cost function
100.
101.     Arguments:
102.     a3 -- post-activation, output of forward propagation
103.     Y -- "true" labels vector, same shape as a3
104.
105.     Returns:
106.     cost - value of the cost function
107.     """
108.     m = Y.shape[1]
109.     k = Y.shape[0]
110.     if k==1: # sigmoid for 2 classes
111.         logprobs = np.multiply(-np.log(a3),Y) + np.multiply(-
np.log(1 - a3), 1 - Y)
112.         cost = 1./m * np.sum(logprobs)
113.
114.     else: #softmax for multiple classes
115.         logprobs = np.multiply(np.log(a3),Y)
116.         cost = -1/m*np.sum(logprobs)
117.         cost = np.squeeze(cost)
118.     return cost
119.
120. def forward_propagation(X, parameters):
121.     """
122.     Implements the forward propagation (and computes the loss) presented in Figur
e 2.
123.
124.     Arguments:
125.     X -- input dataset, of shape (input size, number of examples)
126.     parameters -
- python dictionary containing your parameters "W1", "b1", "W2", "b2", "W3", "b3"
:
127.         W1 -- weight matrix of shape ()
128.         b1 -- bias vector of shape ()
129.         W2 -- weight matrix of shape ()
130.         b2 -- bias vector of shape ()
131.         W3 -- weight matrix of shape ()
132.         b3 -- bias vector of shape ()
133.
134.     Returns:
135.     loss -- the loss function (vanilla logistic loss)
136.     """
137.
138.     # retrieve parameters
139.     W1 = parameters["W1"]
140.     b1 = parameters["b1"]
141.     W2 = parameters["W2"]
142.     b2 = parameters["b2"]
143.     W3 = parameters["W3"]
144.     b3 = parameters["b3"]
145.
146.     # LINEAR -> RELU -> LINEAR -> RELU -> LINEAR -> SIGMOID
147.     z1 = np.dot(W1, X) + b1
148.     a1 = relu(z1)
149.     z2 = np.dot(W2, a1) + b2
150.     a2 = relu(z2)
151.     z3 = np.dot(W3, a2) + b3
152.     a3 = softmax(z3)
153.
154.     cache = (z1, a1, W1, b1, z2, a2, W2, b2, z3, a3, W3, b3)
155.

```

```

156.     return a3, cache
157.
158. def backward_propagation(X, Y, cache):
159.     """
160.     Implement the backward propagation presented in figure 2.
161.
162.     Arguments:
163.     X -- input dataset, of shape (input size, number of examples)
164.     Y -- true "label" vector (containing 0 if cat, 1 if non-cat)
165.     cache -- cache output from forward_propagation()
166.
167.     Returns:
168.     gradients -
169.     - A dictionary with the gradients with respect to each parameter, activation and
170.     pre-activation variables
171.     """
172.     m = X.shape[1]
173.     (z1, a1, W1, b1, z2, a2, W2, b2, z3, a3, W3, b3) = cache
174.     dz3 = 1./m * (a3 - Y)
175.     dW3 = np.dot(dz3, a2.T)
176.     db3 = np.sum(dz3, axis=1, keepdims = True)
177.     da2 = np.dot(W3.T, dz3)
178.     dz2 = np.multiply(da2, np.int64(a2 > 0))
179.     dW2 = np.dot(dz2, a1.T)
180.     db2 = np.sum(dz2, axis=1, keepdims = True)
181.     da1 = np.dot(W2.T, dz2)
182.     dz1 = np.multiply(da1, np.int64(a1 > 0))
183.     dW1 = np.dot(dz1, X.T)
184.     db1 = np.sum(dz1, axis=1, keepdims = True)
185.     gradients = {"dz3": dz3, "dW3": dW3, "db3": db3,
186.                  "da2": da2, "dz2": dz2, "dW2": dW2, "db2": db2,
187.                  "da1": da1, "dz1": dz1, "dW1": dW1, "db1": db1}
188.
189.     return gradients
190.
191. def predict(X, parameters):
192.     """
193.     This function is used to predict the results of a n-layer neural network.
194.
195.     Arguments:
196.     X -- data set of examples you would like to label
197.     parameters -- parameters of the trained model
198.
199.     Returns:
200.     y_esti -- predictions (labels:0,1,2,3,4,,5) for the given dataset X
201.     """
202.     m = X.shape[1]
203.     y_esti = np.zeros((1,m), dtype = np.int)
204.     # Forward propagation
205.     a3, caches = forward_propagation(X, parameters)
206.
207.     for i in range(0, a3.shape[1]):
208.         y_esti[0,i]=np.argmax(a3[:,i])
209.
210.     return y_esti

```

```

215.
216. def load_2D_dataset():
217.     data = scipy.io.loadmat('datasets/data.mat')
218.     train_X = data['X'].T
219.     train_Y = data['y'].T
220.     test_X = data['Xval'].T
221.     test_Y = data['yval'].T
222.
223.     plt.scatter(train_X[0, :], train_X[1, :], c=train_Y.ravel(), s=40, cmap=plt.c
m.Spectral);
224.
225.     return train_X, train_Y, test_X, test_Y
226.
227. def plot_decision_boundary(model, X, y):
228.     # Set min and max values and give it some padding
229.     x_min, x_max = X[0, :].min() - 1, X[0, :].max() + 1
230.     y_min, y_max = X[1, :].min() - 1, X[1, :].max() + 1
231.     h = 0.01
232.     # Generate a grid of points with distance h between them
233.     xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
234.
235.     # Predict the function value for the whole grid
236.     Z = model(np.c_[xx.ravel(), yy.ravel()])
237.     Z = Z.reshape(xx.shape)
238.     # Plot the contour and training examples
239.     plt.contourf(xx, yy, Z, cmap=plt.cm.Spectral)
240.     plt.ylabel('x2')
241.     plt.xlabel('x1')
242.     plt.scatter(X[0, :], X[1, :], c=y.ravel(), cmap=plt.cm.Spectral)
243.     plt.show()
244.
245. def predict_dec(parameters, X):
246.     """
247.     Used for plotting decision boundary.
248.
249.     Arguments:
250.     parameters -- python dictionary containing your parameters
251.     X -- input data of size (m, K)
252.
253.     Returns
254.     predictions -- vector of predictions of our model (red: 0 / blue: 1)
255.     """
256.     # Predict using forward propagation and a classification threshold of 0.5
257.     a3, cache = forward_propagation(X, parameters)
258.     predictions = (a3 > 0.5)
259.     return predictions
260.
261. def load_dataset():
262.     np.random.seed(3)
263.     train_X, train_Y = sklearn.datasets.make_moons(n_samples=300, noise=.2) #300
#0.2
264.     # Visualize the data
265.     plt.scatter(train_X[:, 0], train_X[:, 1], c=train_Y.ravel(), s=40, cmap=plt.c
m.Spectral);
266.     train_X = train_X.T
267.     train_Y = train_Y.reshape((1, train_Y.shape[0]))
268.
269.     return train_X, train_Y
270.
271. ##### for multiple class classification

```

```

272.
273. def load_dataset_signs():
274.     train_dataset = h5py.File('train_signs.h5', "r")
275.     train_set_x_orig = np.array(train_dataset["train_set_x"][:]) # your train set
        features
276.     train_set_y_orig = np.array(train_dataset["train_set_y"][:]) # your train set
        labels
277.
278.     test_dataset = h5py.File('test_signs.h5', "r")
279.     test_set_x_orig = np.array(test_dataset["test_set_x"][:]) # your test set fea
        tures
280.     test_set_y_orig = np.array(test_dataset["test_set_y"][:]) # your test set lab
        els
281.
282.     classes = np.array(test_dataset["list_classes"][:]) # the list of classes
283.
284.     train_set_y_orig = train_set_y_orig.reshape((1, train_set_y_orig.shape[0]))
285.     test_set_y_orig = test_set_y_orig.reshape((1, test_set_y_orig.shape[0]))
286.
287.     return train_set_x_orig, train_set_y_orig, test_set_x_orig, test_set_y_orig,
        classes
288. def convert_to_one_hot(Y, C):
289.     Y = np.eye(C)[Y.reshape(-1)].T
290.     return Y

```

\\machine_learning\\NN_m_class\\ch9.py

```

1.  # -*- coding: utf-8 -*-
2.  """
3.  Created on Tue Jun 11 16:00:23 2019
4.  dnn_opt.py changed to ch9.py
5.  @author: weido
6.  """
7.
8.  import numpy as np
9.  import matplotlib.pyplot as plt
10. import scipy.io
11. import math
12. import sklearn
13. import sklearn.datasets
14.
15. from opt_utils import load_params_and_grads, initialize_parameters, forward_propa
    gation, backward_propagation
16. from opt_utils import compute_cost, predict, predict_dec, plot_decision_boundary,
    load_dataset, load_dataset_signs
17. from opt_utils import convert_to_one_hot
18.
19.
20. #matplotlib inline
21. plt.rcParams['figure.figsize'] = (7.0, 4.0) # set default size of plots
22. plt.rcParams['image.interpolation'] = 'nearest'
23. plt.rcParams['image.cmap'] = 'gray'
24.
25. def update_parameters_with_gd(parameters, grads, learning_rate):
26.     """
27.     Update parameters using one step of gradient descent
28.
29.     Arguments:

```

```

30.     parameters -- python dictionary containing your parameters to be updated:
31.         parameters['W' + str(l)] = w1
32.         parameters['b' + str(l)] = b1
33.     grads -
34.     - python dictionary containing your gradients to update each parameters:
35.         grads['dw' + str(l)] = dw1
36.         grads['db' + str(l)] = db1
37.     learning_rate -- the learning rate, scalar.
38.
39.     Returns:
40.     parameters -- python dictionary containing your updated parameters
41.     """
42.     L = len(parameters) // 2 # number of layers in the neural networks
43.
44.     # Update rule for each parameter
45.     for l in range(L):
46.         ### START CODE HERE ### (approx. 2 lines)
47.         parameters["W" + str(l+1)] = parameters["W" + str(l+1)] - learning_rate*g
48.         parameters["b" + str(l+1)] = parameters["b" + str(l+1)] - learning_rate*g
49.         ### END CODE HERE ###
50.
51.     return parameters
52.
53.
54.
55. def random_mini_batches(X, Y, mini_batch_size = 64, seed = 0):
56.     """
57.     Creates a list of random minibatches from (X, Y)
58.
59.     Arguments:
60.     X -- input data, of shape (input size, number of examples)
61.     Y -
62.     - true "label" vector (1 for blue dot / 0 for red dot), of shape (1, number of ex
63.     amples)
64.     mini_batch_size -- size of the mini-batches, integer
65.
66.     Returns:
67.     mini_batches -- list of synchronous (mini_batch_X, mini_batch_Y)
68.     """
69.     np.random.seed(seed) # To make your "random" minibatches the same
70.     as ours
71.     m = X.shape[1] # number of training examples
72.     mini_batches = []
73.
74.     # Step 1: Shuffle (X, Y)
75.     permutation = list(np.random.permutation(m))
76.     shuffled_X = X[:, permutation]
77.     shuffled_Y = Y[:, permutation].reshape((Y.shape[0],m))
78.
79.     # Step 2: Partition (shuffled_X, shuffled_Y). Minus the end case.
80.     num_complete_minibatches = math.floor(m/mini_batch_size) # number of mini bat
81.     ches of size mini_batch_size in your partitionning
82.     for k in range(0, num_complete_minibatches):
83.         ### START CODE HERE ### (approx. 2 lines)
84.         mini_batch_X = shuffled_X[:, k*mini_batch_size : (k+1)*mini_batch_size]
85.         mini_batch_Y = shuffled_Y[:, k*mini_batch_size : (k+1)*mini_batch_size]
86.         ### END CODE HERE ###

```

```

84.         mini_batch = (mini_batch_X, mini_batch_Y)
85.         mini_batches.append(mini_batch)
86.
87.         # Handling the end case (last mini-batch < mini_batch_size)
88.         if m % mini_batch_size != 0:
89.             ### START CODE HERE ### (approx. 2 lines)
90.             mini_batch_X = shuffled_X[:, num_complete_minibatches*mini_batch_size:]
91.             mini_batch_Y = shuffled_Y[:, num_complete_minibatches*mini_batch_size:]
92.             ### END CODE HERE ###
93.             mini_batch = (mini_batch_X, mini_batch_Y)
94.             mini_batches.append(mini_batch)
95.
96.         return mini_batches
97.
98.
99. def initialize_velocity(parameters):
100.     """
101.     Initializes the velocity as a python dictionary with:
102.         - keys: "dw1", "db1", ..., "dwl", "dbl"
103.         - values: numpy arrays of zeros of the same shape as the corresponding gradients/parameters.
104.     Arguments:
105.         parameters -- python dictionary containing your parameters.
106.             parameters['W' + str(l)] = Wl
107.             parameters['b' + str(l)] = bl
108.
109.     Returns:
110.         v -- python dictionary containing the current velocity.
111.             v['dw' + str(l)] = velocity of dwl
112.             v['db' + str(l)] = velocity of dbl
113.     """
114.
115.     L = len(parameters) // 2 # number of layers in the neural networks
116.     v = {}
117.
118.     # Initialize velocity
119.     for l in range(L):
120.         ### START CODE HERE ### (approx. 2 lines)
121.         v["dw" + str(l+1)] = np.zeros(parameters['W'+str(l+1)].shape)
122.         v["db" + str(l+1)] = np.zeros(parameters['b'+str(l+1)].shape)
123.         ### END CODE HERE ###
124.
125.     return v
126.
127.
128. def update_parameters_with_momentum(parameters, grads, v, beta, learning_rate):
129.     """
130.     Update parameters using Momentum
131.
132.     Arguments:
133.         parameters -- python dictionary containing your parameters:
134.             parameters['W' + str(l)] = Wl
135.             parameters['b' + str(l)] = bl
136.         grads -- python dictionary containing your gradients for each parameters:
137.             grads['dw' + str(l)] = dwl
138.             grads['db' + str(l)] = dbl
139.         v -- python dictionary containing the current velocity:
140.             v['dw' + str(l)] = ...
141.             v['db' + str(l)] = ...
142.         beta -- the momentum hyperparameter, scalar
143.         learning_rate -- the learning rate, scalar

```

```

144.
145.     Returns:
146.     parameters -- python dictionary containing your updated parameters
147.     v -- python dictionary containing your updated velocities
148.     """
149.
150.     L = len(parameters) // 2 # number of layers in the neural networks
151.
152.     # Momentum update for each parameter
153.     for l in range(L):
154.
155.         ### START CODE HERE ### (approx. 4 lines)
156.         # compute velocities
157.         v["dw" + str(l+1)] = beta * v["dw" + str(l+1)] + (1-
beta)*(grads["dw" + str(l+1)])
158.         v["db" + str(l+1)] = beta * v["db" + str(l+1)] + (1-
beta)*(grads["db" + str(l+1)])
159.         # update parameters
160.         parameters["W" + str(l+1)] = parameters["W" + str(l+1)] - learning_rate *
v["dw" + str(l+1)]
161.         parameters["b" + str(l+1)] = parameters["b" + str(l+1)] - learning_rate *
v["db" + str(l+1)]
162.         ### END CODE HERE ###
163.
164.     return parameters, v
165.
166.
167. def initialize_adam(parameters) :
168.     """
169.     Initializes v and s as two python dictionaries with:
170.         - keys: "dw1", "db1", ..., "dWL", "dbL"
171.         - values: numpy arrays of zeros of the same shape as the correspo
nding gradients/parameters.
172.
173.     Arguments:
174.     parameters -- python dictionary containing your parameters.
175.         parameters["W" + str(1)] = W1
176.         parameters["b" + str(1)] = b1
177.
178.     Returns:
179.     v -
180.     - python dictionary that will contain the exponentially weighted average of the g
radient.
181.         v["dw" + str(1)] = ...
182.         v["db" + str(1)] = ...
183.     s -
184.     - python dictionary that will contain the exponentially weighted average of the s
quared gradient.
185.         s["dw" + str(1)] = ...
186.         s["db" + str(1)] = ...
187.
188.     """
189.
190.     L = len(parameters) // 2 # number of layers in the neural networks
191.
192.     v = {}
193.     s = {}
194.
195.     # Initialize v, s. Input: "parameters". Outputs: "v, s".
196.     for l in range(L):
197.         ### START CODE HERE ### (approx. 4 lines)
198.         v["dw" + str(l+1)] = np.zeros(parameters["W" + str(l+1)].shape)

```

```

196.         v["db" + str(l+1)] = np.zeros(parameters['b' + str(l+1)].shape)
197.         s["dw" + str(l+1)] = np.zeros(parameters['w' + str(l+1)].shape)
198.         s["db" + str(l+1)] = np.zeros(parameters['b' + str(l+1)].shape)
199.     ### END CODE HERE ###
200.
201.     return v, s
202.
203. def update_parameters_with_adam(parameters, grads, v, s, t, learning_rate = 0.01,
204.                                beta1 = 0.9, beta2 = 0.999, epsilon = 1e-8):
205.     """
206.     Update parameters using Adam
207.
208.     Arguments:
209.     parameters -- python dictionary containing your parameters:
210.         parameters['w' + str(l)] = w1
211.         parameters['b' + str(l)] = b1
212.     grads -- python dictionary containing your gradients for each parameters:
213.         grads['dw' + str(l)] = dw1
214.         grads['db' + str(l)] = db1
215.     v -- Adam variable, moving average of the first gradient, python dictionary
216.     s --
217.     - Adam variable, moving average of the squared gradient, python dictionary
218.     learning_rate -- the learning rate, scalar.
219.     beta1 -- Exponential decay hyperparameter for the first moment estimates
220.     beta2 -- Exponential decay hyperparameter for the second moment estimates
221.     epsilon -- hyperparameter preventing division by zero in Adam updates
222.
223.     Returns:
224.     parameters -- python dictionary containing your updated parameters
225.     v -- Adam variable, moving average of the first gradient, python dictionary
226.     s --
227.     - Adam variable, moving average of the squared gradient, python dictionary
228.     """
229.     L = len(parameters) // 2 # number of layers in the neural net
230.     works
231.     v_corrected = {} # Initializing first moment estimate
232.     , python dictionary
233.     s_corrected = {} # Initializing second moment estimate
234.     e, python dictionary
235.
236.     # Perform Adam update on all parameters
237.     for l in range(L):
238.         # Moving average of the gradients. Inputs: "v, grads, beta1". Output: "v"
239.         .
240.         ### START CODE HERE ### (approx. 2 lines)
241.         v["dw" + str(l+1)] = beta1 * v['dw' + str(l+1)] + (1-
242.         beta1) * grads['dw' + str(l+1)]
243.         v["db" + str(l+1)] = beta1 * v['db' + str(l+1)] + (1-
244.         beta1) * grads['db' + str(l+1)]
245.         ### END CODE HERE ###
246.
247.         # Compute bias-
248.         corrected first moment estimate. Inputs: "v, beta1, t". Output: "v_corrected".
249.         ### START CODE HERE ### (approx. 2 lines)
250.         v_corrected["dw" + str(l+1)] = v['dw' + str(l+1)] / (1 - np.power(beta1,
251.         t))
252.         v_corrected["db" + str(l+1)] = v['db' + str(l+1)] / (1 - np.power(beta1,
253.         t))
254.         ### END CODE HERE ###

```



```

245.
246.     # Moving average of the squared gradients. Inputs: "s, grads, beta2". Out
put: "s".
247.     ### START CODE HERE ### (approx. 2 lines)
248.     s["dw" + str(l+1)] = beta2 * s['dw' + str(l+1)] + (1-
beta2) * np.power(grads['dw' + str(l+1)], 2)
249.     s["db" + str(l+1)] = beta2 * s['db' + str(l+1)] + (1-
beta2) * np.power(grads['db' + str(l+1)], 2)
250.     ### END CODE HERE ###
251.
252.     # Compute bias-
corrected second raw moment estimate. Inputs: "s, beta2, t". Output: "s_corrected
".
253.     ### START CODE HERE ### (approx. 2 lines)
254.     s_corrected["dw" + str(l+1)] = s['dw' + str(l+1)] / (1 - np.power(beta2,
t))
255.     s_corrected["db" + str(l+1)] = s['db' + str(l+1)] / (1 - np.power(beta2,
t))
256.     ### END CODE HERE ###
257.
258.     # Update parameters. Inputs: "parameters, learning_rate, v_corrected, s_c
orrected, epsilon". Output: "parameters".
259.     ### START CODE HERE ### (approx. 2 lines)
260.     parameters["W" + str(l+1)] = parameters['W' + str(l+1)] - learning_rate *
v_corrected['dw' + str(l+1)] / np.sqrt(s_corrected['dw' + str(l+1)] + epsilon)
261.     parameters["b" + str(l+1)] = parameters['b' + str(l+1)] - learning_rate *
v_corrected['db' + str(l+1)] / np.sqrt(s_corrected['db' + str(l+1)] + epsilon)
262.     ### END CODE HERE ###
263.
264.     return parameters, v, s
265.
266.
267. def model(X, Y, layers_dims, optimizer, learning_rate = 0.0007, mini_batch_size =
64, beta = 0.9,
268.         beta1 = 0.9, beta2 = 0.999, epsilon = 1e-
8, num_epochs = 10000, print_cost = True):
269.     """
270.     3-layer neural network model which can be run in different optimizer modes.
271.
272.     Arguments:
273.     X -- input data, of shape (2, number of examples)
274.     Y -
- true "label" vector (1 for blue dot / 0 for red dot), of shape (1, number of ex
amples)
275.     layers_dims -- python list, containing the size of each layer
276.     learning_rate -- the learning rate, scalar.
277.     mini_batch_size -- the size of a mini batch
278.     beta -- Momentum hyperparameter
279.     beta1 -- Exponential decay hyperparameter for the past gradients estimates
280.     beta2 -
- Exponential decay hyperparameter for the past squared gradients estimates
281.     epsilon -- hyperparameter preventing division by zero in Adam updates
282.     num_epochs -- number of epochs
283.     print_cost -- True to print the cost every 1000 epochs
284.
285.     Returns:
286.     parameters -- python dictionary containing your updated parameters
287.     """
288.
289.     L = len(layers_dims)          # number of layers in the neural networks
290.     costs = []                  # to keep track of the cost

```

```

291.     t = 0                                # initializing the counter required for Adam
      update
292.     seed = 10                            # For grading purposes, so that your "random
" minibatches are the same as ours
293.
294.     # Initialize parameters
295.     parameters = initialize_parameters(layers_dims)
296.     # Initialize the optimizer
297.     if optimizer == "gd":
298.         pass # no initialization required for gradient descent
299.     elif optimizer == "momentum":
300.         v = initialize_velocity(parameters)
301.     elif optimizer == "adam":
302.         v, s = initialize_adam(parameters)
303.
304.     # Optimization loop
305.     for i in range(num_epochs):
306.
307.         # Define the random minibatches. We increment the seed to reshuffle differ-
      ntly the dataset after each epoch
308.         seed = seed + 1
309.         minibatches = random_mini_batches(X, Y, mini_batch_size, seed)
310.
311.         for minibatch in minibatches:
312.
313.             # Select a minibatch
314.             (minibatch_X, minibatch_Y) = minibatch
315.
316.             # Forward propagation
317.             a3, caches = forward_propagation(minibatch_X, parameters)
318.
319.             # Compute cost
320.             cost = compute_cost(a3, minibatch_Y)
321.         # Backward propagation
322.             grads = backward_propagation(minibatch_X, minibatch_Y, caches)
323.
324.             # Update parameters
325.             if optimizer == "gd":
326.                 parameters = update_parameters_with_gd(parameters, grads, learnin
g_rate)
327.             elif optimizer == "momentum":
328.                 parameters, v = update_parameters_with_momentum(parameters, grads
, v, beta, learning_rate)
329.             elif optimizer == "adam":
330.                 t = t + 1 # Adam counter
331.                 parameters, v, s = update_parameters_with_adam(parameters, grads,
v, s,
332.                                                             t, learning_rate,
beta1, beta2, epsilon)
333.
334.         # Print the cost every 1000 epoch
335.         if print_cost and i % 1000 == 0:
336.             #print(t)
337.             print ("Cost after epoch %i: %f" %(i, cost))
338.         # if print_cost and i % 10 == 0:
339.         if print_cost and i % 10 == 0:
340.             costs.append(cost)
341.
342.     # plot the cost
343.     plt.plot(costs)
344.     plt.ylabel('cost')

```

```

345.     plt.xlabel('epochs (per 10)')
346.     plt.title("Learning rate = " + str(learning_rate))
347.     plt.show()
348.
349.     return parameters, costs
350.
351.
352.
353. # train 3-layer model
354. train_X, train_Y = load_dataset()
355.
356. plt.scatter(train_X[0, :], train_X[1, :], c=train_Y[0, :], s=40, cmap=plt.cm.Spectral)
357. plt.xlabel("X[0,:]")
358. plt.ylabel("X[1,:]")
359. plt.title("data visualization, red: y=0, blue: y=1")
360. plt.show()
361.
362. layers_dims = [train_X.shape[0], 5, 2, 2]
363. train_Y_hot = convert_to_one_hot(train_Y, 2)
364. # model(X, Y, layers_dims, optimizer, learning_rate = 0.0007, mini_batch_size
    = 64, beta = 0.9,
365. #         beta1 = 0.9, beta2 = 0.999, epsilon = 1e-
    8, num_epochs = 10000, print_cost = True)
366. #
367.
368. parameters, costs=model(train_X, train_Y_hot, layers_dims, optimizer="momentum",
    learning_rate = 0.1, mini_batch_size = 64, beta = 0.9,
369.         beta1 = 0.9, beta2 = 0.999, epsilon = 1e-
    8, num_epochs = 1000, print_cost = True)
370. # Predict
371. predictions = predict(train_X, parameters)
372. accuracy = np.mean(predictions == train_Y)
373.
374. print("accuracy of train set is "+ str(accuracy))
375. # Predict
376. #predictions = predict(train_X, train_Y, parameters)
377. # Plot decision boundary
378. plt.title("Model with momentum optimization")
379. axes = plt.gca()
380. axes.set_xlim([-1.5,2.5])
381. axes.set_ylim([-1,1.5])
382. plot_decision_boundary(lambda x: predict(x.T, parameters), train_X, train_Y)
383.
384.
385. # multiple class classification
386.
387. np.random.seed(1)
388. # Loading the dataset
389. X_train_orig, Y_train_orig, X_test_orig, Y_test_orig, classes = load_dataset_sign
    s()
390. # Example of a picture
391. index = 0
392. plt.imshow(X_train_orig[index])
393. plt.show()
394. print ("y = " + str(np.squeeze(Y_train_orig[:, index])))
395. # Flatten the training and test images
396. X_train_flatten = X_train_orig.reshape(X_train_orig.shape[0], -1).T
397. X_test_flatten = X_test_orig.reshape(X_test_orig.shape[0], -1).T
398. # Normalize image vectors
399. X_train = X_train_flatten/255.

```

```

400. X_test = X_test_flatten/255.
401. # Convert training and test labels to one hot matrices
402. Y_train = convert_to_one_hot(Y_train_orig, 6)
403. Y_test = convert_to_one_hot(Y_test_orig, 6)
404.
405. print ("number of training examples = " + str(X_train.shape[1]))
406. print ("number of test examples = " + str(X_test.shape[1]))
407. print ("X_train shape: " + str(X_train.shape))
408. print ("Y_train shape: " + str(Y_train.shape))
409. print ("X_test shape: " + str(X_test.shape))
410. print ("Y_test shape: " + str(Y_test.shape))
411.
412. # train 3-layer model
413. layers_dims = [X_train.shape[0], 25, 12, 6]
414. #parameters, costs = model(X_train, Y_train, layers_dims, optimizer="adam", learning_rate = 0.0001, mini_batch_size = 32, num_epochs = 1500)
415. parameters, costs=model(X_train, Y_train, layers_dims, optimizer="adam", learning_rate = 0.0001, mini_batch_size = 32, beta = 0.9,
416.                          beta1 = 0.9, beta2 = 0.999, epsilon = 1e-8, num_epochs = 1500, print_cost = True)
417. # Predict
418. predictions = predict(X_train, parameters)
419. accuracy = np.mean(predictions == Y_train_orig)
420.
421. print("accuracy of train set is "+ str(accuracy))
422.
423. predictions = predict(X_test, parameters)
424. accuracy = np.mean(predictions == Y_test_orig)
425.
426. print("accuracy of test set is "+ str(accuracy))

```

9.4.2 Binary classification for two-feature inputs

We apply the model to classify the data examples shown in Fig.6. To do this, we choose to assign `layers_dims = [train_X.shape[0], 5, 2, 2]`, and then instantiate the model by

```

parameters, costs=model(train_X, train_Y_hot, layers_dims, optimizer="gd", learning_rate
                        = 0.1, mini_batch_size = 64, beta = 0.9, beta1 = 0.9, beta2 = 0.999, epsilon =
                        1e-8, num_epochs = 1000, print_cost = True)

```

Of course, we can tune the hyperparameters by changing their settings.

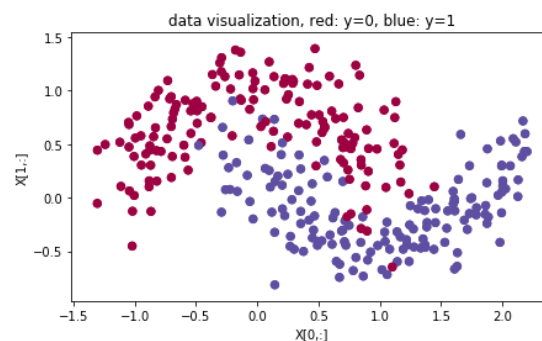


Fig.6 training examples generated by `load_dataset()`

Fig. 7 shows the results using three different optimization methods, based on the above hyperparameter setting. With the current setting, three methods have similar performance. However, a meaningful comparison requires various combinations of hyperparameter values.

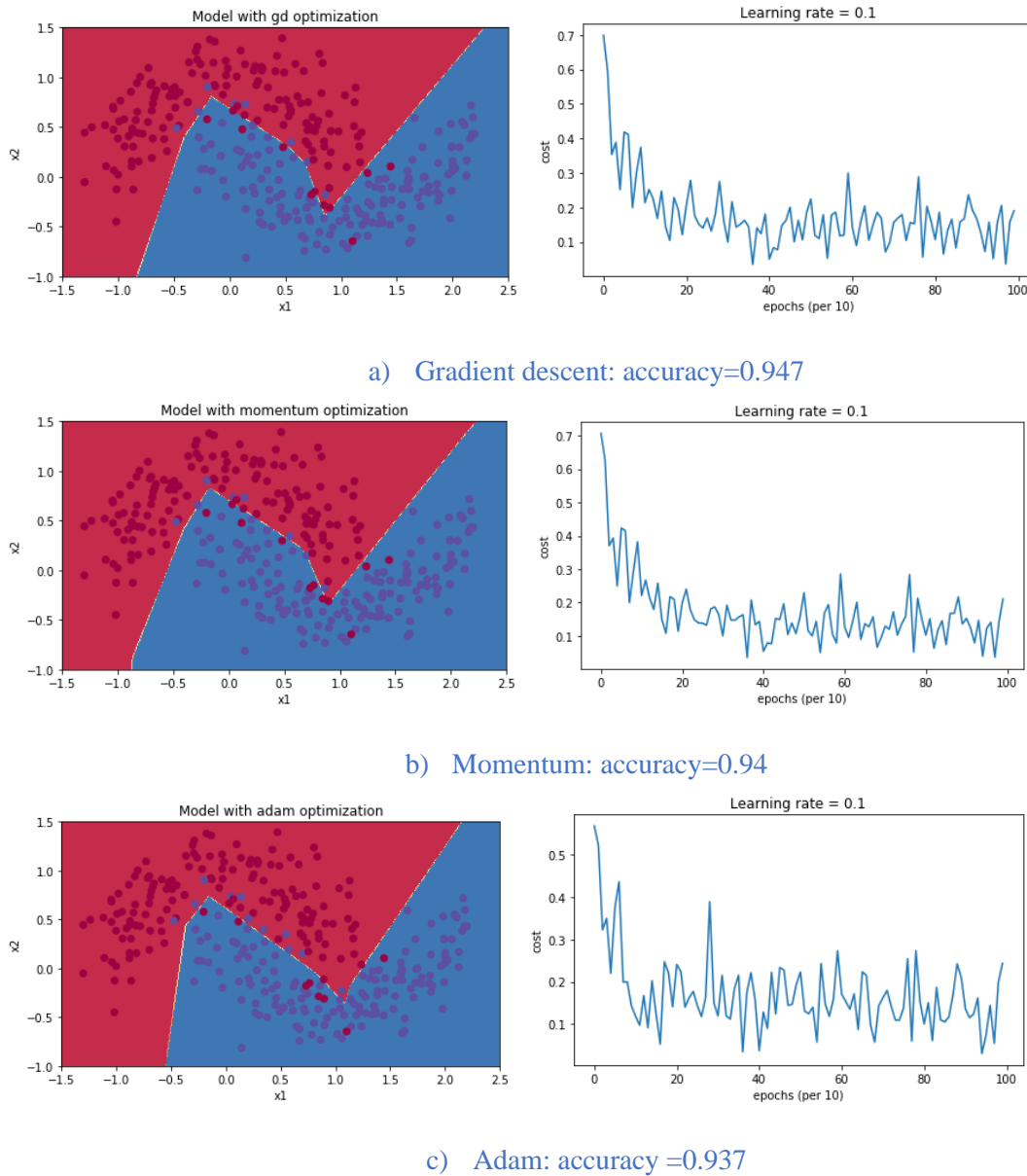


Fig.7 Results of different optimization methods

9.4.3 6-class classification for 12288-feature inputs: finger-signs

Suppose that we use finger-signs to send 6 digits for a communication purpose. Some example signs are shown in Fig.8. Each sign is represented by (64x64)-pixel color picture with a certain combination of finger positions. We will build a neural network to translate a sign to a digit. Thus, this model classifies six classes. The dataset files are train_signs.h5 (1080 images for training) and test_signs.h5 (120 images for testing).

To do this, we instantiate the model by

```
layers_dims = [X_train.shape[0], 25, 12, 6]
parameters, costs=model(X_train, Y_train, layers_dims, optimizer="adam",
learning_rate = 0.0001, mini_batch_size = 32, beta = 0.9, beta1 = 0.9, beta2 =
0.999, epsilon = 1e-8, num_epochs = 1500, print_cost = True)
```

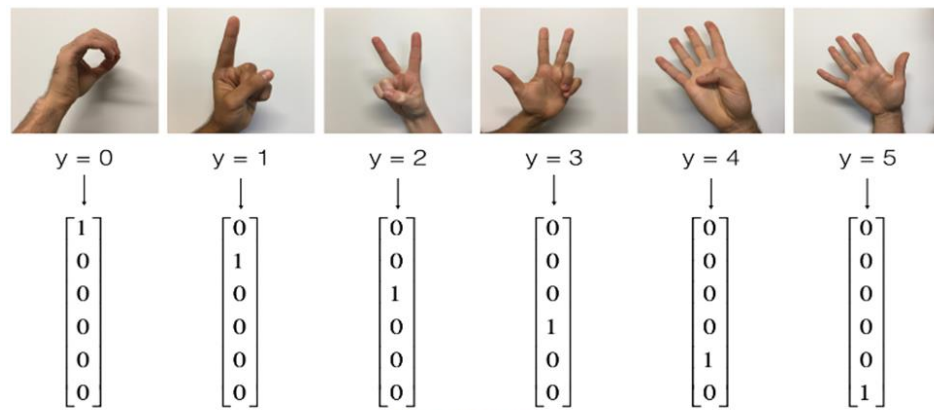


Figure 1: SIGNS dataset

Fig.8 sign-digit mapping

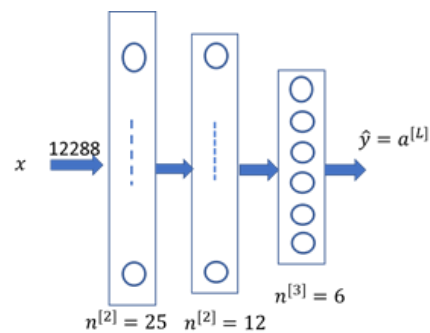
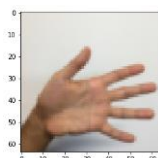


Fig.9 Architecture of neural network for finger sign recognition

By running ch9.py, we got the following results:

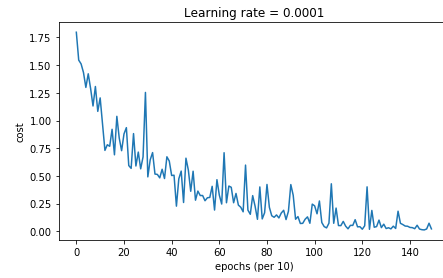


$y = 5$

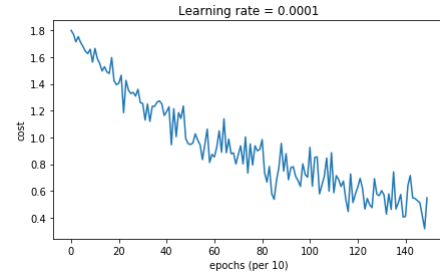
```
number of training examples = 1080
number of test examples = 120
X_train shape: (12288, 1080)
Y_train shape: (6, 1080)
X_test shape: (12288, 120)
Y_test shape: (6, 120)
```

Here we change optimizer and α . The results for some combinations are shown below Fig.10:

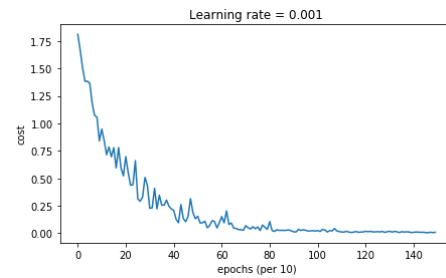
Optimizer="adam", $\alpha=0.0001$, mini-batch=32
 accuracy of train set is 0.9953703703703703
 accuracy of test set is 0.7833333333333333



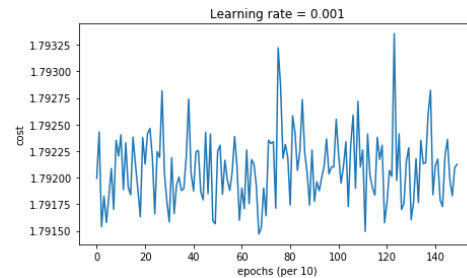
Optimizer="gd", $\alpha=0.0001$, mini-batch=32
 accuracy of train set is 0.887962962962963
 accuracy of test set is 0.8083333333333333



Optimizer="gd", $\alpha=0.001$, mini-batch=32
 accuracy of train set is 1.0
 accuracy of test set is 0.8916666666666667



Optimizer="adam", $\alpha=0.001$, mini-batch=32
 accuracy of train set is 0.16666666666666666
 accuracy of test set is 0.16666666666666666



Note: This setting does not work!

Optimizer = "adam", $\alpha=0.001$, mini-batch=1080 (batch)
 accuracy of train set is 0.9861111111111112
 accuracy of test set is 0.8416666666666667

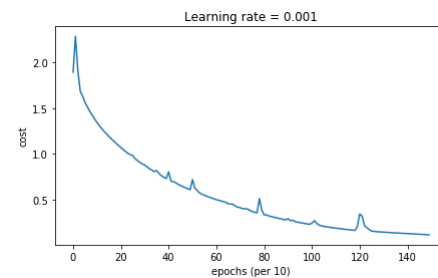


Fig.10 Results for some different settings with fixed arguments: $\beta = 0.9$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 1e-8$, $\text{num_epochs} = 1500$.

Based on the above trials, we realize that Adam may not work if other parameters are not correctly set. Therefore, it is essential to monitor the cost function to make sure the optimization method behaves correctly. In general, the effectiveness of an optimization method may be very sensitive to some hyperparameters. Thus, it is usually necessary to explore the design space of hyperparameters to search for a reasonably satisfied performance.

Summary

To digest the previous chapter (chapter 8), this chapter gives a comprehensive tutorial on mini-batch and Adam optimization through a project. The results of the project show that it is important to plot cost function curve to monitor the behavior of optimizer, and that selection of hyperparameters is essential.

Files:

```
\machine_learning\NN_m_class\opt_utils.py
\machine_learning\NN_m_class\ch9.py
\machine_learning\NN_m_class\train_signs.h5
\machine_learning\NN_m_class\test_signs.h5
```

References

- [1] “Neural networks and deep learning” online course at www.coursera.org
- [2]

Exercises

1. Based on the results shown in Fig.10, the best result has been achieved when `Optimizer="gd"`, $\alpha=0.001$, `mini-batch=32`, (fixed: $\beta = 0.9$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 1e-8$, `num_epochs = 1500`). The resulting training accuracy is 1 and testing accuracy is 0.89. Try different augment assignments to see whether you can find a case works better than this one.
2. Based on the results shown in Fig.10, for most working cases, the training accuracy is higher than testing accuracy by 10%. For example, the case with: `Optimizer="gd"`, $\alpha=0.001$, `mini-batch=32`, (fixed: $\beta = 0.9$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 1e-8$, `num_epochs = 1500`), achieves a training accuracy of 1, and testing accuracy of 0.89. An overfitting probably exists. Please modify `opt_utils.py` and/or `ch9.py` to introduce weight regularization to see if you can improve the testing accuracy.
3. Try different neural networks with different sizes by change `layers_dims` value. Note that the neural network should be still 3-layer.

layers_dims=[12288, n_1, n_2, 6]

The results in Fig.10 are based on *layers_dims=[12288, 25, 12, 6]*. Now you change numbers “25” and “12” to create a different neural network, apply the new neural network to the finger sign dataset with similar setting in Fig.10. Compare your results with Fig.10.