

④

a) In the matrix form, OLS objective is given by,

$$f(\beta) = \min_{\beta} \|y - X\beta\|_2^2$$

$$\begin{aligned}\nabla f(\beta) &= -2X^T y + 2X^T X \beta \\ &= 2X^T (X\beta - y)\end{aligned}$$

Gradient descent formula is given by

$$\begin{aligned}\beta^t &= \beta^{t-1} - \eta \nabla f(\beta^{t-1}) \\ &= \beta^{t-1} - \eta (2X^T (X\beta^{t-1} - y)) \\ &= \beta^{t-1} - 2\eta X^T (X\beta^{t-1} - y)\end{aligned}$$

b) In case of Newton Raphson algorithm, iterative steps, is given by

$$\beta^t = \beta^{t-1} - (\nabla^2 f(\beta^{t-1}))^{-1} \nabla f(\beta^{t-1})$$

$$\begin{aligned}\nabla^2 f(\beta) &= \nabla (2X^T X \beta - 2X^T y) \\ &= 2\nabla (X^T X \beta) \\ &= 2X^T X\end{aligned}$$

In general,  
 $\nabla_{\beta}(A\beta) = A$

$$\begin{aligned}\beta^t &= \beta^{t-1} - (2X^T X)^{-1} (2X^T X \beta^{t-1} - 2X^T y) \\ &= \beta^{t-1} - \frac{1}{2} (X^T X)^{-1} * 2 (X^T X \beta^{t-1} - X^T y) \\ &= \beta^{t-1} - (X^T X)^{-1} (X^T X) \beta^{t-1} + (X^T X)^{-1} X^T y\end{aligned}$$

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

we know  
that  
 $A^{-1}A = I$

$$\boxed{\beta^t = (X^T X)^{-1} X^T y} = \beta_{OLS}$$

Observation  $\rightarrow$  Optimal  $\beta$  is independent of initial  $\beta$  value. It seems like it converges to OLS.