

1003 HW2

Long Chen
lc3424@nyu.edu

February 21st, 2021

Q1

```
1 def feature_normalization(train, test):
2     # Remove constant features, using statistics of training set
3     train = train[:, ~np.all(train[1:] == train[:,-1], axis=0)]
4     test = test[:, ~np.all(test[1:] == test[:,-1], axis=0)]
5
6     # Normalization with min-max
7     max_train = train.max(axis=0)[None, :]
8     min_train = train.min(axis=0)[None, :]
9
10    train_normalized = (train - min_train) / (max_train - min_train)
11    test_normalized = (test - min_train) / (max_train - min_train)
12
13    return train_normalized, test_normalized
```

Q2

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(\mathbf{x}_i) - y_i)^2 \quad (1)$$

$$= \frac{1}{m} \|X\theta - y\|_2^2 \quad (2)$$

Q3

$$J(\theta) = \frac{1}{m} (X\theta - y)^T (X\theta - y) \quad (3)$$

$$= \frac{1}{m} (\theta^T X^T - y^T) (X\theta - y) \quad (4)$$

$$= \frac{1}{m} (\theta^T X^T X\theta - y^T X\theta - \theta^T X^T y + y^T y) \quad (5)$$

$$\nabla J(\theta) = \frac{1}{m} (2X^T X\theta - 2X^T y) \quad (6)$$

$$= \frac{2}{m} X^T (X\theta - y) \quad (7)$$

Q4

$$\theta_{i+1} \leftarrow \theta_i - \eta \nabla J(\theta) \quad (8)$$

Q5

```
1 def compute_square_loss(X, y, theta):
2     m = y.shape[0]
3     y = y.reshape((y.shape[0], 1))
4
5     return (1 / m * (X @ theta - y).T @ (X @ theta - y))[0, 0]
```

Q6

```
1 def compute_square_loss_gradient(X, y, theta):
2     m = y.shape[0]
3     y = y.reshape((-1, 1))
4
5     return ((2 / m) * X.T @ (X @ theta - y)).reshape(-1)}
```

Q7

Gradient checker

```
1 def grad_checker(X, y, theta, epsilon=0.01, tolerance=1e-4):
2     true_gradient = compute_square_loss_gradient(X, y, theta) #
3     The true gradient
4     num_features = theta.shape[0]
5     approx_grad = np.zeros(num_features) # Initialize the gradient
6     we approximate
7     for i in range(num_features):
8         h = np.zeros((num_features, 1))
9         h[i] = 1
10        approx_grad[i] = (compute_square_loss(X, y, theta + epsilon
11        * h)
12                        - compute_square_loss(X, y, theta -
13        epsilon * h)) / (2 * epsilon)
14
15    return np.linalg.norm(approx_grad - true_gradient) <= tolerance
```

Generic gradient checker

```
1 def generic_gradient_checker(X, y, theta, objective_func,
2     gradient_func,
3     epsilon=0.01, tolerance=1e-4):
4     true_gradient = gradient_func(X, y, theta)
5     num_features = theta.shape[0]
```

```

5 approx_grad = np.zeros(num_features)
6 for i in range(num_features):
7     h = np.zeros((num_features, 1))
8     h[i] = 1
9     approx_grad[i] = (objective_func(X, y, theta + epsilon * h)
10                      - objective_func(X, y, theta - epsilon *
11                      h)) / (2 * epsilon)
12 return np.linalg.norm(approx_grad - true_gradient) <= tolerance

```

Q8

```

1 def batch_grad_descent(X, y, alpha=0.1, num_step=1000, grad_check=
  False):
2     num_instances, num_features = X.shape[0], X.shape[1]
3     theta_hist = np.zeros((num_step + 1, num_features)) #
4     # Initialize theta_hist
5     loss_hist = np.zeros(num_step + 1) # Initialize loss_hist
6     theta = np.zeros(num_features) # Initialize theta
7     theta_hist[0] = theta
8     loss_hist[0] = compute_square_loss(X, y, theta.reshape((-1, 1))
9     )
10
11     for i in range(num_step):
12         grad = compute_square_loss_gradient(X, y, theta.reshape
13         ((-1, 1)))
14
15         # Do not update theta if grad_check fails
16         if grad_check and not grad_checker(X, y, theta.reshape((-1,
17         1)))):
18             theta_hist[i + 1, :] = theta
19             loss_hist[i + 1] = compute_square_loss(X, y, theta.
20             reshape((-1, 1)))
21             continue
22
23         theta = theta - alpha * grad.reshape((1, -1))
24         theta_hist[i + 1, :] = theta
25         loss_hist[i + 1] = compute_square_loss(X, y, theta.reshape
26         ((-1, 1)))
27
28     return theta_hist, loss_hist

```

Q9

```

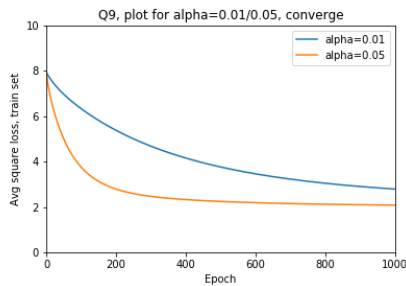
1 plt.ylim([0, 10])
2 plt.xlim([0, 1000])
3 plt.ylabel('Avg square loss, train set')
4 plt.xlabel('Epoch')
5
6 for alpha in [.01, .05]:
7     _, l = batch_grad_descent(X_train, y_train, alpha=alpha,
8     grad_check=False)

```

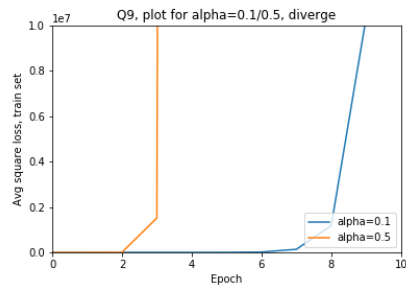
```

8     plt.plot([i for i in range(1.shape[0])], 1, label='alpha={}'.format(alpha))
9
10    plt.legend(loc="upper right")
11    plt.title('Q9, plot for alpha=0.01/0.05, converge')
12    plt.savefig('fig/Q9_small_alpha.png')
13    plt.show()
14
15    plt.ylim([0, 10000000])
16    # plt.yscale('log')
17    plt.xlim([0, 10])
18    plt.ylabel('Avg square loss, train set')
19    plt.xlabel('Epoch')
20
21    for alpha in [.1, .5]:
22        _, l = batch_grad_descent(X_train, y_train, alpha=alpha,
23                                  grad_check=False)
24        plt.plot([i for i in range(1.shape[0])], 1, label='alpha={}'.format(alpha))
25
26    plt.legend(loc="lower right")
27    plt.title('Q9, plot for alpha=0.1/0.5, diverge')
28    plt.savefig('fig/Q9_big_alpha.png')
29    plt.show()

```



(a) Average training loss, small alpha



(b) Average training loss, big alpha

Training with $\alpha = 0.01/0.05$ converges, with $\alpha = 0.05$ converges faster and better final performance (lower average loss) after epoch=1000. Training with $\alpha = 0.1/0.5$ fail to converge.

Q10

```

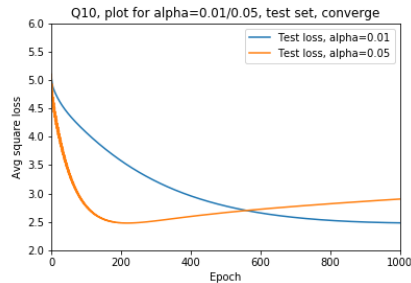
1 plt.ylim([2, 6])
2 plt.xlim([0, 1000])
3 plt.ylabel('Avg square loss')
4 plt.xlabel('Epoch')
5
6 for alpha in [.01, .05]:
7     theta_list, l = batch_grad_descent(X_train, y_train, alpha=
8     alpha, grad_check=False)

```

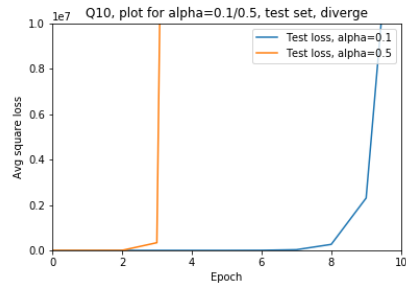
```

8
9     testing_loss = []
10    for i in range(theta_list.shape[0]):
11        testing_loss.append(compute_square_loss(X_test, y_test,
12        theta_list[i].reshape(-1, 1)))
13
14    plt.plot([i for i in range(len(testing_loss))], testing_loss,
15    label='Test loss, alpha={}'.format(alpha))
16
17    plt.legend(loc="upper right")
18    plt.title('Q10, plot for alpha=0.01/0.05, test set, converge')
19    plt.savefig('fig/Q10_small_alpha.png')
20    plt.show()
21
22    plt.ylim([0, 100000000])
23    plt.xlim([0, 10])
24    plt.ylabel('Avg square loss')
25    plt.xlabel('Epoch')
26
27    for alpha in [.1, .5]:
28        theta_list, _ = batch_grad_descent(X_train, y_train, alpha=
29        alpha, grad_check=False)
30
31        testing_loss = []
32        for i in range(theta_list.shape[0]):
33            testing_loss.append(compute_square_loss(X_test, y_test,
34            theta_list[i, :]))
35
36        plt.plot([i for i in range(len(testing_loss))], testing_loss,
37        label='Test loss, alpha={}'.format(alpha))
38
39        plt.legend(loc="upper right")
40        plt.title('Q10, plot for alpha=0.1/0.5, test set, diverge')
41        plt.savefig('fig/Q10_big_alpha.png')
42        plt.show()

```



(a) Average testing loss, small alpha



(b) Average testing loss, big alpha

Training with $\alpha = 0.05$ results in overfitting after around $epoch = 200$. No overfitting found with $\alpha = 0.01$. Still, with $\alpha = 0.1/0.5$, training fail to converge.

Q11

$$J_{\lambda}(\theta) = \frac{1}{m} \|X\theta - y\|_2^2 + \lambda \theta^T \theta \quad (9)$$

$$\nabla J_{\lambda}(\theta) = \frac{2}{m} X^T (X\theta - y) + 2\lambda \theta \quad (10)$$

Q12

```
1 def compute_regularized_square_loss_gradient(X, y, theta,
2       lambda_reg):
3     m = y.shape[0]
4     y = y.reshape((-1, 1))
5
6     return ((2 / m) * X.T @ (X @ theta - y) + 2 * lambda_reg *
7            theta).reshape(-1)
```

Q13

```
1 def regularized_grad_descent(X, y, alpha=0.05, lambda_reg=10 ** -2,
2       num_step=1000):
3     num_instances, num_features = X.shape[0], X.shape[1]
4     theta = np.zeros(num_features) # Initialize theta
5     theta_hist = np.zeros((num_step + 1, num_features)) #
6     Initialize theta_hist
7     loss_hist = np.zeros(num_step + 1) # Initialize loss_hist
8     theta_hist[0, :] = theta
9     loss_hist[0] = compute_square_loss(X, y, theta.reshape((-1, 1))
10    )
11
12    for i in range(num_step):
13        grad = compute_regularized_square_loss_gradient(X, y, theta
14        .reshape((-1, 1)), lambda_reg)
15        theta = theta - alpha * grad.reshape((1, -1))
16        theta_hist[i + 1, :] = theta
17        loss_hist[i + 1] = compute_square_loss(X, y, theta.reshape
18        ((-1, 1)))
19
20    return theta_hist, loss_hist
```

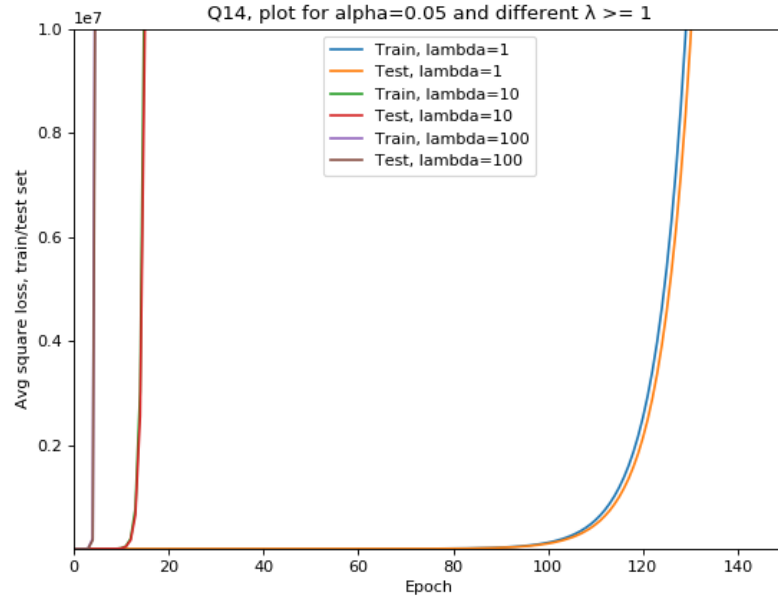
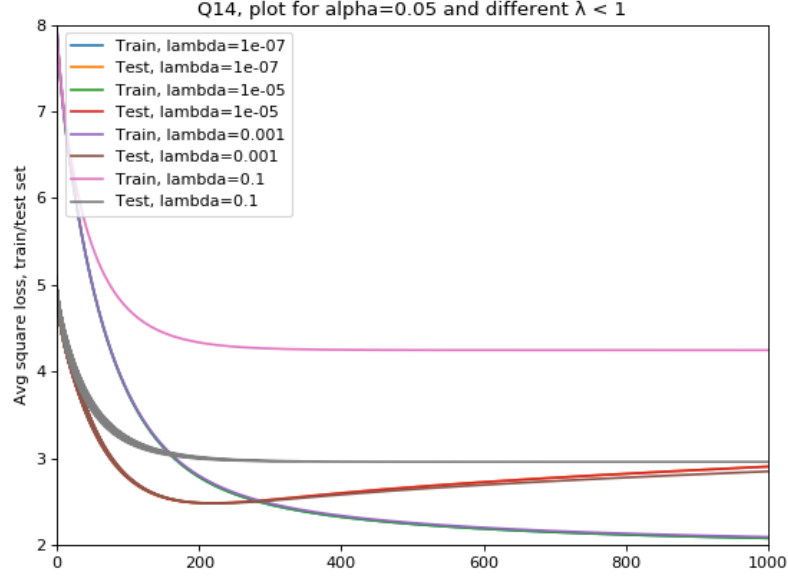
Q14

```
1 from matplotlib.pyplot import figure
2
3 # Lambda < 1
4
5 figure(num=None, figsize=(8, 6), dpi=80, facecolor='w', edgecolor='
6     k')
```

```

6 plt.ylim([2, 8])
7 plt.xlim([0, 1000])
8 plt.ylabel('Avg square loss, train/test set')
9 plt.xlabel('Epoch')
10
11 alpha = 0.05
12 for lambda_reg in [10**-7, 10**-5, 10**-3, 10**-1]:
13     theta_list, training_loss = regularized_grad_descent(X_train,
14     y_train, alpha=alpha, lambda_reg=lambda_reg)
15     plt.plot([i for i in range(1.shape[0])], training_loss, label='
16     Train, lambda={}'.format(lambda_reg))
17
18     testing_loss = []
19     for i in range(theta_list.shape[0]):
20         testing_loss.append(compute_square_loss(X_test, y_test,
21         theta_list[i].reshape((-1, 1))))
22     plt.plot([i for i in range(1.shape[0])], testing_loss, label='
23     Test, lambda={}'.format(lambda_reg))
24
25 plt.legend(loc="upper left")
26 plt.title('Q14, plot for alpha=0.05 and different lambda < 1')
27 plt.savefig('fig/Q14_alpha=0.05_small_lambda.png')
28 plt.show()
29
30 figure(num=None, figsize=(8, 6), dpi=80, facecolor='w', edgecolor='
31     k')
32 plt.ylim([4, 10000000])
33 plt.xlim([0, 150])
34 plt.ylabel('Avg square loss, train/test set')
35 plt.xlabel('Epoch')
36
37 alpha = 0.05
38 for lambda_reg in [1, 10, 100]:
39     theta_list, training_loss = regularized_grad_descent(X_train,
40     y_train, alpha=alpha, lambda_reg=lambda_reg)
41     plt.plot([i for i in range(1.shape[0])], training_loss, label='
42     Train, lambda={}'.format(lambda_reg))
43
44     testing_loss = []
45     for i in range(theta_list.shape[0]):
46         testing_loss.append(compute_square_loss(X_test, y_test,
47         theta_list[i].reshape((-1, 1))))
48     plt.plot([i for i in range(1.shape[0])], testing_loss, label='
49     Test, lambda={}'.format(lambda_reg))
50
51 plt.legend(loc="upper center")
52 plt.title('Q14, plot for alpha=0.05 and different lambda >= 1')
53 plt.savefig('fig/Q14_alpha=0.05_big_lambda.png')
54 plt.show()

```



Choosing $\alpha = 0.05$. With small λ (e.g. $\lambda = 10^{-5}/10^{-7}$), overfitting still presents¹. When $\lambda \geq 1$, training does not converge.

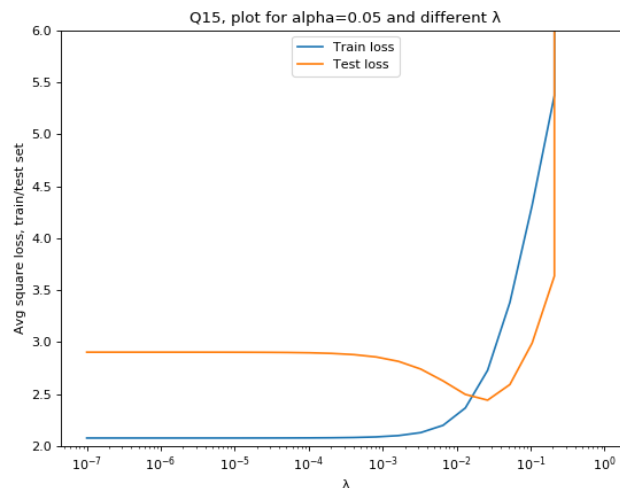
¹For the two small λ , the curve is overlapping on the figure.

Q15

```

1 figure(num=None, figsize=(8, 6), dpi=80, facecolor='w', edgecolor='
  k')
2 plt.ylim([2, 6])
3 plt.xscale('log')
4 plt.ylabel('Avg square loss, train/test set')
5 plt.xlabel('Lambda')
6
7 alpha = 0.05
8 lambda_list = [10**-7]
9 while(lambda_list[-1] < 0.5):
10     lambda_list.append(lambda_list[-1] * 2)
11 lambda_list.append(1)
12
13 plot_list_train = []
14 plot_list_test = []
15
16 for lambda_reg in lambda_list:
17     # Grad descent, training loss
18     theta_list, training_loss = regularized_grad_descent(X_train,
19 y_train, alpha=alpha, lambda_reg=lambda_reg)
20     plot_list_train.append(training_loss[-1])
21     plot_list_test.append(compute_square_loss(X_test, y_test,
22 theta_list[-1].reshape((-1, 1))))
23
24 plt.plot(lambda_list, plot_list_train, label='Train loss')
25 plt.plot(lambda_list, plot_list_test, label='Test loss')
26
27 plt.legend(loc="upper center")
28 plt.title('Q15, plot for alpha=0.05 and different lambda')
29 plt.savefig('fig/Q15_alpha=0.05.png')
30 plt.show()

```



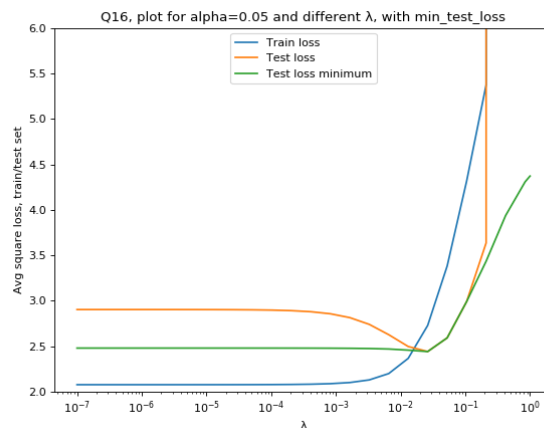
Choose $\lambda \approx 5^{-1}$, the λ of lowest average square loss with test set.

Q16

```

1 figure(num=None, figsize=(8, 6), dpi=80, facecolor='w', edgecolor='
  k')
2 plt.ylim([2, 6])
3 plt.xscale('log')
4 plt.ylabel('Avg square loss, train/test set')
5 plt.xlabel('Lambda')
6
7 alpha = 0.05
8 lambda_list = [10 ** -7]
9 while (lambda_list[-1] < 0.5):
10     lambda_list.append(lambda_list[-1] * 2)
11 lambda_list.append(1)
12 plot_list_train = []
13 plot_list_test = []
14 plot_list_test_min = []
15
16 for lambda_reg in lambda_list:
17     theta_list, training_loss = regularized_grad_descent(X_train,
18 y_train, alpha=alpha, lambda_reg=lambda_reg)
19     plot_list_train.append(training_loss[-1])
20     plot_list_test.append(compute_square_loss(X_test, y_test,
21 theta_list[-1].reshape((-1, 1))))
22     test_min_cur = float('inf')
23     for theta in theta_list:
24         test_min_cur = min(test_min_cur, compute_square_loss(X_test
25 , y_test, theta.reshape((-1, 1))))
26     plot_list_test_min.append(test_min_cur)
27
28 plt.plot(lambda_list, plot_list_train, label='Train loss')
29 plt.plot(lambda_list, plot_list_test, label='Test loss')
30 plt.plot(lambda_list, plot_list_test_min, label='Test loss min')
31 plt.legend(loc="upper center")
32 plt.title('Q16, plot for alpha=0.05 and different lambda, with
33 min_test_loss')
34 plt.savefig('fig/Q16_alpha=0.05.png')
35 plt.show()

```



Still choosing the same lambda as in Q15.

Q17

In practice, we don't want to early stop as soon as an increase in testing loss is observed because we might have fluctuating testing set performances. One way to consider is to set a threshold and compare current epoch testing loss with the minimum testing loss previously achieved. Should the threshold being breached, we will stop the training process and return the optimal θ .

Q18*

$$J_\lambda(\theta) = \frac{1}{m} \sum_{i=1}^m (\theta^T x_i - y_i)^2 + \lambda \theta^T \theta \quad (11)$$

$$= \frac{1}{m} \sum_{i=1}^m [(\theta^T x_i - y_i)^2 + \lambda \theta^T \theta] \quad (12)$$

Therefore,

$$f_i(\theta) = (\theta^T x_i - y_i)^2 + \lambda \theta^T \theta \quad (13)$$

Q19*

$$\nabla J_\lambda(\theta) = \frac{2}{m} \sum_{i=1}^m [x_i(x_i^T \theta - y_i)] + 2\lambda \theta \quad (14)$$

$$\nabla_\theta f_i(\theta) = 2x_i(\theta^T x_i - y_i) + 2\lambda \theta \quad (15)$$

$$E[\nabla_\theta f_i(\theta)] = \frac{1}{m} \sum_{i=1}^m (2x_i(\theta^T x_i - y_i) + 2\lambda \theta) \quad (16)$$

$$\begin{aligned} &= \frac{2}{m} \sum_{i=1}^m (x_i(\theta^T x_i - y_i) + \lambda \theta) \\ &= \nabla J_\lambda(\theta) \end{aligned} \quad (17)$$

Q20*

$$\theta_{i+1} \leftarrow \theta_i - \eta \nabla_\theta f_i(\theta) \quad (18)$$

Q21*

Q23

$$L(\theta) = \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-y_i h_{\theta,b}(x_i)}) \quad (19)$$

If $y_i = 1$,

$$(11) = \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-h_{\theta,b}(x_i)}) \quad (20)$$

$$= \frac{1}{2m} \sum_{i=1}^m (1 + y_i) \log(1 + e^{-y_i h_{\theta,b}(x_i)}) \quad (21)$$

If $y_i = -1$,

$$(11) = \frac{1}{m} \sum_{i=1}^m \log(1 + e^{h_{\theta,b}(x_i)}) \quad (22)$$

$$= \frac{1}{2m} \sum_{i=1}^m (1 - y_i) \log(1 + e^{y_i h_{\theta,b}(x_i)}) \quad (23)$$

When $y_i = 1$, eq.(23) = 0. When $y_i = -1$, eq.(21) = 0. Since $y_i \in \{-1, 1\}$:

$$L(\theta) = \frac{1}{2m} \sum_{i=1}^m (1 + y_i) \log(1 + e^{-h_{\theta,b}(x_i)}) + (1 - y) \log(1 + e^{h_{\theta,b}(x_i)}) \quad (24)$$

Q24

$$L(\theta) = \frac{1}{2m} \sum_{i=1}^m [(1 + y_i) \log(1 + e^{-h_{\theta,b}(x_i)}) + (1 - y) \log(1 + e^{h_{\theta,b}(x_i)})] + \alpha \sum_{i=1}^{784} |\theta_i| \quad (25)$$

Q25

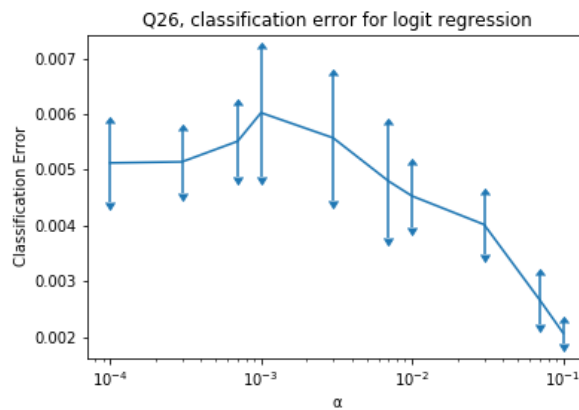
```
1 def classification_error(clf, X, y):
2     return np.sum(clf.predict(X) != y) / y.shape[0]
```

Q26

```

1 from collections import defaultdict
2
3 X_train, y_train = sub_sample(100, X_train, y_train)
4 q26_res = defaultdict(list)
5
6 for alpha in [0.0001, 0.0003, 0.0007, 0.001, 0.003, 0.007, 0.01,
7               0.03, 0.07, 0.1]:
8     for _ in range(10):
9         clf = SGDClassifier(loss='log', max_iter=1000,
10                             tol=1e-3,
11                             penalty='l1', alpha=alpha,
12                             learning_rate='invscaling',
13                             power_t=0.5,
14                             eta0=0.01,
15                             verbose=1)
16         clf.fit(X_train, y_train)
17         q26_res[alpha].append(classification_error(clf, X_test,
18                                                    y_test))
19
20 q26_plot = [(alpha, np.mean(val), np.std(val)) for alpha, val in
21             q26_res.items()]
22
23 plt.xscale('log')
24 plt.ylabel('Classification Error')
25 plt.xlabel('alpha')
26 plt.errorbar([x[0] for x in q26_plot], [x[1] for x in q26_plot],
27              yerr=[x[2] for x in q26_plot], uplims = True, lolims = True,)
28 plt.title('Q26, classification error for logit regression')
29 plt.savefig('fig/Q26.png')
30 plt.show()

```



Q27

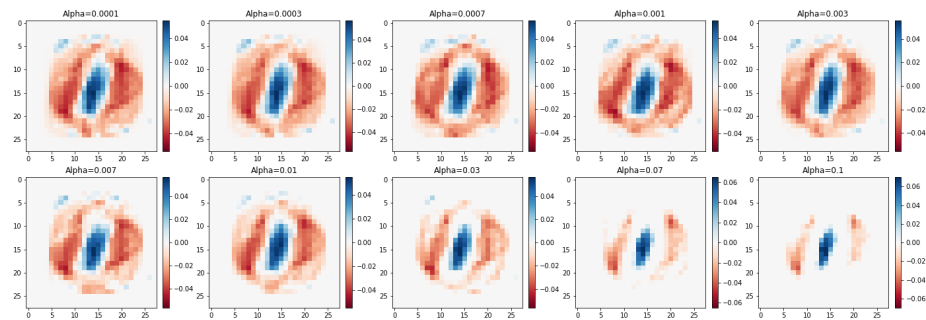
The randomness comes from the stochastic sampling of training data in SGD classifier.

Q28

$\alpha = 0.1$, with lowest mean classification error.

Q29

```
1 X_train, y_train = sub_sample(100, X_train, y_train)
2
3 fig, axs = plt.subplots(2, 5, figsize=(24, 8))
4
5 for idx, alpha in enumerate([0.0001, 0.0003, 0.0007, 0.001, 0.003,
6                               0.007, 0.01, 0.03, 0.07, 0.1]):
7     clf = SGDClassifier(loss='log', max_iter=1000,
8                           tol=1e-3,
9                           penalty='l1', alpha=alpha,
10                          learning_rate='invscaling',
11                          power_t=0.5,
12                          eta0=0.01,
13                          verbose=1)
14
15     clf.fit(X_train, y_train)
16
17     theta = clf.coef_.reshape((28, 28))
18     scale = np.abs(theta).max()
19     im = axs[idx//5, idx%5].imshow(theta, cmap=plt.cm.RdBu, vmax=
20                                     scale, vmin=-scale)
21     axs[idx//5, idx%5].title.set_text("Alpha={}".format(alpha))
22     plt.colorbar(im, ax=axs[idx//5, idx%5], fraction=0.046, pad
23                 =0.04)
24
25 plt.savefig('fig/Q29.png')
```



Q30

A larger regularization term "penalizes" high dimension. As we include regularization term in loss function, larger regularization term will force the training algorithm to return a lower θ . As we can see from the figure of Q29, the higher the regularization term, the smaller the coefficients.