

## Feature Normalization

When feature values differ greatly, we can get much slower rates of convergence of gradient-based algorithms. Furthermore, when we start using regularization (introduced in a later problem), features with larger values are treated as “more important”, which is not usually what you want. One common approach to feature normalization is perform an affine transformation (i.e. shift and rescale) on each feature so that all feature values in the training set are in  $[0; 1]$ . Each feature gets its own transformation. We then apply the same transformations to each feature on the test set. It’s important that the transformation is “learned” on the training set, and then applied to the test set. It is possible that some transformed test set values will lie outside the  $[0; 1]$  interval. Modify function *feature\_normalization* to normalize all the features to  $[0; 1]$ . (Can you use numpy’s “broadcasting” here?) Note that a feature with constant value cannot be normalized in this way. Your function should discard features that are constant in the training set.

```

1 def feature_normalization(train, test):
2     """Rescale the data so that each feature in the training set is in
3     the interval [0,1], and apply the same transformations to the test
4     set, using the statistics computed on the training set.
5
6     Args:
7         train – training set, a 2D numpy array of size (num_instances, num_features)
8         test – test set, a 2D numpy array of size (num_instances, num_features)
9
10    Returns:
11        train_normalized – training set after normalization
12        test_normalized – test set after normalization
13    """
14    # TODO
15    train_min = np.min(train, axis = 0)
16    train_max = np.max(train, axis = 0)
17    train_range = train_max - train_min
18    train_normalized = (train - train_min)/train_range
19    test_normalized = (test - train_min)/train_range
20    return train_normalized, test_normalized

```

Function 1: feature normalization

To do the feature normalization, we need to calculate the range of the numpy array, so compute the *train\_range* and then

$$(train - train\_min)/train\_range$$

$$(test - train\_min)/train\_range$$

to finish the normalized.

## Gradient Descent Setup

In linear regression, we consider the hypothesis space of linear functions  $h_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$ , where

$$h_\theta(x) = \theta^T x$$

for  $\theta, x \in \mathbb{R}^d$ , and we choose  $\theta$  that minimizes the following "average square loss" objective function:

$$j(\theta) = 1/m * \sum_{i=1}^m (h_\theta(x_i) - y_i)^2$$

where  $(x_1, y_1), \dots, (x_m, y_m) \in \mathbb{R}^d \rightarrow \mathbb{R}$  is our training data. While this formulation of linear regression is very convenient, it's more standard to use a hypothesis space of "affine" functions:

$$h_{\theta,b}(x) = \theta^T x + b$$

which allows a "bias" or nonzero intercept term. The standard way to achieve this, while still maintaining the convenience of the first representation, is to add an extra dimension to  $x$  that is always a fixed value, such as 1. You should convince yourself that this is equivalent. We'll assume this representation, and thus we'll actually take  $\theta, x \in \mathbb{R}^d$ .

1. Let  $X \in \mathbb{R}^{m \times (d+1)}$  be the **design matrix**, where the  $i$ th row of  $X$  is  $x_i$ . Let  $y = (y_1, \dots, y_m)^T \in \mathbb{R}^m$  be the "response". Write the objective function  $J(\theta)$  as a matrix/vector expression, without using an explicit summation sign.

$$\begin{aligned} j(\theta) &= 1/m * \sum_{i=1}^m (h_\theta(x_i) - y_i)^2 \\ &= \frac{1}{m} (X\theta - y)^T (X\theta - y) \end{aligned}$$

Or we can also write:

$$\begin{aligned} j(\theta) &= 1/2m * \sum_{i=1}^m (h_\theta(x_i) - y_i)^2 \\ &= \frac{1}{2m} (X\theta - y)^T (X\theta - y) \end{aligned}$$

The  $1/m$  is to "average" the squared error over the number of components so that the number of components doesn't affect the function. The latter function  $2J$  may seem more "natural," but the factor of 2 does not matter when optimizing. Can have better comparison for across model.

2. Write down an expression for the gradient of  $J$  (again, as a matrix/vector expression, without using an explicit summation sign).

$$j(\theta) = \frac{\partial(\frac{1}{m}(X\theta - y)^T(X\theta - y))}{\partial(X\theta - y)} \frac{\partial(X\theta - y)}{\partial\theta} = \frac{2}{m}(X\theta - y)^T X$$

Or we can also write:

$$j(\theta) = \frac{\partial(\frac{1}{2m}(X\theta - y)^T(X\theta - y))}{\partial(X\theta - y)} \frac{\partial(X\theta - y)}{\partial\theta} = \frac{1}{m}(X\theta - y)^T X$$

3. Use the gradient to write down an approximate expression for  $J(\theta + \eta h) - J(\theta)$ . The gradient at point  $\theta$  is the best linear approximation of  $J$  at that point.

$$J(\theta + \eta h) - J(\theta) \approx \nabla J(\theta)^T(\theta + \eta h - \theta) = \eta \nabla J(\theta)h$$

4. Write down the expression for updating  $\theta$  in the gradient descent algorithm. Let  $\eta$  be the step size.

$$\theta_{i+1} = \theta_i - \eta \nabla J(\theta), \text{ where } \eta > 0$$

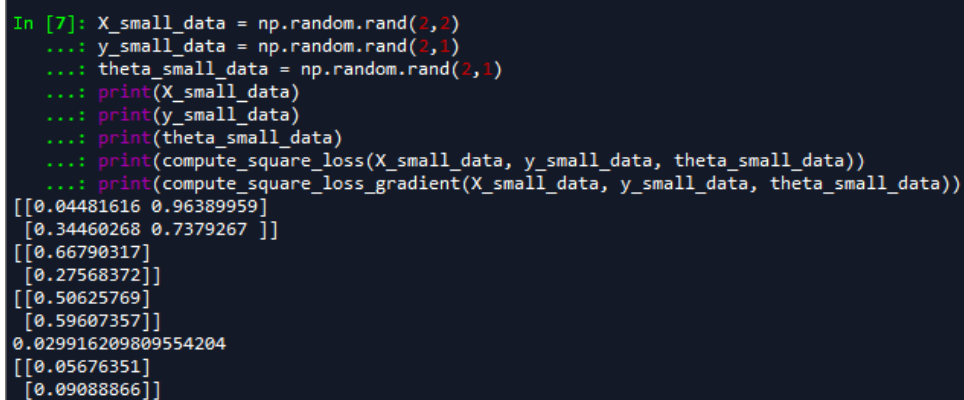
## 5. Modify the function compute square loss

```

1 def compute_square_loss(X, y, theta):
2     """
3     Given a set of X, y, theta, compute the average square loss for predicting y with X*theta.
4     Args:
5         X – the feature vector, 2D numpy array of size (num_instances, num_features)
6         y – the label vector, 1D numpy array of size (num_instances)
7         theta – the parameter vector, 1D array of size (num_features)
8
9     Returns:
10        loss – the average square loss, scalar
11    """
12    loss = 0 #Initialize the average square loss
13    #TODO
14    y_predict = np.dot(X, theta)
15    loss = np.dot(np.power(np.subtract(y_predict, y), 2), np.ones(X.shape[0])) / (2 * X.shape[0])
16    return loss
17 # test with small dataset
18 X_small_data = np.random.rand(2,2)
19 y_small_data = np.random.rand(2,1)
20 theta_small_data = np.random.rand(2,1)
21 print(X_small_data)
22 print(y_small_data)
23 print(theta_small_data)
24 print(compute_square_loss(X_small_data, y_small_data, theta_small_data))

```

Function 2: compute square loss



```

In [7]: X_small_data = np.random.rand(2,2)
...: y_small_data = np.random.rand(2,1)
...: theta_small_data = np.random.rand(2,1)
...: print(X_small_data)
...: print(y_small_data)
...: print(theta_small_data)
...: print(compute_square_loss(X_small_data, y_small_data, theta_small_data))
...: print(compute_square_loss_gradient(X_small_data, y_small_data, theta_small_data))
[[0.04481616 0.96389959]
 [0.34460268 0.7379267 ]]
[[0.66790317]
 [0.27568372]]
[[0.50625769]
 [0.59607357]]
0.029916209809554204
[[0.05676351]
 [0.09088866]]

```

Figure 1: small dataset to test square loss

In order to figure out whether this compute square loss function work, I used the *X\_small\_data* for the training data, which only two data, and *y\_small\_data* for the class. And the result is right.

## 6. Modify the function compute square loss gradient

```

1 def compute_square_loss_gradient(X, y, theta):
2     """
3     Compute the gradient of the average square loss (as defined in compute_square_loss), at the point
4     ↪ theta.
5     Args:
6     X – the feature vector, 2D numpy array of size (num_instances, num_features)
7     y – the label vector, 1D numpy array of size (num_instances)
8     theta – the parameter vector, 1D numpy array of size (num_features)
9     Returns:
10    grad – gradient vector, 1D numpy array of size (num_features)
11    """
12    #TODO
13    y_predict = np.dot(X, theta)
14    loss_gradient = np.dot(np.subtract(y_predict, y), X) / X.shape[0]
15    return loss_gradient
16
17 theta = np.ones(X_train.shape[1])
18 loss = compute_square_loss(X_train, y_train, theta)
19 loss_gradient = compute_square_loss_gradient(X_train, y_train, theta)
20 print(loss)
21 print(loss_gradient)
22 # test with small dataset
23 X_small_data = np.random.rand(2,2)
24 y_small_data = np.random.rand(2,1)
25 theta_small_data = np.random.rand(2,1)
26 print(X_small_data)
27 print(y_small_data)
28 print(theta_small_data)
29 print(compute_square_loss(X_small_data, y_small_data, theta_small_data))
30 print(compute_square_loss_gradient(X_small_data, y_small_data, theta_small_data))

```

Function 3: compute square loss gradient and run both functions

```

In [7]: X_small_data = np.random.rand(2,2)
...: y_small_data = np.random.rand(2,1)
...: theta_small_data = np.random.rand(2,1)
...: print(X_small_data)
...: print(y_small_data)
...: print(theta_small_data)
...: print(compute_square_loss(X_small_data, y_small_data, theta_small_data))
...: print(compute_square_loss_gradient(X_small_data, y_small_data, theta_small_data))
[[0.04481616 0.96389959]
 [0.34460268 0.7379267 ]]
[[0.66790317]
 [0.27568372]]
[[0.50625769]
 [0.59607357]]
0.029916209809554204
[[0.05676351]
 [0.09088866]]

```

Figure 2: small dataset to test square loss and loss gradient

Name	Type	Size	Value
x	float64	(200, 48)	[[ 0.00000000e+00 0.00000000e+00 0.00000000e+00 ... -7.81824199e-01 ...
x_test	float64	(100, 49)	[[1. 1. 1. ... 0.21229954 0.21229954 1. ...
x_train	float64	(100, 49)	[[1. 1. 1. ... 0.13241261 0.13241261 1. ...
df	DataFrame	(200, 49)	Column names: x0, x1, x2, x3, x4, x5, x6, x7, x8, x9, x10, x11, x12, x ...
loss	float64	1	441.4342610243893
loss_gradient	float64	(49,)	[26.07013205 25.58326502 25.05054258 ... 16.85255853 1... 26.2 ...
theta	float64	(49,)	[1. 1. 1. ... 1. 1. 1.]
x_test	float64	(100, 48)	[[ 1. 1. 1. ... -0.61692371 ... -3 ...
x_train	float64	(100, 48)	[[ 1. 1. 1. ... -0.67897469 ... -3 ...
y	float64	(200,)	[-1.37657523 0.87878245 1.10870068 ... 3.37935567 ... 3.9 ...
y_test	float64	(100,)	[-1.07263709 -1.91105204 -1.07036941 ... -2.95765686 -... 4.2 ...
y_train	float64	(100,)	[-2.79297326 4.73657935 1.43620799 ... -3.4037511 ...

Figure 3: import data

```

loading the dataset
Split into Train and Test
Scaling all to [0, 1]
441.4342610243893
[26.07013205 25.58326502 25.05054258 24.74081471 24.24546864 23.79891226
 22.90863314 22.90863314 21.90026135 20.0569049 18.36416632 17.3017451
 14.5002761 12.45800123 8.21040281 7.34556074 5.24340575 0.90770525
 23.14742097 23.14742097 23.14742097 21.48832491 21.48832491 21.48832491
 19.8972414 19.8972414 19.8972414 19.15895467 19.15895467 19.15895467
 18.74898619 18.74898619 18.74898619 11.85229611 11.85229611 11.85229611
 14.39885467 14.39885467 14.39885467 15.83930209 15.83930209 15.83930209
 16.49353617 16.49353617 16.49353617 16.85255853 16.85255853 16.85255853
 26.23768143]

```

Figure 4: loss and gradient loss result

In order to figure out whether this compute square loss function work, I used the  $X\_small\_data$  for the training data, which only two data, and  $y\_small\_data$  for the class same with the last question. And the result is right.

## Batch Gradient Descent

At the end of the skeleton code, the data is loaded, split into a training and test set, and normalized. We'll now finish the job of running regression on the training set. Later on we'll plot the results together with SGD results.

### 1. Complete batch gradient descent.

```

1 def batch_grad_descent(X, y, alpha=0.1, num_step=1000, grad_check=False):
2     """
3     In this question you will implement batch gradient descent to
4     minimize the average square loss objective.
5
6     Args:
7         X – the feature vector, 2D numpy array of size (num_instances, num_features)
8         y – the label vector, 1D numpy array of size (num_instances)
9         alpha – step size in gradient descent
10        num_step – number of steps to run
11        grad_check – a boolean value indicating whether checking the gradient when updating
12
13    Returns:
14        theta_hist – the history of parameter vector, 2D numpy array of size (num_step+1,
15                        ↪ num_features)
16                        for instance, theta in step 0 should be theta_hist[0], theta in step (num_step) is
17                        ↪ theta_hist[-1]
18        loss_hist – the history of average square loss on the data, 1D numpy array, (num_step+1)
19    """
20    num_instances, num_features = X.shape[0], X.shape[1]
21    theta_hist = np.zeros((num_step+1, num_features)) #Initialize theta_hist
22    loss_hist = np.zeros(num_step+1) #Initialize loss_hist
23    theta = np.zeros(num_features) #Initialize theta
24    #TODO
25    loss_hist[0] = compute_square_loss(X, y, theta) # initial loss
26    theta_hist[0] = theta # initial theta
27    for i in range(num_step):
28        if grad_check:
29            grad = grad_checker(X, y, theta.T)
30            print('Grade check:{0}'.format(grad))
31
32        loss_gradient = compute_square_loss_gradient(X, y, theta.T)
33        theta = theta - alpha*loss_gradient.T
34        theta_hist[i+1] = theta
35        loss = compute_square_loss(X, y, theta.T)
36        loss_hist[i+1] = loss
37
38    return theta_hist, loss_hist

```

Function 4: batch gradient descent function

loss_return	float64	(1001, 49)	[[ 0. 0. 0. ... 0. ... 0. ...
theta	float64	(49,)	[1. 1. 1. ... 1. 1. 1.]
theta_return	float64	(1001,)	[3.98075917 3.91325976 3.85188424 ... 1.03900245 1.03892598 1.03884969 ...

Figure 5: data result for batch gradient algo

```

loss return: [[ 0. 0. 0. ... 0. 0.
 0.
 [-0.01259103 0.00750755 0.02025476 ... 0.01444064 0.01444064
 -0.01848617]
 [-0.05580238 -0.01564154 0.00998855 ... 0.00770169 0.00770169
 -0.06757517]
 ...
 [-1.5671494 0.64936919 1.43509585 ... -0.00771457 -0.00771457
 -1.16289895]
 [-1.56727897 0.649472 1.43460544 ... -0.00782174 -0.00782174
 -1.16264481]
 [-1.56740818 0.64957502 1.43411528 ... -0.00792888 -0.00792888
 -1.16239102]]
theta return: [3.98075917 3.91325976 3.85188424 ... 1.03900245 1.03892598
1.03884969]

```

Figure 6: result for batch gradient



# Stochastic Gradient Descent

October 8, 2019

2. Now let's experiment with the step size. Note that if the step size is too large, gradient descent may not converge. Starting with a step-size of 0.1, try various different fixed step sizes to see which converges most quickly and/or which diverge. As a minimum, try step sizes 0.5, 0.1, .05, and .01. Plot the average square loss as a function of the number of steps for each step size. Briefly summarize your findings.

```
1 def step_size_plot(X, y):  
2     alphas = [0.01, 0.05, 0.1, 0.5]  
3     color = ["red", "green", "blue", "yellow"]  
4     for i, alpha in enumerate(alphas):  
5         theta_hist, loss_hist = batch_grad_descent(X, y, alpha, num_step=1000)  
6         plt.plot(loss_hist, label=f"Alpha={alpha}", color = color[i])  
7     plt.yscale('log')  
8     plt.xlabel('Steps')  
9     plt.ylabel('Loss')  
10    plt.legend()  
11    plt.show()
```

Function 5: plot for log loss with different alpha

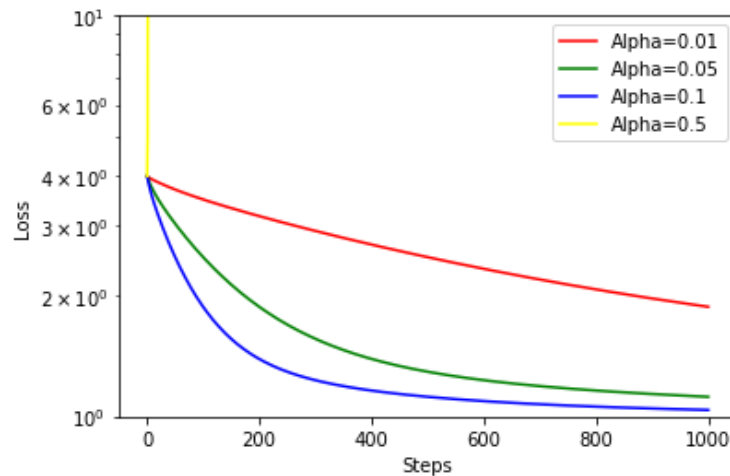


Figure 7: plot for log loss with different alpha

The solution we can conclude from the plot is the larger the alpha, the converge will faster.

3. Implement backtracking line search. How does it compare to the best fixed step-size you found in terms of number of steps? In terms of time? How does the extra time to run backtracking line search at each step compare to the time it takes to compute the gradient? (You can also compare the operation counts.)

```

1 def backtracking_line_search(X, y, theta):
2     sigma = 0.01
3     beta = 0.5
4     k = 0
5     alpha = 1 #step_size
6     while(True):
7         loss_gradient = compute_square_loss_gradient(X,y,theta)
8         theta_next = theta - alpha * loss_gradient
9         if (compute_square_loss(X,y,theta) - compute_square_loss(X,y,theta_next) >= sigma * alpha
10             ↪ * np.dot(loss_gradient.T,loss_gradient)):
11             break
12         else:
13             alpha = beta * alpha
14             theta = theta_next
15             k = k+1
16     print(alpha)
17 def step_size_plot(X, y):
18     alphas = [0.01, 0.05, 0.1, 0.5, 0.0625]
19     color = ["red", "green", "blue", "yellow", "pink"]
20     for i,alpha in enumerate(alphas):
21         start = time.time()
22         theta_hist, loss_hist = batch_grad_descent(X, y, alpha, num_step=1000)
23         end = time.time()
24         print("the" + str(i) + "alpha test cost:" + str(end - start))
25         plt.plot(loss_hist, label=f"Alpha={alpha}", color = color[i])
26     plt.yscale('log')
27     plt.xlabel('Steps')
28     plt.ylabel('Loss')
29     plt.legend()
30     plt.show()

```

Function 6: backtracking line search

```

In [90]: backtracking_line_search(X_train,y_train,np.ones(X_train.shape[1]))
0.0625

```

Figure 8: alpha after backtracking line search

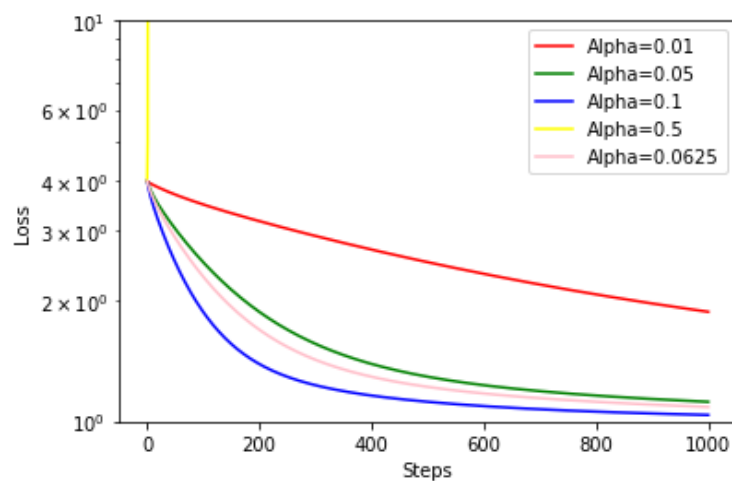


Figure 9: backtracking result is 0.0625 in pink line

```
the0alpha test cost:0.01694035530090332
the1alpha test cost:0.015965700149536133
the2alpha test cost:0.015947580337524414
the3alpha test cost:0.016998291015625
the4alpha test cost:0.016953468322753906
```

Figure 10: alpha after backtracking line search

The pink line is for backtracking result, that one converges in less steps, but only slightly different after 100 steps. And when I use `time.time()` to record the time for each different alphas, there are not too big difference between them.

## Stochastic Gradient Descent

When the training data set is very large, evaluating the gradient of the objective function can take a long time, since it requires looking at each training example to take a single gradient step. When the objective function takes the form of an average of many values, such as

$$J(\theta) = 1/m \sum_{i=1}^m f_i(\theta)$$

1. Show that the objective function  $j(\theta) = 1/m * \sum_{i=1}^m (h_\theta(x_i) - y_i)^2$  can be written in the form  $j(\theta) = 1/m * \sum_{i=1}^m f_i(\theta)$  by giving an expression for  $f_i(\theta)$  that makes the two expressions equivalent.

I think the function should be  $j(\theta) = 1/m * \sum_{i=1}^m (h_\theta(x_i) - y_i)^2 + \lambda \theta^T \theta$ .

$$\begin{aligned} j(\theta) &= 1/m * \sum_{i=1}^m (h_\theta(x_i) - y_i)^2 + \lambda \theta^T \theta \\ &= 1/m * \sum_{i=1}^m [(h_\theta(x_i) - y_i)^2 + \lambda \theta^T \theta] &= (h_\theta(x_i) - y_i)^2 + \lambda \theta^T \theta \end{aligned}$$

Or if we don't contain  $\lambda \theta^T \theta$ , then the answer is  $(h_\theta(x_i) - y_i)^2$ .

2. Show that the stochastic gradient

$$\mathbb{E}_i \nabla f_i(x) = \sum_{j=1}^m \mathbb{P}(i = j) \nabla f_j(\theta) = 1/m \sum_{i=1}^m \nabla f_i(x) = \nabla J(x)$$

3. Write down the update rule for  $\theta$  in SGD for the linear regression objective function.  
 because:  $\nabla f_i(\theta) = 2(h_\theta(x_i) - y_i)x_i + 2\lambda\theta$

$$\theta_{t+1} = \theta_t - \eta [2/m(\theta_t^T x_i - y_i)x_i + 2\lambda\theta_t]$$

Or we can write and there is no matter 1 or 2

$$\theta_{t+1} = \theta_t - \eta (2/m(\theta_t^T x_i - y_i)x_i)$$

## 4. Implement stochastic grad descent.

```

1 def stochastic_grad_descent(X, y, alpha=0.01, lambda_reg = 0.01, num_epoch=1000):
2     """
3     In this question you will implement stochastic gradient descent with regularization term
4
5     Args:
6         X – the feature vector, 2D numpy array of size (num_instances, num_features)
7         y – the label vector, 1D numpy array of size (num_instances)
8         alpha – string or float, step size in gradient descent
9             NOTE: In SGD, it's not a good idea to use a fixed step size. Usually it's set to 1/sqrt(t) or
10                 ↪ 1/t
11                 if alpha is a float, then the step size in every step is the float.
12                 if alpha == "1/sqrt(t)", alpha = 1/sqrt(t).
13                 if alpha == "1/t", alpha = 1/t.
14         num_epoch – number of epochs to go through the whole training set
15
16     Returns:
17         theta_hist – the history of parameter vector, 3D numpy array of size (num_epoch,
18                     ↪ num_instances, num_features)
19                     for instance, theta in epoch 0 should be theta_hist[0], theta in epoch (num_epoch) is
20                     ↪ theta_hist[-1]
21         loss_hist – the history of loss function vector, 2D numpy array of size (num_epoch,
22                     ↪ num_instances)
23     """
24     num_instances, num_features = X.shape[0], X.shape[1]
25     theta = np.ones(num_features) #Initialize theta
26     theta_hist = np.zeros((num_epoch, num_instances, num_features)) #Initialize theta_hist
27     loss_hist = np.zeros((num_epoch, num_instances)) #Initialize loss_hist
28     #TODO
29     step_size = 0
30     for i in range(num_epoch):
31
32         for j in range(num_instances):
33             theta_hist[i,j] = theta
34             regulariztion_loss = lambda_reg * np.dot(theta.T, theta)
35             loss_hist[i,j] = compute_square_loss(X, y, theta) + regulariztion_loss
36             gradient_part = 2*(np.dot(theta.T, X[j]) - y[j])*X[j] + 2*lambda_reg*theta
37
38             if alpha == "1/sqrt(t)":
39                 theta = theta - 0.1/step_size * gradient_part
40             elif alpha == "1/t":
41                 theta = theta - 0.1/np.sqrt(step_size)*gradient_part
42             else:
43                 theta = theta - alpha*gradient_part
44             step_size += 1
45     return theta_hist, loss_hist

```

Function 7: stochastic gradient descent

loss_gradient	float64	(49,)	[26.07013205 25.58326502 25.05054258 ... 16.85255853 16.85255853 26.2 ...
loss_hist	float64	(1000, 100)	[[441.92426102 396.7787976 18.63538062 ... 3.30285244 3.30271564 ...
loss_return	float64	(1001, 49)	[[ 0. 0. 0. ... 0. 0. 0. ...
theta	float64	(49,)	[1. 1. 1. ... 1. 1. 1.]
theta_hist	float64	(1000, 100, 49)	[[[ 1. 1. 1. ... 1. 1. 1. ...
theta_return	float64	(1001,)	[3.98075917 3.91325976 3.85188424 ... 1.03900245 1.03892598 1.03884969 ...

Figure 11: SGD output

5. In this case we are investigating the convergence rate of the optimization algorithm with different step size schedules, thus we're interested in the value of the objective function, which includes the regularization term. Sometimes the initial step size  $C$  for  $C/t$  and  $C/\sqrt{t}$  is too aggressive and will get you into a part of parameter space from which you can't recover. Try reducing  $C$  to counter this problem.

I plot the data for different alpha,  $[1/t, 1/\sqrt{t}, 0.0005, 0.001, 0.01]$ , and when the step size equal to  $0.1/\sqrt{t}$  performs the best. but the step size equal to  $0.1/t$  does not converge. ( Because I define the step\_size to  $0.1/\sqrt{t}$  if the enter alpha is  $1/\sqrt{t}$ . Same with  $1/t$ )

```

1  def test_SGD(X, y):
2      alphas = ["1/t", "1/sqrt(t)", 0.0005, 0.001, 0.01]
3      fig = plt.figure(figsize = (20,8))
4      plt.subplot(222)
5      for alpha in alphas:
6          [theta_hist, loss_hist] = stochastic_grad_descent(X, y, alpha, num_epoch = 5)
7          plt.plot(np.log(loss_hist.ravel()), label = 'alpha:' + str(alpha))
8      plt.legend()

```

Function 8: plot the SGD with different alpha

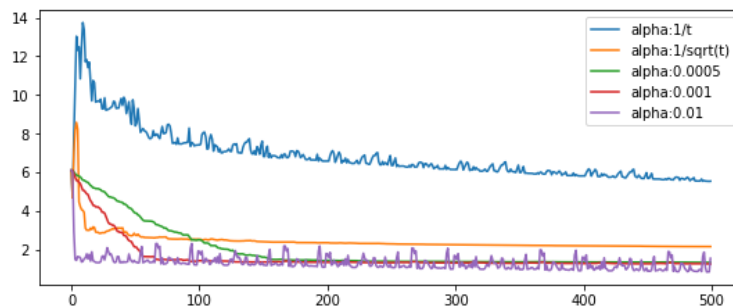


Figure 12: SGD plot