

## Лабораторная работа №3 “Переобучение и регуляризация”

Долматович Алина, 858641

Набор данных ex3data1.mat представляет собой файл формата \*.mat (т.е. сохраненного из Matlab). Набор содержит две переменные X (изменения уровня воды) и y (объем воды, вытекающий из дамбы). По переменной X необходимо предсказать y. Данные разделены на три выборки: обучающая выборка (X, y), по которой определяются параметры модели; валидационная выборка (Xval, yval), на которой настраивается коэффициент регуляризации; контрольная выборка (Xtest, ytest), на которой оценивается качество построенной модели.

```
In [1]: 1 import pandas
        2 from scipy.io import loadmat
        3 import matplotlib.pyplot as pyplot
        4 import numpy as np
        5 import scipy.optimize as optimize
```

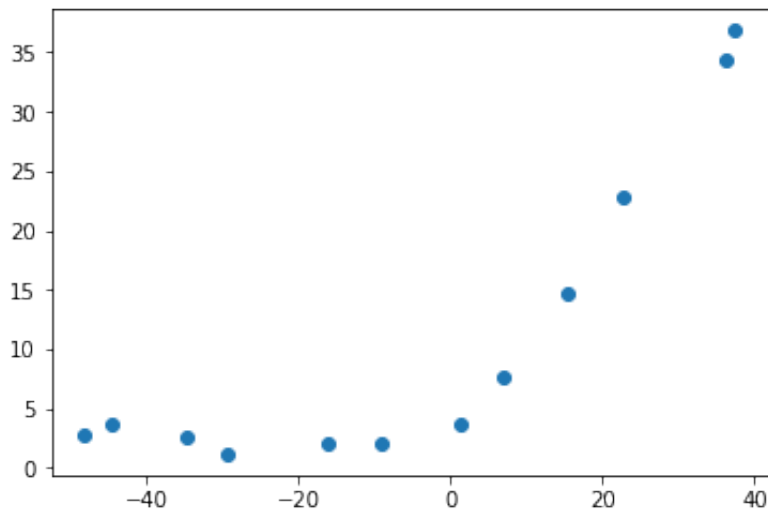
Загрузите данные ex3data1.mat из файла.

```
In [2]: 1 matData = loadmat('ex3data1.mat')
        2 # print(matData)
        3
        4 x = matData["X"]
        5 y = matData["y"].squeeze()
        6
        7 xValidation = matData["Xval"]
        8 yValidation = matData["yval"].squeeze()
        9
       10 xTest = matData["Xtest"]
       11 yTest = matData["ytest"].squeeze()
       12
       13 print(x.shape, y.shape)
       14 print(xValidation.shape, yValidation.shape)
       15 print(xTest.shape, yTest.shape)
       16
```

```
((12, 1), (12,))
((21, 1), (21,))
((21, 1), (21,))
```

Постройте график, где по осям откладываются X и y из обучающей выборки.

```
In [3]: 1 pyplot.scatter(x, y)
        2 pyplot.show()
```



Реализуйте функцию стоимости потерь для линейной регрессии с L2-регуляризацией.

```
In [4]: 1 def h(theta, x):
        2     return np.dot(theta, x.T)
        3
        4 def costWithL2(theta, x, y, lambda):
        5     prediction = h(theta, x)
        6     difference = prediction - y.squeeze()
        7     sqareDifference = sum(difference ** 2)
        8     l2Parameter = lambda * sum(theta ** 2)
        9     return (sqareDifference + l2Parameter) / (2*len(y))
        10
        11 def getExtendedData(x, y):
        12     extendedX = np.hstack((np.ones((len(x), 1)), x))
        13     theta = np.ones(extendedX.shape[1])
        14     return extendedX, theta
        15
        16
        17 extendedX, theta = getExtendedData(x, y)
        18 lambda = 1
        19 cost = costWithL2(theta, extendedX, y, lambda)
        20 print(cost)
```

304.034858887

Реализуйте функцию градиентного спуска для линейной регрессии с L2-регуляризацией.

```
In [5]: 1 def gradientDescentWithL2(theta, x, y, lambda):
2         predictions = h(theta, x)
3         gradient = np.dot(x.T, (predictions - y))
4         regularization = lambda * theta
5         return (gradient + regularization) / len(y)
6
7 gradient = gradientDescentWithL2(theta, extendedX, y, lambda)
8 print(gradient)
```

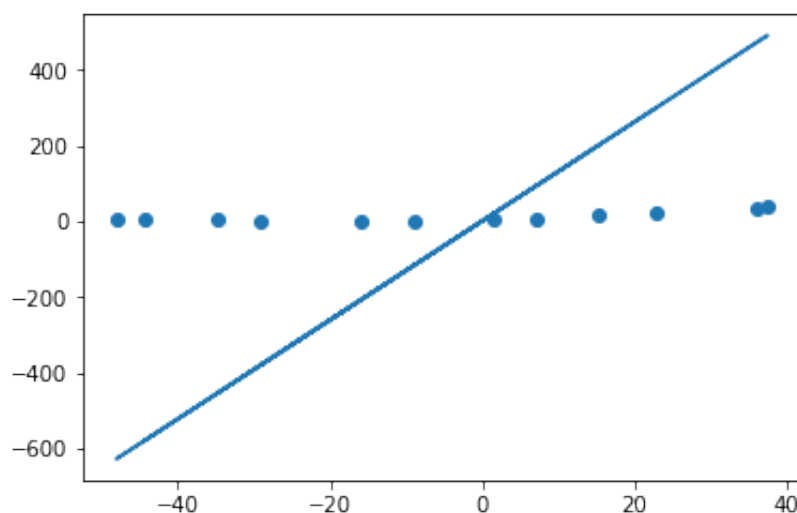
```
[ -15.21968234  598.25074417]
```

Постройте модель линейной регрессии с коэффициентом регуляризации 0 и постройте график полученной функции совместно с графиком из пункта 2. Почему регуляризация в данном случае не работает?

```
In [6]: 1 theta = optimize.minimize(costWithL2, theta, (extendedX, y, 0), me
2         print(theta)
```

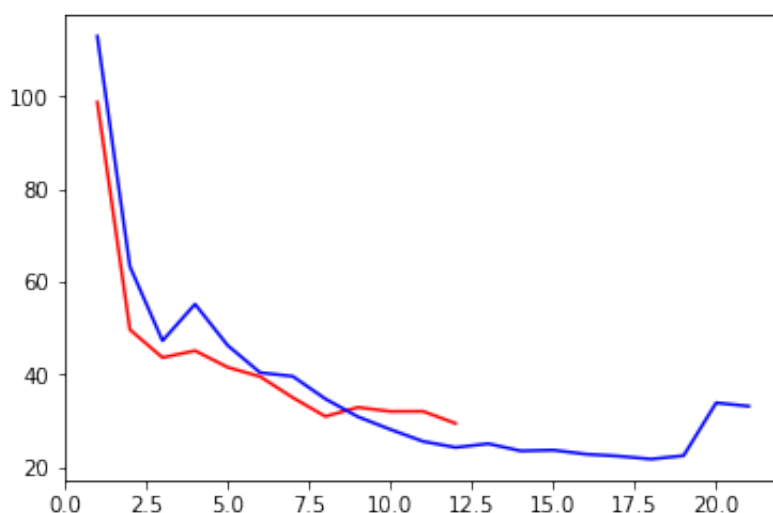
```
[ 13.08773879  0.36777194]
```

```
In [7]: 1 def lineValue(x, theta):
2         return x*theta[0] + theta[1]
3
4 pyplot.scatter(x, y)
5 pyplot.plot(x, lineValue(x, theta))
6 # pyplot.axis((-60,40,0,50))
7 pyplot.show()
```



Постройте график процесса обучения (learning curves) для обучающей и валидационной выборки. По оси абсцисс откладывается число элементов из обучающей выборки, а по оси ординат - ошибка (значение функции потерь) для обучающей выборки (первая кривая) и валидационной выборки (вторая кривая). Какой вывод можно сделать по построенному графику?

```
In [8]: 1 def plotLearningCurves(theta, x, y, lambda, color="r"):
2         counts = []
3         costs = []
4
5         for iterationIndex in range(len(x)):
6             sliceValue = iterationIndex+1
7             counts.append(sliceValue)
8             theta = optimize.minimize(costWithL2, theta, (x, y, 0), me
9             cost = costWithL2(theta, x[:sliceValue], y[:sliceValue], l
10            costs.append(cost)
11
12         pyplot.plot(counts, costs, c=color)
13
14     extendedXValidation, thetaValidation = getExtendedData(xValidation
15
16     plotLearningCurves(theta, extendedX, y, lambda, "r")
17     plotLearningCurves(thetaValidation, extendedXValidation, yValidati
18
19     pyplot.show()
```



Реализуйте функцию добавления  $p - 1$  новых признаков в обучающую выборку ( $X_2, X_3, X_4, \dots, X_p$ ).

In [9]:

```
1 def addNewParameters(x, count):
2     exist = x.squeeze()
3     newParameters = np.zeros((len(x), count-1))
4     return np.hstack((x, newParameters))
5
6 newX = addNewParameters(x, 2)
7 print(newX)
```

```
[[-15.93675813  0.          ]
 [-29.15297922  0.          ]
 [ 36.18954863  0.          ]
 [ 37.49218733  0.          ]
 [-48.05882945  0.          ]
 [ -8.94145794  0.          ]
 [ 15.30779289  0.          ]
 [-34.70626581  0.          ]
 [  1.38915437  0.          ]
 [-44.38375985  0.          ]
 [  7.01350208  0.          ]
 [ 22.76274892  0.          ]]
```

Поскольку в данной задаче будет использован полином высокой степени, то необходимо перед обучением произвести нормализацию признаков.

```
In [10]: 1 #Xnorm = (X - Xmin) / (Xmax - Xmin)
2
3 def normalize(x):
4     m, n = x.shape
5
6     for columnIndex in range(n):
7         column = x[:,columnIndex]
8         minValue, maxValue = min(column), max(column)
9         denominator = maxValue - minValue if (maxValue - minValue)
10         normalizeColumn = (column - minValue) / denominator
11         x[:,columnIndex] = normalizeColumn
12     return x
13
14 normalize(newX)
15
```

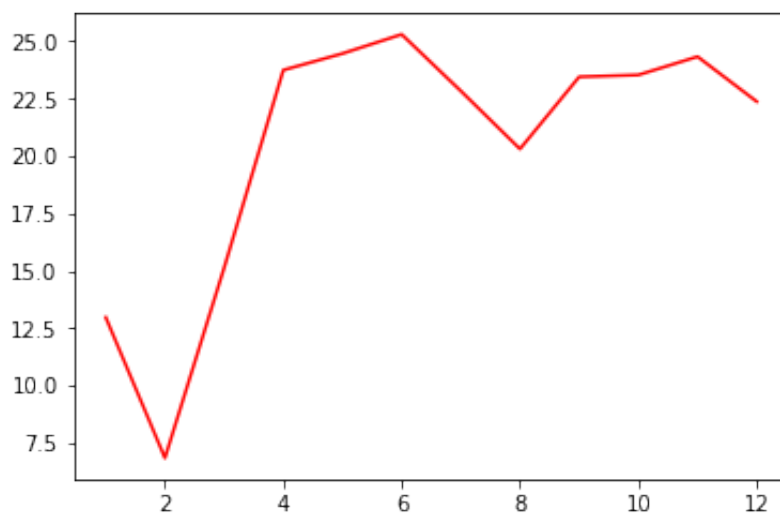
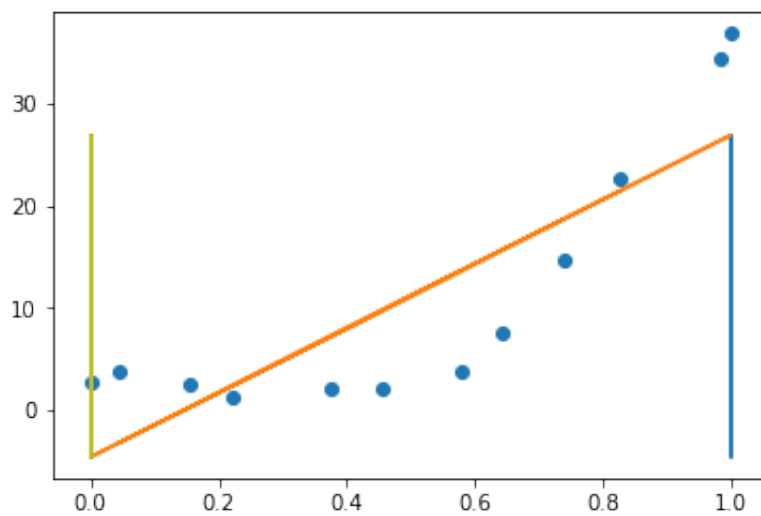
```
Out[10]: array([[ 0.3754727 ,  0.          ],
 [ 0.22098919,  0.          ],
 [ 0.98477355,  0.          ],
 [ 1.          ,  0.          ],
 [ 0.          ,  0.          ],
 [ 0.45724029,  0.          ],
 [ 0.74068813,  0.          ],
 [ 0.15607721,  0.          ],
 [ 0.57799411,  0.          ],
 [ 0.04295764,  0.          ],
 [ 0.64373673,  0.          ],
 [ 0.8278286 ,  0.          ]])
```

Обучите модель с коэффициентом регуляризации 0 и  $p = 8$ .

```
In [11]: 1 def model(x, lmbda=0, p=8):
2     newX = addNewParameters(x, 8)
3     normalize(newX)
4     extendedNewX, newTheta = getExtendedData(newX, y)
5
6     # gradient = gradientDescentWithL2(newTheta, extendedNewX, y,
7     # print(gradient)
8
9     newTheta = optimize.minimize(costWithL2, newTheta, (extendedNe
10     return extendedNewX, newTheta
11
12 newX, theta = model(x)
```

Постройте график модели, совмещенный с обучающей выборкой, а также график процесса обучения. Какой вывод можно сделать в данном случае?

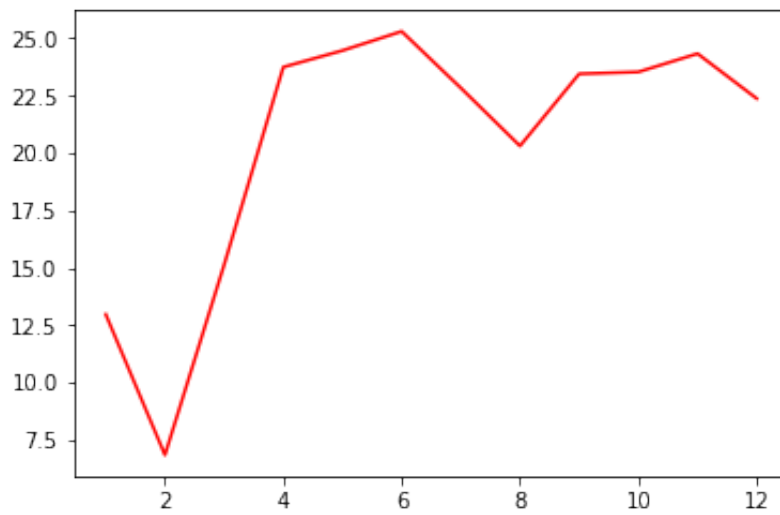
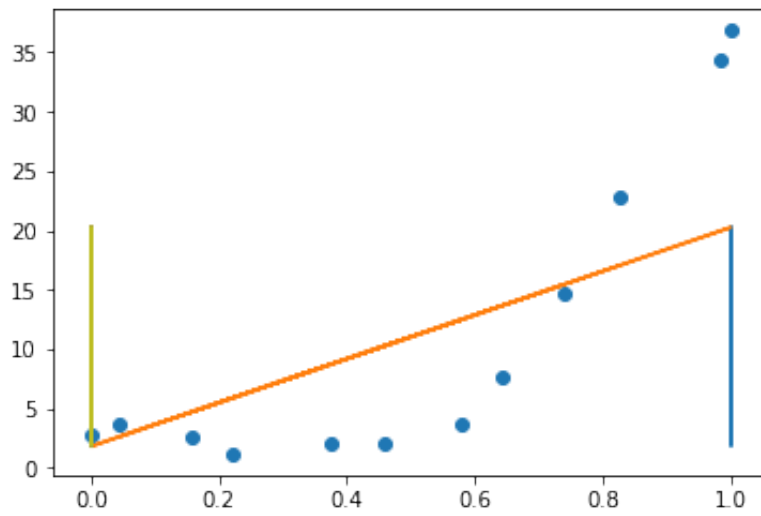
```
In [12]: 1 def plotModelAndLearningCurves(theta, x, y, lambda=0):
2         pyplot.scatter(x[:, 1], y)
3         pyplot.plot(x, h(theta, x))
4         pyplot.show()
5
6         plotLearningCurves(theta, x, y, lambda)
7         pyplot.show()
8
9 plotModelAndLearningCurves(theta, newX, y)
```



Постройте графики из пункта 10 для моделей с коэффициентами регуляризации 1 и 100. Какие выводы можно сделать?

In [13]:

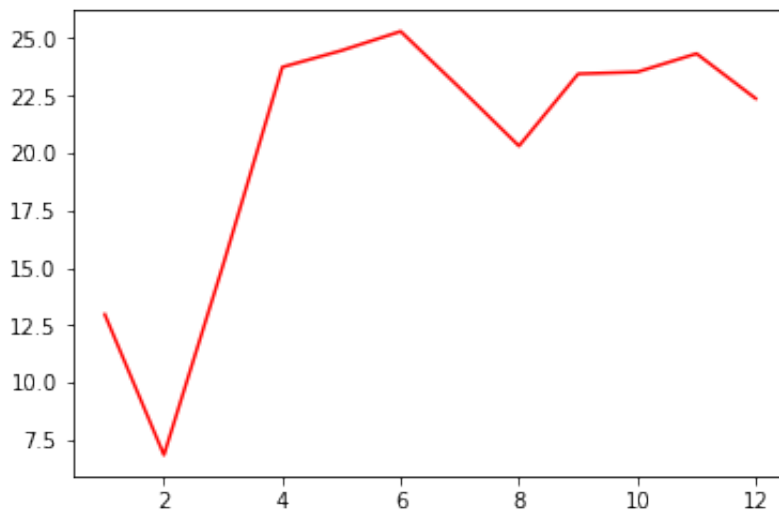
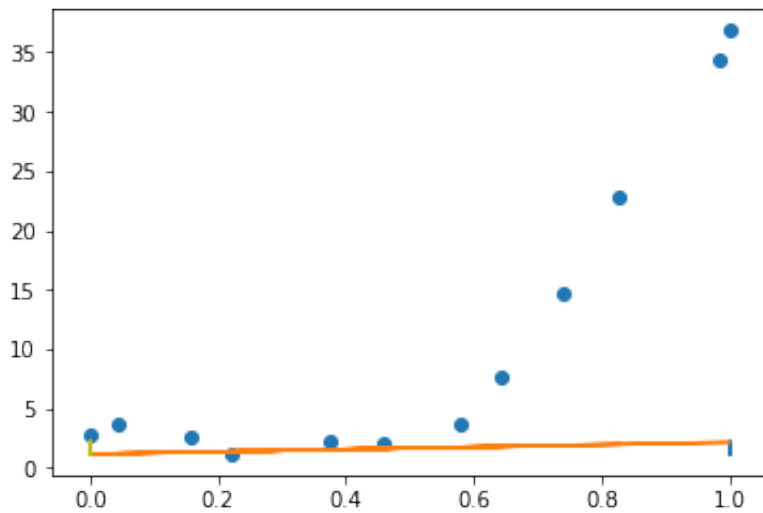
```
1 newX, theta = model(x, lambda=1)
2 plotModelAndLearningCurves(theta, newX, y)
```





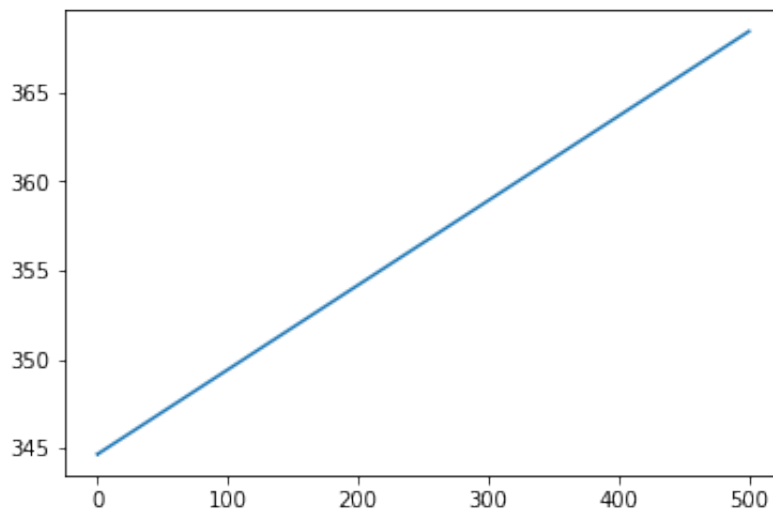
In [14]:

```
1 newX, theta = model(x, lambda=100)
2 plotModelAndLearningCurves(theta, newX, y)
```



С помощью валидационной выборки подберите коэффициент регуляризации, который позволяет достичь наименьшей ошибки. Процесс подбора отразите с помощью графика (графиков).

```
In [15]: 1 lambdaValues = [500, 250, 100, 1, 0.5, 0, 0.01, 0.0001, 0.0000001]
2
3 costs = []
4
5 for lambda in lambdaValues:
6     x, theta = getExtendedData(xValidation, yValidation)
7     cost = costWithL2(theta, x, yValidation, lambda)
8     costs.append(cost)
9
10 pyplot.plot(lambdaValues, costs)
11 pyplot.show()
```



Вычислите ошибку (потерю) на контрольной выборке.

```
In [16]: 1 xResult, thetaResult = getExtendedData(xTest, yTest)
2 thetaResult = optimize.minimize(costWithL2, thetaResult, (xResult,
3 cost = costWithL2(thetaResult, xResult, yTest, lambda)
4 print(cost)
5
```

30.1205372978

**Вывод**

Переобучение — негативное явление, возникающее, когда алгоритм обучения вырабатывает предсказания, которые слишком близко или точно соответствуют конкретному набору данных и поэтому не подходят для применения алгоритма к дополнительным данным или будущим наблюдениям

Возможные решения при переобучении:

- Увеличение количества данных в наборе;
- Уменьшение количества параметров модели;
- Добавление регуляризации / увеличение коэффициента регуляризации.