

Machine Learning Homework 1

Hunter Carroll

January 2024

Problem 1.a

For the multiple regression model:

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

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

(a) Show that the residual sum of squares,

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

can be written in matrix form as,

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

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

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

My solution:

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

Let \mathbf{y} be a column vector of size $N \times 1$,

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

Let $\boldsymbol{\beta}$ be a column vector of size $(n + 1) \times 1$,

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{bmatrix}$$

Finally, let X be a matrix of size $N \times (n + 1)$,

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

Then,

$$X\boldsymbol{\beta} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{N1} & x_{N2} & \dots & x_{Nn} \end{bmatrix} \cdot \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{bmatrix} = \begin{bmatrix} \beta_0 + x_{11}\beta_1 + x_{12}\beta_2 + \dots + x_{1n}\beta_n \\ \beta_0 + x_{21}\beta_1 + x_{22}\beta_2 + \dots + x_{2n}\beta_n \\ \vdots \\ \beta_0 + x_{N1}\beta_1 + x_{N2}\beta_2 + \dots + x_{Nn}\beta_n \end{bmatrix}$$

Note then the result of $X\boldsymbol{\beta}$ is a column vector of size $N \times 1$ containing

our predicted values. Let $\mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_N \end{bmatrix}$ be the $N \times 1$ column vector containing

the residuals that is, the actual values represented by column vector \mathbf{y} and predicted values represented as the column vector $X\boldsymbol{\beta}$ such that the residuals are the difference between the two column vectors,

$$\mathbf{u} = \mathbf{y} - X\boldsymbol{\beta}$$

Therefore, we can rewrite the Residual Sum of Squares in matrix form as follows,

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

Thus, $\mathbf{u}^T \cdot \mathbf{u}$ is the sum of the squared elements of \mathbf{u} (Residuals).

Problem 1.b

Show that the partial derivative of RSS with respect to β is,

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

Use the following matrix calculus results:

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

If A is symmetric

My solution:

Given $RSS(\beta) = (\mathbf{y} - X\beta)^T(\mathbf{y} - X\beta)$ we can expand this to get,

$$= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T X \beta - \beta^T X^T \mathbf{y} + \beta^T X^T X \beta$$

also note that $\mathbf{y}^T X \beta = \beta^T X^T \mathbf{y}$ this is because we have that both are equivalent scalars that are $\mathbf{y}^T X \beta = (1 \times N)(N \times (n+1))((n+1) \times 1)$ which results in a row vector $(1 \times (n+1))$ times column vector $((n+1) \times 1)$ and thus a scalar and taking the transpose of a scalar is the scalar itself. Implying that $\beta^T X^T \mathbf{y}$ is an equivalent expression. We can now simplify the $RSS(\beta)$ to,

$$RSS(\beta) = \mathbf{y}^T \mathbf{y} - 2\beta^T X^T \mathbf{y} + \beta^T X^T X \beta$$

When taking the partial derivative of the $RSS(\beta)$ the first term ($\mathbf{y}^T \mathbf{y}$) goes away since with respect to β this is considered a constant. The second term ($-2\beta^T X^T \mathbf{y}$) when differentiated with respects to β results in $-2X^T \mathbf{y}$. Finally, the third term we know that $X^T X$ will result in a symmetric $(n+1) \times (n+1)$ matrix thus we can apply the fact that,

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

giving,

$$\frac{\partial \beta^T X^T X \beta}{\partial \beta} = 2X^T X \beta$$

Therefore,

$$\frac{\partial}{\partial \beta} (\mathbf{y}^T \mathbf{y}) = 0 \tag{1}$$

$$\frac{\partial}{\partial \beta} (-2\beta^T X^T \mathbf{y}) = -2X^T \mathbf{y} \tag{2}$$

$$\frac{\partial}{\partial \beta} (\beta^T X^T X \beta) = 2X^T X \beta \tag{3}$$

Combining the above results we find,

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

Problem 1.c

Show that the least squares solution (normal equation):

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

minimizes RSS .

My solution:

First, we set the partial derivative of the Residual Sum of Squares with respect to beta equal to zero to find the value β that minimizes the RSS ,

$$\frac{\partial RSS(\beta)}{\partial \beta} = 0$$

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

Then we can factor a 2 out,

$$2(-X^T \mathbf{y} + X^T X \beta) = 0$$

Resulting in,

$$-X^T \mathbf{y} + X^T X \beta = 0$$

Move $-X^T \mathbf{y}$ to the right hand side,

$$X^T X \beta = X^T \mathbf{y}$$

Given $X^T X$ is symmetric and the determinant is non-zero, that is the columns are not linearly dependent ensuring full rank, then we can multiply each side by the inverse of $(X^T X)$,

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

However, a matrix multiplied by its inverse is simply the identity matrix. Therefore, the β that minimizes the RSS is,

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

as desired.

Problem 2

Consider using Ridge Regression for modeling. Use the following form of cost function:

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

Problem 2.a

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

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

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

My solution:

Let β be a column vector of size $(n+1) \times 1$,

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{bmatrix}$$

Let A be a square matrix of form $(n+1) \times (n+1)$,

$$A = \begin{bmatrix} 0 & 0 & \cdots & 0 \\ 0 & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & 1 \end{bmatrix}$$

Then note the term $\alpha \sum_{i=1}^n \beta_i^2$ will serve as a regularization term ensuring that we keep model weights small. Further, the reason we have 0 in the top-left cell of A is to exclude β_0 in the regularization term given our sum starts at $i = 1$. We know that β is an $(n+1) \times 1$ column vector. Thus, we want to find the value of the sum of our β^2 excluding β_0 . In order to do this we will transpose β resulting in a $1 \times (n+1)$ row vector resulting in the following when multiplied by A ,

$$[\beta_0 \ \beta_1 \ \dots \ \beta_n] \begin{bmatrix} 0 & 0 & \cdots & 0 \\ 0 & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & 1 \end{bmatrix} = [0 \ \beta_1 \ \beta_2 \ \dots \ \beta_n]$$

We now have that $\beta^T A$ is a $1 \times (n+1)$ row vector excluding β_0 which is now compatible to multiply with β the $(n+1) \times 1$ column vector. The result is a scalar that represents $\sum_{i=1}^n \beta_i^2$,

$$[0 \ \beta_1 \ \beta_2 \ \dots \ \beta_n] \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_n \end{bmatrix} = 0 \cdot \beta_0 + \beta_1^2 + \dots + \beta_n^2$$

Therefore, we can represent $\sum_{i=1}^n \beta_i^2$ in matrix form as $\beta^T A \beta$.

Problem 2.b

Show the closed-form solution for the Ridge Regression is given by,

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

My solution:

Given Residual Sum of Squares in matrix form is,

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

Then the partial derivative of $RSS(\beta)$ with respects to β is,

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

Now, Ridge Regression is similar to RSS but now we add a regularization term $\alpha \sum_{i=1}^n \beta_i^2$ which can be written as $\alpha \beta^T A \beta$ where α is some constant, then we have that Ridge Regression in matrix form is,

$$J(\beta) = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T X \beta - \beta^T X^T \mathbf{y} + \beta^T X^T X \beta + \alpha \beta^T A \beta$$

Taking the partial derivative of $J(\beta)$ with respects to β and using the fact that A is symmetric,

$$\frac{\partial J(\beta)}{\partial \beta} = -2X^T \mathbf{y} + 2X^T X \beta + 2\alpha A \beta$$

Set this equal to zero to find $\hat{\beta}$,

$$-2X^T\mathbf{y} + 2X^TX\beta + 2\alpha A\beta = 0$$

Factor out 2 and multiply it over,

$$-X^T\mathbf{y} + X^TX\beta + \alpha A\beta = 0$$

Move $-X^T\mathbf{y}$ over,

$$X^TX\beta + \alpha A\beta = X^T\mathbf{y}$$

Isolate β ,

$$\beta(X^TX + \alpha A) = X^T\mathbf{y}$$

Apply the inverse of $(X^TX + \alpha A)$ to both sides. Then $J(\beta)$ is minimized for,

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