Problem 1

For the multiple regression model:

$$f(\boldsymbol{x}) = \beta_0 + \sum_{j=1}^n \beta_j x_j$$

where x_1, \ldots, x_n are inputs and β_0, \ldots, β_n are coefficients of the model. Suppose the training set include N samples: $(\boldsymbol{x}_1, y_1), (\boldsymbol{x}_2, y_2), \ldots, (\boldsymbol{x}_N, y_N)$, where $\boldsymbol{x}_i = (x_{i1}, \ldots, x_{in})^T$.

(a) Show that the residual sum of squares

$$RSS(\boldsymbol{\beta}) = \sum_{i=1}^{N} \left(y_i - \beta_0 - \sum_{j=1}^{n} \beta_j x_{ij} \right)^2$$
 (1)

can be written in matrix form as

$$RSS(\boldsymbol{\beta}) = (\boldsymbol{y} - X\boldsymbol{\beta})^T (\boldsymbol{y} - X\boldsymbol{\beta})$$

where $\mathbf{y} = (y_1, y_2, \dots, y_N)^T$, $\boldsymbol{\beta} = (\beta_0, \dots, \beta_n)^T$, and X is the matrix:

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1n} \\ 1 & x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{N1} & x_{N2} & \cdots & x_{Nn} \end{bmatrix}$$

(b) Show that the partial derivative of RSS with respect to β is:

$$\frac{\partial \text{RSS}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = -2X^T \boldsymbol{y} + 2X^T X \boldsymbol{\beta}$$

Use the following matrix calculus results:

$$\frac{\partial \boldsymbol{x}^T A \boldsymbol{x}}{\partial \boldsymbol{x}} = 2A\boldsymbol{x}$$

if A is symmetric.

(c) Show that the least squares solution (normal equation):

$$\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \boldsymbol{y}$$

minimizes RSS.

1) (a)
$$f(x) = \beta_0 + \sum_{j=1}^{n} \beta_j X_j$$

$$X = \begin{bmatrix} 1 & X_{1,1} & X_{1,2} & \cdots & X_{1,n} \\ 1 & X_{2,1} & X_{2,2} & \cdots & X_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{N,1} & X_{N,2} & \cdots & X_{N,n} \end{bmatrix}$$

$$RSS(\beta) = \sum_{i=1}^{N} \left(y_i - \beta_0 - \sum_{j=1}^{n} \beta_j x_{i,j} \right)^2$$

$$-f(x)$$

$$RSS(\beta) = \sum_{i=1}^{N} (y_i - (\beta_0 + \sum_{j=1}^{n} \beta_j X_{i,j}))^2$$

$$X\beta = \beta_0 + \sum_{j=1}^{n} \beta_j X_j$$

$$RSS(\beta) = \sum_{i=1}^{N} (Y_i - X_i)^2$$

$$\Gamma = (Y_i - X_i) \rightarrow \Gamma^2 = \Gamma^T \Gamma$$

$$RSS(\beta) = (Y - x\beta)^{2}$$

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

(b)
$$\frac{\partial x^T A x}{\partial x} = 2A x$$

$$(y - x\beta)^{T}(y - x\beta) = y^{T}y - \beta^{T}x^{T}y - y^{T}x\beta + \beta^{T}x^{T}x\beta$$

$$= y^{T}y - 2\beta^{T}x^{T}y + \beta^{T}x^{T}x\beta$$

$$\frac{\partial RSS(\beta)}{\partial \beta} = -2x^{T}y + 2x^{T}x\beta$$

(C)

$$-2x^{T}y + 2x^{T}x\beta = 0$$

$$\Rightarrow -x^{T}y + x^{T}x\beta = 0$$

$$\therefore \times^{T} \times \beta = \times^{T} y$$

$$\Rightarrow (X^{T} \times)^{-1} \times^{T} \times \beta = (X^{T} \times)^{-1} \times^{T} y$$
Thus $\hat{\beta} = (X^{T} \times)^{-1} \times^{T} y$

Problem 2

Consider using Ridge Regression for modeling. Use the following form of cost function (equivalent to Eqn. 5):

$$J(\beta) = \sum_{i=1}^{N} \left(y_i - \beta_0 - \sum_{j=1}^{n} \beta_j x_{ij} \right)^2 + \alpha \sum_{i=1}^{n} \beta_i^2$$

where all the notations from Problem 1 are still valid here.

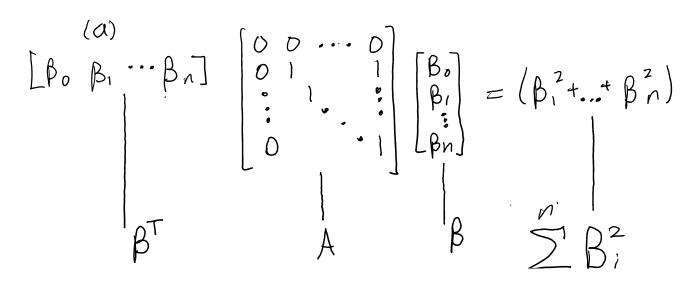
(a) Show that $\sum_{i=1}^{n} \beta_i^2$ can be written in matrix form as:

$$\sum_{i=1}^{n} \beta_i^2 = \boldsymbol{\beta}^T A \boldsymbol{\beta}$$

where A is the $(n+1) \times (n+1)$ identity matrix except with a 0 in the top-left cell.

(b) Show that the closed-form solution for the Ridge Regression is given by

$$\hat{\boldsymbol{\beta}} = (X^T X + \alpha A)^{-1} X^T \boldsymbol{y}$$



$$\beta^{T} \qquad A \qquad \beta \qquad \sum_{i=1}^{n} \beta_{i}^{2}$$

$$\beta^{T}A\beta = \sum_{i=1}^{n} \beta_{i}^{2}$$

$$(b) \qquad \beta = (X^{T}X + \alpha A)^{-1} X^{T}y$$

$$\sigma(\beta) = \sum_{i=1}^{n} (y_{i} - \beta_{0} - \sum_{i=1}^{n} \beta_{i} \times y_{i})^{2} + \alpha \sum_{i=1}^{n} \beta_{i}^{2}$$

$$= (y - x\beta^{T})(y - x\beta) + \alpha \beta^{T}A\beta$$

$$= y^{T}y - \beta^{T}X^{T}y - y^{T}x\beta + \beta^{T}X^{T}x\beta + \alpha \beta^{T}A\beta$$

$$= y^{T}y - \beta^{T}X^{T}y + \beta^{T}X^{T}x\beta + \alpha \beta^{T}A\beta$$

$$\beta(\sigma(\beta) = y^{T}y - 2\beta^{T}X^{T}y + \beta^{T}x^{T}x\beta + \alpha \beta^{T}A\beta)$$

$$\frac{\partial \sigma(\beta)}{\partial \beta} = -2 \times y + 2 \times x + 2 \times x$$

Problem 3 (MATH 5388 Only)

Consider the Ridge Regression problem where the cost function is given in Problem 2. Show that the problem is equivalent to the problem

$$\hat{\beta}^{c} = \arg\min_{\beta^{c}} \left\{ \sum_{i=1}^{N} \left(y_{i} - \beta_{0}^{c} - \sum_{j=1}^{n} (x_{ij} - \bar{x}_{j}) \beta_{j}^{c} \right)^{2} + \alpha \sum_{j=1}^{n} (\beta_{j}^{c})^{2} \right\}$$

where \bar{x}_{i} denotes the mean of the jth feature of all samples:

$$\bar{x}_j = \frac{1}{N} \sum_{i=1}^{N} x_{ij}$$

Give the correspondence between β^c and the original β in the cost function $J(\beta)$ in Problem 2.

$$\mathcal{J}(\beta) = \sum_{i=1}^{n} (y_i - \beta_0 - \sum_{j=1}^{n} \beta_j \times ij)^2 + \alpha \sum_{i=1}^{n} \beta_i^2$$
to center data:

Let $\times ij = \times ij - \times j$

Then $\beta = \overline{y}$

Then
$$\beta_0 = \overline{y}$$

lecause $\beta_0 = \beta_0 + \overline{\beta}_0 \times \overline{y} = \overline{y} \times \overline{y}$
 $\beta_0 = \beta_0$
 $\beta_0 = \beta_0$