

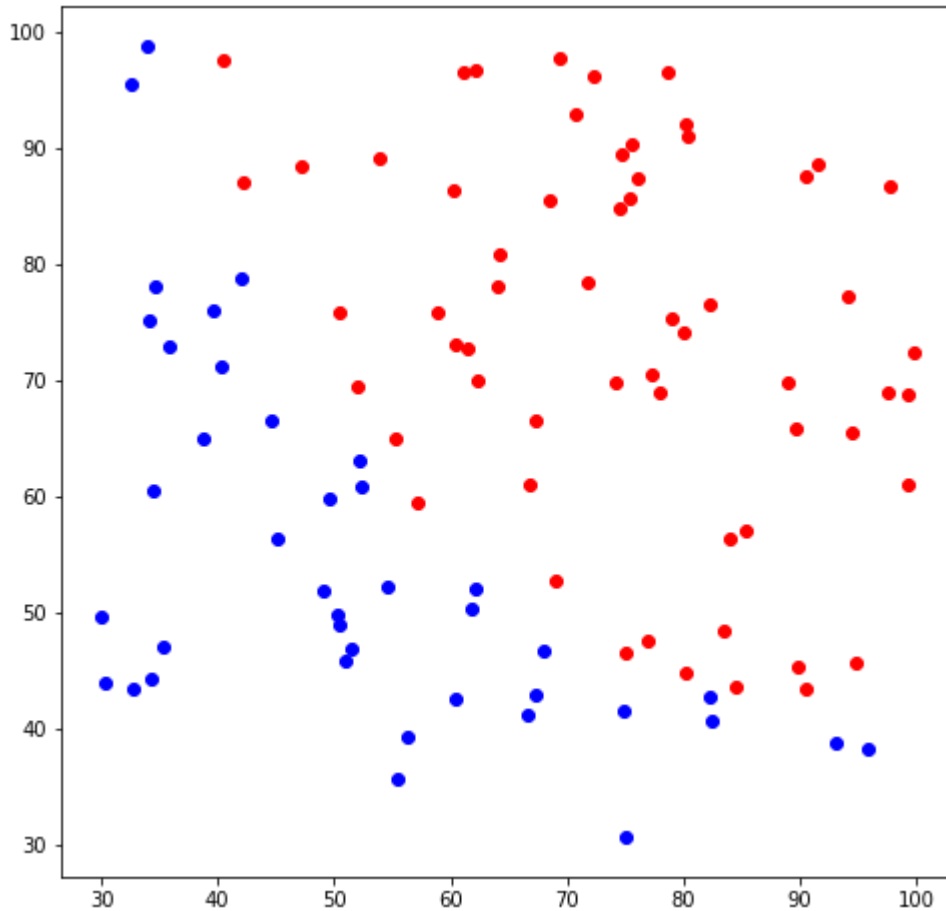
▼ 1. Load data

- data.txt 파일을 동일 디렉터리에 위치시키고 데이터를 읽었습니다.
- 읽은 데이터들을 visualize할때 label이 0인 data는 blue, 1인 data는 red로 시각화했습니다.

```

1 import csv
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import math
5 '''
6 1. Data
7 - Config
8 - Load Data
9 '''
10
11 # Config
12 learning_rate=0.0002014
13 opt_threshold=0.00000002748
14 t0, t1, t2=-10., -5., 10
15 theta=np.array([[t0, t1, t2]])
16
17 # Load data
18 path = "drive/My Drive/Classroom/Machine Learning (2) 2020-1\
19 /class-MachineLearning/assignment05/"
20 train=[]
21 test=[]
22
23 data = np.genfromtxt(path+"data.txt", delimiter=',')
24
25 x = data[:, 0]
26 y = data[:, 1]
27 m = len(x)
28
29 label = data[:, 2]
30
31 x_label0 = x[label == 0]
32 x_label1 = x[label == 1]
33
34 y_label0 = y[label == 0]
35 y_label1 = y[label == 1]
36 temp=np.ones(m)
37
38 input_data=np.hstack((temp.reshape(-1, 1), x.reshape(-1, 1), y.reshape(-1, 1)))
39
40 plt.figure(figsize=(8, 8))
41 plt.scatter(x_label0, y_label0, alpha=1, c='b')
42 plt.scatter(x_label1, y_label1, alpha=1, c='r')
43 plt.show()

```



▼ 2. Optimize Logistic regression

- train data에 대한 hypothesis, objective function, gradient descent를 정의했습니다.
- 연산과정들을 matrix로 빠르게 계산하기 위해 np.dot함수를 사용했습니다.
- 각 theta(feature)들에 맞는 update함수를 정의함으로써 학습동안 각각의 theta들이 모두 update되도록
- 수렴하는 기준을 0.00000002748로 설정하고 iteration에 대한 loss변화량이 이 기준값보다 작으면 수렴했다 판단했습니다.
- Objective function에서 log0으로 인한 -inf값 방지를 위해 동일한 notation이지만 inf를 생성하지 않도록 식을 수정했습니다.

```

1 '''
2 2. Hypothesis
3 - Define hypothesis and sigmoid hypothesis.
4
5 3. Objective function
6 - Define objective function, modify notation for J.
7
8 4. Gradient Descent
9 - Define derivations.
10 - Update theta0, theta1, theta2.
'''

```

```

11 '''
12
13 # Define hypothesis
14 def sigmoid_h() :
15     return 1/(1+np.exp(-z))
16
17 def hypothesis() :
18     return input_data.dot(theta.T).flatten()
19
20 z=hypothesis()
21 h=sigmoid_h()
22
23 # Define objective function
24 def objective_function() :
25     positive_z=np.array(list(map(lambda x : max(x, 0), z)))
26     J = np.sum(positive_z - z*label + np.log(1+np.exp(-np.abs(z)))) / m
27     return J
28
29 # Define derivations
30 def dev() :
31     dev0=np.sum(h-label) / m
32     dev1=np.sum((h-label)*x) / m
33     dev2=np.sum((h-label)*y) / m
34     dev=[dev0, dev1, dev2]
35     return dev
36
37 # Update theta0, theta1, theta2, theta3
38 def gradient_descent() :
39     return (theta[0][0]-(learning_rate*dev0), theta[0][1]-(learning_rate*dev1),\
40            theta[0][2]-(learning_rate*dev2))
41
42 J_list=[]
43 theta0_list=[]
44 theta1_list=[]
45 theta2_list=[]
46 h_list=[]
47 # Train
48 count=0
49 for i in range(20000) :
50     h_list.append(h)
51     count+=1
52     J=objective_function()
53
54     J_list.append(J)
55
56     theta0_list.append(theta[0][0])
57     theta1_list.append(theta[0][1])
58     theta2_list.append(theta[0][2])
59
60
61     dev0, dev1, dev2 = dev()
62     temp0, temp1, temp2 = theta[0]
63     theta[0][0], theta[0][1], theta[0][2] = gradient_descent()
64     z=hypothesis()
65     h=sigmoid_h()

```

```

66
67     if len(J_list)>=2 and \
68         abs(J_list[-2] - J_list[-1]) <opt_threshold :
69         break
70 print("theta0 : ", theta[0][0])
71 print("theta1 : ", theta[0][1])
72 print("theta2 : ", theta[0][2])
73 print("loss : ", J)

```

```

↳ theta0 : -10.173494868641017
   theta1 :  0.0868318479067422
   theta2 :  0.08059519104611516
   loss :   0.27049902948840765

```

▼ 3. Model parameters

- update를 할때마다 변하는 theta0, theta1, theta2, theta3의 값들을 따로 리스트에 저장해뒀었습니다.
- 각 iteration마다 위 theta들이 각각 어떻게 변하는지 시각화했습니다.

```

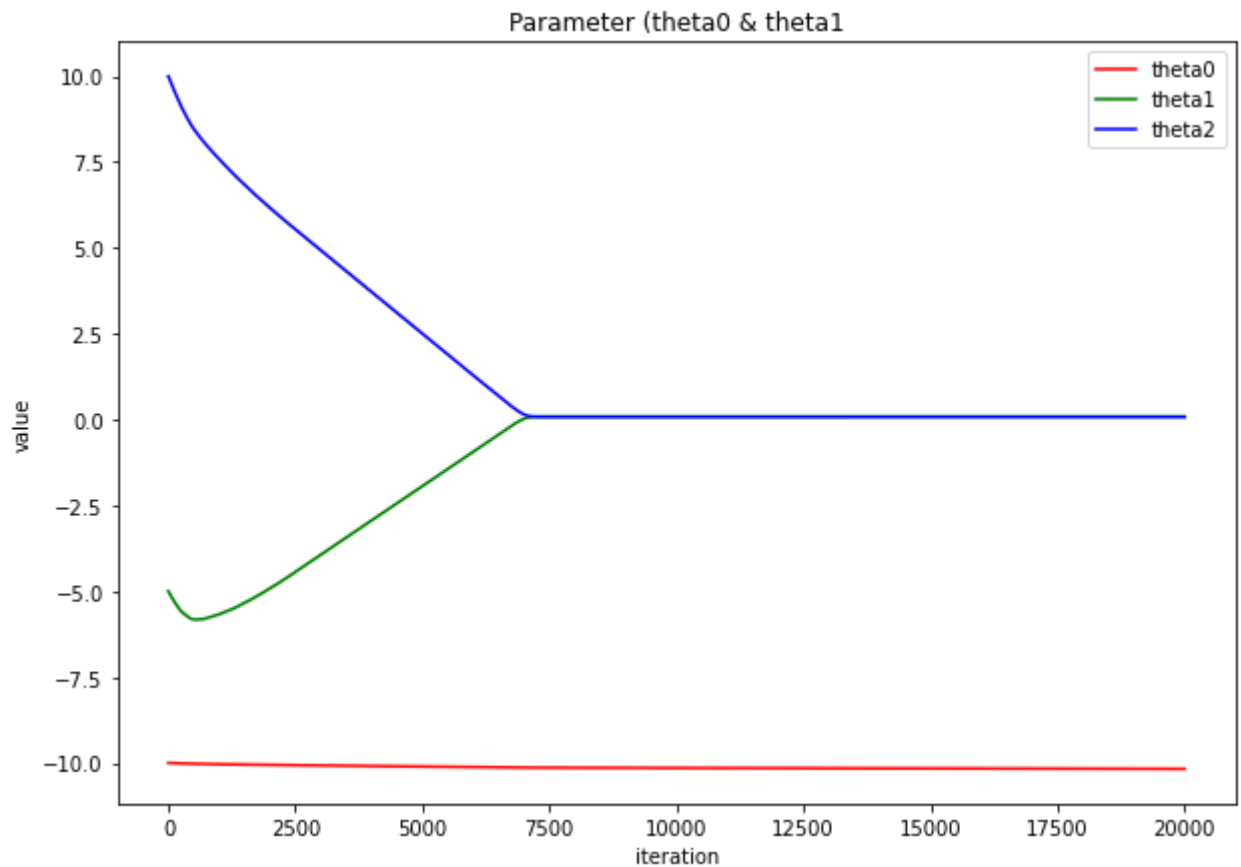
1 '''
2 Visualize theta
3 '''
4 plt.figure(figsize=(10, 7))
5 plt.plot(theta0_list, color='red', label='theta0')
6 plt.plot(theta1_list, color='green', label='theta1')
7 plt.plot(theta2_list, color='blue', label='theta2')
8 plt.legend(loc='upper right')
9 plt.xlabel("iteration")
10 #plt.xlim(left=-1000)
11 plt.ylabel("value")
12 plt.title("Parameter (theta0 & theta1)")
13 plt.show()

```

```

↳

```



▼ 4. Training error

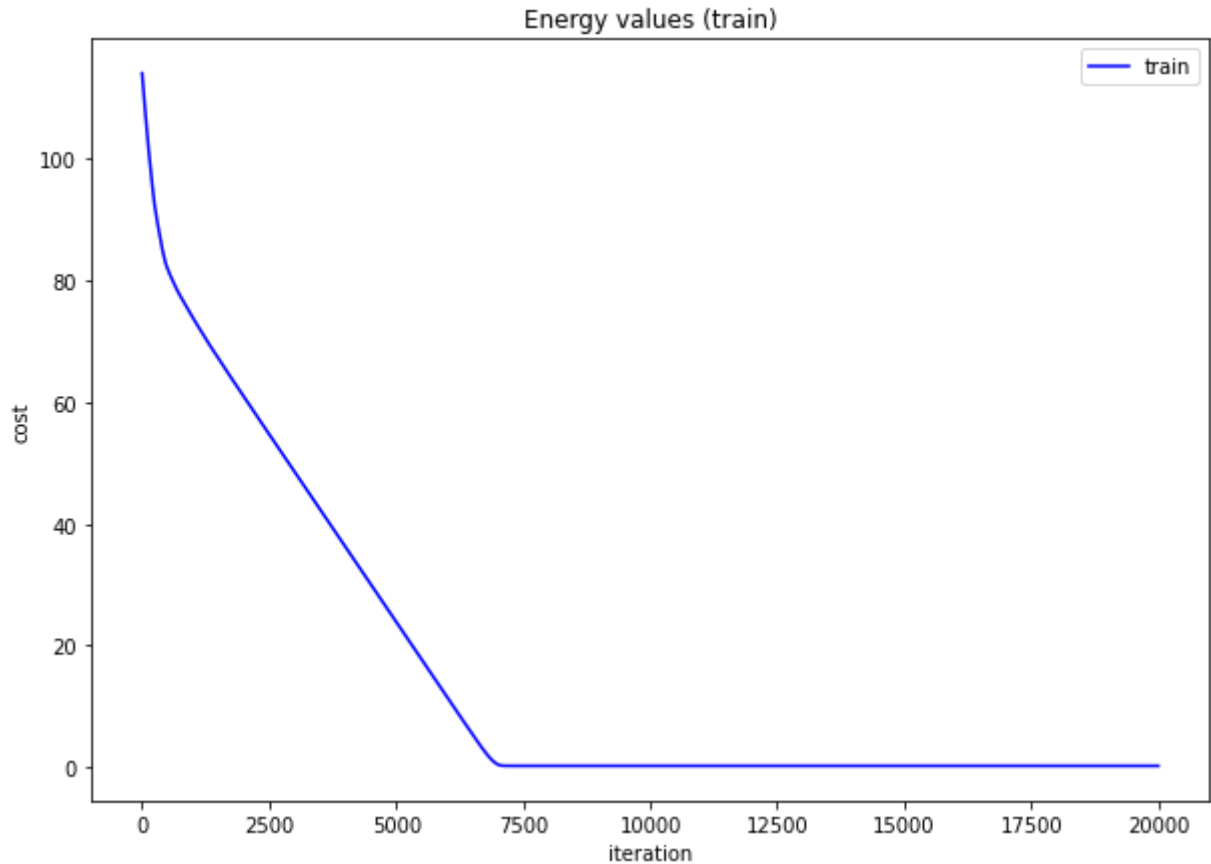
- training data에 대해서 update를 할때마다 train_J_list에 그동안의 loss 값들을 넣었습니다.
- 각 iteration마다 loss값이 어떻게 변하는지 시각화했습니다.

```

1 '''
2 Visualize Loss
3 '''
4 plt.figure(figsize=(10, 7))
5 plt.plot(J_list, color='blue', label='train')
6 plt.xlabel("iteration")
7 plt.ylabel("cost")
8 plt.legend(loc='upper right')
9 plt.title("Energy values (train)")
10 plt.show()

```





▼ 5. Visualize classifier

- 그래프에서 x와 y에 대한 범위는 30~100으로 설정하고 해당 범위 내에서 0.5씩 증가하면서 각각 141개의 값을 갖게됩니다.
- meshgrid를 통해서 x와 y값들을 좌표평면에 대응시킬 수 있는 grid matrix로 만듭니다.
- 각 x, y 좌표에 따른 hypothesis값을 구하는 함수를 정의하고 0에 가까울수록 blue, 1에 가까울수록 red로 plot했습니다.

```

1 '''
2 Visualize classifier
3 '''
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import matplotlib.cm as cm
7 def H(x, y) :
8     H=theta[0][0] + theta[0][1] * x + theta[0][2] * y
9     return 1/(1+np.exp(-H))
10
11 grid_x=np.linspace(start=30, stop=100, num=141)
12 grid_y=np.linspace(start=30, stop=100, num=141)
13 # Construct meshgrid between x and y
14 grid_x, grid_y=np.meshgrid(grid_x, grid_y)
15

```

```
16
17 fig = plt.figure(figsize=(12, 8))
18
19 grid_H=np.array([H(x, y) for (x, y) in \
20                 zip(np.ravel(grid_x), np.ravel(grid_y))])
21 grid_H = grid_H.reshape(grid_x.shape)
22
23 plt.scatter(grid_x, grid_y, c=grid_H, cmap='jet', s=3)
24 plt.scatter(x_label0, y_label0, alpha=1, c='b')
25 plt.scatter(x_label1, y_label1, alpha=1, c='r')
26 plt.show()
```

