

1a)

$$RSS(\theta) = (Y - X\theta)^T (Y - X\theta)$$

$$= (Y^T - \theta^T X^T) (Y - X\theta)$$

$$= Y^T Y - \theta^T X^T Y - Y^T X \theta + \theta^T X^T X \theta$$

$$\frac{\partial}{\partial \theta} (RSS(\theta)) = -X^T Y - X^T Y + (X^T X + (X^T X)^T) \theta$$

$$= -2X^T Y + 2X^T X \theta$$

$$= 0$$

$$\theta_{\text{critical}} = (X^T X)^{-1} X^T Y$$

1b)

$$\text{Ridge}_\lambda(\theta) = \text{RSS}(\theta) + \lambda^2 \theta^T \theta$$

$$\frac{\partial}{\partial \theta} (\text{Ridge}_\lambda \theta) = \frac{\partial}{\partial \theta} (\text{RSS}(\theta)) + \lambda^2 \frac{\partial}{\partial \theta} (\theta^T \theta)$$

$$= -2X^T Y + 2X^T X \theta + 2\lambda^2 \theta$$

$$= -2X^T Y + 2(X^T X + \lambda^2 I) \theta$$

$$= 0$$

$$\Rightarrow \theta_{\text{critical}} = (X^T X + \lambda^2 I)^{-1} X^T Y$$

3a)

$$J(\vec{\theta}, \vec{x}) = \sum_{i=1}^n w^{(i)} \left(\vec{\theta}^T \vec{x}^{(i)} - y^{(i)} \right)^2$$

$$= (Y - X\theta)^T W (Y - X\theta)$$

where,

$$Y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

 $n \times 1$

$$X = \begin{bmatrix} 1 & x^{(1)(1)} & x^{(1)(2)} & \dots & x^{(1)(d)} \\ 1 & x^{(2)(1)} & x^{(2)(2)} & \dots & x^{(2)(d)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x^{(n)(1)} & x^{(n)(2)} & \dots & x^{(n)(d)} \end{bmatrix}$$

 $n \times (d+1)$

$$\theta = \begin{bmatrix} \theta^{(0)} \\ \theta^{(1)} \\ \theta^{(2)} \\ \vdots \\ \theta^{(d)} \end{bmatrix}$$

 $(d+1) \times 1$

$$W = \begin{bmatrix} w^{(1)} & & & 0 \\ & w^{(2)} & & \\ & & \ddots & \\ 0 & & & w^{(n)} \end{bmatrix}$$

 $n \times n$

$$J(\vec{\theta}, \vec{x}) = (Y^T - \theta^T x^T) (WY - WX\theta)$$

$$= Y^T W Y - Y^T W X \theta - \theta^T x^T W Y + \theta^T x^T W X \theta$$

(a)

$$\begin{aligned} \frac{\partial}{\partial \theta} (J(\vec{\theta}, \vec{x})) &= 2 X^T W X \theta - 2 X^T W Y \\ &= 2 X^T W (X\theta - Y) \end{aligned}$$

(b)

$$\begin{aligned} \frac{\partial^2}{\partial \theta^2} (J(\vec{\theta}, \vec{x})) &= 2 (X^T W X)^T \\ &= 2 X^T W^T X \\ &= 2 X^T W X \end{aligned}$$

(c)

$$\begin{aligned}\theta^{(t+1)} &= \theta^{(t)} - 2\eta X^T W (X \theta^{(t)} - Y) \\ &= (I_n - 2\eta X^T W X) \theta^{(t)} + 2\eta X^T W Y\end{aligned}$$

(d)

$$\begin{aligned}\theta^{(t+1)} &= \theta^{(t)} - \frac{1}{2} (X^T W X)^{-1} 2 X^T W (X \theta^{(t)} - Y) \\ &= \theta^{(t)} - \theta^{(t)} + (X^T W X)^{-1} X^T W Y \\ &= (X^T W X)^{-1} X^T W Y\end{aligned}$$

④
(1)

$$f(x) = \beta_0 + \beta_1 \sin(x) + \beta_2 \cos(x)$$

$$X = \begin{bmatrix} 1 & \sin(x_1) & \cos(x_1) \\ 1 & \sin(x_2) & \cos(x_2) \\ \vdots & \vdots & \vdots \\ 1 & \sin(x_n) & \cos(x_n) \end{bmatrix}$$

$\sin(x)$ and $\cos(x)$ are linearly independent. Therefore, we have design matrix with linearly independent columns. Hence, we can use least squares method.

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$\hat{y} = X\beta$$

$$\beta = (X^T X)^{-1} X^T Y$$

(2)

$$g(x) = \beta_0 + \sin(\beta_2 x) + \cos(\beta_2 x)$$

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}, \quad \text{RSS} = \sum_{i=1}^n (g(x^{(i)}) - y^{(i)})^2$$

$$\beta^{(t+1)} = \beta^{(t)} - \eta \left(\frac{\partial (\text{RSS})}{\partial \beta} \right)^{(t)}$$

learning rate

$$\frac{\partial (\text{RSS})}{\partial \beta} = \begin{bmatrix} \frac{\partial (\text{RSS})}{\partial \beta_0} \\ \frac{\partial (\text{RSS})}{\partial \beta_1} \\ \frac{\partial (\text{RSS})}{\partial \beta_2} \end{bmatrix}$$

$$\frac{\partial (RSS)}{\partial \beta_j} = \sum_{i=1}^n 2(g(x^{(i)}) - y^{(i)}) \cdot \frac{\partial (g(x^{(i)}))}{\partial \beta_j}$$

$$\frac{\partial (g(x^{(i)}))}{\partial \beta_0} = 1$$

$$\frac{\partial (g(x^{(i)}))}{\partial \beta_1} = x^{(i)} \cos(\beta_1 x^{(i)})$$

$$\frac{\partial (g(x^{(i)}))}{\partial \beta_2} = -x^{(i)} \sin(\beta_2 x^{(i)})$$

Substitute the
3 terms
in above equation
to get
derivatives
with β_j
for $j=0, 1, 2$

Homework2 - Sai Nikhil

February 16, 2022

1 Problem 2

```
[1]: import numpy as np
import matplotlib.pyplot as plt
```

```
[2]: x = np.array([1.2, 3.2, 5.1, 3.5, 2.6]).reshape((-1, 1))
y = np.array([7.8, 1.2, 6.4, 2.6, 8.1]).reshape((-1, 1))
```

```
[3]: def normal_equation(x, y, lmd=0):
    n, d = x.shape[0], x.shape[1]
    ones = np.ones(n).reshape(-1, 1)
    x = np.hstack((ones, x))
    return np.linalg.inv(x.T @ x + (lmd**2) * np.eye(d + 1)) @ x.T @ y
```

```
[4]: lmd = [0, 1, 10]
theta = []
print("2a)")
for i, l in enumerate(lmd):
    theta.append(normal_equation(x, y, l))
    print(f"Equation for lambda = {l} is y = {theta[i][1][0]} * x + \u2192{theta[i][0][0]}")
```

2a)

Equation for lambda = 0 is $y = -0.6766317887394105 * x + 7.331091180866961$

Equation for lambda = 1 is $y = 0.47491248541423553 * x + 3.1152275379229852$

Equation for lambda = 10 is $y = 0.4671668474177275 * x + 0.1791637826693662$

```
[5]: x1 = np.linspace(1, 6, 100).reshape((-1, 1))
n, d = x1.shape[0], x1.shape[1]
ones = np.ones(n).reshape(-1, 1)
x1 = np.hstack((ones, x1))
y1 = x1 @ theta[0]

x2 = np.linspace(1, 6, 100).reshape((-1, 1))
n, d = x2.shape[0], x2.shape[1]
ones = np.ones(n).reshape(-1, 1)
x2 = np.hstack((ones, x2))
y2 = x2 @ theta[1]
```

```

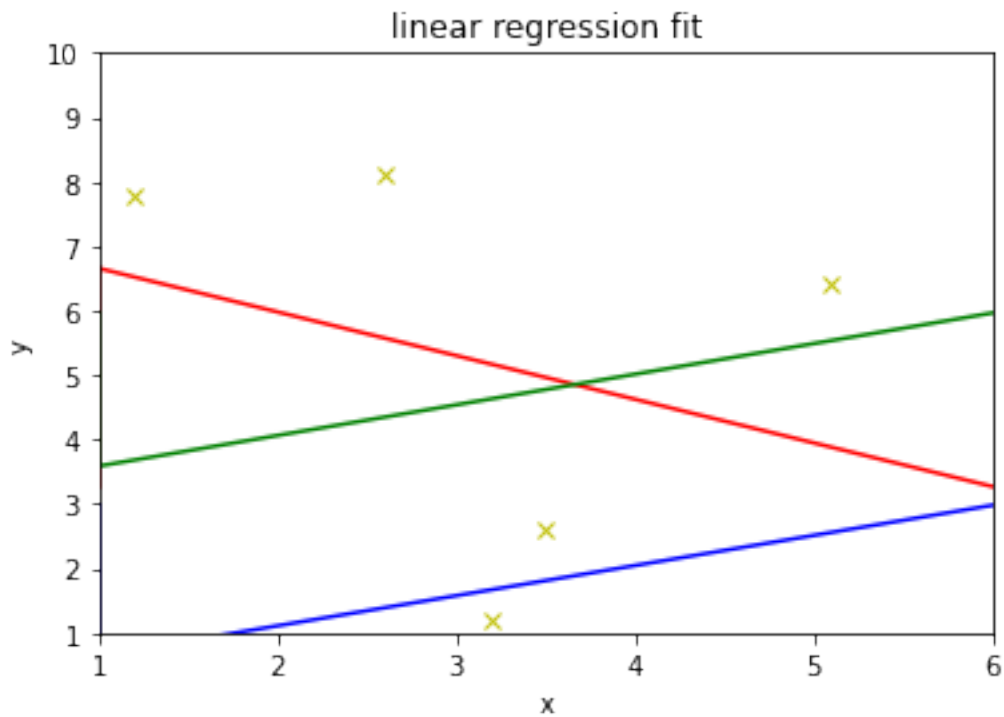
x3 = np.linspace(1, 6, 100).reshape((-1, 1))
n, d = x3.shape[0], x3.shape[1]
ones = np.ones(n).reshape(-1, 1)
x3 = np.hstack((ones, x3))
y3 = x3 @ theta[2]

```

```

plt.plot(x1, y1, '-r')
plt.plot(x2, y2, '-g')
plt.plot(x3, y3, '-b')
plt.plot(x, y, 'yx')
plt.axis([1, 6, 1, 10])
plt.xlabel('x')
plt.ylabel('y')
plt.title('linear regression fit')
plt.show()

```



1.1 2b

From the above plots, it is clear that when $\lambda = 0$ (or no regularization), it is overfitting the outliers (Red line). When $\lambda = 1$, it is just fitting the data (Green line) and not sensitive to outliers. When $\lambda = 10$, there is too much of regularization. Hence, it is underfitting the data (Blue line).

2 Problem 4

```
[6]: def get_gradient(beta, g, x, y):  
    del_g = np.array([x * np.cos(beta[1][0] * x), -x * np.sin(beta[2][0] * x)]).  
    ↪reshape((10, 2))  
    n = x.shape[0]  
    del_g = np.hstack((del_g, np.ones(n).reshape(-1, 1)))  
    return np.sum(2 * (g(beta, x) - y) * del_g, axis=0).reshape((-1, 1))
```

```
[7]: def g(beta, x):  
    return beta[0][0] + np.sin(beta[1][0] * x) + np.cos(beta[2][0] * x)
```

```
[8]: def RSS(beta, g, x, y):  
    return np.sum((g(beta, x) - y)**2)
```

```
[9]: def gradient_descent(g, x, y, lr=0.01, iterations=100, threshold=0.001):  
    beta = np.array([np.random.rand(), np.random.rand(), np.random.rand()]).  
    ↪reshape((-1, 1))  
    beta_prev = beta  
  
    losses = []  
    for i in range(iterations):  
        gradient = get_gradient(beta_prev, g, x, y)  
        beta = beta_prev - lr * gradient  
        if np.all(np.abs(beta - beta_prev) < threshold):  
            print(f"Convergence found at iteration {i}.")  
            break  
        beta_prev = beta  
        losses.append(RSS(beta, g, x, y))  
  
    return losses, beta
```

```
[10]: x = np.array([0, 2, 4, 6, 8, 10, 12, 14, 16, 18])  
y = np.array([2.85, 1.5, 0.49, 1.57, 1.9, 0.6, 0.38, 2.33, 1.65, 0.3])  
x = x.reshape((-1, 1))  
y = y.reshape((-1, 1))
```

```
[11]: losses, beta = gradient_descent(g, x, y, lr=0.0001, iterations=50, threshold=0.  
    ↪00001)
```

```
[12]: plt.plot(losses, 'rx')  
plt.xlabel('iterations')  
plt.ylabel('losses')  
plt.title('iterations vs loss plot')  
plt.show()
```



[]: