

ECE521 A1

Baiwu Zhang 1001127540

Daiqing Li 1000795357

January 2017

1

1.1

1.1.1

Construct a data set as follow, alternatively place two nodes from different class one after the other. A graphic illustration of the data set is shown in Figure 1. The data should be large enough to avoid corner cases at the edges. Accuracy with different K values is in Table ??.



Figure 1: The data set for KNN with periodical accuracy.

Table 1: K vs Accuracy

K	1	2	3	4	5	6	7	8	9	...
Acc	100	50	0	50	100	50	0	50	100	(period = 4)

Analysis:

$K = 1$: Every data point chooses itself. Accuracy is 100%.

$K = 2$: Every data point chooses itself and its opposite color neighbour. A tie is broke by randomly pick one. Accuracy is 50%.

$k = 3$: Other than the two points that already been picked, each data point picks the third nearest node, which is in the opposite color. Accuracy is 0%.

$k = 4$: Other than the above three points, two closest points to the data point have the same distance and same color. Pick either one, a tie is created and it will have 50% accuracy.

$k = 5$: Pick the one not been picked in last step, which has the same color as the data point. Accuracy is 100%. The pattern starts to repeat.

1.1.2

$$\begin{aligned}
\text{var}\left(\frac{\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|_2^2}{E[\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|_2^2]}\right) &= \text{var}\left(\frac{\|\mathbf{d}\|_2^2}{E[\|\mathbf{d}\|_2^2]}\right) \\
&= \frac{\text{var}(\|\mathbf{d}\|_2^2)}{E[\|\mathbf{d}\|_2^2]^2} \\
&= \frac{E[\|\mathbf{d}\|_2^4] - E[\|\mathbf{d}\|_2^2]^2}{E[\|\mathbf{d}\|_2^2]^2} \\
&= \frac{E[\|\mathbf{d}\|_2^4]}{E[\|\mathbf{d}\|_2^2]^2} - 1 \\
&= \frac{E[(\sum_{i=1}^N d_i^2)^2]}{E[\sum_{i=1}^N d_i^2]^2} - 1 \\
&= \frac{E[\sum_{i=1}^N \sum_{j=1}^N (d_i^2 d_j^2)]}{\sum_{i=1}^N \sum_{j=1}^N E[d_i^2] E[d_j^2]} - 1 \\
&= \frac{\sum_{i=1}^N \sum_{j=1}^N E[d_i^2 d_j^2]}{\sum_{i=1}^N \sum_{j=1}^N E[d_i^2] E[d_j^2]} - 1 \\
&= \frac{\sum_{i=1}^N E[d_i^4] + \sum_{i=1}^N \sum_{j=1, j \neq i}^N (E[d_i^2] E[d_j^2])}{\sum_{i=1}^N E[d_i^2]^2 + \sum_{i=1}^N \sum_{j=1, j \neq i}^N (E[d_i^2] E[d_j^2])} - 1 \\
&= \frac{N * 12\sigma^4 + \sum_{i=1}^N \sum_{j=1, j \neq i}^N 4\sigma^4}{N * 4\sigma^4 + \sum_{i=1}^N \sum_{j=1, j \neq i}^N 4\sigma^4} - 1 \\
&= \frac{N * 12\sigma^4 + (N^2 - N)4\sigma^4}{N * 4\sigma^4 + (N^2 - N)4\sigma^4} - 1 \\
&= \frac{3N + N^2 - N}{N + N^2 - N} - 1 \\
&= \frac{N + 2}{N} - 1
\end{aligned} \tag{1}$$

1.2

1.2.1

$$\begin{aligned}
\|\mathbf{x}^{(m)} - \mathbf{x}^*\|_2^2 &= \sum_{n=1}^N (x_n^{(m)} - x_n^*)^2 \\
&= \sum_{n=1}^N (x_n^{(m)2} + x_n^{*2}) - 2 \sum_{n=1}^N x_n^{(m)} x_n^* \\
&= 2C - 2\mathbf{x}^{(m)T} \mathbf{x}^* \\
&= \text{constant} + \text{constant} * -\mathbf{x}^{(m)T} \mathbf{x}^*
\end{aligned} \tag{2}$$

From equation (2), we can conclude that as long as we rank $-\mathbf{x}^{(m)T} \mathbf{x}^*$ for every input, we can the nearest neighbour for \mathbf{x}^* .

1.2.2

Please see notebook for the coding and testing.

1.3

1.3.1

Please see code section for our implementation.

1.3.2

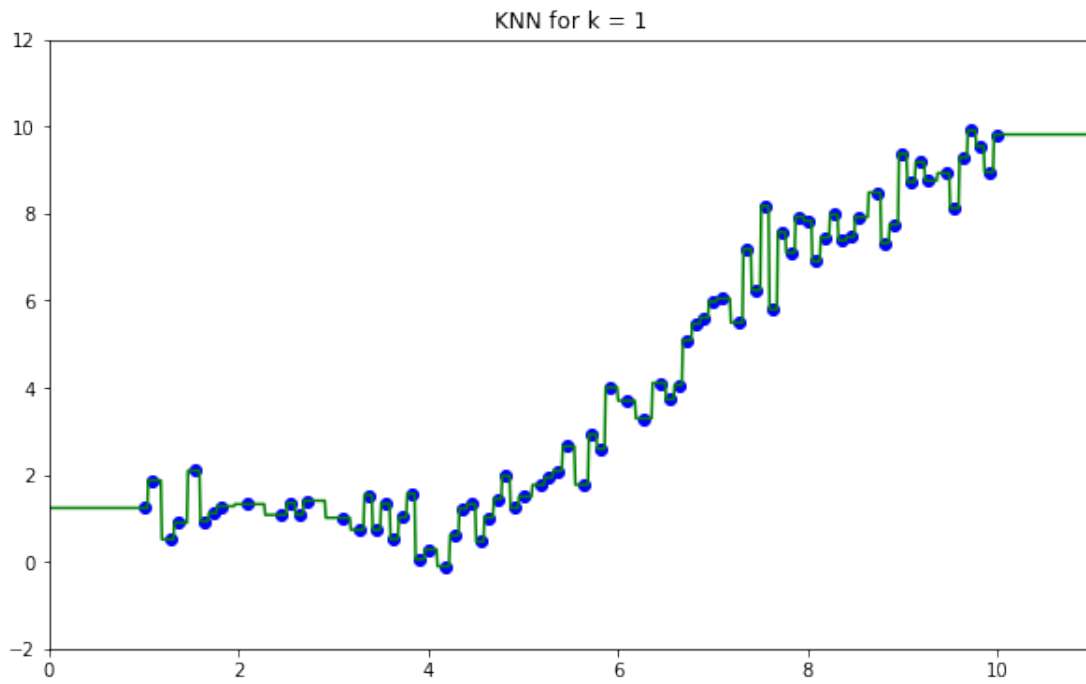


Figure 2: k: 1, trainMSE: 0.00, validMSE: 5.43, testMSE: 6.22

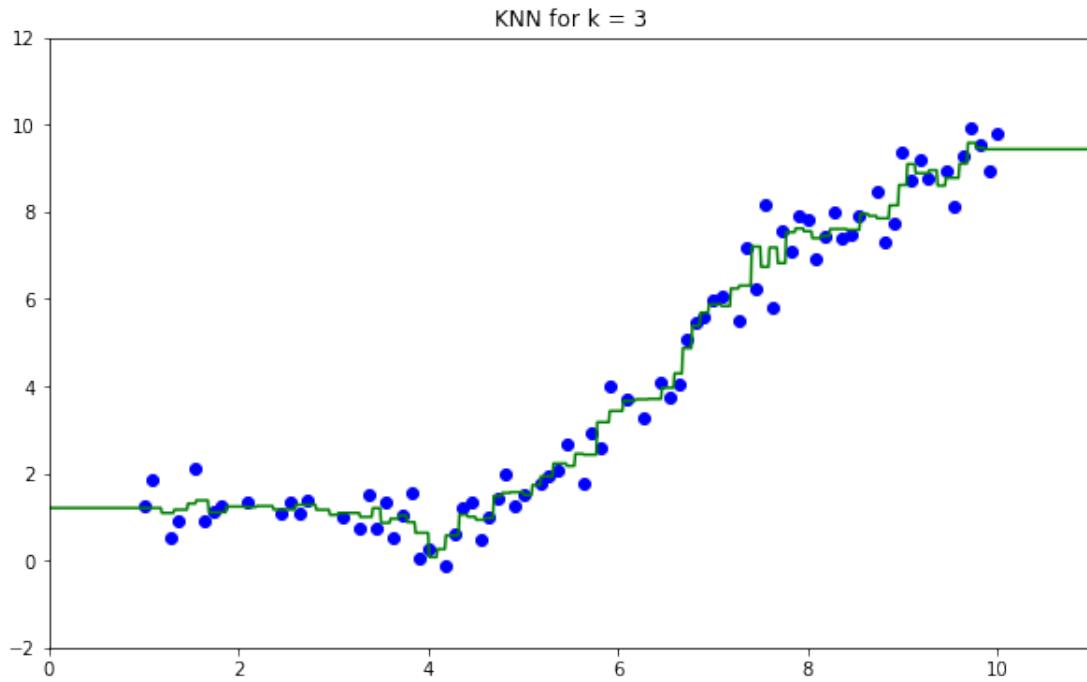


Figure 3: k: 3, trainMSE: 16.84, validMSE: 6.53, testMSE: 2.90

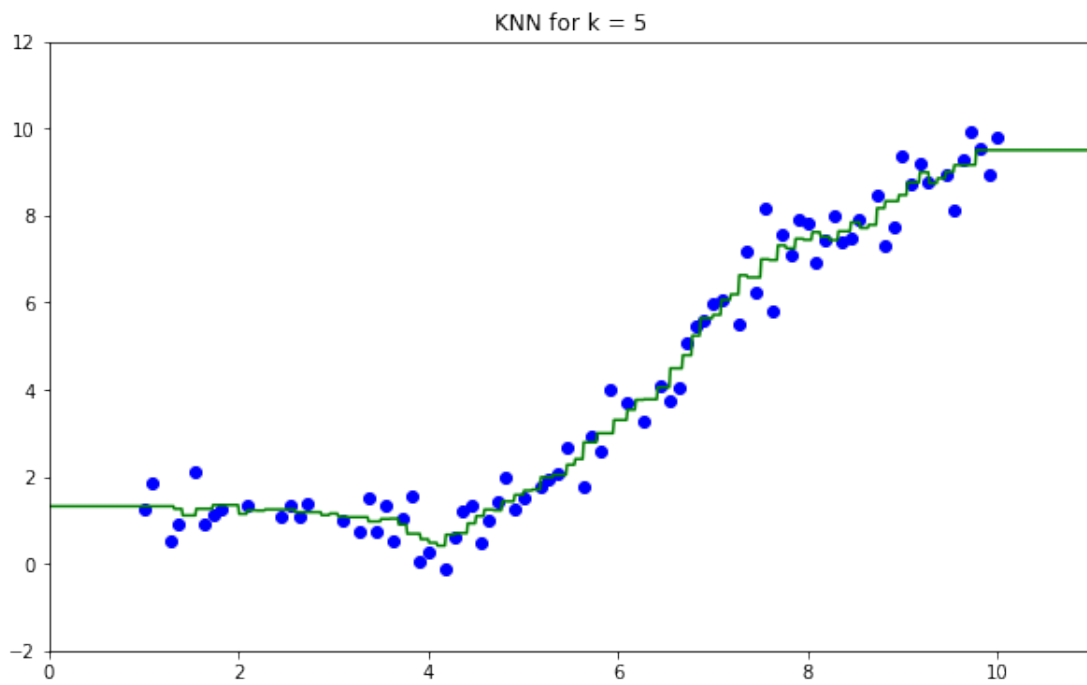


Figure 4: k: 5, trainMSE: 18.97, validMSE: 6.21, testMSE: 3.57

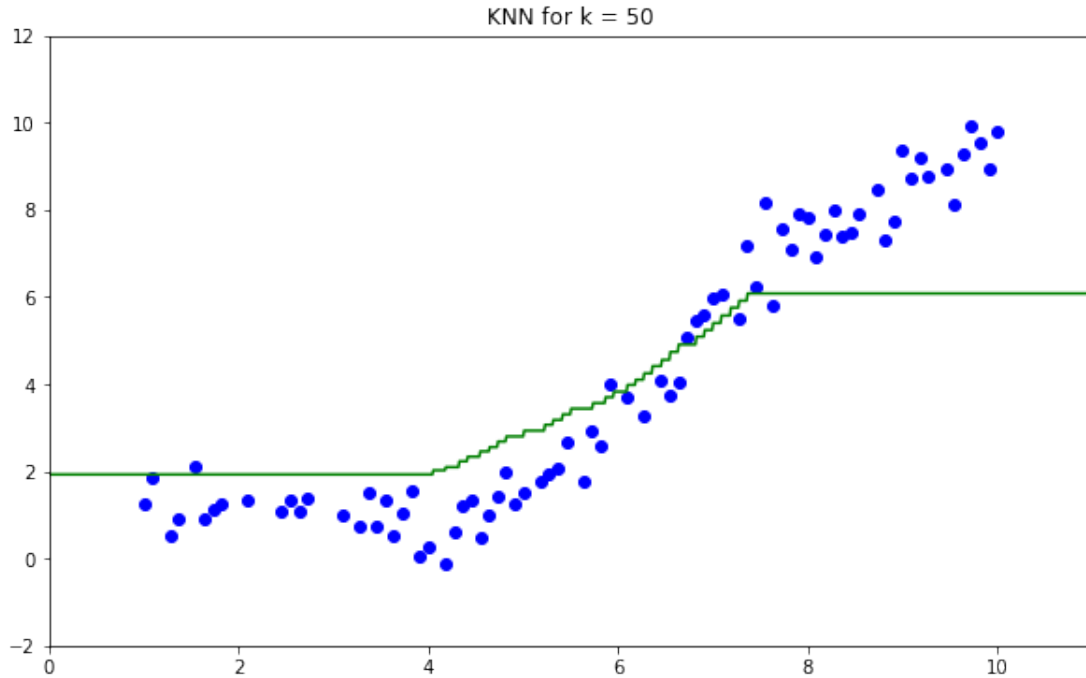


Figure 5: k: 50, trainMSE: 199.68, validMSE: 24.57, testMSE: 14.14

From the calculation, we find the best k value based on validation error is $k = 1$. From the plotting, $k = 1$ fits every training points perfectly, but it is overfitting. For $k = 3$, it is more smooth, but still a little bit overfitting. For $k = 50$, the prediction has a great offset comparing to the training points. We think $k = 5$ is a better choice because it is more general and smooth.

1.4

1.4.1

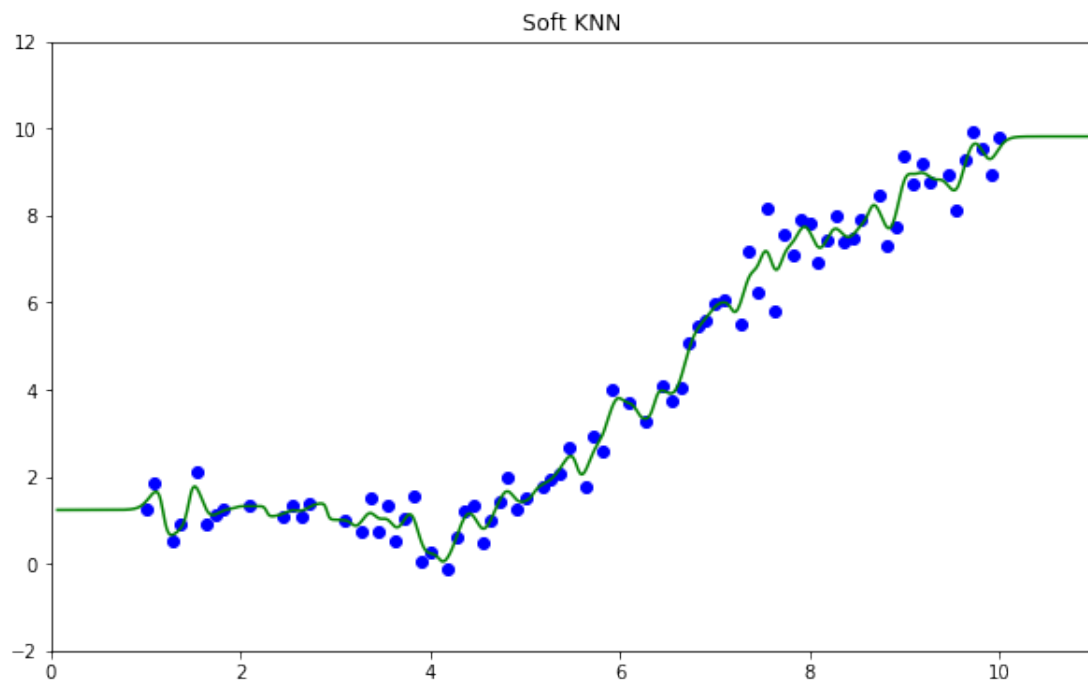


Figure 6: Soft KNN

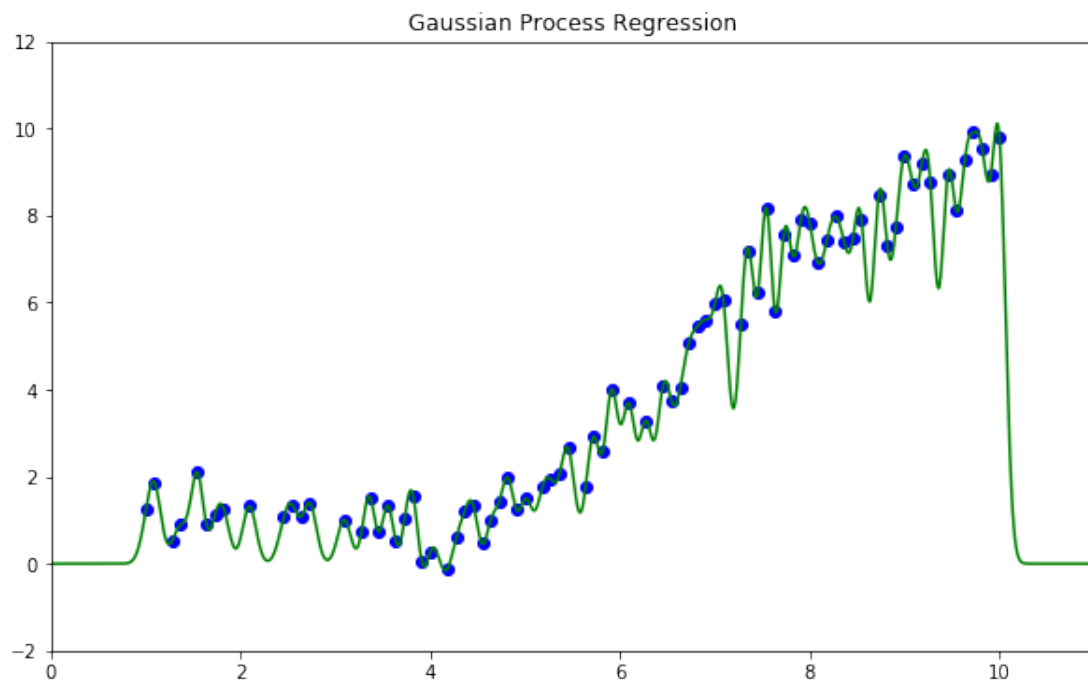


Figure 7: Gaussian Process Regression

The soft KNN doesn't overfit the data, and it gives a smooth curve compared to regular KNN. The Gaussian process regression overfits the data. But it gives confident predictions when the testing point is close to any of the training points. In the region where no data point exists, it gives predictions close to its prior, which is zero in this case.

1.4.2

Let $z = y^* + Ay_{train}$ where $A = -\Sigma_{y^*y_{train}}\Sigma_{y_{train}y_{train}}^{-1}$.
Now we can write

$$\begin{aligned} cov(z, y_{train}) &= cov(y^*, y_{train}) + cov(Ay_{train}, y_{train}) \\ &= \Sigma_{y^*y_{train}} + Avar(y_{train}) \\ &= \Sigma_{y^*y_{train}} - \Sigma_{y^*y_{train}}\Sigma_{y_{train}y_{train}}^{-1}\Sigma_{y_{train}y_{train}} \\ &= 0 \end{aligned} \tag{3}$$

Therefore z and y_{train} are uncorrelated and, since they are jointly normal, they are independent. We know $E(z) = \mu_{y^*} + A\mu_{y_{train}} = 0$. Therefore,

$$\begin{aligned} \mu^* &= E(y^*|y_{train}) = E(z - Ay_{train}|y_{train}) \\ &= E(z|y_{train}) - E(Ay_{train}|y_{train}) \\ &= E(z) - Ay_{train} \\ &= -Ay_{train} \\ &= \Sigma_{y_{train}y^*}^T \Sigma_{y_{train}y_{train}}^{-1} y_{train} \end{aligned} \tag{4}$$

Variance:

$$\begin{aligned} \Sigma^* &= var(y^*|y_{train}) = var(z - Ay_{train}|y_{train}) \\ &= var(z|y_{train}) + var(Ay_{train}|y_{train}) - Acov(z, -y_{train}) - cov(z, -y_{train})A' \\ &= \Sigma_{y^*y^*} + \Sigma_{y^*y_{train}}\Sigma_{y_{train}y_{train}}^{-1}\Sigma_{y_{train}y_{train}}\Sigma_{y_{train}y_{train}}^{-1}\Sigma_{y_{train}y^*} - 2\Sigma_{y^*y_{train}}\Sigma_{y_{train}y_{train}}^{-1}\Sigma_{y_{train}y^*} \\ &= \Sigma_{y^*y^*} + \Sigma_{y^*y_{train}}\Sigma_{y_{train}y_{train}}^{-1}\Sigma_{y_{train}y^*} - 2\Sigma_{y^*y_{train}}\Sigma_{y_{train}y_{train}}^{-1}\Sigma_{y_{train}y^*} \\ &= \Sigma_{y^*y^*} - \Sigma_{y_{train}y^*}^T \Sigma_{y_{train}y_{train}}^{-1} \Sigma_{y_{train}y^*} \end{aligned} \tag{5}$$

2

2.1

2.1.1

Let $L(\mathbf{W}) = \sum_{m=1}^M \frac{1}{2M} \left\| \mathbf{W}^T \mathbf{x}^m + b - y^m \right\|_2^2 + \frac{\lambda}{2} \|\mathbf{W}\|_2^2$.

For \mathbf{W} :

$$\begin{aligned}
L(t\mathbf{W}_1 + (1-t)\mathbf{W}_2) &= \sum_{m=1}^M \frac{1}{2M} \left\| (t\mathbf{W}_1 + (1-t)\mathbf{W}_2)^T \mathbf{x}^m + b - y^m \right\|_2^2 + \frac{\lambda}{2} \|t\mathbf{W}_1 + (1-t)\mathbf{W}_2\|_2^2 \\
&= \sum_{m=1}^M \frac{1}{2M} \left\| t\mathbf{W}_1^T \mathbf{x}^m + (1-t)\mathbf{W}_2^T \mathbf{x}^m + b - y^m \right\|_2^2 + \frac{\lambda}{2} \|t\mathbf{W}_1 + (1-t)\mathbf{W}_2\|_2^2 \\
&= \sum_{m=1}^M \frac{1}{2M} \left\| t(\mathbf{W}_1^T \mathbf{x}^m + b - y^m) + (1-t)(\mathbf{W}_2^T \mathbf{x}^m + b - y^m) \right\|_2^2 + \frac{\lambda}{2} \|t\mathbf{W}_1 + (1-t)\mathbf{W}_2\|_2^2 \\
&\leq \sum_{m=1}^M \frac{1}{2M} \left\| t(\mathbf{W}_1^T \mathbf{x}^m + b - y^m) \right\|_2^2 + \frac{\lambda}{2} \|t\mathbf{W}_1\|_2^2 + \\
&\quad \sum_{m=1}^M \frac{1}{2M} \left\| (1-t)(\mathbf{W}_2^T \mathbf{x}^m + b - y^m) \right\|_2^2 + \frac{\lambda}{2} \|(1-t)\mathbf{W}_2\|_2^2 \\
&\leq t \sum_{m=1}^M \frac{1}{2M} \left\| \mathbf{W}_1^T \mathbf{x}^m + b - y^m \right\|_2^2 + t \frac{\lambda}{2} \|\mathbf{W}_1\|_2^2 + \\
&\quad (1-t) \sum_{m=1}^M \frac{1}{2M} \left\| \mathbf{W}_2^T \mathbf{x}^m + b - y^m \right\|_2^2 + (1-t) \frac{\lambda}{2} \|(1-t)\mathbf{W}_2\|_2^2 \\
&\leq tL(\mathbf{W}_1) + (1-t)L(\mathbf{W}_2)
\end{aligned} \tag{6}$$

For b :

$$\begin{aligned}
L(tb_1 + (1-t)b_2) &= \sum_{m=1}^M \frac{1}{2M} \left\| \mathbf{W}^T \mathbf{x}^m + (tb_1 + (1-t)b_2) - y^m \right\|_2^2 + \frac{\lambda}{2} \|\mathbf{W}\|_2^2 \\
&= \sum_{m=1}^M \frac{1}{2M} \left\| t(\mathbf{W}^T \mathbf{x}^m + b_1 - y^m) + (1-t)(\mathbf{W}^T \mathbf{x}^m + b_2 - y^m) \right\|_2^2 + \frac{\lambda}{2} \|\mathbf{W}\|_2^2 \\
&\leq \sum_{m=1}^M \frac{1}{2M} \left\| t(\mathbf{W}^T \mathbf{x}^m + b_1 - y^m) \right\|_2^2 + \sum_{m=1}^M \frac{1}{2M} \left\| (1-t)(\mathbf{W}^T \mathbf{x}^m + b_2 - y^m) \right\|_2^2 + \frac{\lambda}{2} \|\mathbf{W}\|_2^2 \\
&\leq t \sum_{m=1}^M \frac{1}{2M} \left\| \mathbf{W}^T \mathbf{x}^m + b_1 - y^m \right\|_2^2 + t \frac{\lambda}{2} \|\mathbf{W}\|_2^2 \\
&\quad + (1-t) \sum_{m=1}^M \frac{1}{2M} \left\| \mathbf{W}^T \mathbf{x}^m + b_2 - y^m \right\|_2^2 + (1-t) \frac{\lambda}{2} \|\mathbf{W}\|_2^2 \\
&\leq tL(b_1) + (1-t)L(b_2)
\end{aligned} \tag{7}$$

2.1.2

Nth dimension of W , W_n , will be scaled by $\frac{1}{\alpha}$ while other dimensions will not change. It is obvious that change in one dimension is orthogonal to all other dimensions, thus no effect on

other dimensions of the weights. Scaling one dimension is essentially spread out all the data points on one dimension (if $\alpha > 1$), and the weight is subsequently scaled by $\frac{1}{\alpha}$.

If we assume shifting happens after scaling, bias will decrease by $\beta * W_n$ where W_n is the new one. This offset happens because shifting a line in one dimension is equivalent to shift the line in another dimension with the offset scaled by the gradient.

Plug it into the inner part of loss equation:

$$\begin{aligned}
\|W^{*T} X^{*(m)} + b^* - y^{(m)}\| &= \left\| \sum_{i=1, i \neq n}^N W_i X_i^{(m)} + W_n^* (\alpha * X_n^{(m)} + \beta) + (b - \beta * W_n^*) - y^{(m)} \right\| \\
&= \left\| \sum_{i=1, i \neq n}^N W_i X_i^{(m)} + W_n * X_n^{(m)} + W_n^* * \beta + b - \beta * W_n^* - y^{(m)} \right\| \quad (8) \\
&= \left\| \sum_{i=1}^N W_i X_i^{(m)} + b - y^{(m)} \right\| \\
&= \|W^{*T} X^{*(m)} + b^* - y^{(m)}\|
\end{aligned}$$

We can see the loss stays the same. We can reconfirm our choice above since in order for the loss to be minimum, the derivative of the loss need to equal to 0, which is equivalent for the inner term of the loss function to be zero.

2.1.3

Since we want to find the minimum loss, which does not depend on what method we use to find it, whether regularized or not regularized. Also because linear regression has a exact solution that achieves minimum loss, W and b should be the same as the ones in 2.1.2.

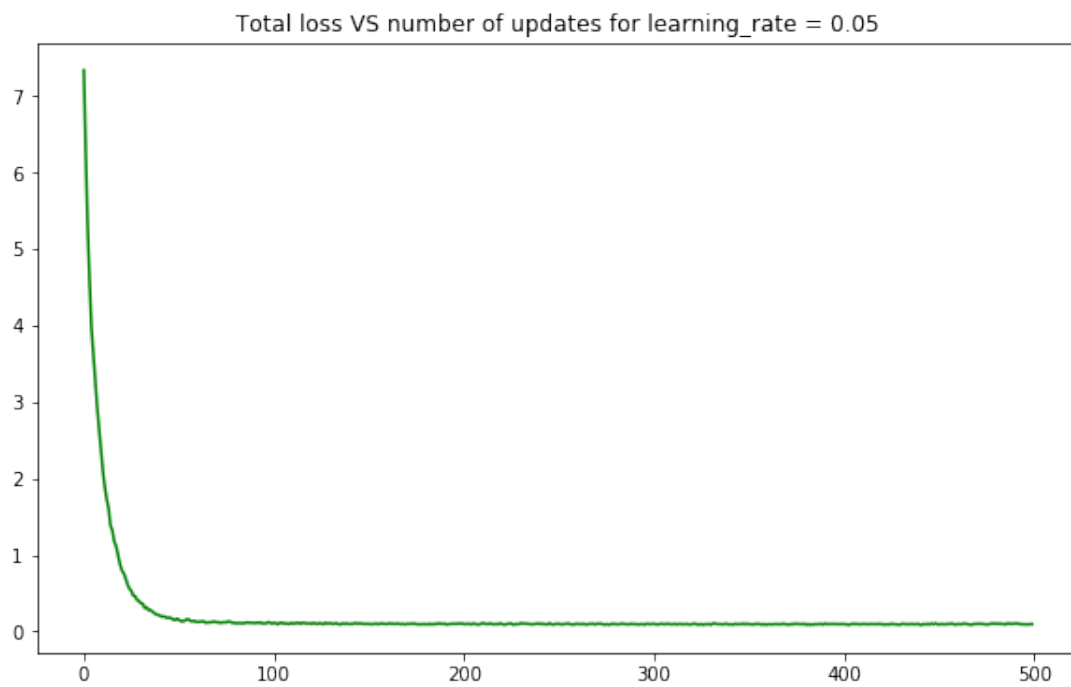
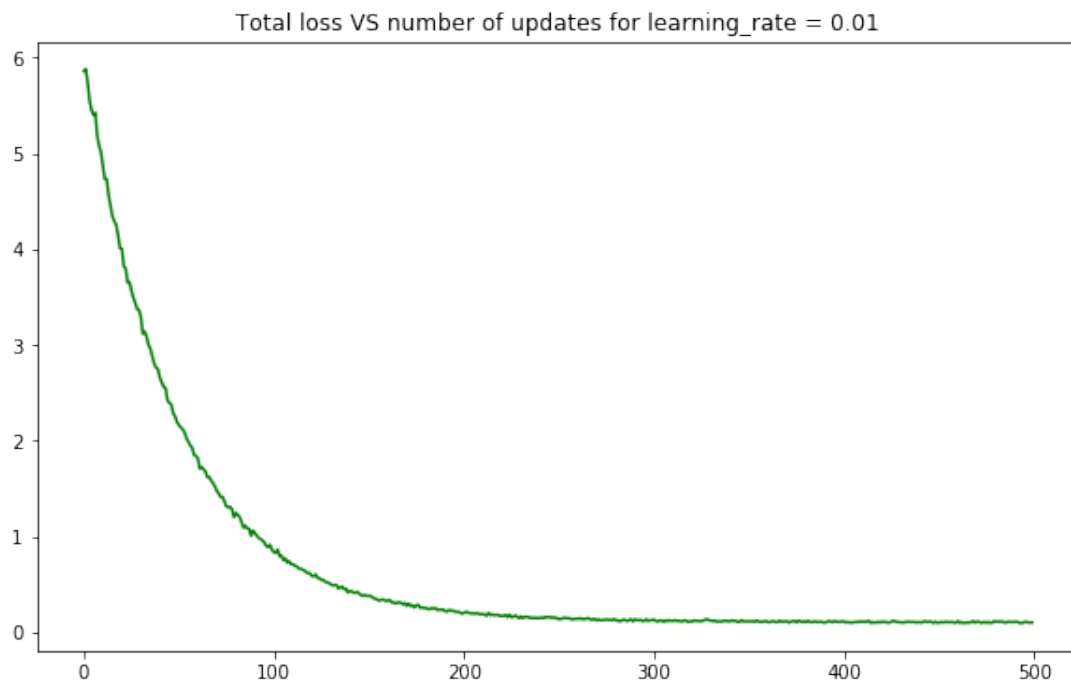
However, the loss will be smaller since one dimension of W is scaled by $\frac{1}{\alpha}$.

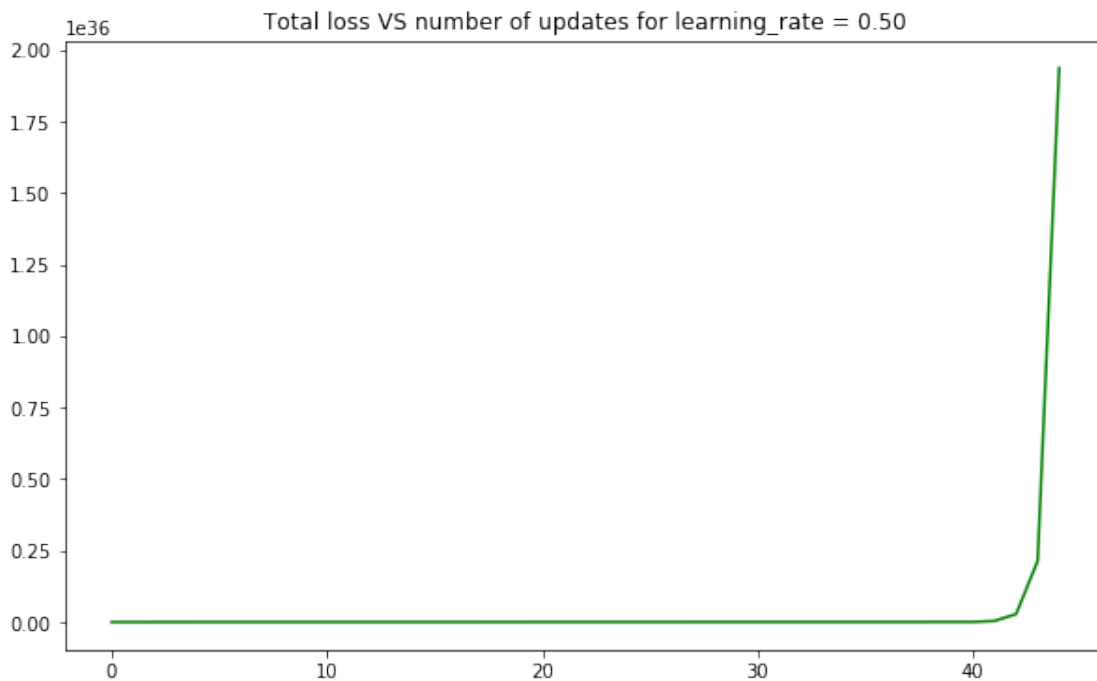
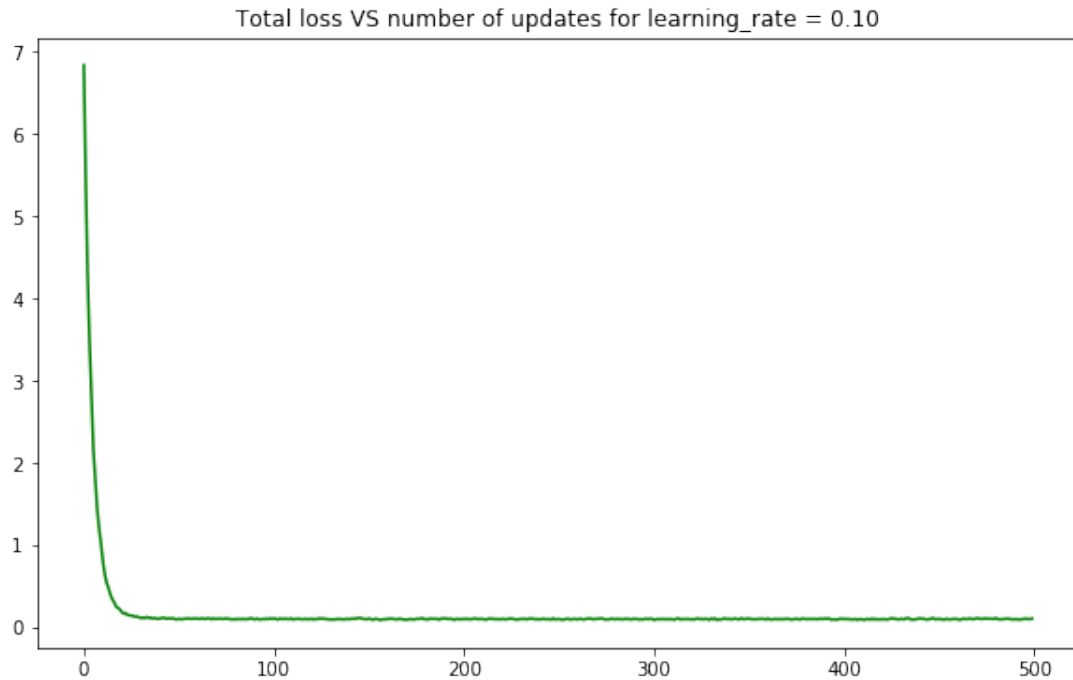
2.1.4

We can use on-hot encoding for every class, and obtain D classifiers for D classes in training set. To classify a new data point, we will feed it into all classifiers and pick the maximum score as the result.

2.2

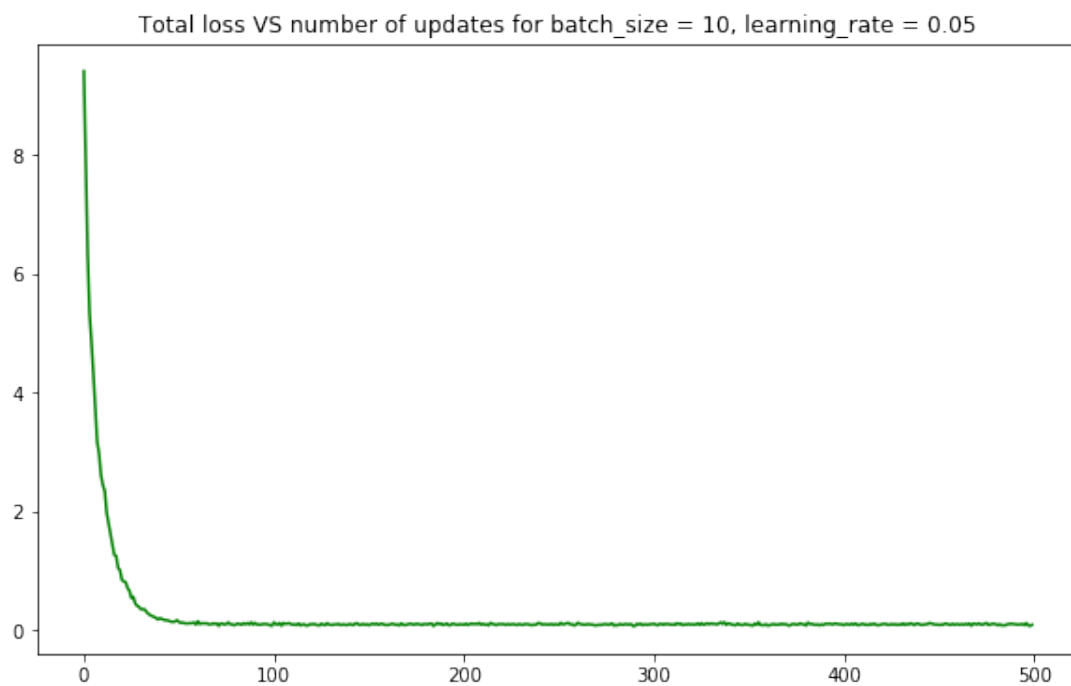
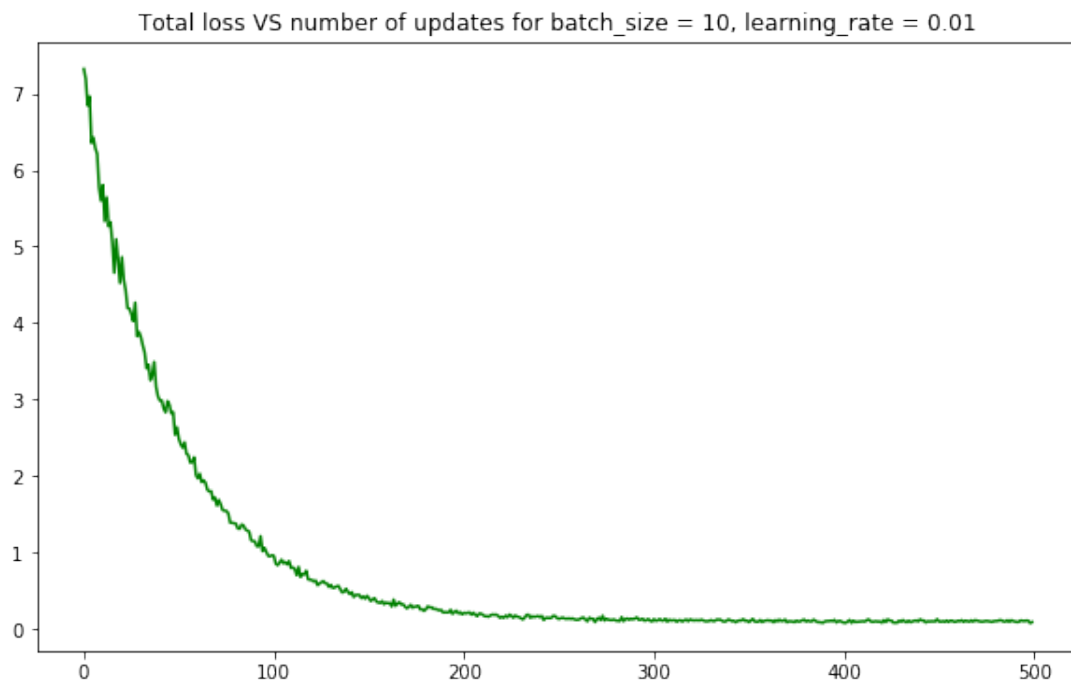
2.2.1

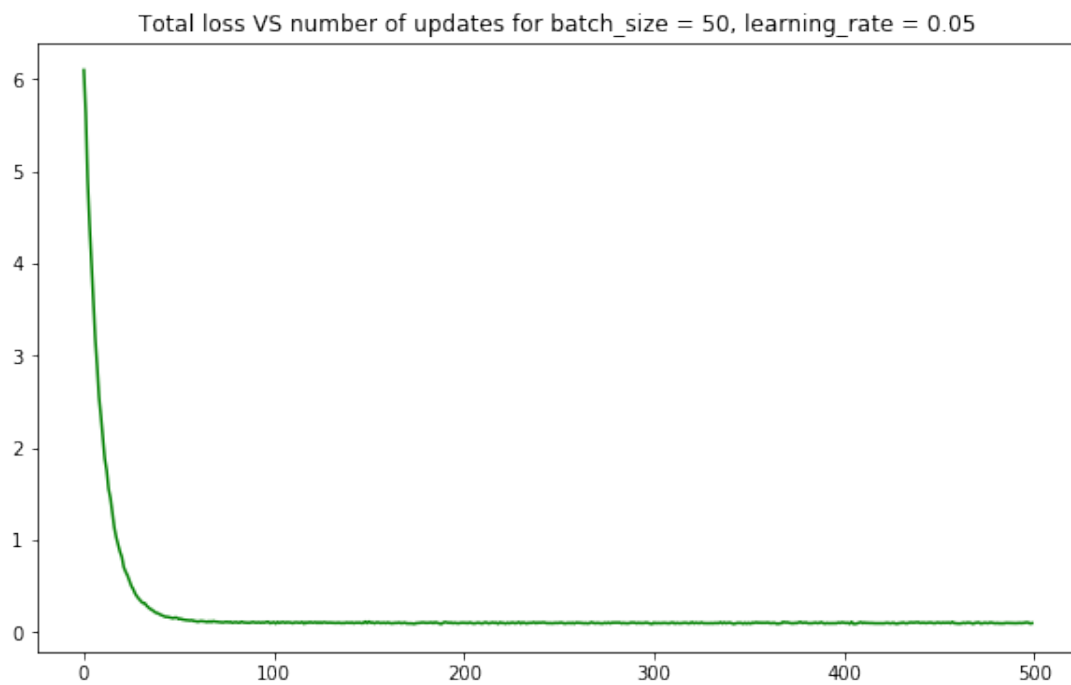
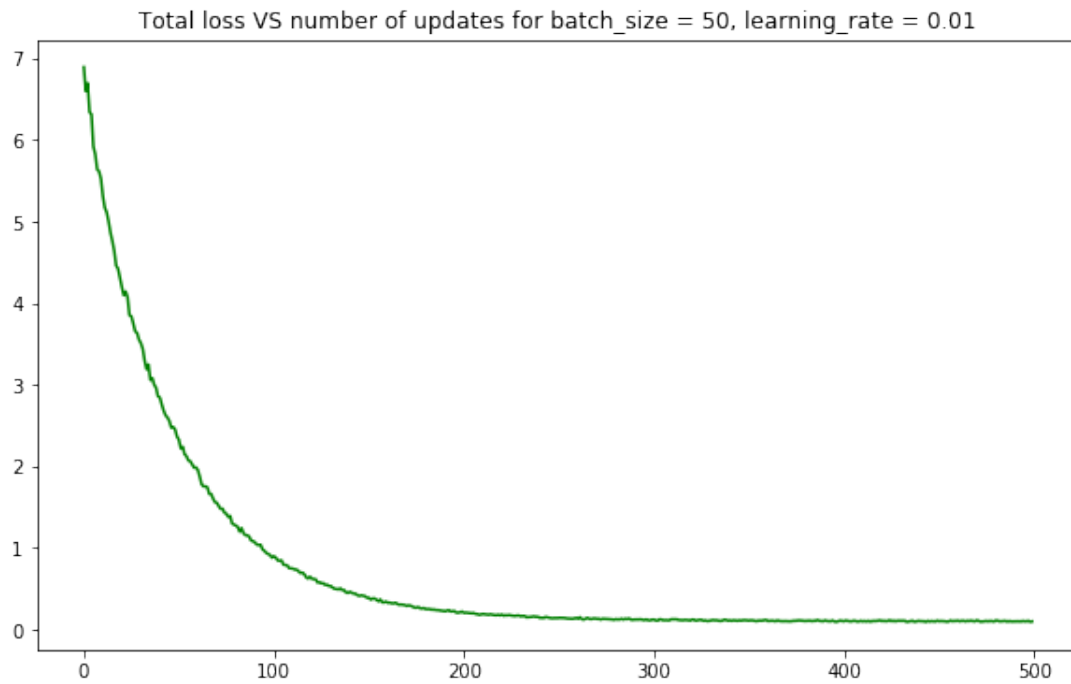


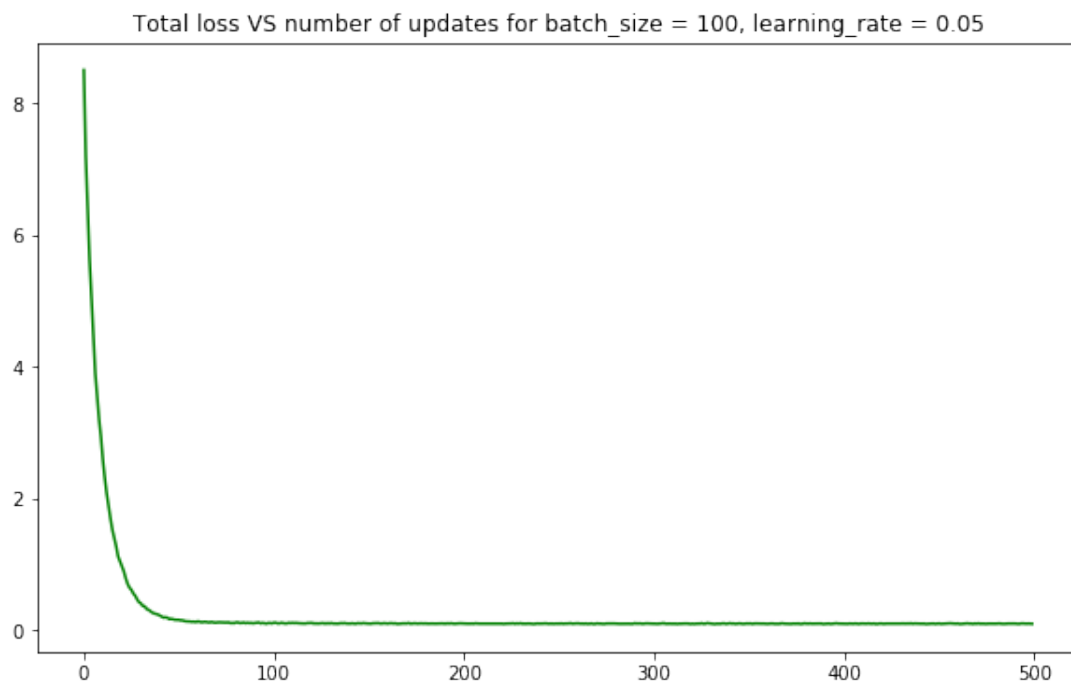
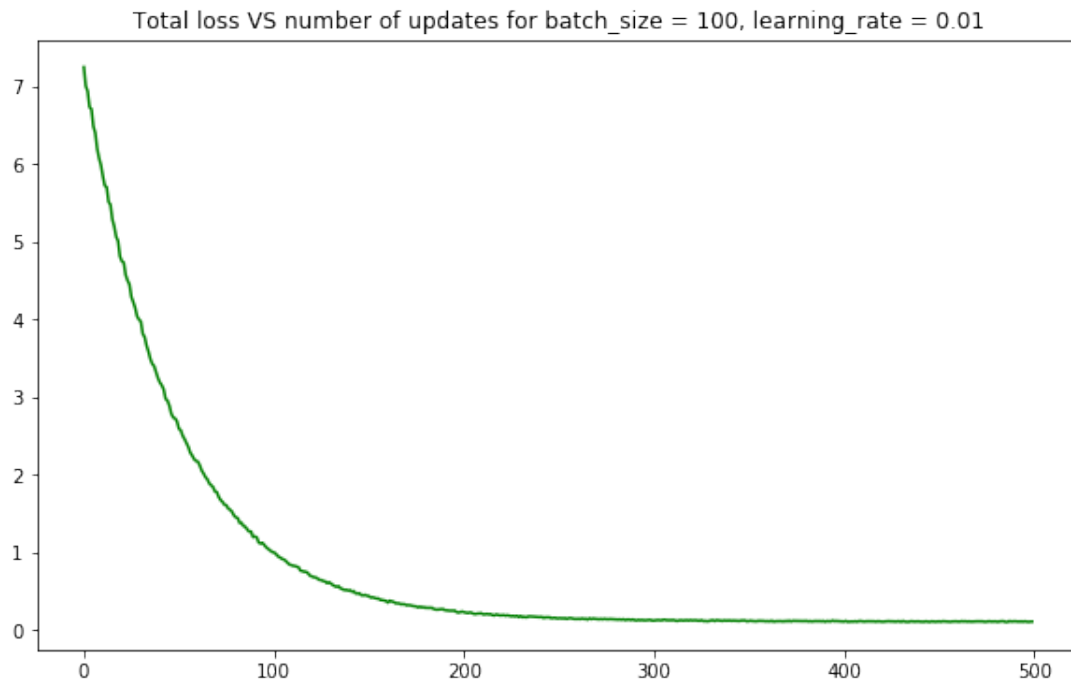


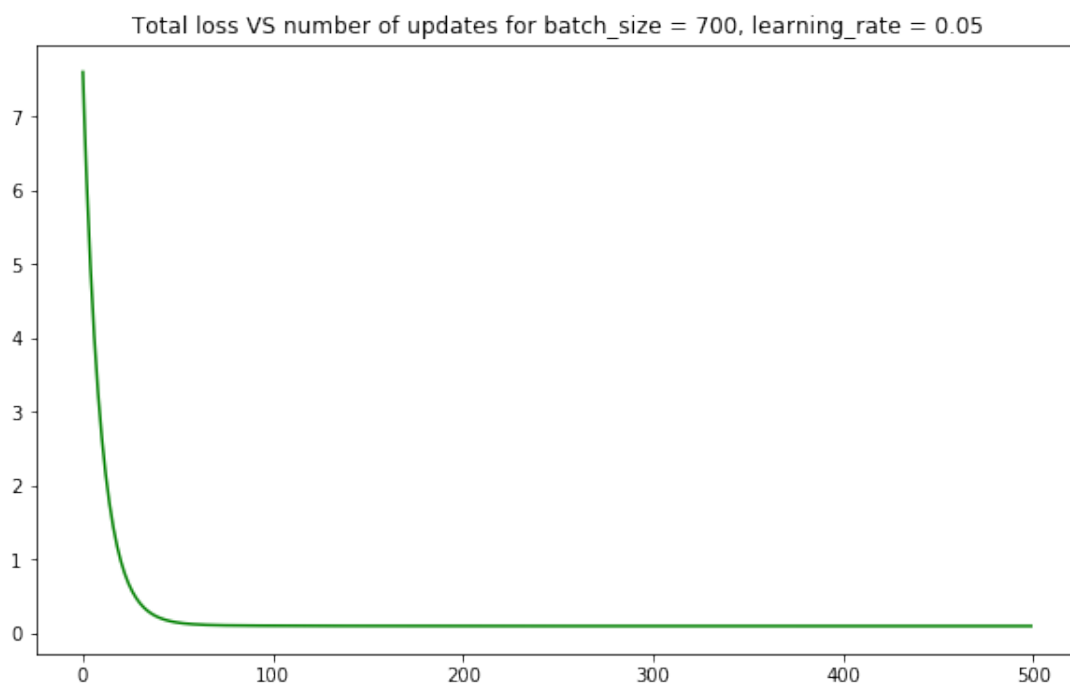
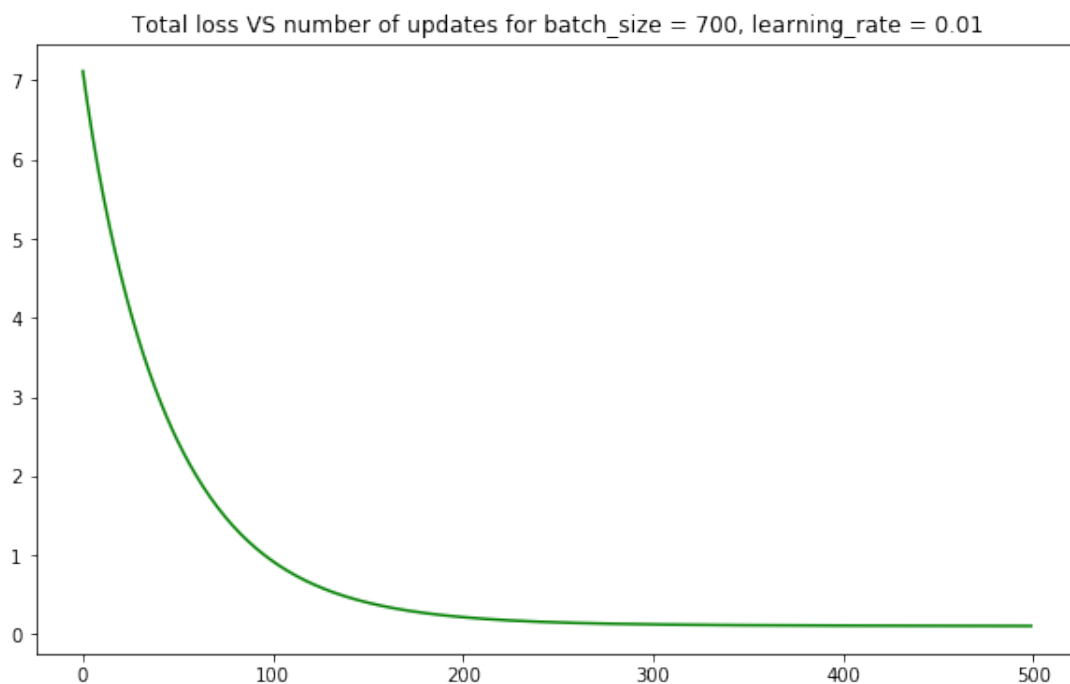
In this problem, we fix batch size to be 50 and decay rate to be 1 and try different learning rate 0.01, 0.05, 0.1, 0.5. From the plot we can see increase learning rate will speed up convergence rate, and we obtain the best learning rate 0.1. However, as learning rate increases, it might diverge like the case of learning rate 0.5.

2.2.2



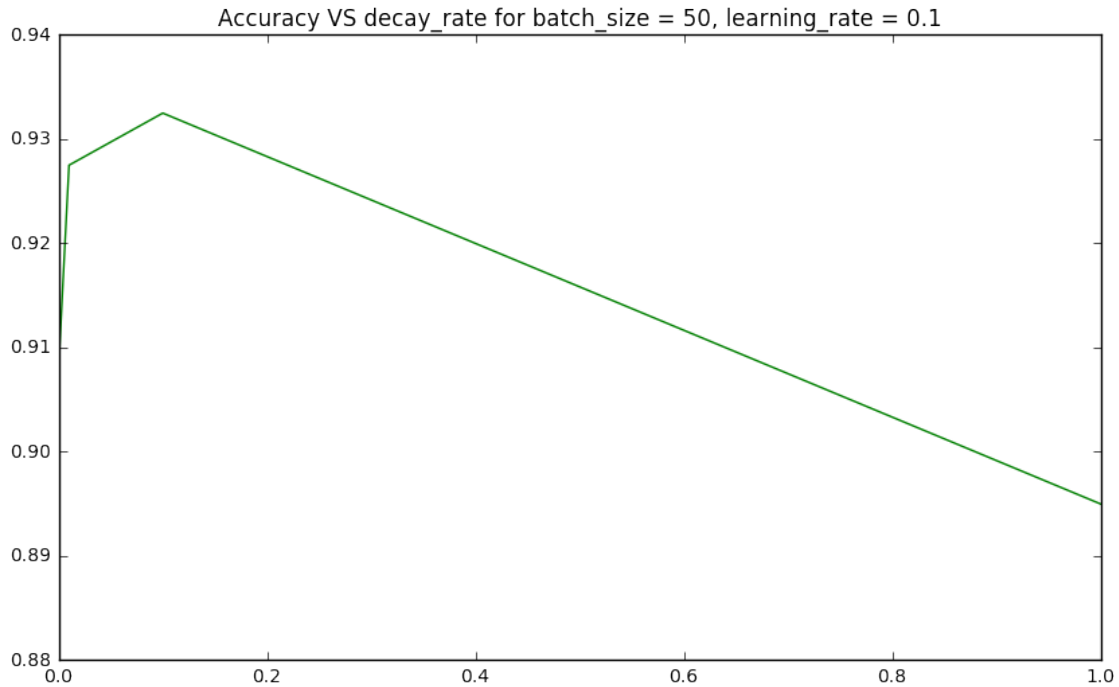






In this problem we try different batch size and tune the learning rate respectively. We observe that setting learning rate 0.1 for all different batch size obtains the best convergence time overall. And we found that batch size has linear relation with training time. Considering the trade off between training time and randomization, we think batch size 50 is a good choice.

2.2.3



From the experiment, we found that the best accuracy in valid set is 0.94 when the decay rate is 0. Although decay rate 0 obtains good accuracy in valid set, it performs bad in test set with accuracy less than 90%. It encounters overfitting problem when we do not assign any regularization. However, when decay rate is set to 1, we can see its accuracy in test set also decreases dramatically. It is a problem of under fitting. The best decay rate is about 0.1 which obtains best accuracy in test set about 93%. The reason why we use validation set to tune hyper-parameter is to make sure we do not overfit the model. If we use training set to tune hyper-parameter, the weight will be adjusted to minimize the error. However, we cannot test if the model is general enough for some data set that is not used for training. Validation set is such a data set for testing if the training is overfit. If the accuracy in training set is increasing, but the accuracy in validation set remains the set, we know it might be overfitting.

3 Code

```
1 import tensorflow as tf
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5 def pairwise_dist(x, z):
6     z = tf.transpose(z)
7     return tf.squared_difference(x, z)
8
9 def get_respon_mat(m, k):
10     # get sorted index
11     values, indices = tf.nn.top_k(-m, k, sorted=True)
12
13     # build up [[1,2], [1,3]] index form
14     indices_pair = tf.tile(tf.reshape(tf.range(0, tf.shape(m)[0]), [-1,1]), [1,k])
15
```



```

16     concated = tf.concat(2, [tf.reshape(indices_pair, [-1,k,1]),
17                             tf.reshape(indices, [-1,k,1])])
18
19     concated = tf.reshape(concated, [-1,2])
20
21     # return dense matrix
22     value = 1.0 / k
23     res = tf.sparse_to_dense(sparse_indices=concated, output_shape=[tf.shape(m)[0],
24                             tf.shape(m)[1]], sparse_values=value, validate_indices=False)
25     return res
26
27 # testing for get_respon_mat function
28 train = tf.placeholder(tf.float32)
29 X = tf.placeholder(tf.float32)
30 target = tf.placeholder(tf.float32)
31 test_input = np.linspace(0.0, 11.0, num = 1000)
32 test_input = np.expand_dims(test_input, axis=0)
33 test_input = np.transpose(test_input)
34 r = get_respon_mat(pairwise_dist(X, train), 2)
35 y_hat = tf.matmul(tf.transpose(target), tf.transpose(r))
36 sess.run(init)
37 predict = sess.run(y_hat, feed_dict={X: test_input, train: trainData, target:
    trainTarget})
38
39 plt.plot(trainData, trainTarget, 'bo', test_input, np.transpose(predict), 'g-')
40 plt.axis([0,11, -2, 12])
41 plt.title("KNN for k = 2")
42 plt.show()

```

Listing 1: 1.3.1

```

43 import tensorflow as tf
44 import matplotlib.pyplot as plt
45 import numpy as np
46 kList = [1,3,5,50]
47 train = tf.placeholder(tf.float32)
48 X = tf.placeholder(tf.float32)
49 Y = tf.placeholder(tf.float32)
50 target = tf.placeholder(tf.float32)
51 test_input = np.linspace(0.0, 11.0, num = 1000)
52 test_input = np.expand_dims(test_input, axis=0)
53 test_input = np.transpose(test_input)
54 plot_input = tf.constant(test_input)
55 best_k = 0
56 min_err = float("inf")
57 sess.run(init)
58 for k in kList:
59     # Error definition
60     r = get_respon_mat(pairwise_dist(X, train), k)
61     y_hat = tf.matmul(tf.transpose(target), tf.transpose(r))
62     meanSquaredError = tf.reduce_mean(tf.reduce_sum(tf.square(y_hat -
63         tf.transpose(Y)), reduction_indices=[1]))
64     trainMSE = sess.run(meanSquaredError, feed_dict={X: trainData, Y: trainTarget,
        train: trainData, target: trainTarget})
65     validMSE = sess.run(meanSquaredError, feed_dict={X: validData, Y: validTarget,
        train: trainData, target: trainTarget})
66     testMSE = sess.run(meanSquaredError, feed_dict={X: testData, Y: testTarget,
        train: trainData, target: trainTarget})
67     plot_predict = sess.run(y_hat, feed_dict={X: test_input, train: trainData,
        target: trainTarget})
68     if validMSE < min_err:
69         min_err = validMSE
70         best_k = k
71     print("k: %2d, trainMSE: %0.2f, validMSE: %0.2f, testMSE: %0.2f"%(k, trainMSE,
        validMSE, testMSE))
72     plt.plot(trainData, trainTarget, 'bo', test_input, np.transpose(plot_predict),
73         'g-')
74     plt.axis([0,11, -2, 12])
75
76

```

```

77     plt.title("KNN for k = %d"%(k))
78     plt.show()
79     print("Best k value based on validation err is: %d"%(best_k))

```

Listing 2: 1.3.2

```

80 import tensorflow as tf
81 import matplotlib.pyplot as plt
82 import numpy as np
83
84 def get_soft_knn_res_mat(x, z, lamb):
85     z = tf.transpose(z)
86     Kxx = tf.exp(tf.cast(tf.squared_difference(x, z), tf.float32) * -lamb)
87     K_norm = tf.expand_dims(tf.reduce_sum(Kxx, reduction_indices=1), 1)
88     return tf.div(Kxx, K_norm)
89
90 def get_Gaussian_knn_res_mat(x, z, lamb):
91     z = tf.transpose(z)
92     KxX = tf.exp(tf.cast(tf.squared_difference(x, z), tf.float32) * -lamb)
93     KXX = tf.exp(tf.cast(tf.squared_difference(x, tf.transpose(x)), tf.float32)
94                 * -lamb)
95     KXX_inverse = tf.matrix_inverse(KXX)
96     return tf.matmul(KXX_inverse, KxX)
97
98
99 myLambda = 100
100 train = tf.placeholder(tf.float32)
101 X = tf.placeholder(tf.float32)
102 target = tf.placeholder(tf.float32)
103 test_input = np.linspace(0.0, 11.0, num = 1000)
104 test_input = np.expand_dims(test_input, axis=0)
105 test_input = np.transpose(test_input)
106
107 # r = get_Gaussian_knn_res_mat(train, X, myLambda)
108 r = get_soft_knn_res_mat(X, train, myLambda)
109
110 y_hat = tf.matmul(tf.transpose(target), tf.transpose(r))
111 sess.run(init)
112 predict = sess.run(y_hat, feed_dict={X: test_input, train: trainData, target:
    trainTarget})
113
114 plt.plot(trainData, trainTarget, 'bo', test_input, np.transpose(predict), 'g-')
115 plt.axis([0,11, -2, 12])
116 plt.title("Soft KNN")
117 plt.show()

```

Listing 3: 1.4.1

```

118 import tensorflow as tf
119 import matplotlib.pyplot as plt
120 import numpy as np
121
122 with np.load("../data/tinymnist.npz") as data :
123     trainData, trainTarget = data ["x"], data["y"]
124     validData, validTarget = data ["x_valid"], data ["y_valid"]
125     testData, testTarget = data ["x_test"], data ["y_test"]
126
127 def buildGraph(decay_rate, learning_rate):
128     # Variable creation
129     W = tf.Variable(tf.truncated_normal(shape=[64,1], stddev=0.5), name='weights')
130     b = tf.Variable(0.0, name='biases')
131     X = tf.placeholder(tf.float32, [None, 64], name='input_x')
132     y_target = tf.placeholder(tf.float32, [None,1], name='target_y')
133
134     # Graph definition
135     y_predicted = tf.matmul(X,W) + b
136
137     # Error definition

```

```

138     error = 0.5 *
139         tf.reduce_mean(tf.cast(tf.reduce_sum(tf.cast(tf.square(y_predicted -
140             y_target), tf.float32), reduction_indices=1,
141             name='squared_error'),
142             tf.float32), name='mean_squared_error')
143     + 0.5 * decay_rate * tf.reduce_sum(tf.cast(tf.square(W), tf.float32))
144
145     # Training mechanism
146     optimizer = tf.train.GradientDescentOptimizer(learning_rate = learning_rate)
147     train = optimizer.minimize(loss=error)
148     return W, b, X, y_target, y_predicted, error, train
149
150 def getRandomBatch(trainData, trainTarget, size):
151     idx = np.random.choice(trainData.shape[0], size, replace=False)
152     return trainData[idx,:], trainTarget[idx,:]

```

Listing 4: 2.2

```

154 #2.1 fix decay_rate = 1 and tune learning_rate
155 decay_rate = 1
156 batch_size = 50
157 learning_rate_list = [0.01, 0.05, 0.1, 0.5]
158 for learning_rate in learning_rate_list:
159     W, b, X, y_target, y_predicted, error, train = buildGraph(decay_rate,
160         learning_rate)
161     init = tf.global_variables_initializer()
162     sess = tf.InteractiveSession()
163     sess.run(init)
164     loss_recorder = np.array([])
165     for itr in range(500):
166         batch_xs, batch_ys = getRandomBatch(trainData, trainTarget, batch_size)
167         loss, _ = sess.run([error, train], feed_dict={X: batch_xs, y_target:
168             batch_ys})
169         loss_recorder = np.append(loss_recorder, loss)
170     plt.plot(np.arange(500), loss_recorder, 'g')
171     #plt.axis([0,2000, 0, 2])
172     plt.title("Total loss VS number of updates for learning_rate = %0.2f"%(
173         learning_rate))
174     plt.show()

```

Listing 5: 2.2.1

```

172 #2.2 fix decay_rate = 1 and tune learning_rate and batch_size
173 decay_rate = 1
174 batch_size_list = [10,50,100,700]
175 learning_rate_list = [0.01, 0.05, 0.1, 0.5]
176 for batch_size in batch_size_list:
177     for learning_rate in learning_rate_list:
178         W, b, X, y_target, y_predicted, error, train = buildGraph(decay_rate,
179             learning_rate)
180         init = tf.global_variables_initializer()
181         sess = tf.InteractiveSession()
182         sess.run(init)
183         loss_recorder = np.array([])
184         for itr in range(500):
185             batch_xs, batch_ys = getRandomBatch(trainData, trainTarget, batch_size)
186             loss, _ = sess.run([error, train], feed_dict={X: batch_xs, y_target:
187                 batch_ys})
188             loss_recorder = np.append(loss_recorder, loss)
189             plt.plot(np.arange(500), loss_recorder, 'g')
190             #plt.axis([0,2000, 0, 2])
191             plt.title("Total loss VS number of updates for batch_size = %d,
192                 learning_rate = %0.2f"%(batch_size, learning_rate))
193             plt.show()

```

Listing 6: 2.2.2

```

191 #2.3 fix batch_size = 50 and learning_rate = 0.1 and tune decay_rate
192 decay_rate_list = [0., 0.0001, 0.001, 0.01, 0.1, 1.]
193 batch_size = 50
194 learning_rate = 0.1
195 best_decay_rate = 0
196 best_acc = 0
197 test_acc_recorder = np.array([])
198 for decay_rate in decay_rate_list:
199     W, b, X, y_target, y_predicted, error, train = buildGraph(decay_rate,
200     learning_rate)
201     correct_prediction = tf.equal(y_target, (tf.sign(y_predicted - 0.5) + 1)/2)
202     accuracy = tf.reduce_mean(tf.cast(correct_prediction, tf.float32))
203     init = tf.global_variables_initializer()
204     sess = tf.InteractiveSession()
205     sess.run(init)
206     for itr in range(500):
207         batch_xs, batch_ys = getRandomBatch(trainData, trainTarget, batch_size)
208         sess.run(train, feed_dict={X: batch_xs, y_target: batch_ys})
209     # use valid data set accuracy to tune decay rate
210     cur_acc = sess.run(accuracy, feed_dict={X: validData, y_target: validTarget})
211     if cur_acc > best_acc:
212         best_acc = cur_acc
213         best_decay_rate = decay_rate
214     # use test set to test accuracy
215     test_acc_recorder = np.append(test_acc_recorder, sess.run(accuracy, \
216     feed_dict={X:
217     testData, y_target: testTarget}))
218 print("Best accuracy in valid set is: %0.2f, decay_rate is: %0.2f"%(best_acc,
219     best_decay_rate))
220 plt.plot(decay_rate_list, test_acc_recorder, 'g')
221 #plt.axis([0,2000, 0, 2])
222 plt.title("Accuracy VS decay_rate for batch_size = 50, learning_rate = 0.1")
223 plt.show()

```

Listing 7: 2.2.3