

# DS-GA 1003 HW1

Yiyan Chen

January 2018

## 1 Introduction

## 2 Linear Regression

### 2.1 Feature Normalization

---

```
1 def feature_normalization(train, test):
2     train_normalized = np.zeros((train.shape[0],train.shape[1]))
3     test_normalized = np.zeros((test.shape[0],test.shape[1]))
4     for idx in range(train.shape[1]):
5         mean_val = np.mean(train[:,idx], axis = 0)
6         max_val = np.max(train[:,idx], axis = 0)
7         min_val = np.min(train[:,idx], axis = 0)
8
9         if max_val != min_val:
10             train_normalized[:,idx] = (train[:,idx] - mean_val)/(max_val - min_val)
11             test_normalized[:,idx] = (test[:,idx] - mean_val)/(max_val - min_val)
12         else:
13             train_normalized[:,idx] = train[:,idx]
14             test_normalized[:,idx] = test[:,idx]
15
16     return (train_normalized, test_normalized)
```

---

## 2.2 Gradient Descent Setup

### 2.2.1

$$\begin{aligned}e_{\theta} &= y - X\theta \\J(\theta) &= \frac{1}{m} \sum e_i^2 \\&= \frac{1}{m} e_{\theta}^T e_{\theta} \\&= \frac{1}{m} (y - X\theta)^T (y - X\theta) \\&= \frac{1}{m} (y^T y - y^T X\theta - \theta^T X^T y + \theta^T X^T X\theta) \\&= \frac{1}{m} (y^T y - 2y^T X\theta + \theta^T X^T X\theta)\end{aligned}$$

### 2.2.2

$$\nabla_{\theta} J(\theta) = \frac{2}{m}(-X^T y + X^T X \theta)$$

### 2.2.3

$$\begin{aligned} J(\theta + \eta h) - J(\theta) &= \frac{1}{m} [y^T y - 2y^T X(\theta + \eta h) + (\theta + \eta h)^T X^T X(\theta + \eta h) - y^T y + 2y^T X\theta - \theta^T X^T X\theta] \\ &= \frac{1}{m} [-2\eta y^T Xh + 2\eta h^T X^T X\theta + \eta^2 h^T X^T Xh] \\ &= \eta h^T \nabla_{\theta} J(\theta) + \frac{\eta^2 h^T X^T Xh}{m} \\ &\approx \eta h^T \nabla_{\theta} J(\theta) \end{aligned}$$

#### 2.2.4

$$\begin{aligned}\theta_{k+1} &= \theta_k - \eta \nabla_{\theta} J(\theta_k) \\ &= \theta_k - \frac{2\eta}{m} (-X^T y + X^T X \theta)\end{aligned}$$

### 2.2.5

---

```
1 def compute_square_loss(X, y, theta):
2     loss = 0 #initialize the square_loss
3     TODO
4     m = X.shape[0]
5     return 1/m*(np.matmul(y.transpose(),y)-2*np.matmul(np.matmul(y.transpose(),X),theta)
6               +np.matmul(theta.transpose(),np.matmul(X.transpose(), np.matmul(X, theta))))
```

---

### 2.2.6

---

```
1 def compute_square_loss_gradient(X, y, theta):
2     m = X.shape[0]
3     return 2/m*(-np.matmul(y.transpose(), X) + np.matmul(X.transpose(), np.matmul(X, theta)))
```

---

## 2.3 (OPTIONAL)Gradient Checker

---

```
1 def grad_checker(X, y, theta, epsilon=0.01, tolerance=1e-4):
2     true_gradient = compute_square_loss_gradient(X, y, theta) #the true gradient
3     num_features = theta.shape[0]
4     approx_grad = np.zeros(num_features) #Initialize the gradient we approximate
5     TODO
6     for i in range(num_features):
7         e_i = np.zeros(num_features)
8         e_i[i] = 1
9         approx_grad[i] = (compute_square_loss(X, y, theta+epsilon*e_i)
10                          -compute_square_loss(X, y, theta-epsilon*e_i))/(2*epsilon)
11     distance = np.linalg.norm(approx_grad-true_gradient.transpose())
12     if distance > tolerance:
13         return False
14     return True
```

---

```
1 def generic_gradient_checker(X, y, theta, objective_func, gradient_func, epsilon=0.01, tolerance=1e-4):
2     TODO
3     true_gradient = gradient_func(X,y,theta)
4     num_features = theta.shape[0]
5     approx_grad = np.zeros(num_features) #Initialize the gradient we approximate
6     for i in range(num_features):
7         e_i = np.zeros(num_features)
8         e_i[i] = 1
9         approx_grad[i] = (objective_func(X, y, theta+epsilon*e_i)
10                          -objective_func(X, y, theta-epsilon*e_i))/(2*epsilon)
11     distance = np.linalg.norm(approx_grad-true_gradient.transpose())
12     if distance > tolerance:
13         return False
14     return True
```

---



## 2.4 Batch Gradient Descent

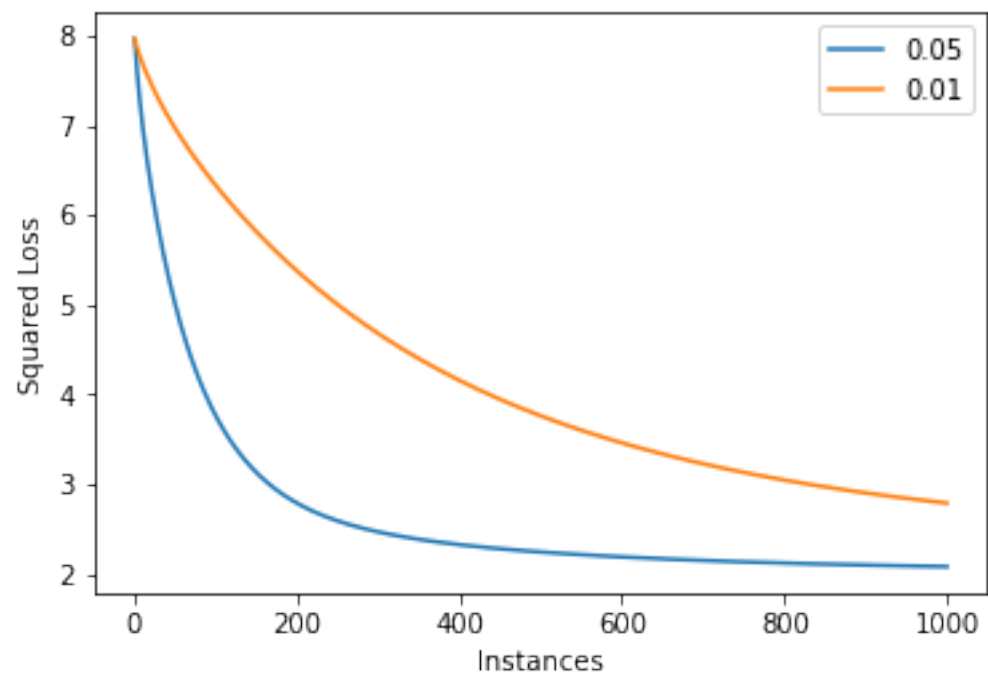
### 2.4.1

---

```
1 def batch_grad_descent(X, y, alpha=0.1, num_iter=1000, check_gradient=False):
2     num_instances, num_features = X.shape[0], X.shape[1]
3     theta_hist = np.zeros((num_iter+1, num_features)) #Initialize theta_hist
4     loss_hist = np.zeros(num_iter+1) #initialize loss_hist
5     theta = np.zeros(num_features) #initialize theta
6     #TODO
7     theta_hist[0] = theta
8     loss_hist[0] = compute_square_loss(X,y,theta)
9     for i in range(num_iter):
10         if check_gradient and (not grad_checker(X,y,theta_hist[i])):
11             return False
12         elif (not check_gradient) or (check_gradient and grad_checker(X,y,theta_hist[i])):
13             grad = compute_square_loss_gradient(X,y,theta_hist[i])
14             loss = compute_square_loss(X,y,theta_hist[i])
15             theta_hist[i+1] = theta_hist[i]-alpha*grad
16             loss_hist[i+1] = loss
17     if check_gradient and (not grad_checker(X,y,theta_hist[num_iter])):
18         return False
19     return (theta_hist, loss_hist)
```

---

### 2.4.2



### **2.4.3 OPTIONAL**

## 2.5 Ridge Regression

### 2.5.1

$$\nabla_{\theta} J(\theta) = \frac{2}{m}(-X^T y + X^T X \theta) + 2\lambda \theta$$

$$\begin{aligned}\theta_{k+1} &= \theta_k - \eta \nabla_{\theta} J(\theta) \\ &= \theta_k - \frac{2\eta}{m}(-X^T y + X^T X \theta) - 2\eta \lambda \theta\end{aligned}$$

### 2.5.2

---

```
1 def compute_regularized_square_loss_gradient(X, y, theta, lambda_reg=1):
2     m = X.shape[0]
3     return 2/m*(-np.matmul(X.transpose(), y)
4         + np.matmul(X.transpose(), np.matmul(X,theta)))+2*lambda_reg*theta
```

---

### 2.5.3

---

```
1 def regularized_grad_descent(X, y, alpha=0.1, lambda_reg=1, num_iter=1000):
2     (num_instances, num_features) = X.shape
3     theta = np.zeros(num_features) #Initialize theta
4     theta_hist = np.zeros((num_iter+1, num_features)) #Initialize theta_hist
5     loss_hist = np.zeros(num_iter+1) #Initialize loss_hist
6     #TODO
7     theta_hist[0] = theta
8     loss_hist[0] = compute_square_loss(X,y,theta)
9     for i in range(num_iter):
10         grad = compute_regularized_square_loss_gradient(X, y, theta_hist[i], lambda_reg)
11         theta_hist[i+1] = theta_hist[i]-alpha*grad
12         loss = compute_square_loss(X,y,theta_hist[i+1])
13         loss_hist[i+1] = loss
14     return (theta_hist, loss_hist)
```

---

#### 2.5.4

$$\theta_{k+1} = \theta_k + \frac{2\eta}{m} X^T y - (\frac{2\eta}{m} X^T X \theta + 2\eta\lambda)\theta$$

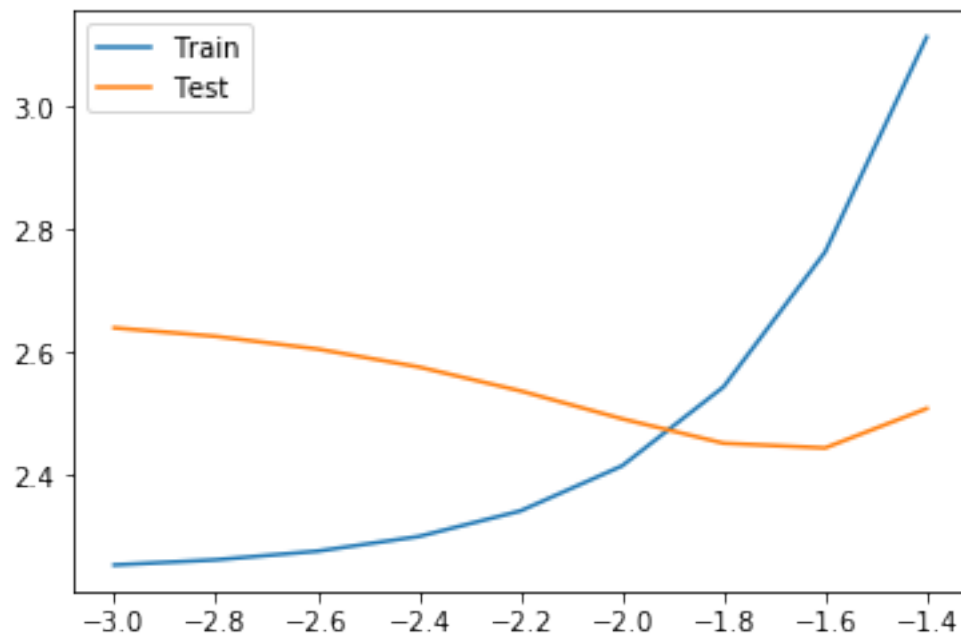
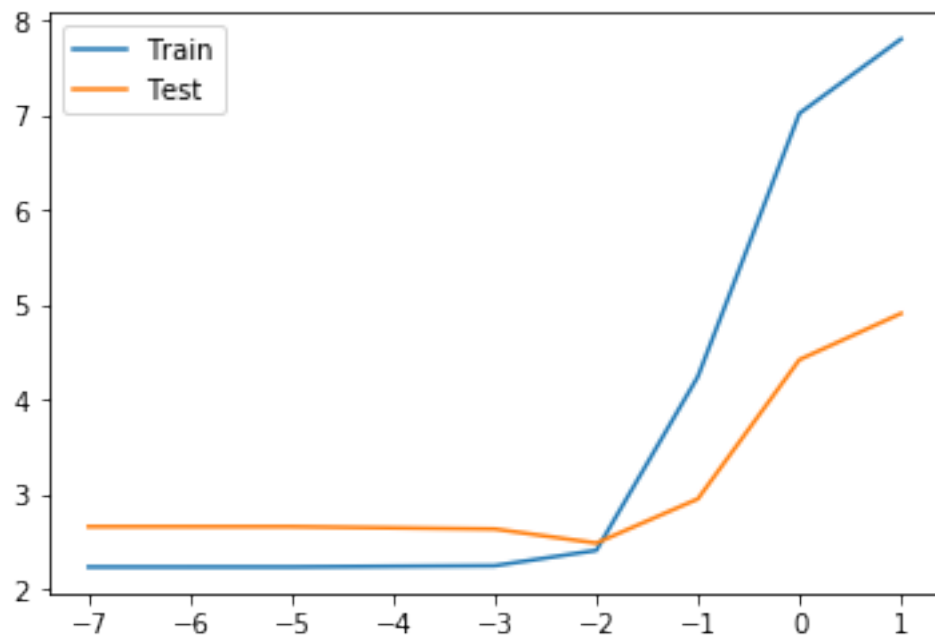
Making B larger means the  $\theta$  is replaced by  $B\theta$ . When the  $B\theta$  increases, the new  $\theta$  gets small. The bigger  $B\theta$ , the smaller new  $\theta$  is. The effective regularization on the bias term, i.e.  $\theta$  decreases. We make regularization weaker when B is larger. But B shouldn't be too large to avoid the convergence of loss function.

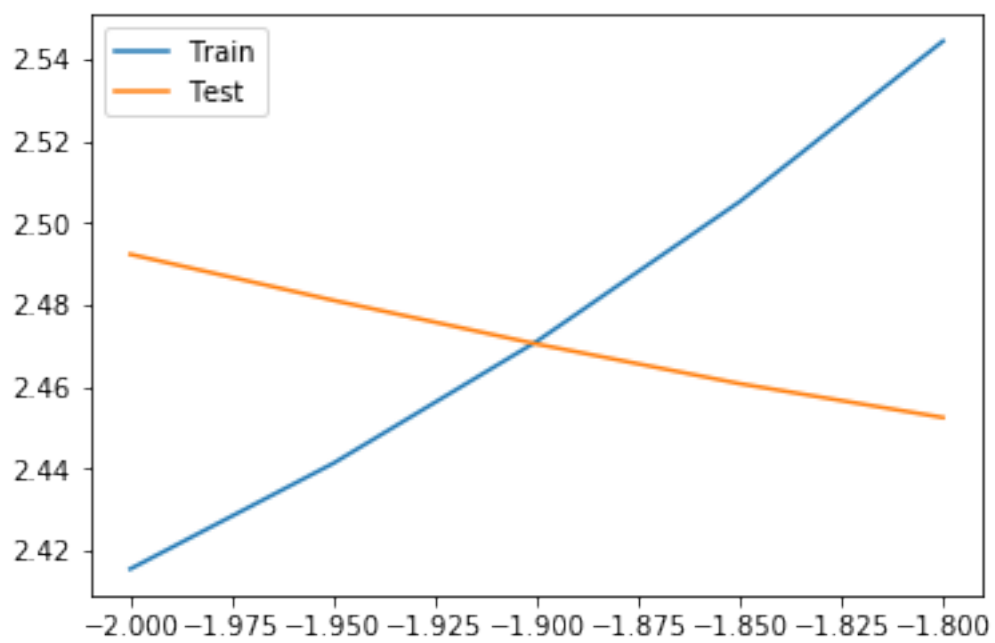
### **2.5.5    OPTIONAL**



### **2.5.6 OPTIONAL**

2.5.7





### 2.5.8

The  $\theta$  should be  $10^{-1.9}$  which is the intersection of train and test square loss.  
Reason: it minimizes the test error

## 2.6 Stochastic Gradient Descent

### 2.6.1

$$\begin{aligned} J(\theta) &= \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2 + \lambda \theta^T \theta \\ &= \frac{1}{m} \left[ \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2 + m \lambda \theta^T \theta \right] \\ &= \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x_i) - y_i]^2 + \lambda \theta^T \theta \\ f_i(\theta) &= (h_{\theta}(x_i) - y_i)^2 + \lambda \theta^T \theta \end{aligned}$$

### 2.6.2

$$\begin{aligned} E(\nabla f_i(\theta)) &= \sum (\nabla f_i(\theta) p(x_i)) \\ &= \sum (\nabla f_i(\theta) \frac{1}{m}) \\ &= \frac{1}{m} \sum \nabla f_i(\theta) \\ &= \nabla J(\theta) \end{aligned}$$

### 2.6.3

$$\begin{aligned}\theta_{k+1} &= \theta_k - \eta \nabla_{\theta} f_k(\theta) \\ &= \theta_k - \eta \frac{d((h_{\theta}(x_k) - y_k)^2 + \lambda \theta_k^T \theta_k)}{d\theta_k} \\ &= \theta_k - \eta \frac{d((\theta_k^T x_k - y_k)^2 + \lambda \theta_k^T \theta_k)}{d\theta_k} \\ &= \theta_k - 2\eta((\theta_k^T x_k - y_k)x_k + \lambda \theta_k)\end{aligned}$$

## 2.6.4

---

```
1 def compute_sgd_loss(X, y, theta, lambda_reg):
2     return (np.dot(theta,X)-y)**2 + lambda_reg*np.dot(theta,theta)
3 def compute_sgd_gradient(X, y, theta, lambda_reg):
4     return 2*(np.dot(theta, X)-y)*X + 2 * lambda_reg*theta
```

---

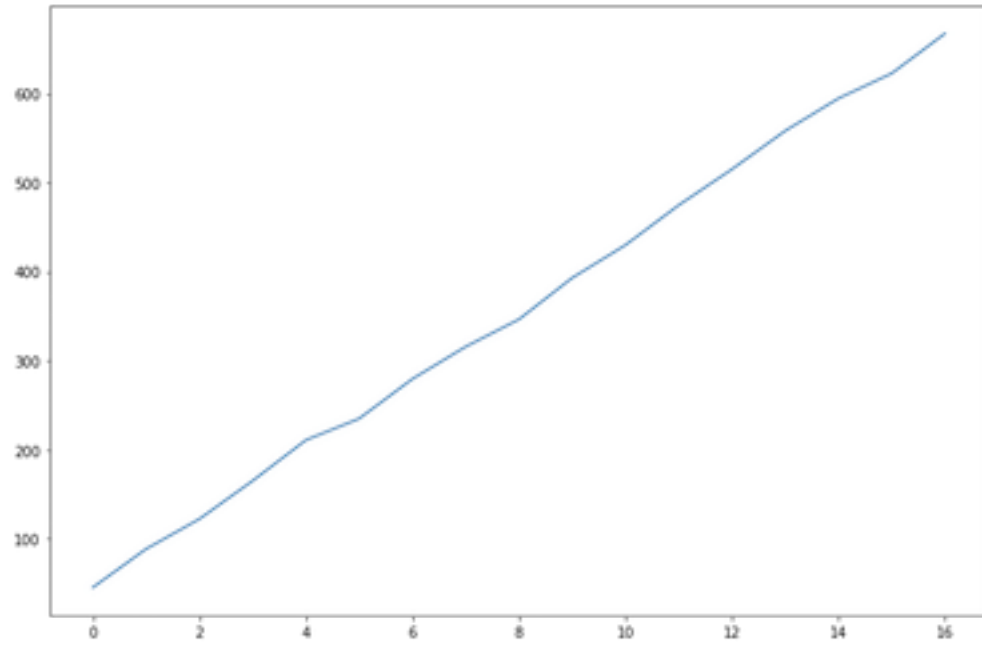
```
1 def stochastic_grad_descent(X, y, alpha=0.1, lambda_reg=1, num_iter=1000):
2     num_instances, num_features = X.shape[0], X.shape[1]
3     theta = np.ones(num_features) #Initialize theta
4
5
6     theta_hist = np.zeros((num_iter, num_instances, num_features)) #Initialize theta_hist
7     loss_hist = np.zeros((num_iter, num_instances)) #Initialize loss_hist
8     #TODO
9     # grad0 = compute_square_loss_gradient(X[row], y, theta_hist[i][row]) + 2*lambda_reg*theta_hist[i][row]
10    step_size = 10
11
12    for i in range(num_iter):
13        change_idx = np.arange(num_instances)
14        np.random.shuffle(change_idx)
15        for row in change_idx:
16            theta_hist[i][row] = theta
17            loss_hist[i][row] = compute_sgd_loss(X[row], y[row], theta, lambda_reg)
18            grad = compute_sgd_gradient(X[row], y[row], theta, lambda_reg)
19            if isinstance(alpha, float):
20                theta = theta - alpha*grad
21            elif alpha == '1/sqrt(t)':
22                theta = theta - 1./np.sqrt(step_size)*grad
23                step_size += 1
24            elif alpha == '1/t':
25                theta = theta - 1./step_size*grad
26                step_size += 1
27            else:
28                theta = theta - theta/(1+step_size*lambda_reg*theta)*grad
29                step_size += 1
30    return (theta_hist, loss_hist)
```

---

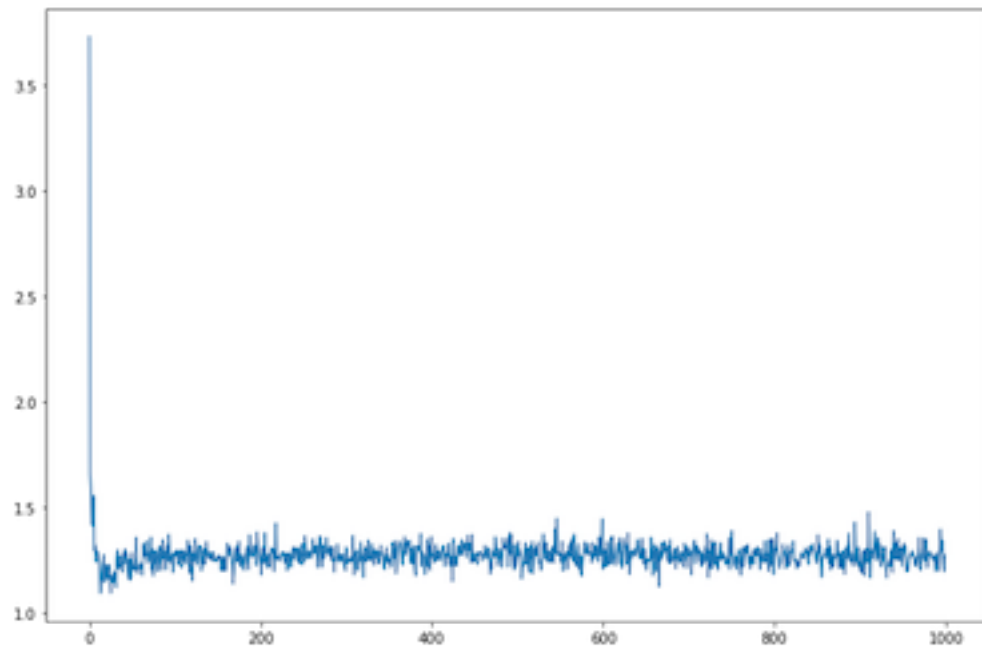


### 2.6.5

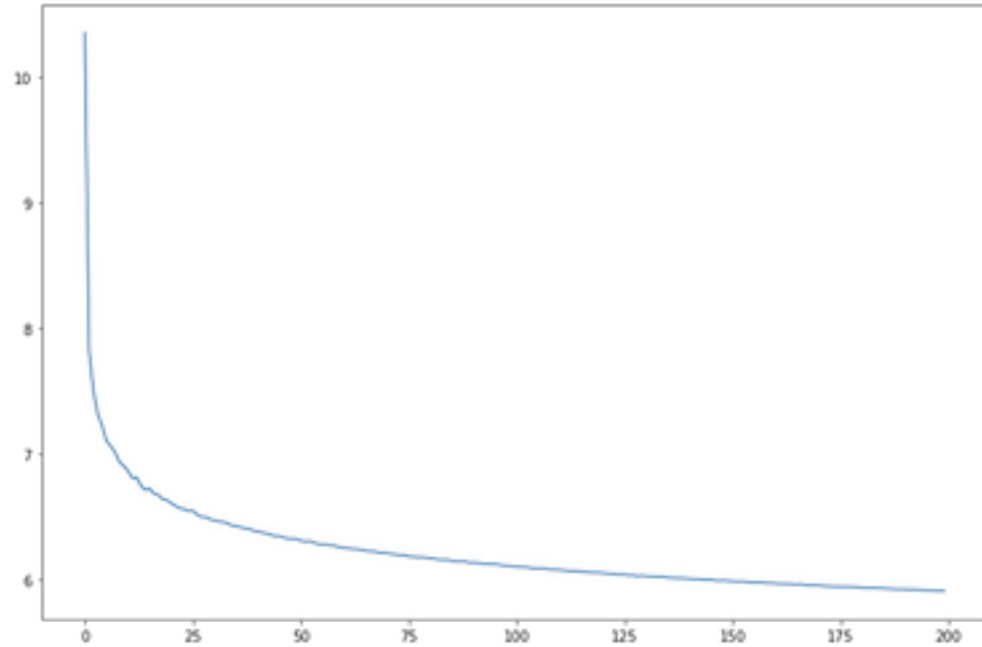
The stochastic gradient descent plot with fixed step size as 0.05:



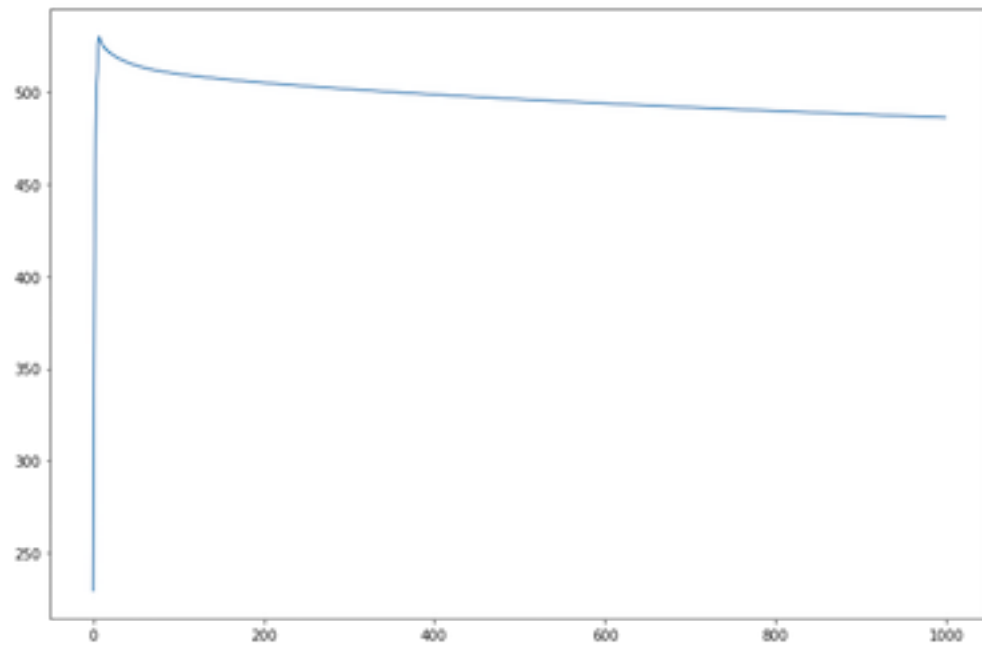
The stochastic gradient descent plot with fixed step size as 0.005:



When the step size is 0.05, the stochastic gradient descent cannot converge.  
The stochastic gradient descent plot with step size as  $1/t$ :



The stochastic gradient descent plot with step size as  $1/\sqrt{t}$ :



When step size is  $1/\sqrt{t}$ , the stochastic gradient descent cannot converge.

### **2.6.6 OPTIONAL**

### 3 Risk Minimization

#### 3.1 Square Loss

##### 3.1.1

$$\text{Risk} = R(a) = E[(a - y)^2]$$

$$\begin{aligned} E[(a - y)^2] &= E[(a - E(y) + E(y) - y)^2] \\ &= E[(a - E(y))^2 + 2(a - E(y))(E(y) - y) + (E(y) - y)^2] \end{aligned}$$

$$\begin{aligned} E[2(a - E(y))(E(y) - y)] &= 2E[aE(y) - ya - E(y)^2 + yE(y)] \\ &= 2[E(a)E(y) - E(y)E(a) - E(E(y))^2 + E(y)E(y)] \\ &= 2[E(a)E(y) - E(y)E(a) - E(y)^2 + E(y)E(y)] \\ &= 0 \end{aligned}$$

$$E[(a - y)^2] = E[a - E(y)]^2 + E[(E(y) - y)^2]$$

Since the second term in the expectation function is independent of  $a$ :

$$a^* = E(y)$$

Bayes risk:

$$\begin{aligned} R(a^*) &= E[(a^* - y)^2] \\ &= E[(E(y) - y)^2] \\ &= E(y^2 - 2yE(y) + (E(y))^2) \\ &= E(y^2) - 2E(y)E(y) + (E(y))^2 \\ &= E(y^2) - (E(y))^2 \\ &= \text{Var}(y) \end{aligned}$$

### 3.1.2 (a)

$$a^* = E(y|x)$$

### 3.1.3 (b)

$$\begin{aligned} E[(a - y)^2|x] &= E[(a - E(y|x) + E(y|x) - y)^2|x] \\ &= E(a - E(y|x))^2 + 2E[(a - E(y|x))(E(y|x) - y)|x] + E(E(y|x) - y)^2|x \end{aligned}$$

$$\begin{aligned} E[(a - E(y|x))(E(y|x) - y)|x] &= (a - E(y|x))(E(y|x) - E(y|x)) \\ &= 0 \end{aligned}$$

$$\begin{aligned} E[(a - y)^2|x] &= E[a - E(y|x)^2|x] + E[(E(y|x) - y)^2|x] \\ &= E[a - E(y|x)^2|x] + E[(a^* - y)^2|x] \\ &= E[f(x) - E(y|x)^2|x] + E[(f(x)^* - y)^2|x] \\ E[(f(x) - y)^2|x] &= E[f(x) - E(y|x)^2|x] + E[(f(x)^* - y)^2|x] \end{aligned}$$

$$E[(f(x)^* - y)^2|x] \leq E[(f(x) - y)^2|x]$$

$$\begin{aligned} E[E[(f(x)^* - y)^2|x]] &\leq E[E[(f(x) - y)^2|x]] \\ E[(f(x)^* - y)^2] &\leq E[(f(x) - y)^2] \end{aligned}$$

## 3.2 OPTIONAL