

$$1.1 \quad \frac{1}{N} ((w^T x + b) - y) \quad \nabla \text{mse}(w, b) = \begin{pmatrix} \frac{1}{N} \sum (w^T x + b - y) x \\ \frac{1}{N} \sum (w^T x + b - y) \end{pmatrix}$$

$$2.1 \quad a) \hat{Y} = XW$$

$$b) \text{mse} = \frac{1}{2N} \sum_{i=1}^N (\hat{y}_i - y_i)^2 = \frac{1}{2N} (XW - Y)^T (XW - Y)$$

$$\frac{\partial \text{mse}}{\partial W} = \frac{1}{2N} \frac{\partial (XW - Y)^T (XW - Y)}{\partial W}$$

$$= \frac{1}{2N} \left(\frac{\partial}{\partial W} ((XW)^T - Y^T) (XW - Y) \right)$$

$$= \frac{1}{2N} \left(\frac{\partial}{\partial W} (W^T X^T X W - W^T X^T Y - Y^T X W + Y^T Y) \right)$$

$$= \frac{1}{2N} \left(2X^T X W - X^T Y - \frac{\partial}{\partial W} (Y^T X W)^T \right)$$

$$= \frac{1}{2N} \left(2X^T X W - X^T Y - \frac{\partial}{\partial W} (W^T X^T Y) \right)$$

$$= \frac{1}{2N} (2X^T X W - X^T Y - X^T Y) = \frac{X^T X W - X^T Y}{N}$$

$$c) \text{mse} = \frac{1}{2N} \sum_{i=1}^N (\hat{y}_i - y_i)^2 + \lambda \|W\|^2$$

$$= \text{mse}_{OLS} + \lambda W^T W$$

$$\frac{\partial \text{mse}}{\partial W} = \frac{\partial}{\partial W} \text{mse}_{OLS} + \lambda \frac{\partial}{\partial W} W^T W$$

$$= \frac{X^T X W - X^T Y}{N} + 2\lambda W$$

3.1 Let $R = \begin{bmatrix} r_1 & 0 & 0 & 0 & \dots & 0 \\ 0 & r_2 & 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & r_i & \dots & 0 \\ 0 & 0 & 0 & \dots & \dots & r_n \end{bmatrix}_{n \times n}$

$$E(W) = \frac{1}{2N} \sum_{i=1}^N r_i^2 (y_i - w^T x_i)^2$$

$$= \frac{1}{2N} (RY - RXW)^T (RY - RXW)$$

$$\frac{\partial E(W)}{\partial W} = \frac{1}{2N} \frac{\partial}{\partial W} (RY^T R^T RY - Y^T R^T R X W - W^T X^T R^T R Y + W^T X^T R^T R X W)$$

$$= \frac{1}{2N} \left(\frac{\partial}{\partial W} (-Y^T R^T R X W)^T - X^T R^T R Y + 2X^T R^2 X W \right)$$

$$= \frac{1}{2N} \left(\frac{\partial}{\partial W} (-W^T X^T R^T R Y^T) - X^T R^2 Y + 2X^T R^2 X W \right)$$

$$= \frac{1}{2N} (-X^T R^2 Y - X^T R^2 Y + 2X^T R^2 X W)$$

$$= \frac{1}{N} (X^T R^2 X W - X^T R^2 Y) = 0$$

$$\Rightarrow W = (X^T R^2 X)^{-1} (X^T R^2 Y)$$

4.2 The closed form of OLS doesn't have a solution if the columns of X are linearly ~~in~~dependent. Yes, the gradient descent converges to a sol. in that case, a closed form solution doesn't exist.

The modification I did to the data was to convert X to $I(4 \times 4)$

4.1 The program was not executing properly because, it was calculating $(X^T X)^{-1}$ for an X , whose columns are linearly dependent. So, I modified the data in such a way that the columns are linearly independent and no. of columns \leq no. of rows.

$|X^T X| \neq 0$ iff X has ~~non~~ full column rank (given no. of columns \leq no. of rows)

If X has full column rank and $|X^T X| = 0$,

Assume $X^T X v = 0 \Rightarrow v^T X^T X v = 0 \Rightarrow (Xv)^T Xv = 0$

$\Rightarrow \|Xv\|^2 = 0 \Rightarrow Xv = 0$, a contradiction

If $|X^T X| \neq 0$ and X has non-full column rank,

Assume $Xv = 0 \Rightarrow X^T X v = 0 \Rightarrow |X^T X| = 0$, a contradiction

Hence, proved.

If Let no. of columns be n and no. of rows be m .

If $n > m$,

~~rank~~ col. rank $(X) \leq n$

col. rank $(X^T) \leq m$

Since col. rank $(X^T X) \leq \min(\text{col. rank}(X^T), \text{col. rank}(X))$
 $\leq \min(m, n)$

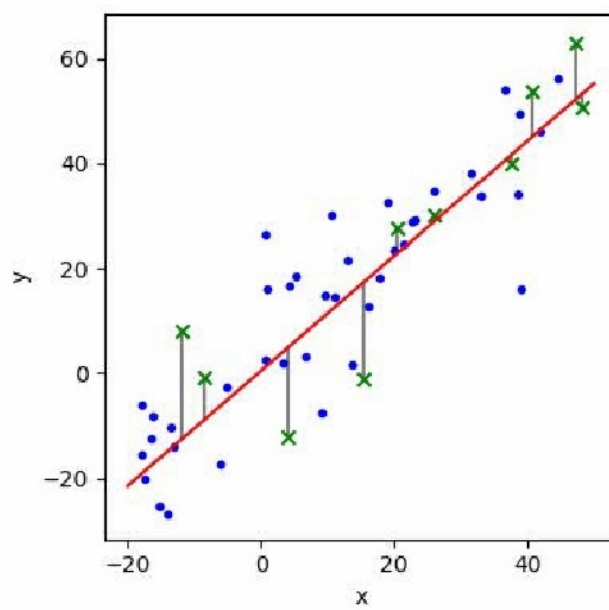
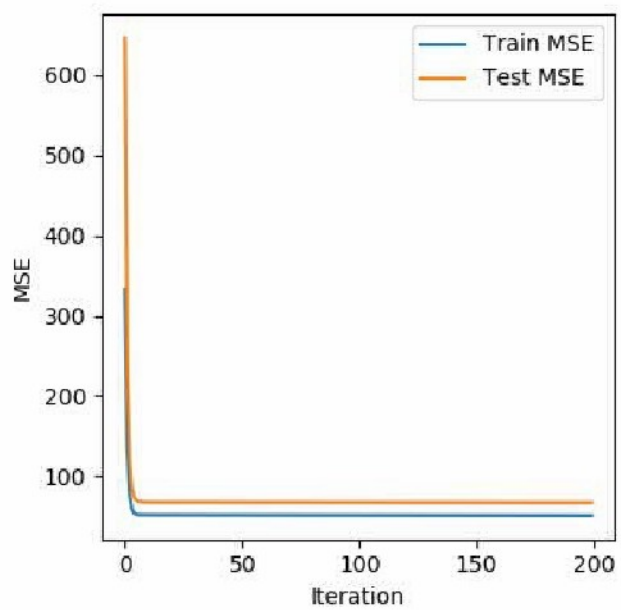
$\leq m < n \Rightarrow X^T X$ is non-invertible

4.2 (continued) Also, for OLS to ^{not} have a closed form sol. no. of col. $(X) \geq$ no. of rows (X)

If there exists an x_0 such that, $X^T X x_0 = X^T Y$, then there exist infinite solutions of the form $(x_0 + \lambda v)$, where λ is any scalar and v belongs to null space of $X^T X$ (i.e. $X^T X v = 0$), as we are discussing the case when $X^T X$ is not invertible. Hence, gradient descent converges to any of these solutions.

If a minima to which gradient descent converges exists, it should satisfy the eq. $X^T X w_0 = X^T Y$, where w_0 is a point of minima ($\because X^T X w_0 = X^T Y$ is obtained by equating derivative to zero, which is the condition for minima). Hence, we can say that multiple minima exist with same mse value (see above paragraph) and gradient descent, by the nature of the algorithm, converges to any one of them.

1.2 d) The blue points denote training data and the green points denote test data and the red line denotes the best fit line. The line best fits on the training data, as it is trained on the training data, but also fits the test data well. Some points have positive error and some have negative error, but the overall error is minimized as we used the square of the error as our measure.



The vertical distance between the points and the line denote $\hat{y} - y$, where \hat{y} is predicted by the line and y is the original value. As the line passes through the middle of the data, it also matches our intuition and prior knowledge that the mean reduces mean square error. But, we can also see that ~~that~~ the line is ~~also~~ influenced by outliers, since we treat all training data equally.