# Assignment 2

1.

**1.1** $P(t=1|x) = \dfrac{P(x|t=1)p(t=1)}{p(x|t=1)p(t=1)+p(x|t=0)p(t=0)}$  (Using Baye's Rule)

$$= \dfrac{1}{1+\dfrac{P(x|t=0)P(t=0)}{P(x|t=1)P(t=1)}} = \dfrac{1}{1+\dfrac{\prod_{i=1}^{D}(2\pi\sigma_i^2)^{-\frac{1}{2}}\exp\left(\sum_{i=1}^{D}\dfrac{(x_i-\mu_{i0})^2}{-2\sigma_i^2}\right)(1-\alpha)}{\prod_{i=1}^{D}(2\pi\sigma_i^2)^{-\frac{1}{2}}\exp\left(\sum_{i=1}^{D}\dfrac{(x_i-\mu_{i1})^2}{-2\sigma_i^2}\right)\alpha}}$$

$$= \dfrac{1}{1+\exp\left(\sum_{i=1}^{D}\dfrac{(x_i-\mu_{i1})^2-(x_i-\mu_{i0})^2}{2\sigma_i^2}\right)\dfrac{1-\alpha}{\alpha}}$$

$$= \dfrac{1}{1+\exp\left(\sum_{i=1}^{D}\dfrac{2(\mu_{i0}-\mu_{i1})x_i+\mu_{i1}^2-\mu_{i0}^2}{2\sigma_i^2}+\ln\dfrac{1-\alpha}{\alpha}\right)}$$

$$= \dfrac{1}{1+\exp\left(\sum_{i=1}^{D}\dfrac{\mu_{i0}-\mu_{i1}}{\sigma_i^2}x_i+\sum_{i=1}^{D}\dfrac{\mu_{i1}^2-\mu_{i0}^2}{2\sigma_i^2}+\ln\dfrac{1-\alpha}{\alpha}\right)}$$

$$= \dfrac{1}{1+\exp\left(-\sum_{i=1}^{D}\dfrac{\mu_{i1}-\mu_{i0}}{\sigma_i^2}x_i-\left(\sum_{i=1}^{D}\dfrac{\mu_{i0}^2-\mu_{i1}^2}{2\sigma_i^2}+\ln\dfrac{\alpha}{1-\alpha}\right)\right)}$$

Let $W_i = \dfrac{\mu_{i1}-\mu_{i0}}{\sigma_i^2}$, $b = \sum_{i=1}^{D}\dfrac{\mu_{i0}^2-\mu_{i1}^2}{2\sigma_i^2}+\ln\dfrac{\alpha}{1-\alpha}$

**1.2** Let $\theta_i = x_i^T w + b$, $i \in [1,N]$.

$$L = \prod_{i=1}^{N}\left(\dfrac{1}{1+e^{-\theta_i}}\right)^{t_i}\left(1-\dfrac{1}{1+e^{-\theta_i}}\right)^{1-t_i}$$

$$= \prod_{i=1}^{N}\left(\dfrac{1}{1+e^{-\theta_i}}\right)^{t_i}\left(\dfrac{1}{1+e^{\theta_i}}\right)^{1-t_i}$$

$$l = \log L = \sum_{i=1}^{N} t_i \log\dfrac{1}{1+e^{-\theta_i}} + (1-t_i)\log\dfrac{1}{1+e^{\theta_i}}$$

$$= \sum_{i=1}^{N} t_i\left(\log\dfrac{1}{1+e^{-\theta_i}} - \log\dfrac{1}{1+e^{\theta_i}}\right) + \log\dfrac{1}{1+e^{\theta_i}}$$

$$= \sum_{i=1}^{N} t_i \log\left(\dfrac{1+e^{\theta_i}}{1+e^{-\theta_i}}\right) + \log\dfrac{1}{1+e^{\theta_i}}$$

$$= \sum_{i=1}^{N} t_i \log e^{\theta_i} - \log(1+e^{\theta_i}) = \sum_{i=1}^{N} t_i \theta_i - \log(1+e^{\theta_i})$$

$$l(w,b) = -l = \sum_{i=1}^{N} \log(1+e^{\theta_i}) - t_i \theta_i$$

$$\frac{dl}{dw_j} = \frac{dl}{d\theta_i} \frac{d\theta_i}{dw_j} = \left(\frac{1}{1+e^{-\theta_i}} - t_i\right) x_{ij},$$

write this as matrix:

$$\frac{dl}{dw} = X^T \left( \begin{pmatrix} \frac{1}{1+e^{-\theta_1}} \\ \frac{1}{1+e^{-\theta_2}} \\ \vdots \\ \frac{1}{1+e^{-\theta_N}} \end{pmatrix} - t \right)$$

$$\frac{dl}{db} = \begin{pmatrix} \frac{1}{1+e^{-\theta_1}} \\ \frac{1}{1+e^{-\theta_2}} \\ \vdots \\ \frac{1}{1+e^{-\theta_N}} \end{pmatrix} - t$$

<u>1.3</u> $p(w,b | D,\lambda) = \dfrac{p(w,b,D|\lambda)}{p(D|\lambda)}$

$$\alpha \, p(w,b,D|\lambda) = C \cdot p(w,b|\lambda) p(D|\lambda)$$

$$= C \cdot \prod_{i=1}^{N} \left(\frac{1}{1+e^{-x_i w - b}}\right)^{t_N} \left(1 - \frac{1}{1+e^{-x_i w - b}}\right)^{1-t_i} (2\pi \tfrac{1}{\lambda})^{-\frac{N}{2}} \exp\left(-\frac{\lambda}{2} w^T w\right)$$
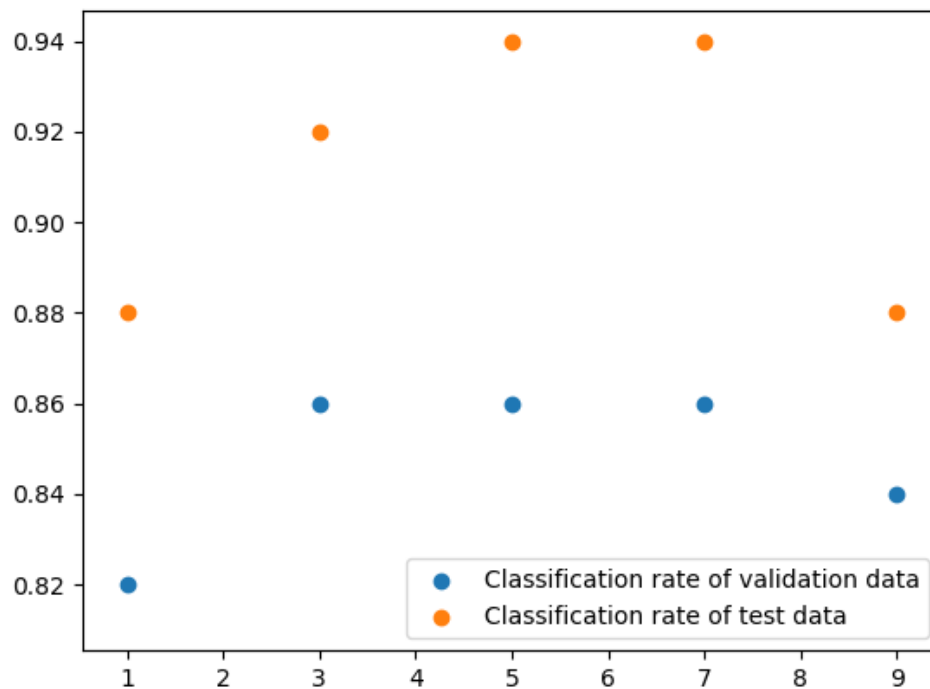
$$= C_1 \cdot \prod_{i=1}^{N} \left(\frac{1}{1+e^{-x_i w - b}}\right)^{t_i} \left(1 - \frac{1}{1+e^{-x_i w - b}}\right)^{1-t_i} \exp\left(-\frac{\lambda}{2} w^T w\right)$$

$$l_{post}(w,b) = -\left(\log C_1 - l(w,b) - \frac{\lambda}{2} w^T w\right)$$
$$= l(w,b) + \frac{\lambda}{2} w^T w - \log C_1$$
$$= l(w,b) + \frac{\lambda}{2} \sum_{j=1}^{D} w_j^2 + C_2$$

2.1



The classification rate on validation set has the trend of increasing then decreasing. We choose k = 5 according to the classification rate performance on validation set. K = 3 and k = 7 both have classification rate 86% as high as k = 5. The test performance corresponds to validation set generally.

Code of 2.1

```python
import numpy as np
import matplotlib.pyplot as plt
from run_knn import run_knn
from logistic_regression_template import run_logistic_regression

# (a)
train = np.load("mnist_train.npz")
# print(train.files)
train_data = train['train_inputs']
# print(train_data.shape)
train_labels = train['train_targets']

valid = np.load("mnist_valid.npz")
# print(valid.files)
valid_data = valid['valid_inputs']
# print(valid_data.shape)
valid_labels = valid['valid_targets']

test = np.load("mnist_test.npz")
# print(test.files)
test_data = test['test_inputs']
# print(test_data.shape)
test_labels = test['test_targets']


classification_rate_validation = []
classification_rate_test = []
for k in [1, 3, 5, 7, 9]:
```

```
22     # print(test_data.shape)
23     test_labels = test['test_targets']
24
25
26     classification_rate_validation = []
27     classification_rate_test = []
28     for k in [1, 3, 5, 7, 9]:
29         vk_labels = run_knn(k, train_data, train_labels, valid_data)
30         tk_labels = run_knn(k, train_data, train_labels, test_data)
31         classification_rate_validation.append(np.count_nonzero(vk_labels == valid_labels) / len(val
32         classification_rate_test.append(np.count_nonzero(tk_labels == test_labels) / len(test_label
33
34     print(classification_rate_validation)
35     print(classification_rate_test)
36
37     plt.scatter(np.array([1, 3, 5, 7, 9]), classification_rate_validation)
38     plt.scatter(np.array([1, 3, 5, 7, 9]), classification_rate_test)
39     plt.legend(['Classification rate of validation data', 'Classification rate of test data'])
40     plt.show()
41
42
43
44
```
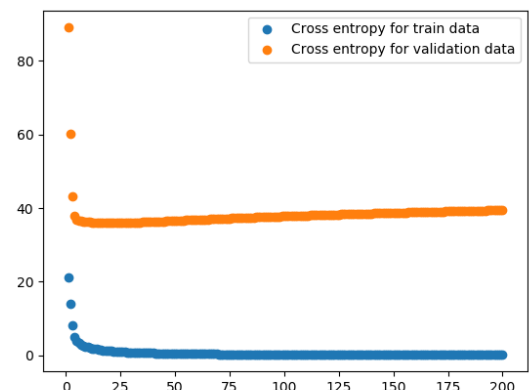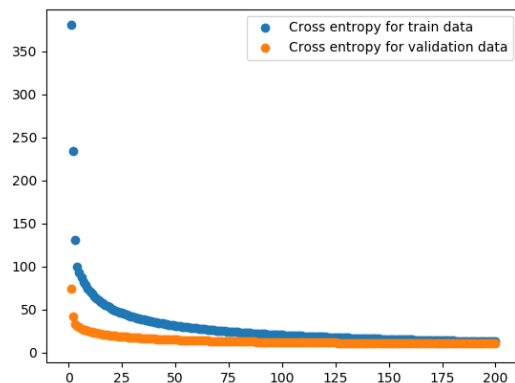
## 2.2

The best hyperparameter I found for this model is learning rate = 0.1, num_interations
= 200. Final cross entropy and classification error on the training, validation and test
sets are:

```
TRAIN CE:12.710649916810699 TRAIN FRAC:100.0   VALID CE:10.610897864366496   VALID FRAC:88.0
TEST CE:10.452382270510103 TEST FRAC:92.0
```



The left plot is the cross entropy change for mnist_train set and validation set. The right
pot is the cross entropy change for minst_train_small set and validation set. With the
best hyperparameter set I choose, the two plot remains relatively stable in each time I
run the code.

# Code for 2.2

```python
def logistic_predict(weights, data):
    """
    Compute the probabilities predicted by the logistic classifier.

    Note: N is the number of examples and
          M is the number of features per example.

    Inputs:
        weights:    (M+1) x 1 vector of weights, where the last element
                    corresponds to the bias (intercepts).
        data:       N x M data matrix where each row corresponds
                    to one data point.
    Outputs:
        y:          :N x 1 vector of probabilities. This is the output of the classifier.
    """
    intercept = np.ones((data.shape[0], 1))
    x = np.append(data, intercept, axis=1)
    log_odds = np.dot(x, weights)
    y = sigmoid(log_odds)
    return y


def evaluate(targets, y):
    """
    Compute evaluation metrics.
    Inputs:
        targets : N x 1 vector of targets.
```

```
        targets : N x 1 vector of targets.
        y       : N x 1 vector of probabilities.
    Outputs:
        ce             : (scalar) Cross entropy. CE(p, q) = E_p[-log q]. Here we want to compute
        frac_correct : (scalar) Fraction of inputs classified correctly.
    """
    ce = -(np.dot(targets.T, np.log(y)) + np.dot((1 - targets).T, np.log(1 - y)))
    y = (y >= 0.5).astype(int)
    frac_correct = np.count_nonzero(y == targets) / len(targets)
    return ce, frac_correct


def logistic(weights, data, targets, hyperparameters):
    """
    Calculate negative log likelihood and its derivatives with respect to weights.
    Also return the predictions.

    Note: N is the number of examples and
          M is the number of features per example.

    Inputs:
        weights:    (M+1) x 1 vector of weights, where the last element
                    corresponds to bias (intercepts).
        data:       N x M data matrix where each row corresponds
                    to one data point.
        targets:    N x 1 vector of targets class probabilities.
        hyperparameters: The hyperparameters dictionary.
```
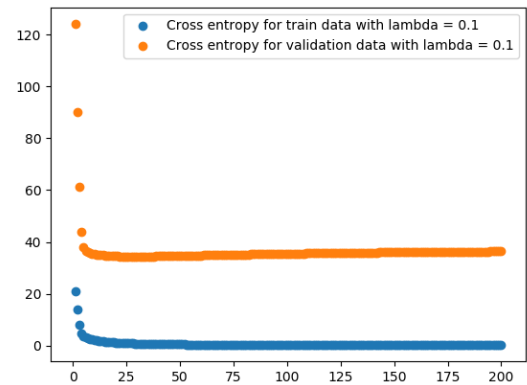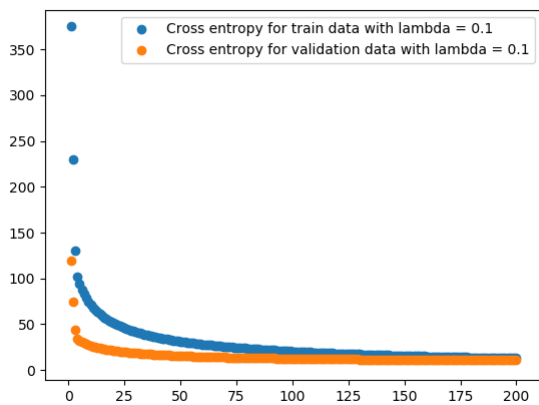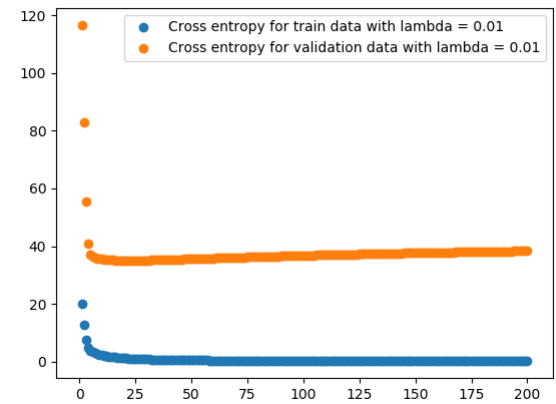
```
51              M is the number of features per example.
52
53          Inputs:
54              weights:    (M+1) x 1 vector of weights, where the last element
55                          corresponds to bias (intercepts).
56              data:       N x M data matrix where each row corresponds
57                          to one data point.
58              targets:    N x 1 vector of targets class probabilities.
59              hyperparameters: The hyperparameters dictionary.
60
61          Outputs:
62              f:          The sum of the loss over all data points. This is the objective that we want
63              df:         (M+1) x 1 vector of derivative of f w.r.t. weights.
64              y:          N x 1 vector of probabilities.
65          """
66
67          y = logistic_predict(weights, data)
68          intercept = np.ones((data.shape[0], 1))
69          x = np.append(data, intercept, axis=1)
70          log_odds = np.dot(x, weights)
71          f = np.sum(np.log(1+np.exp(log_odds)) - np.multiply(log_odds, targets))
72          df = np.dot(x.T, y - targets)
73          return f, df, y
74
75
76      def logistic_pen(weights, data, targets, hyperparameters):
77          """
78          Calculate negative log likelihood and its derivatives with respect to weights.
79          Also return the predictions
```
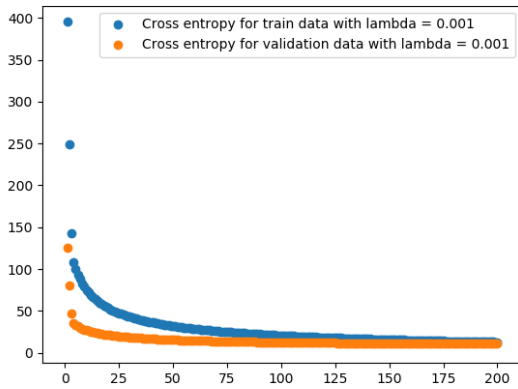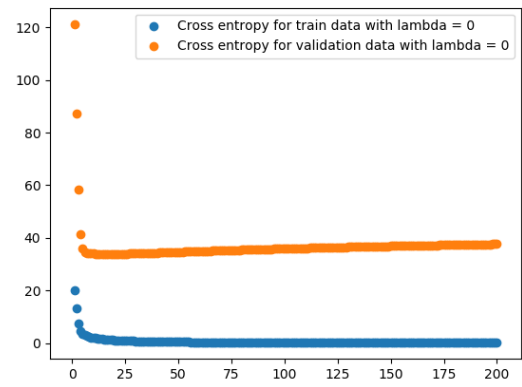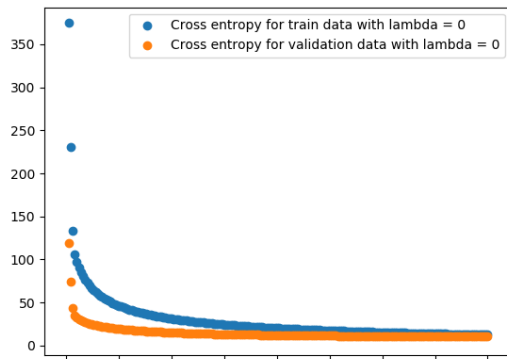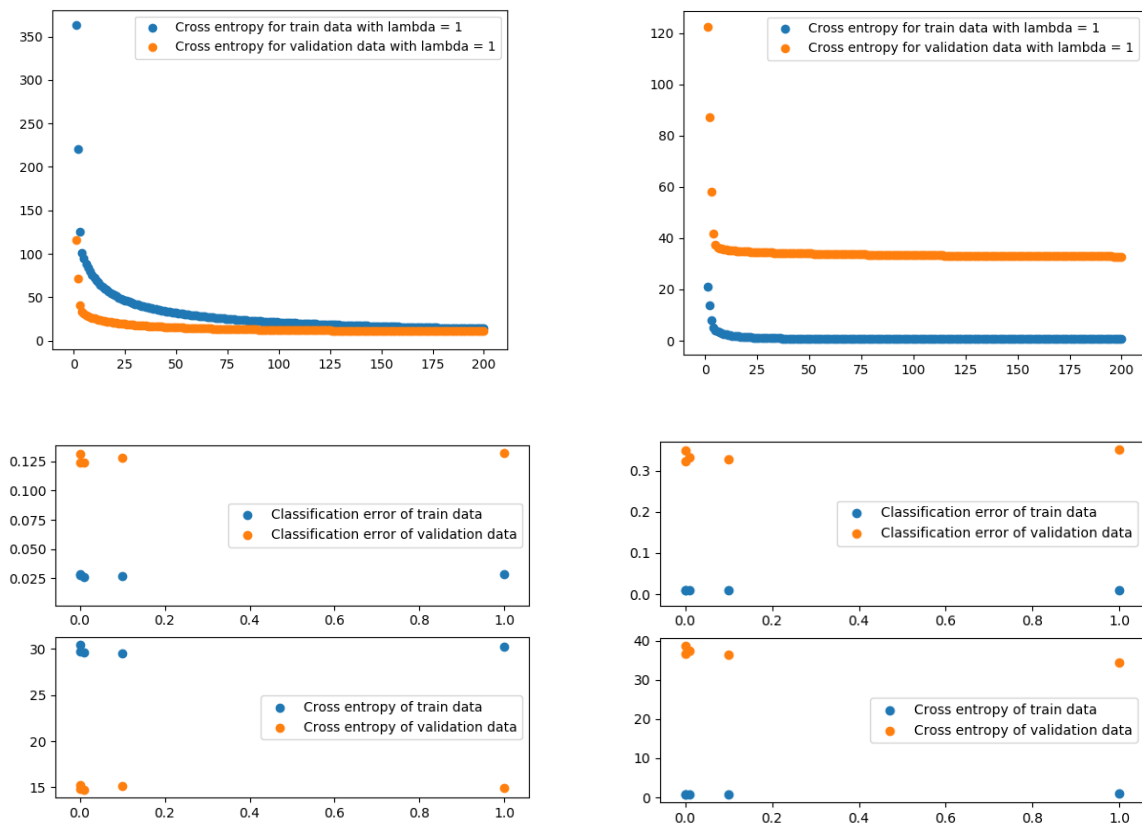
logistic()

72:5   LF÷ UTF-8÷

```
123
124             # print some stats
125             # print("ITERATION:{}  TRAIN NLOGL:{}  TRAIN CE:{} "
126             #       "TRAIN FRAC:{}  VALID NLOGL:{}  VALID CE:{}  VALID FRAC:{}".format(
127             #        t + 1, f / N, cross_entropy_train, frac_correct_train * 100, f_valid / valid_inp
128             #        cross_entropy_valid, frac_correct_valid * 100))
129             print("ITERATION:{}  TEST NLOGL:{}  TEST CE:{} "
130                   "TEST FRAC:{}".format(
131                   t + 1, f_test[0][0] / test_targets.shape[0],
132                   cross_entropy_test[0][0], frac_correct_test * 100))
133
134         plt.scatter(np.arange(1, hyperparameters['num_iterations'] + 1), TRAIN_CE)
135         plt.scatter(np.arange(1, hyperparameters['num_iterations'] + 1), VALID_CE)
136         plt.legend(['Cross entropy for train data with lambda = ' + str(lmbda),
137                     'Cross entropy for validation data with lambda = ' + str(lmbda)])
138         plt.show()
139
```

## 2.3

Here are several plots:

The plots on the left are cross entropy changes each iteration of train and validation set with different $\lambda$ using training set minst_train, plots on the right are generated in the same procedure as left plots but using training set minst_small.

The last two plots both containing 2 subplots representing the average change over iterations of cross entropy and classification error change against different values of $\lambda$. Both cross entropy and classification error have the trend of increasing first then decreasing. This is reasonable as when we increase the penalty $\lambda$, we may fix some collinearity issue in the input set. However, when the penalty is large, our estimator could be severely biased, so that classification error will increase. We choose $\lambda = 0.001$ at last, the test classification error is 8%, with 92% accuracy.

ITERATION:200  TEST NLOGL:0.20692403247235058  TEST CE:10.341793538031553  TEST FRAC:92.0
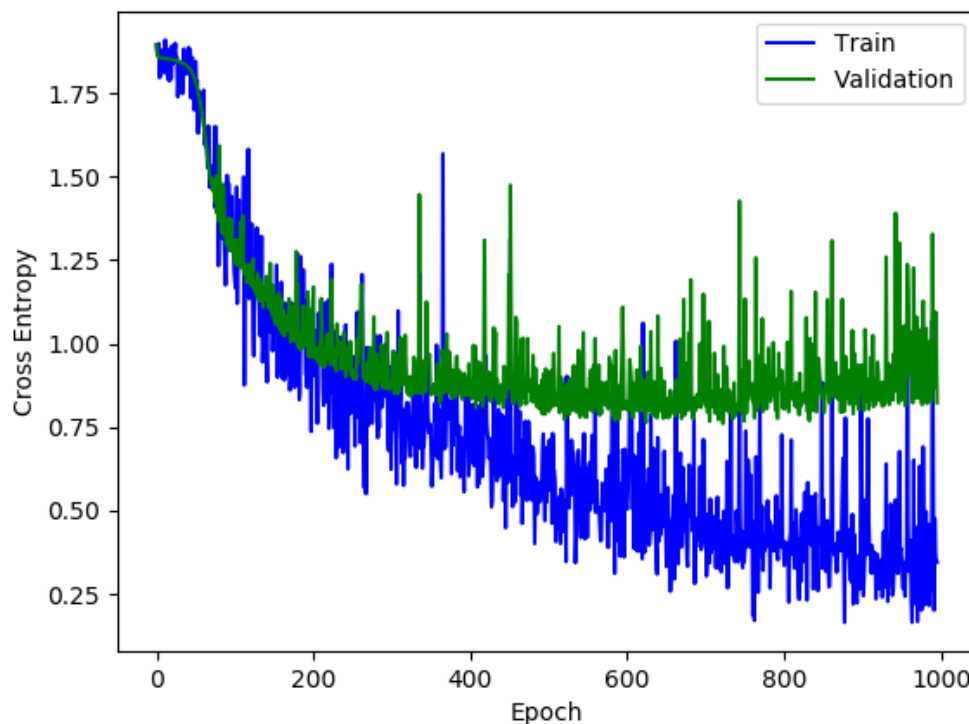
Finally, we see that there is no significant difference in adding the penalty. We think this is probably because the data values in 28*28 pixels are kind of random distributed, so columns in input matrix are quite linearly independent so that not much regularization is needed to modify this model.
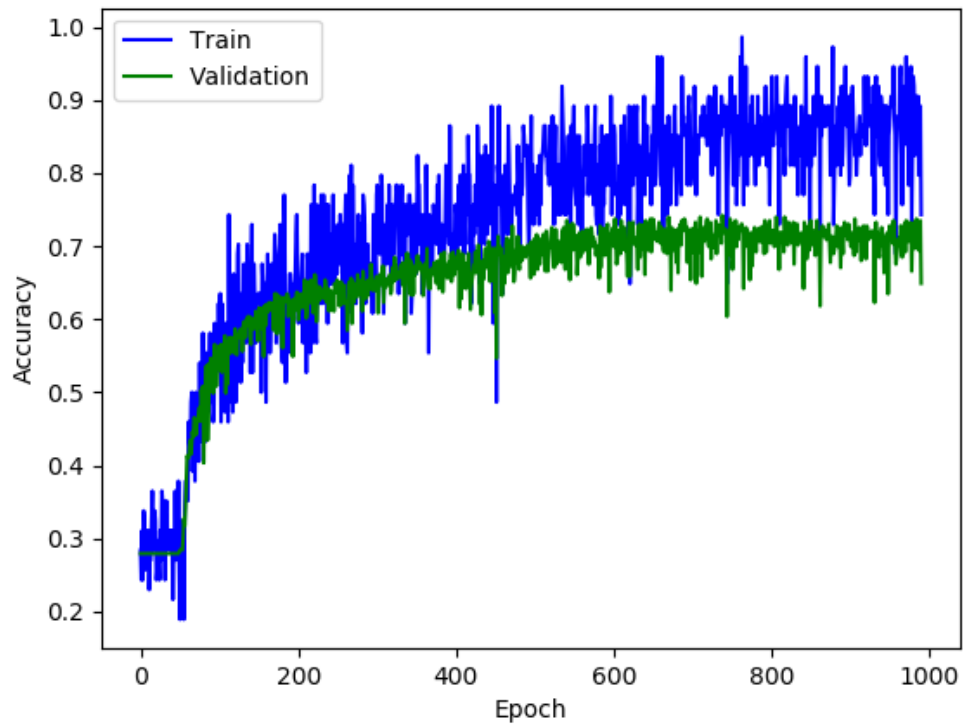
Code for 2.3:

```
1   from logistic_regression_template import *
2   import matplotlib.pyplot as plt
3
4   list_train_error_avg = []
5   list_train_ce_avg = []
6   list_valid_error_avg = []
7   list_valid_ce_avg = []
8   list_lmbda = [0, 0.001, 0.01, 0.1, 1]
9   for lmbda in list_lmbda:
10      train_error_avg, train_ce_avg,\
11          valid_error_avg, valid_ce_avg = run_pen_logistic_regression(lmbda)
12      list_train_error_avg.append(train_error_avg)
13      list_train_ce_avg.append(train_ce_avg)
14      list_valid_error_avg.append(valid_error_avg)
15      list_valid_ce_avg.append(valid_ce_avg)
16
17  fig, (ax1, ax2) = plt.subplots(2)
18  ax1.scatter(np.array(list_lmbda), list_train_error_avg)
19  ax1.scatter(np.array(list_lmbda), list_valid_error_avg)
20  ax1.legend(['Classification error of train data', 'Classification error of validation data'])
21  ax2.scatter(np.array(list_lmbda), list_train_ce_avg)
22  ax2.scatter(np.array(list_lmbda), list_valid_ce_avg)
23  ax2.legend(['Cross entropy of train data', 'Cross entropy of validation data'])
24  plt.show()
25
26  run_pen_logistic_regression(0.001)
```

3.1

We get the two plots indicating cross entropy and accuracy with random hyperparameters:
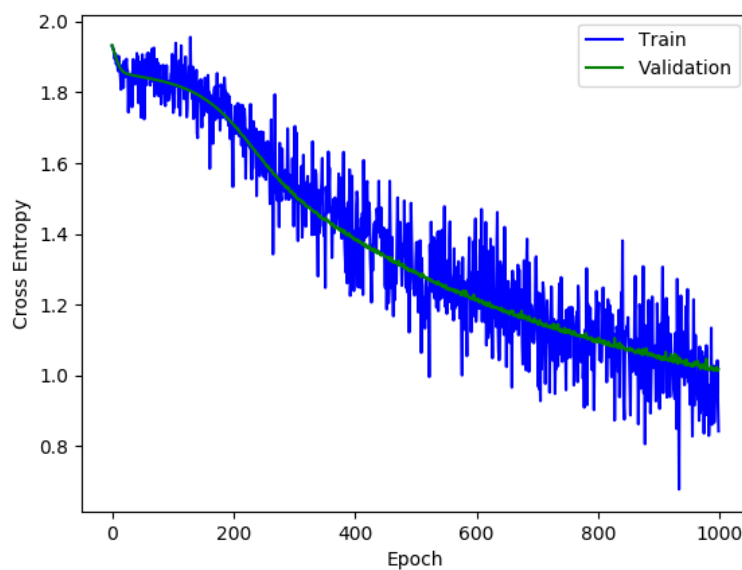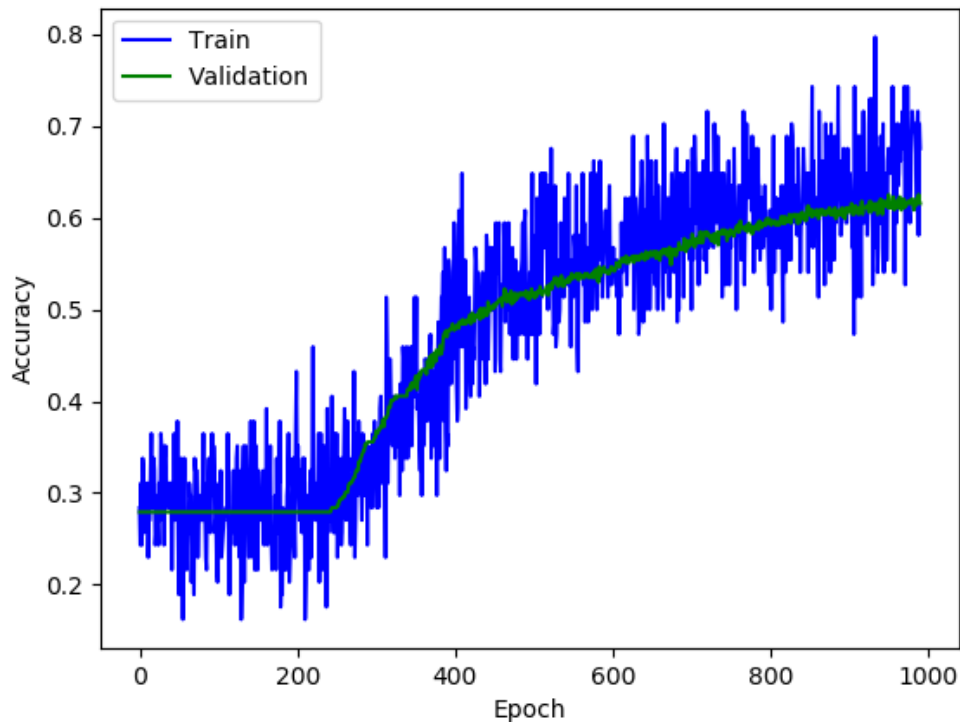
Both training and validation set has the same trend: increasing accuracy, decreasing cross entropy.
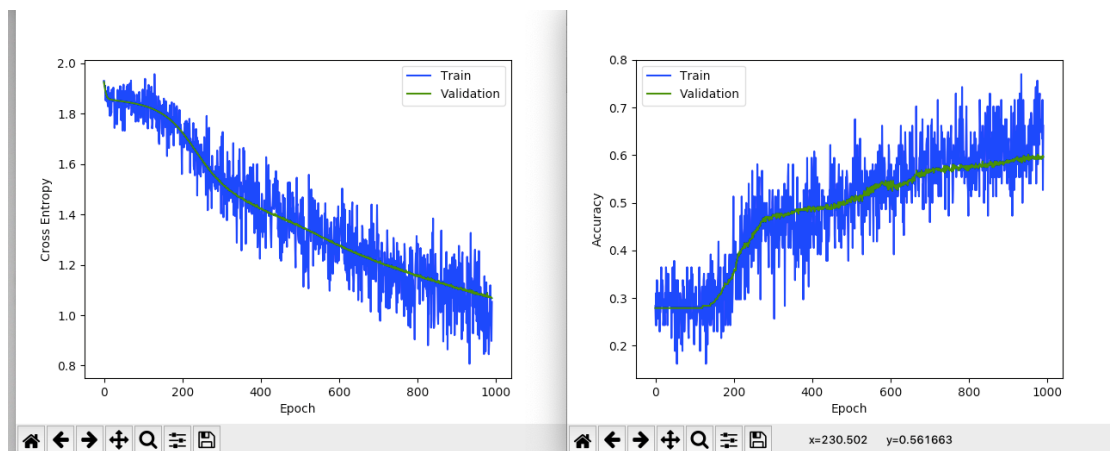
3.2
After experimenting, we find that as learning and increases, accuracy increases. Plot for learning rate = 0.001:

It seems that momentum = 0.6 works best:



Statistics for the above plots:
CE: Train 1.03181 Validation 1.06957 Test 1.09550
Acc: Train 0.62893 Validation 0.59189 Test 0.59740
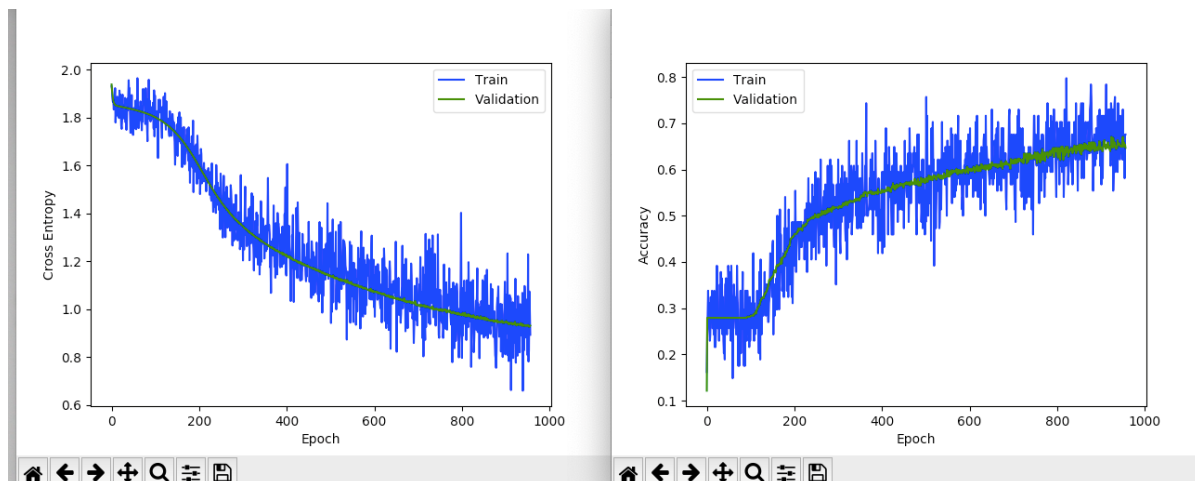
We found that batch size increases, the accuracy decreases and the cross entropy increases. Therefore, we choose batch size = 100.

3.3
Due to time pressure, we only tested a few choices for number of hidden layers and this choice seems to work better. Thus, we guess that with increasing hidden layers, the accuracy should increase, but should be a upper bound for this relation.

HIDDEN [32,64]:



Statistics for the above plots:
CE: Train 0.98503 Validation 1.05586 Test 1.05219
Acc: Train 0.64878 Validation 0.60382 Test 0.60519

3.4
Some faces may have ambiguous face expressions, so it may be hard to classify them into the 7 categories. We try to show those face photos from the original data set.

Code for 3.1-3.4:

```python
def AffineBackward(grad_y, h, w):
    """Computes gradients of affine transformation.
    hint: you may need the matrix transpose np.dot(A,B).T = np.dot(B,A) and (A.T).T = A

    Args:
        grad_y: gradient from last layer
        h: inputs from the hidden layer
        w: weights

    Returns:
        grad_h: Gradients wrt. the inputs/hidden layer.
        grad_w: Gradients wrt. the weights.
        grad_b: Gradients wrt. the biases.
    """
    ###########################
    # Insert your code here.
    grad_w = np.dot(h.T, grad_y)
    grad_h = np.dot(grad_y, w.T)
    grad_b = np.sum(grad_y, axis=0).T
    return grad_h, grad_w, grad_b
    ###########################
```

```python
104
105
106   def ReLUBackward(grad_h, z):
107       """Computes gradients of the ReLU activation function wrt. the unactivated inputs.
108
109       Returns:
110           grad_z: Gradients wrt. the hidden state prior to activation.
111       """
112       ############################
113       # Insert your code here.
114       grad_z = np.zeros(grad_h.shape)
115       for i in range(grad_h.shape[0]):
116           for j in range(grad_h.shape[1]):
117               if z[i][j] > 0:
118                   grad_z[i][j] = grad_h[i][j]
119       return grad_z
120       ############################
121
184   def NNUpdate(model, eps, momentum):
185       """Update NN weights.
186
187       Args:
188           model:    Dictionary of all the weights.
189           eps:      Learning rate.
190           momentum: Momentum.
191       """
192       ############################
193       # Insert your code here.
194       # Update the weights.
195       model['W1'] = model['W1'] - eps * (momentum * model['W1'] +
196                                          (1 - momentum) * model['dE_dW1'])
197       model['W2'] = model['W2'] - eps * (momentum * model['W2'] +
198                                          (1 - momentum) * model['dE_dW2'])
199       model['W3'] = model['W3'] - eps * (momentum * model['W3'] +
200                                          (1 - momentum) * model['dE_dW3'])
201       model['b1'] = model['b1'] - eps * (momentum * model['b1'] +
202                                          (1 - momentum) * model['dE_db1'])
203       model['b2'] = model['b2'] - eps * (momentum * model['b2'] +
204                                          (1 - momentum) * model['dE_db2'])
205       model['b3'] = model['b3'] - eps * (momentum * model['b3'] +
206                                          (1 - momentum) * model['dE_db3'])
207       ############################
208
```