

Лабораторная работа №1 “Линейная регрессия”

Долматович Алина, 858641

Набор данных ex1data1.txt представляет собой текстовый файл, содержащий информацию о населении городов (первое число в строке) и прибыли ресторана, достигнутой в этом городе (второе число в строке). Отрицательное значение прибыли означает, что в данном городе ресторан терпит убытки.

Набор данных ex1data2.txt представляет собой текстовый файл, содержащий информацию о площади дома в квадратных футах (первое число в строке), количестве комнат в доме (второе число в строке) и стоимости дома (третье число).

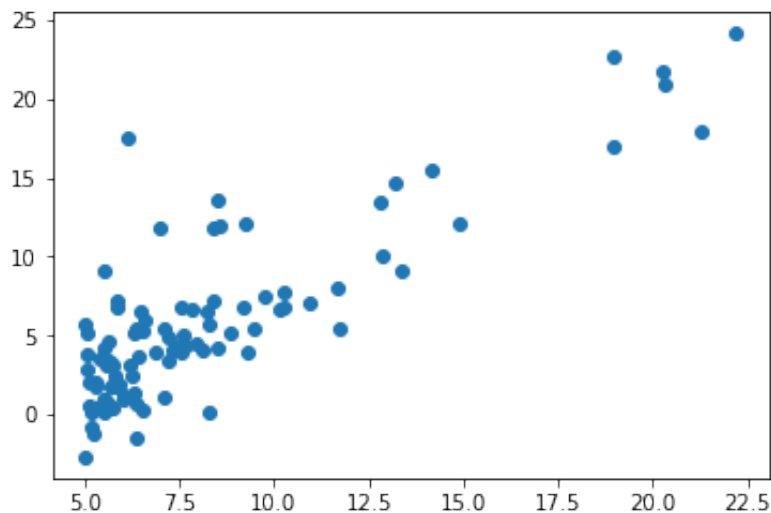
```
In [1]: 1 import pandas
        2 import matplotlib.pyplot as mp
        3 import numpy as np
        4 import mpl_toolkits.mplot3d
        5 from matplotlib import cm
        6 import time
```

Загрузите набор данных ex1data1.txt из текстового файла.

```
In [2]: 1 def getData(fileName):
        2     return pandas.read_csv(fileName, sep = ",", header=None)
        3
        4 data = getData('ex1data1.txt').sort_values(by=0)
        5 # print(data)
```

Постройте график зависимости прибыли ресторана от населения города, в котором он расположен.

```
In [3]: 1 population = list(data[0])
2 profit = list(data[1])
3
4 def defaultPlot():
5     mp.scatter(population, profit)
6
7 # mp.plot([4, 25], [0, 27], 'r') #  $y = 1,286x - 5,143$ 
8 defaultPlot()
9 mp.show()
```



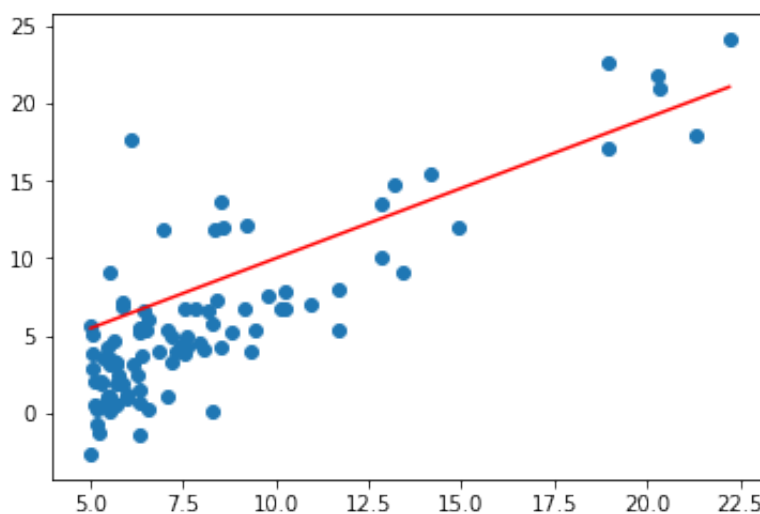
Реализуйте функцию потерь $J(\theta)$ для набора данных ex1data1.txt.

```
In [4]: 1 def MSECostFunction(y, theta):
2     def predictedValue(x):
3         return theta[0]*x + theta[1]
4
5     def lossValue(existent, predicted):
6         return (existent - predicted)
7
8     hypothesis = map(predictedValue, y)
9     lossValue = map(lossValue, y, hypothesis)
10    costValue = sum(map(lambda x: x ** 2, lossValue)) / (2*len(y))
11    return hypothesis, lossValue, costValue
12
13 theta = [0.89710913, 0.89710913]
14 _, _, cost = MSECostFunction(population, theta)
15 print(cost)
```

0.0801101386171

Реализуйте функцию градиентного спуска для выбора параметров модели. Постройте полученную модель (функцию) совместно с графиком из пункта 2.

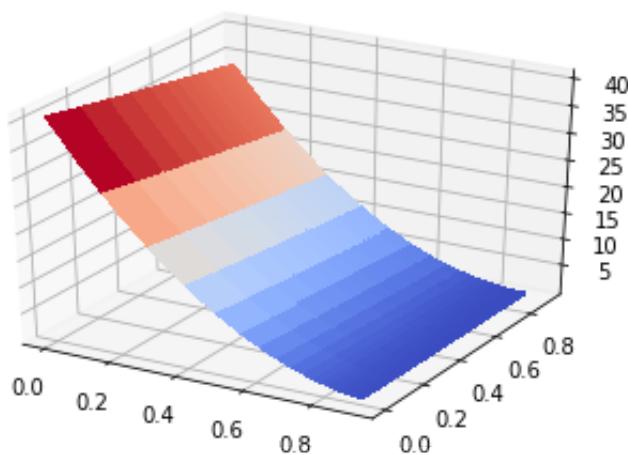
```
In [5]: 1 theta1 = []
2 theta2 = []
3
4 def gradientDescent(xValue, yValue, theta=[0, 0], alpha=0.001):
5     theta = np.array(theta)
6     xValue = np.array(xValue)
7     yValue = np.array(yValue)
8
9     optimalTheta, optimalCost = theta, 10000000
10
11     for i in range(100):
12         hypothesis, lossValue, costValue = MSECostFunction(xValue,
13             theta1.append(theta[0])
14             theta2.append(theta[1])
15
16         gradient = np.dot(xValue, lossValue) / len(xValue)
17         if optimalCost > costValue:
18             optimalTheta, optimalCost = theta, costValue
19             theta = theta + alpha * gradient
20
21     return optimalTheta, optimalCost
22
23
24 theta, cost = gradientDescent(population, profit)
25
26 hypothesis, _, _ = MSECostFunction(population, theta)
27 mp.plot(population, hypothesis, 'r')
28 defaultPlot()
29 mp.show()
```



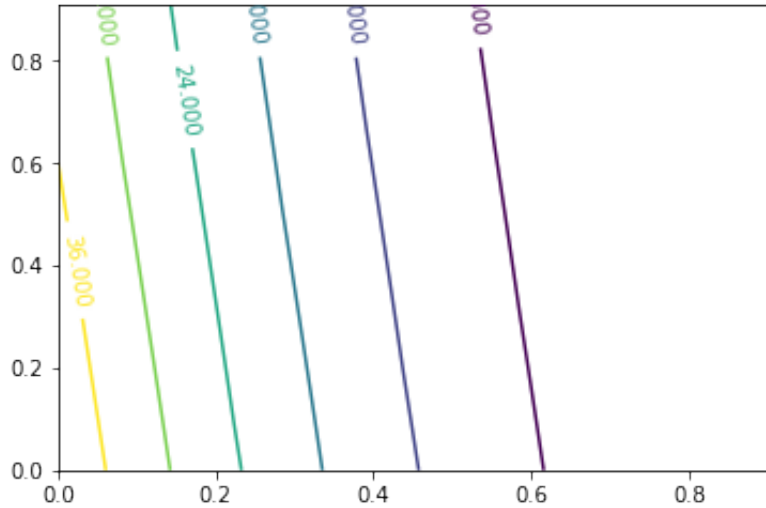
Постройте трехмерный график зависимости функции потерь от параметров модели (θ_0 и θ_1) как в виде поверхности, так и в виде изолиний (contour plot).

```
In [6]: 1 def getPlotData(theta1, theta2):
2         X = np.array(theta1)
3         Y = np.array(theta2)
4         X, Y = np.meshgrid(X, Y)
5         Z = np.zeros((len(X), len(Y)))
6
7         for rowIndex in range(len(X)):
8             for columnIndex in range(len(X[rowIndex])):
9                 theta1 = X[rowIndex][columnIndex]
10                theta2 = Y[rowIndex][columnIndex]
11                _, _, cost = MSECostFunction(population, [theta1, thet
12                Z[rowIndex][columnIndex] = cost
13
14         return X, Y, Z
```

```
In [7]: 1 figure = mp.figure()
2         ax = figure.gca(projection='3d')
3         startTimeNonVectorized = time.time()
4         X, Y, Z = getPlotData(theta1, theta2)
5         finishTimeNonVectorized = time.time()
6         surface = ax.plot_surface(X.copy(), Y.copy(), Z.copy(), cmap=cm.co
7                               linewidth=0, antialiased=False)
8         mp.show()
```



```
In [8]: 1 fig, ax = mp.subplots()
2 CS = ax.contour(X.copy(), Y.copy(), Z.copy())
3 ax.clabel(CS, inline=1, fontsize=10)
4
5 mp.show()
```



Загрузите набор данных ex1data2.txt из текстового файла.

```
In [9]: 1 data2 = getData('ex1data2.txt')
2 # print(data2)
```

Произведите нормализацию признаков. Повлияло ли это на скорость сходимости градиентного спуска? Ответ дайте в виде графика.

In [10]:

```
1 #Xnorm = (X - Xmin) / (Xmax - Xmin)
2
3 sqareList = list(data2[0])
4 roomList = list(data2[1])
5 priceList = list(data2[2])
6
7 def normalize(list):
8     normalizeList = []
9     xMax, xMin = max(list), min(list)
10    for item in list:
11        normalizeValue = float(item - xMin)/float(xMax - xMin)
12        normalizeList.append(normalizeValue)
13    return normalizeList
14
15 normalizedPrice = normalize(priceList)
16 normalizedRoom = normalize(roomList)
17 normalizedSquare = normalize(sqareList)
```

Реализуйте функции потерь $J(\theta)$ и градиентного спуска для случая многомерной линейной регрессии с использованием векторизации

In [13]:

```

1 def vectorizedCostFunction(X, Y, theta):
2     temp = (X.dot(theta) - Y)
3     return (1. / (2 * len(X))) * np.dot(temp.T, temp)
4
5 def vectorizedGradientDescent(X, Y, theta, iterations=500, alpha=0.01):
6     optimalTheta, optimalCost = theta, float("inf")
7     thetas = [theta]
8     costs = []
9     for i in range(iterations):
10         value = ((X.dot(theta) - Y).T).dot(X).T
11         a = alpha * (1. / len(Y)) * value
12         theta -= alpha * (1. / len(Y)) * value
13         thetas.append(theta)
14         cost = vectorizedCostFunction(X, Y, theta)
15         costs.append(cost)
16         if cost < optimalCost:
17             optimalCost = cost
18             optimalTheta = theta
19     return np.array(thetas), np.array(costs)
20
21 X = np.array(data2.iloc[:, 0:2])
22 Y = np.array(data2[2])
23 theta = [165, 2]
24 # cost = vectorizedCostFunction(X, Y, theta)
25 # print(cost)
26 # optimalTheta, optimalCost, optimalThetas = vectorizedGradientDescent(X, Y, theta)
27 optimalThetas, optimalCosts = vectorizedGradientDescent(X, Y, theta)
28
29 normalizedX = np.column_stack([normalizedSquare, normalizedRoom])
30 normalizedY = np.array(normalizedPrice)
31 # normalizedCost = vectorizedCostFunction(normalizedX, normalizedY, theta)
32 # print(normalizedCost)
33 # normalizedOptimalTheta, normalizedOptimalCost, normalizedOptimalThetas = vectorizedGradientDescent(normalizedX, normalizedY, theta)
34 normalizedOptimalThetas, normalizedOptimalCosts = vectorizedGradientDescent(normalizedX, normalizedY, theta)

```

In [14]:

```

1 def getPlotData(xArray, yArray, theta):
2     X = np.array(theta[:,0])
3     Y = np.array(theta[:,1])
4     X, Y = np.meshgrid(X, Y)
5     Z = np.zeros((len(X), len(Y)))
6
7     for rowIndex in range(len(X)):
8         for columnIndex in range(len(X[rowIndex])):
9             theta1 = X[rowIndex][columnIndex]
10            theta2 = Y[rowIndex][columnIndex]
11            cost = vectorizedCostFunction(xArray, yArray, [theta1, theta2])
12            Z[rowIndex][columnIndex] = cost

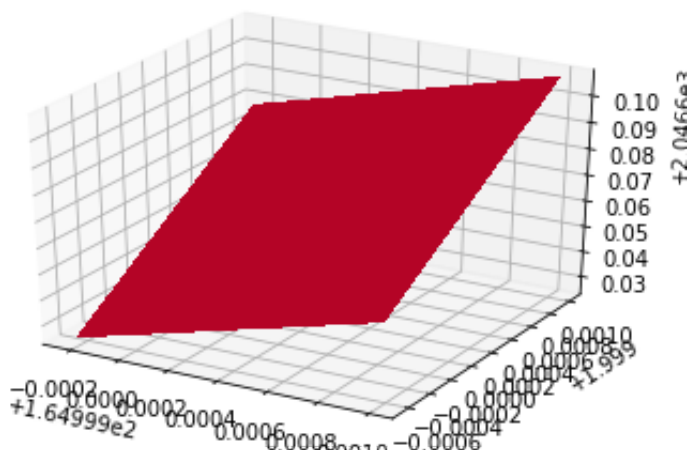
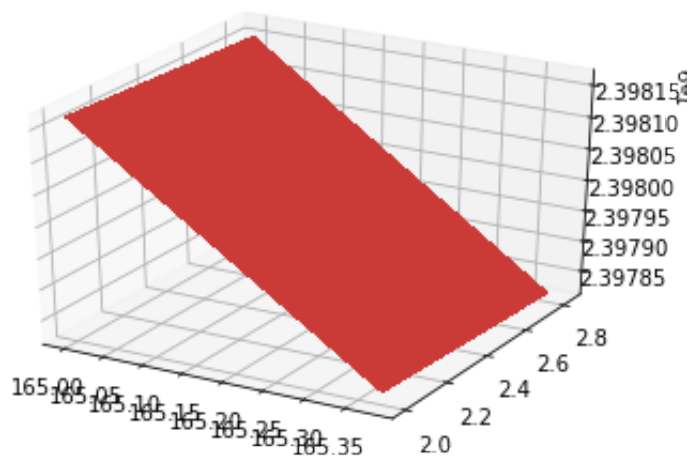
```



```

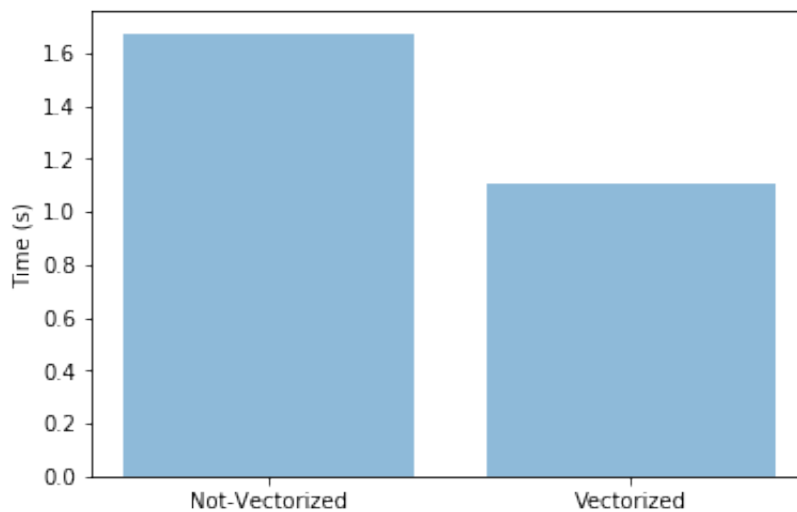
12         Z[rowIndex][columnIndex] = cost
13
14     return X, Y, Z
15
16     startTimeVectorized = time.time()
17     X1, Y1, Z1 = getPlotData(X, Y, optimalThetas)
18     finishTimeVectorized = time.time()
19     X2, Y2, Z2 = getPlotData(normalizedX, normalizedY, normalizedOptimalThetas)
20
21     # Plot the surface.
22     figure = mp.figure()
23     ax = figure.gca(projection='3d')
24     surface = ax.plot_surface(X1.copy(), Y1.copy(), Z1.copy(), cmap=cm.coolwarm,
25                             linewidth=0, antialiased=False)
26     figure = mp.figure()
27     ax = figure.gca(projection='3d')
28     surface = ax.plot_surface(X2.copy(), Y2.copy(), Z2.copy(), cmap=cm.coolwarm,
29                             linewidth=0, antialiased=False)
30
31     mp.show()

```



Покажите, что векторизация дает прирост производительности.

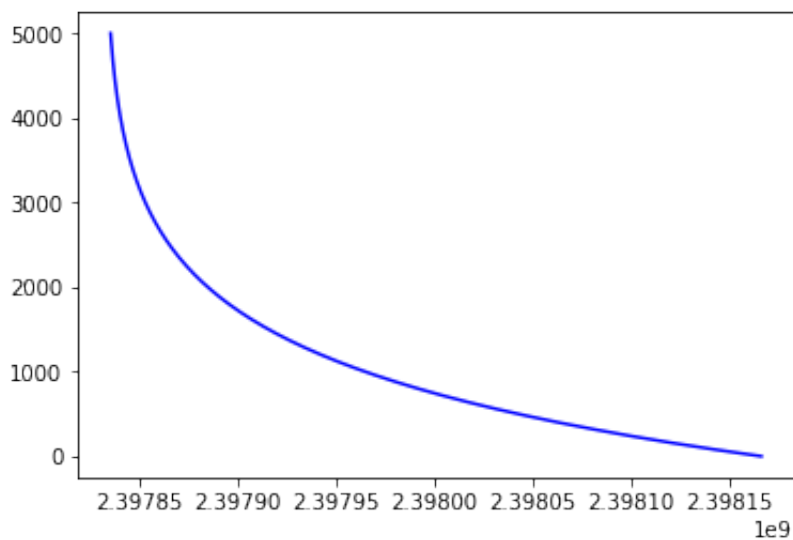
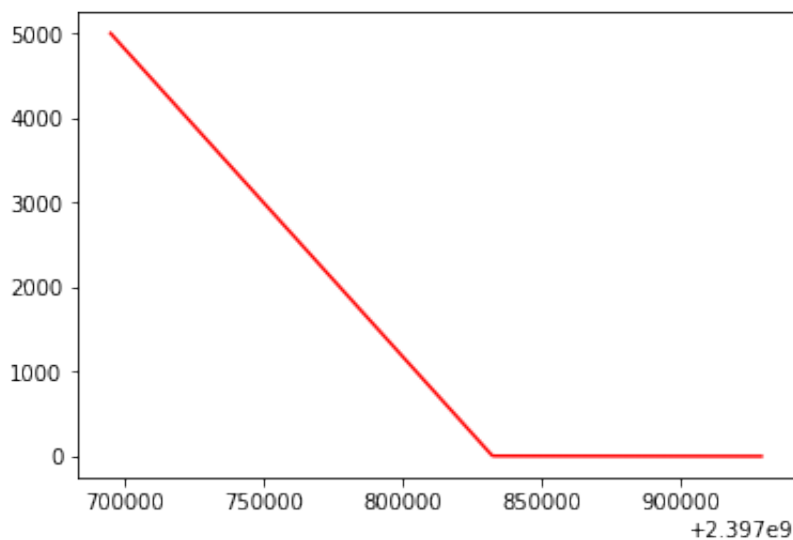
```
In [15]: objects = ('Not-Vectorized', 'Vectorized')
y2pos = np.arange(len(objects))
performance = [finishTimeNonVectorized-startTimeNonVectorized, finishTimeVectorized-startTimeVectorized]
mp.bar(y_pos, performance, align='center', alpha=0.5)
mp.xticks(y_pos, objects)
mp.ylabel('Time (s)')
mp.show()
```



Попробуйте изменить параметр α (коэффициент обучения). Как при этом изменяется график функции потерь в зависимости от числа итераций градиентного спуска? Результат изобразите в качестве графика.

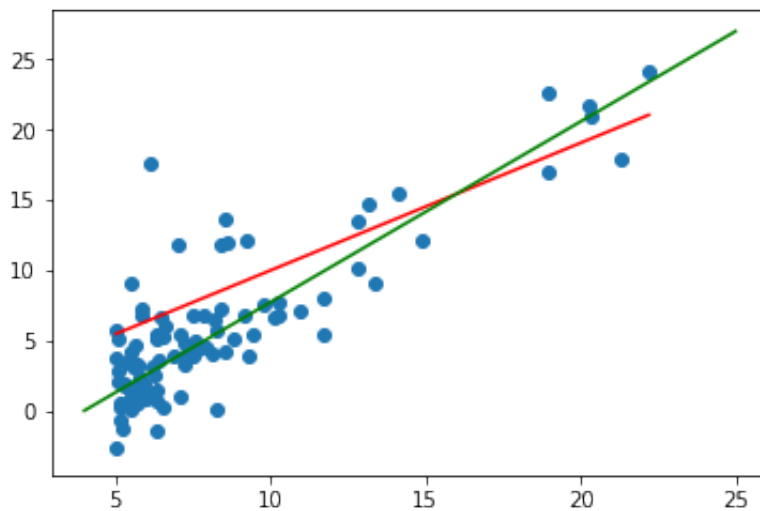
In [18]:

```
1 _, optimalCosts1 = vectorizedGradientDescent(X, Y, theta, 5000, 0.)
2 _, optimalCosts2 = vectorizedGradientDescent(X, Y, theta, 5000, 0.)
3
4 y1 = np.array(range(len(optimalCosts1)))
5 y2 = np.array(range(len(optimalCosts2)))
6
7 mp.plot(optimalCosts1, y1, 'r')
8 mp.show()
9 mp.plot(optimalCosts2, y2, 'b')
10 mp.show()
```



Постройте модель, используя аналитическое решение, которое может быть получено методом наименьших квадратов. Сравните результаты данной модели с моделью, полученной с помощью градиентного спуска.

```
In [19]: 1 mp.plot(population, hypothesis, 'r')
2 mp.plot([4, 25], [0, 27], 'g') #  $y = 1,286x - 5,143$ 
3 defaultPlot()
4 mp.show()
```



Вывод

Линейная регрессия — модель зависимости переменной x от одной или нескольких других переменных (факторов, регрессоров, независимых переменных) с линейной функцией зависимости.

Линейная регрессия относится к задаче определения «линии наилучшего соответствия» через набор точек данных и стала простым предшественником нелинейных методов, которые используют для обучения нейронных сетей.

Цель линейной регрессии — поиск линии, которая наилучшим образом соответствует этим точкам.

```
In [ ]: 1
```