

Exercise 1

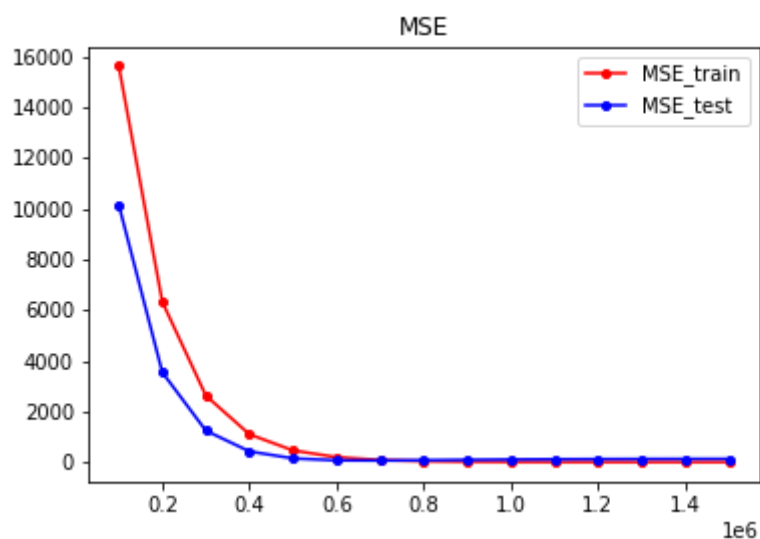
(a)

该数据集中一共有两个自变量，因此需要训练的参数有三个，此处设它们分别为 $[w_1, w_2, w_3]$ ，其中 w_1 为第一个自变量的权重， w_2 为第二个自变量的权重， w_3 为常量（偏置）。因为使用梯度下降法更新参数，需要手动设置的超参数有每次更新时的步幅大小，此处设为 α 。

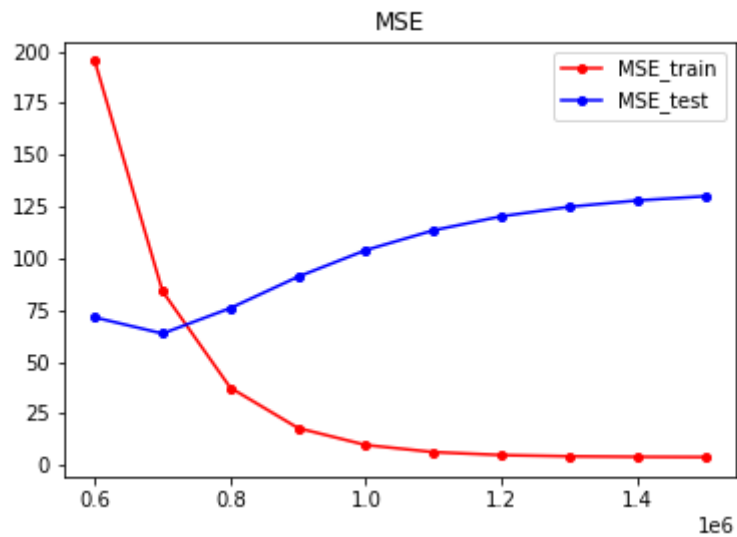
实现细节：

开始的时候，我未对数据进行预处理，直接进行线性回归，预测房价，后来计算出来的 MSE 过大，超出了 python float 能表示的精度范围，对参数 w 进行输出也发现 w 的数值很大。于是我对数据进行了输出，发现第一个自变量（房的面积）和第二个自变量（对应房子距离双鸭山职业技术学院的距离）的数值不在一个数量级上，这可能是导致 MSE 和 w 过大的原因，于是我对数据进行了归一化处理。处理后的结果是，不再出现以上现象，MSE 和 w 也在迭代过程中逐渐收敛。

实验结果：



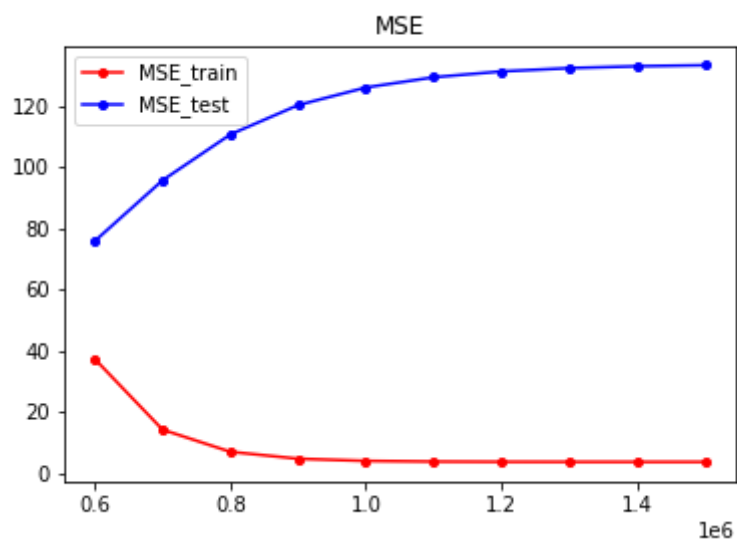
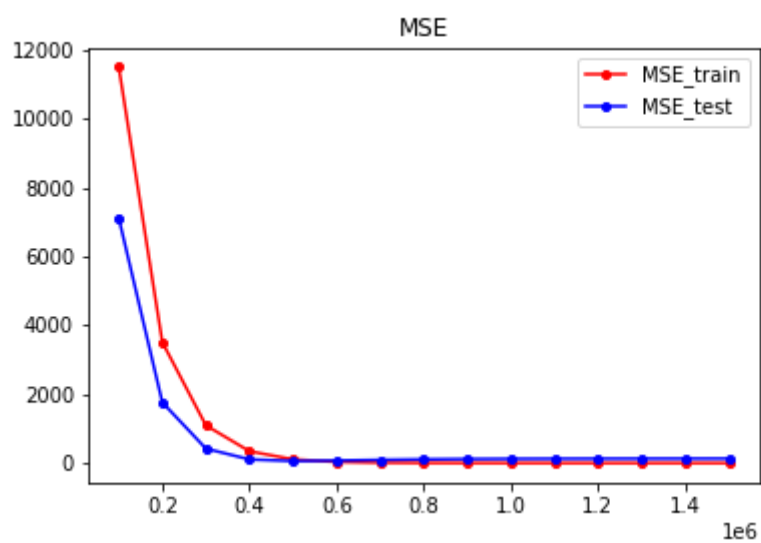
可以看到在 600000 次迭代前，训练集和测试集上的 MSE 都是下降的，但是训练集的 MSE 大于测试集的 MSE。600000 次迭代后，两个 MSE 的范围接近，于是对大于 600000 次迭代后得到的 MSE 进行二次绘图，结果如下：



可以看到继续训练训练集上的 MSE 继续下降，但是约在 700000 次迭代后，测试集上的 MSE 不降反升，说明此时模型已经出现了过拟合现象。

(b)

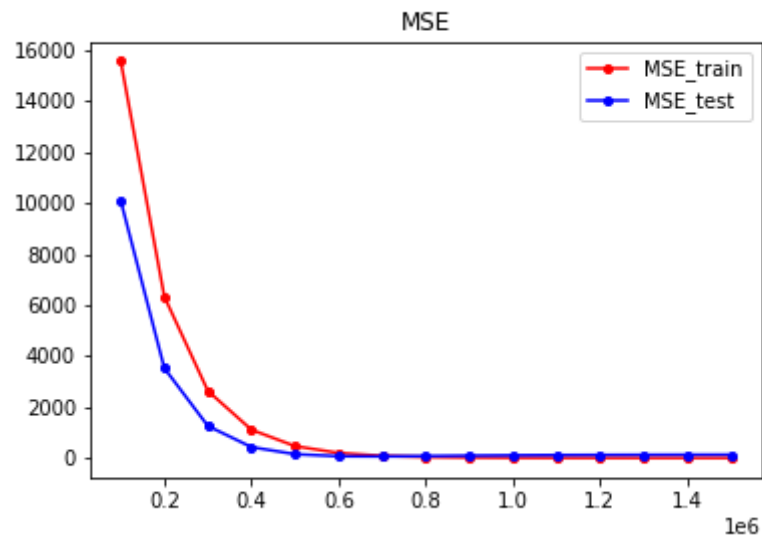
将学习率由 (a) 中的 0.00015 改为 0.0002 后可以发现



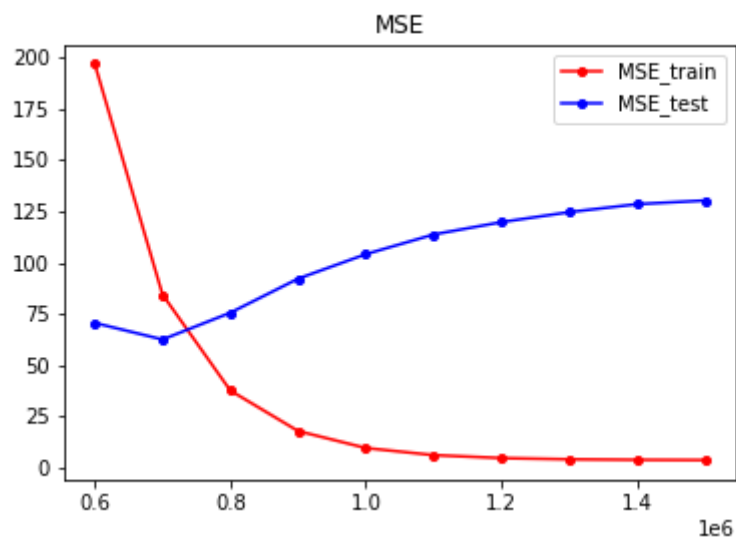
在前 500000 次迭代中训练集和测试集的 MSE 都在变小，在 500000 次迭代后训练集的 MSE 继续变小，但是测试集的 MSE 开始变大。并且可以发现，学习率为 0.00015 的 MSE 比学习率为 0.0002 的 MSE 下降得快，说明学习率的选取会影响模型的最终表现。

(c)

此处采用了随机梯度下降的方法。我可能可以用梯度下降法得到最优结果，但是不能保证。采用随机梯度下降法，查看训练集和测试集上的 MSE。学习率为 0.00015。迭代次数为 1500000。



可以看到在 600000 次迭代前，训练集和测试集上的 MSE 都是下降的，但是训练集的 MSE 大于测试集的 MSE。600000 次迭代后，两个 MSE 的范围接近，于是对大于 600000 次迭代后得到的 MSE 进行二次绘图，结果如下：



可以看到继续训练训练集上的 MSE 继续下降，但是约在 700000 次迭代后，测试集上的 MSE 不降反升，说明此时模型已经出现了过拟合现象。

对比同样参数下的梯度下降法，即题 (a) 中数据，可以发现两者相差甚微，但是从更新时间上看，随机梯度下降比梯度下降快，因此随机梯度下降是梯度下降的一个不错的替代方法。但两者均无法保证可以取得全局最优解。

代码细节:

读取数据的代码如下:

```
1 def openreadtxt(file_name):
2     data = []
3     file = open(file_name, 'r') #打开文件
4     file_data = file.readlines() #读取所有行
5     for row in file_data:
6         tmp_list = row.split(' ') #按' '切分每行的数据
7         tmp_list[-1] = tmp_list[-1].replace('\n', '') #去掉换行符
8         for i in range(len(tmp_list)): #将字符串转化为浮点数
9             tmp_list[i] = float(tmp_list[i])
10        data.append(tmp_list) #将每行数据插入data中
11    return data
```

对数据做归一化处理函数实现如下:

```
1 #归一化处理
2 def autoNorm(data, mins, maxs): #传入一个矩阵
3     ranges = maxs - mins #最大值列表 - 最小值列表 = 差值列表
4     normData = np.zeros(np.shape(data)) #生成一个与 data矩阵同规格的normData
5     #全0矩阵, 用于装归一化后的数据
6     row = data.shape[0] #返回 data矩阵的行数
7     normData = data - np.tile(mins, (row, 1)) #data矩阵每一列数据都减去每一列的最小值
8     normData = normData / np.tile(ranges, (row, 1)) #data矩阵每一列数据都除去每一
9     #列的差值 (差值 = 某列的最大值 - 某列最小值)
10    return normData
```

归一化的公式是 $data' = \frac{data - data_{min}}{data_{max} - data_{min}}$. 测试集上的归一化处理使用的是训练集上的最小值和最大值。

对每个输入自变量x, 预测 y 的函数实现如下:

```
1 def pred(x, w):
2     pred_y = 0
3     for i in range(len(w)):
4         if(i == len(w)-1): #w[i]是偏置
5             pred_y += w[i]
6         else:
7             pred_y += (x[i] * w[i])
8     #print(pred_y)
9     return pred_y
```

对输入一系列 X, 预测一系列 Y 的函数实现如下:

```
1 def pred_Y(X, w):
2     Y_pred = []
3     for i in range(len(X)):
4         y = pred(X[i], w)
5         Y_pred.append(y)
6     return Y_pred
```

随机梯度下降更新参数实现如下:

```

1 def SGD(Y, X, w, alpha):
2     '''
3     w: 参数
4     alpha: 学习率
5     '''
6     i = randint(0, len(Y)-1)
7     for j in range(len(w)):
8         if( j == len(w) - 1 ):
9             w[j] = w[j] + alpha * (Y[i] - pred(X[i],w))
10        else:
11            w[j] = w[j] + alpha * (Y[i] - pred(X[i],w)) * X[i][j]

```

梯度下降实现如下:

```

1 def GD(Y, X, w, alpha):
2     sum0 = 0
3     sum1 = 0
4     sum2 = 0
5     for i in range(len(X)):
6         y_pred = pred(X[i],w)
7         sum0 += (Y[i] - y_pred) * X[i][0]
8         sum1 += (Y[i] - y_pred) * X[i][1]
9         sum2 += (Y[i] - y_pred)
10    w[0] += alpha * (sum0 / len(X))
11    w[1] += alpha * (sum1 / len(X))
12    w[2] += alpha * (sum2 / len(X))

```

计算均方误差函数实现如下:

```

1 def MSE(Y_pred, Y):
2     MSE = 0
3     #for i in range(len(Y)):
4     #     #print(round((Y_pred[i] - Y[i]),2))
5     #     tmp = pow((Y_pred[i] - Y[i]),2)
6     #     MSE += tmp
7     #print(MSE)
8     #MSE = MSE/len(Y)
9     MSE = metrics.mean_squared_error(np.array(Y), np.array(Y_pred))
10    return MSE

```

绘图代码如下:

```

1 xx =
    [100000,200000,300000,400000,500000,600000,700000,800000,900000,1000000,11000
    00,1200000,1300000,1400000,1500000]
2 plt.plot(xx, MSE_train,color = 'red',marker='o',markersize=4,
    label='MSE_train')
3 plt.plot(xx,MSE_test, color = 'blue',marker='o',markersize=4, label=
    'MSE_test')
4 plt.legend()
5 plt.title("MSE")
6 plt.savefig('MSE.png')

```

2、

(a)

因为是二分类问题，可使用 sigmoid 函数作为逻辑函数，表示样本属于某个类别的概率。

$$\text{Sigmoid}(y) = \frac{1}{1 + e^{-y}}$$

用 $\text{Sigmoid}(y)$ 表示样本属于正类的概率，用 $1 - \text{Sigmoid}(y)$ 表示样本属于负类的概率。

条件似然函数可写成：

$$\begin{aligned} l(w) &= \sum_l y^l \ln P(y^l = 1 | x^l, w) + (1 - y^l) \ln P(y^l = 0 | x^l, w) \\ &= \sum_l y^l (w_0 + \sum_{i=1}^p w_i x_i^l) - \ln(1 + \exp(w_0 + \sum_{i=1}^p w_i x_i^l)) \end{aligned}$$

其中

$$P(y^l | x^l, w) = \text{Sigmoid}(\sum_{i=1}^p w_i x_i^l)$$

(b)

$$\begin{aligned} \frac{\partial}{\partial w_i} l(w) &= \sum_l y^l x_i^l - \frac{\sum_{i=1}^p x_i^l}{1 + \exp(w_0 + \sum_{i=1}^p w_i x_i^l)} \\ &= \sum_l x_i^l (y^l - \hat{P}(y^l = 1 | x^l, w)) \end{aligned}$$

同理，对 w_0 求导可得

$$\frac{\partial}{\partial w_0} l(w) = \sum_l (y^l - \hat{P}(y^l = 1 | x^l, w))$$

(c)

使用梯度下降法更新参数。

$$w_i = w_i + \alpha \sum_l x_i^l (y^l - \hat{P}(y^l = 1 | x^l, w))$$

其中 α 为学习率。

代码实现如下：

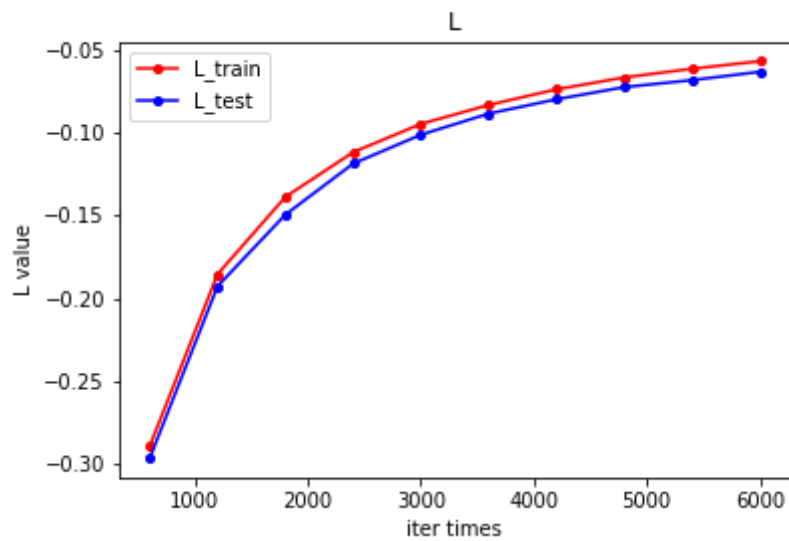
```
1  def GD(Y, x, w, alpha):
2      sum0 = 0
3      sum1 = 0
4      sum2 = 0
5      sum3 = 0
6      sum4 = 0
7      sum5 = 0
8      sum6 = 0
9      for i in range(len(x)):
10         y_pred = pred(x[i], w)
11         sum0 += (Y[i] - y_pred) * x[i][0]
12         sum1 += (Y[i] - y_pred) * x[i][1]
13         sum2 += (Y[i] - y_pred) * x[i][2]
```

```

14     sum3 += (Y[i] - y_pred) * X[i][3]
15     sum4 += (Y[i] - y_pred) * X[i][4]
16     sum5 += (Y[i] - y_pred) * X[i][5]
17     sum6 += (Y[i] - y_pred)
18     w[0] += alpha * (sum0 )
19     w[1] += alpha * (sum1 )
20     w[2] += alpha * (sum2 )
21     w[3] += alpha * (sum3 )
22     w[4] += alpha * (sum4 )
23     w[5] += alpha * (sum5 )
24     w[6] += alpha * (sum6 )

```

实验结果：



将迭代次数调至 6001，学习率调为 0.02，将结果绘制成图，其中纵坐标为 (a) 中的条件对数似然函数，可以看到随着训练次数增加，L 值在变大。

(d)

将上题的参数视为最优，在测试集上进行预测和对结果进行统计。

实现代码如下：

```

1  #在测试集上统计结果
2  test_pred = pred_Y(X_test,w)
3  sum = 0
4  for i in range(len(test_pred)):
5      tmp = 0
6      if(test_pred[i] > 0.5):
7          tmp = 1
8      sum += abs(tmp - Y_test[i])
9  print(sum)

```

结果显示为 0，即没有样本被分错类。

(e)

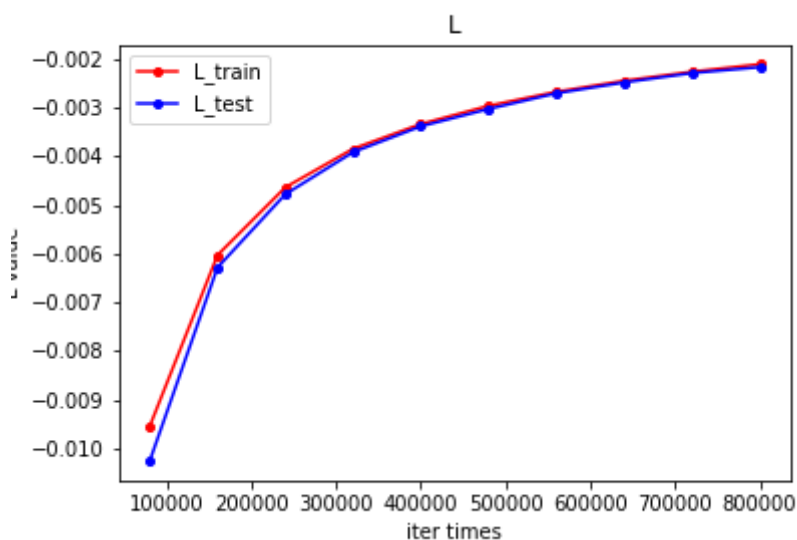
使用梯度下降法更新参数。

实现代码如下：

```
1 def SGD(Y, x, w, alpha):
2     '''
3     w: 参数
4     alpha: 学习率
5     '''
6     i = randint(0, len(Y)-1)
7     for j in range(len(w)):
8         if( j == len(w) - 1 ):
9             w[j] = w[j] + alpha * (Y[i] - sigmoid(pred(x[i],w)))
10        else:
11            w[j] = w[j] + alpha * (Y[i] - sigmoid(pred(x[i],w))) * x[i][j]
```

经过多次测试，我将学习率 α 设置为 0.02，将迭代次数设置为 500001。

实验结果如下图：



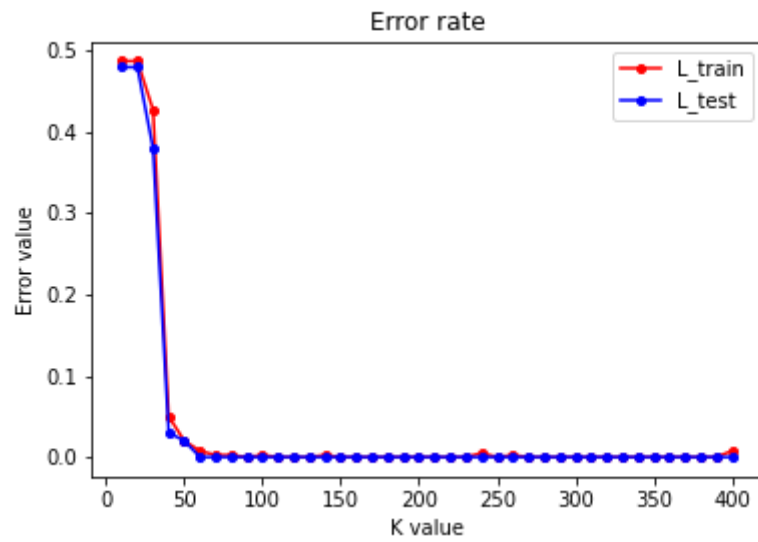
可以看到训练集和测试集上结果都变好，训练集上的结果略好于测试集上的结果。迭代收敛需要400000次左右。

(f)

错误率的计算函数：

```
1 def error(y_pred,y):
2     sum = 0
3     for i in range(len(y_pred)):
4         tmp = 0
5         if(y_pred[i] > 0.5):
6             tmp = 1
7         sum += abs(tmp - y[i])
8     return sum / len(y)
```

设置训练迭代次数为 100 次，学习率为 0.02，实验结果如下：



可以看到训练集和测试集上的错误率变化比较一致，当 k 小于 50 时，有较高的错误率，当 k 大于 50 后，基本能预测准确。

注：第2题的代码详细见 `Exercise 2.ipynb`。