$X$ data matrix

N rows (samples)

K+1 columns (feature)

$K+1^{st}$ column is $\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$

$Y$ target data $N \times 1$ column vector

Goal: find $K+1 \times 1$ vector $M = \begin{pmatrix} M_1 \\ \vdots \\ M_{K+1} \end{pmatrix}$

so that

$$E = \| Y - XM \|^2 \quad \text{is minimized.}$$

$$MSE(M_1, \ldots, M_{K+1}) = \frac{1}{N} E = \frac{1}{N} \sum_{i=1}^{N} \left( y_i - \underbrace{\sum_{j=1}^{K+1} x_{ij} M_j}_{} \right)^2$$

K+1 variables $M_1, \ldots, M_{K+1}$ $\qquad \| Y - Xu \|^2$

$$\nabla E = \begin{bmatrix} \partial E / \partial M_1 \\ \partial E / \partial M_2 \\ \vdots \\ \partial E / \partial M_{K+1} \end{bmatrix} = 0$$

$A = (a_{ij}) \quad r \times S$
$B = (b_{ij}) \quad S \times t$

$\qquad C_{ij} = \sum_{k=1}^{S} a_{ik} b_{kj}$

$C = (C_{ij}) \quad r \times t$

$\qquad A = (a_{ij})$
$\qquad A^T = (a_{ji}) \quad S \times r \text{ matrix}$

$$A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \qquad A^T = \begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix}$$

**Proposition:** The gradient of $MSE(M) = E$ is given by

*(handwritten, top right)* $E = \|Y - XM\|^2$

$X \quad N*(K+1) \quad M \quad (K+1)*1$
$Y \quad N*1$

$$\nabla E = \begin{bmatrix} \frac{\partial}{\partial M_1}E \\ \frac{\partial}{\partial M_2}E \\ \vdots \\ \frac{\partial}{\partial M_{\text{KM}}}E \end{bmatrix} = -2X^\mathsf{T}Y + 2X^\mathsf{T}XM \tag{5}$$

*(handwritten)* $= -2\begin{bmatrix} X^\mathsf{T}Y - X^\mathsf{T}XM \end{bmatrix}$

$(K+1)\times 1 \text{ matrix}$

where $X^\mathsf{T}$ is the transpose of $X$.

**Proof:** First, remember that the $ij$ entry of $X^\mathsf{T}$ is the $ji$ entry of $X$. Also, we will use the notation $X[j,:]$ to mean the $j^{th}$ row of $X$ and $X[:,i]$ to mean the $i^{th}$ column of $X$. (This is copied from the Python programming language; the ':' means that index runs over all possibilities).

Since

*(handwritten, right)* $2(Y_j - \sum_{s=1}^{k+1} X_{js}M_s)[-X_{jt}]$

$$E = \sum_{j=1}^{N}(Y_j - \sum_{s=1}^{k+1}X_{js}M_s)^2$$

*(handwritten)* $X_{jt}M_t$

$1 \le t \le K+1$

we compute:

$$\frac{\partial}{\partial M_t}E = -2\sum_{j=1}^{N}X_{jt}(Y_j - \sum_{s=1}^{k+1}X_{js}M_s)$$

*(handwritten, right)* $\sum X_{tj}^\mathsf{T}Y_j$

$$= -2(\sum_{j=1}^{N}Y_jX_{jt} - \sum_{j=1}^{N}\sum_{s=1}^{k+1}X_{jt}X_{js}M_s)$$

*(handwritten)* $t^{th} \left( x_{t1} \; x_{t2} \cdots \right) \begin{pmatrix} Y_1 \\ \vdots \\ Y_N \end{pmatrix}$

$X^\mathsf{T}$

$$= -2(\sum_{j=1}^{N}X_{tj}^\mathsf{T}Y_j - \sum_{j=1}^{N}\sum_{s=1}^{k+1}X_{tj}^\mathsf{T}X_{js}M_s) \tag{6}$$

$$= -2(X^\mathsf{T}[t,:]Y - \sum_{s=1}^{k+1}\sum_{j=1}^{N}X_{tj}^\mathsf{T}X_{js}M_s)$$

*(handwritten, left)* $X^\mathsf{T}[t,:]$
$t^{th}$ row of $X$

$$= -2(X^\mathsf{T}[t,:]Y - \sum_{s=1}^{k+1}(X^\mathsf{T}X)_{ts}M_s)$$

*(handwritten, right)* $t^{th}$ entry of $\nabla E$ is this formula

$$= -2(X^\mathsf{T}[t,:]Y - (X^\mathsf{T}X)[t,:]M)$$

Stacking up the different rows to make $E$ yields the desired formula.

$$\nabla E = -2(X^T Y \; \phi - X^T X M) = 0$$

$$\underbrace{X^T X}_{K+1 \times K+1} \overset{\underset{K+1 \times 1}{\uparrow}}{M} = X^T Y \qquad {}_{K+1 \times 1}$$

Assume

$D = X^T X \quad (K+1) \times (K+1)$  $\underline{D \text{ is invertible}}$

$$M = D^{-1} X^T Y$$

$$Y_{predicted} = XM$$

$$E = \|Y - XM\|^2 = \|Y - Y_{predicted}\|^2$$

$$Y_{predicted} = \underbrace{X \overset{K+1 \times K+1}{D^{-1}} \overset{K+1 \times N}{X^T} Y}_{N \times 1} \qquad \underset{N \times M}{}$$

predicted

$M \Leftrightarrow m, b$

$y = mx + b$

$Y \; N \times 1$