3)

Given $\quad y = \beta_0 + \beta_1 x + \epsilon$

where $\epsilon$ is noise and $N(\epsilon : 0, \sigma_e^2)$ i.e. errors (noise)

follows a normal distribution.

$$\hat{y}_i = \beta_0 + \beta_1 x_i \quad \Rightarrow \quad y_i - \hat{y}_i = \epsilon$$

Function of $\epsilon$, $\quad f(\epsilon_i) = \dfrac{1}{\sqrt{2\pi\sigma^2}} e\left(-\dfrac{1}{2\pi\sigma^2}(\epsilon_i)^2\right)$

The likelihood is given by,

$$L = \dfrac{1}{(2\pi\sigma^2)^{-N/2}} e^{-\left(\frac{1}{2\sigma^2}\sum\epsilon^2\right)}$$

Log likelihood is,

$$\ell = -\dfrac{N}{2}\ln(2\pi) - \dfrac{N}{2}\ln(\sigma^2) - \dfrac{1}{2\sigma^2}\sum_{i=1}^{N}\epsilon_i^2$$

$$\ell = -\dfrac{N}{2}\ln(2\pi\sigma^2) - \dfrac{1}{2\sigma^2}\sum_{i=1}^{N}\epsilon_i^2$$

where $N$ is number of datapoints

So, for the likelihood to be maximum, we choose,

$$\beta_{MLE} = \arg\max_{\beta_1} - \sum_{i=1}^{N}\epsilon_i^2$$

In other words,

$$\beta_{MLE} = \arg\min_{\beta_1} \dfrac{1}{N}\sum_{i=1}^{N}\epsilon_i^2$$

or $\quad \beta_{MLE} = \arg\min_{\beta_1} \text{MSE}$

which essentially means that if we minimize MSE, we get maximum likelihood. Conversely, it is clear that maximization of likelihood leads to minimized MSE.