$X$ data matrix — $N$ rows, 1 for each sample

$K+1$ columns

$K$ features

Last column is $\begin{pmatrix} \vdots \\ \vdots \end{pmatrix}$

$Y$ target values — $N \times 1$ column vector

$M$ — $(K+1) \times 1$ column vector — $M_{K+1}$ is the intercept

Minimize $MSE = \frac{1}{N} \| Y - XM \|^2$

$Y - XM$ is an $N \times 1$ vector

$$Y - XM = \left( y_i - \sum_{j=1}^{K+1} x_{ij} m_j \right)$$

$$\| Y - XM \|^2 = \frac{1}{N} \sum_{i=1}^{N} \left( y_i - \sum_{j=1}^{K+1} x_{ij} m_j \right)^2$$

← minimize this

$m_j$ are variables

$A = (a_{ij})$ which is $r \times s$

$B = (b_{ij})$ which is $s \times t$

$C = AB$ is $r \times t$

$$c_{ik} = \sum_{j=1}^{s} a_{ij} b_{jk}$$

$A^T = (a_{ji})$

$s \times r$ matrix

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

*N.*

$E = ||Y - XM||^2$

$$\nabla E = \begin{bmatrix} \frac{\partial}{\partial M_1} E \\ \frac{\partial}{\partial M_2} E \\ \vdots \\ \frac{\partial}{\partial m_{M_{k+1}}} E \end{bmatrix} = -2X^\mathsf{T}Y + 2X^\mathsf{T}XM \qquad (5)$$

$$= -2 \left[ \underset{(k+1)\times 1}{X^\mathsf{T}Y} - \underset{(k+1)\times 1}{X^\mathsf{T}XM} \right]$$

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

Minimize MSE
by taking its $\frac{\partial}{\partial M_i}$
$i = 1, \ldots, k+1$
and setting $= 0$.

$X \quad N \times (k+1)$
$X^\mathsf{T} \quad (k+1) \times N$
$Y \quad N \times 1$
$X^\mathsf{T}X = (k+1) \times (k+1)$
$M \quad (k+1) \times 1$

**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

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

we compute:

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

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

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

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

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

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

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

$X^\mathsf{T}[t, :] \quad t^{th}$ row of $X^\mathsf{T}$

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

$$\nabla E = -2 \left[ X^\mathsf{T}Y - X^\mathsf{T}XM \right].$$

$$\nabla E = 0$$

$$X^\mathsf{T}Y = X^\mathsf{T}XM$$
$(k+1)\times N \quad N\times 1 \qquad (k+1)\times(k+1) \quad (k+1)\times 1$

$$D = X^\mathsf{T}X \qquad (k+1)\times(k+1)$$

Suppose $D$ is an invertible matrix

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