# Assignment 1

## Question 2

1. Show that ridge regression can be rewritten as a non regularized linear regression

$$\min_{w \in \mathbb{R}^a, b \in \mathbb{R}} \quad \frac{1}{2n} \| Xw + b1 - y \|_2^2 + \lambda \| w \|_2^2$$

$$= \frac{1}{2n} \left\{ \| Xw + b1 - y \|^2 + (\sqrt{2n\lambda})^2 \| w \|_2^2 \right\}$$

$$= \frac{1}{2n} \left( \underbrace{(Xw + b1 - y)^T (Xw + b1 - y)}_{\text{dot product}} + \underbrace{\sqrt{2n\lambda} w^T (\sqrt{2n\lambda} w)}_{\text{dot product}} \right)$$

This is similar to a dot product

$$= \frac{1}{2n} \left( \begin{bmatrix} Xw + b1 - y \\ \sqrt{2n\lambda} \, I_d \, w \end{bmatrix}^T \begin{bmatrix} Xw + b1 - y \\ \sqrt{2n\lambda} \, I_d \, w \end{bmatrix} \right)$$

$$= \frac{1}{2n} \left\| \begin{bmatrix} Xw + b1 - y \\ \sqrt{2n\lambda} \, I \, w + 0_d \end{bmatrix} \right\|_2^2 \quad \text{since } w \text{ is of size } d$$

$$= \frac{1}{2n} \left\| \begin{bmatrix} Xw + b1 \\ \sqrt{2n\lambda} \, w \end{bmatrix} - \begin{bmatrix} y \\ 0_d \end{bmatrix} \right\|_2^2$$

$Iw = w$ since $I$ is identity matrix of size $d$

$$= \frac{1}{2n} \left\| \begin{bmatrix} Xw + b1_n \\ \sqrt{2n\lambda} \, I_d \, w \end{bmatrix} - \begin{bmatrix} y \\ 0_d \end{bmatrix} \right\|_2^2$$

$1 \rightarrow 1_n$ because $y \in \mathbb{R}^n$

$$= \frac{1}{2n} \left\| \begin{bmatrix} X & 1_n \\ \sqrt{2n\lambda} \, I_d & 0_d \end{bmatrix} \begin{bmatrix} w \\ b \end{bmatrix} - \begin{bmatrix} y \\ 0_d \end{bmatrix} \right\|_2^2$$

$$\therefore \text{Proved}$$

2.

$$\frac{\partial}{\partial w} \frac{1}{2n} \|Xw + b1 - y\|_2^2 + \lambda \|w\|_2^2$$

$$\text{1st term} = \frac{1}{2n}\left((Xw + b1 - y)^T (Xw + b1 - y)\right)$$

$$= \frac{1}{2n}(w^TX^T + 1^Tb - y^T)(Xw + b1 - y)$$

$$= \frac{1}{2n}(w^TX^TXw + 1^Tb\,Xw - y^TXw + w^TX^Tb1 + b1^Tb$$
$$\quad - y^Tb1 - w^TX^Ty - 1^Tb\,y + y^Ty)$$

Now derive w.r.t $w$

$$= \frac{1}{2n}\left(\frac{\partial}{\partial w}\left(w^TX^TXw + b^TXw - y^TXw + w^TX^Tb - w^TX^Ty\right.\right.$$

$$(b^TXw)^T = w^TX^Tb \qquad (y^TXw)^T = w^TX^Ty$$

$$= \frac{1}{2n}\left(\frac{\partial}{\partial w}\left(w^TX^TXw + 2b\,X^Tw\,w^TX^Tb - 2w^TX^Ty\right.\right.$$

$$= \frac{1}{2n}\left(\frac{\partial}{\partial w}\left(w^TX^TXw + 2(X^Tb)^Tw - 2(X^Ty)^Tw\right)\right)$$

$$= \frac{1}{2n}\left(\frac{\partial(w^TX^TXw)}{\partial w} + 2X^Tb - 2X^Ty\right)$$

$$\frac{\partial}{\partial w}(w^TX^TXw) = w^T(\cancel{X^TX} + (X^TX)^T)$$

Use product rule of matrix partial derivative

$$= \frac{1}{2n}(2X^TXw + 2X^Tb - 2X^Ty)$$

$$= \frac{1}{n}X^T(Xw + b - y)$$

2nd term $\Rightarrow \dfrac{\partial}{\partial w}\lambda\|w\|_2^2 = 2\lambda w$

Combining terms

$$= \frac{1}{n} X^T (Xw + b1 - y) + 2\lambda w$$

$$\therefore \text{Proved}$$

$$\frac{\partial}{\partial b} \frac{1}{2n} \|Xw + b1 - y\|_2^2 + \lambda \|w\|_2^2$$

$$= \frac{1}{2n} \frac{\partial}{\partial b} \|Xw + b1 - y\|_2^2$$

$$= \frac{1}{2n} \frac{\partial}{\partial b} (Xw + b1 - y)^T (Xw + b1 - y)$$

$$= \frac{1}{2n} \frac{\partial}{\partial b} (w^T X^T + b^T - y^T)(Xw + b - y)$$

$$= \frac{1}{2n} \frac{\partial}{\partial b} (w^T X^T Xw + b^T Xw - y^T Xw + w^T X^T b + b^T b$$
$$- y^T b - w^T X^T y - b^T y + y^T y)$$

$$= \frac{1}{2n} \frac{\partial}{\partial b} (2b^T Xw + b^T b - 2y^T b)$$

$$= \frac{1}{2n} \frac{\partial}{\partial b} (2(Xw)^T b + 1^T b^2 - 2yb^T)$$

$$= \frac{1}{2n} \frac{\partial}{\partial b} (2b1^T Xw + 1^T b^2 - 2b1^T y)$$

$$= \frac{1}{n} 1^T (Xw + b - y)$$

$$\therefore \text{Proved}$$

Question 2
3. Algorithm Submitted
4. Algorithm submitted

5. Comparing run time of 2 approaches
   For $\lambda = 10$
   Runtime of closed form approach = 0.0003
   Runtime of gradient descent approach = 0.442
   $\lambda = 0$
   Closed form Training error = 0.062297
            Training loss = 0.062297
            Test error = 0.528986

   $\lambda = 10$
   Closed form Training error = 0.41589
            Training loss = 0.45481
            Test error = 0.44393

   $\lambda = 0$
   Gradient Descent Training error = 0.062595
            Training loss = 0.062595
            Test error = 0.050952

   $\lambda = 10$
   Gradient Descent Training error = 0.41642
            Training loss = 0.45481
            Testing error = 0.44419

   I believe $\lambda = 10$, closed form is a better approach
   because it not only has a faster runtime but
   also a smaller testing error

   $\therefore$ For this dataset, $\lambda = 10$ closed form approach
       is better.