**Q1.** Given a simple linear regression model $Y = \beta_0 + \beta_1 X$.

(a) $\hat{\beta}_i$ is the least squares estimates of $\beta_i$, $i = 0, 1$. A constant $c$ is added to each response. We want to know: the new least squares estimates of $\beta_i$.

$Y = f(X) + c$   $Y = \beta_0 + \beta_i X + c$ for $i = 0, 1$

So we define the residual sum of squares (RSS) as:

$RSS = c^2 + c^2 + c^2 + \cdots + c^2 = nc^2$

$= (y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^2 + (y_2 - \hat{\beta}_0 - \hat{\beta}_1 x_2)^2 + \cdots + (y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n)^2$

According to the note: $\hat{\beta}_1 = \dfrac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2}$     $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$

Since a constant $c$ is added to each response:

$\hat{\beta}_i$: $i = 1$, $\hat{\beta}_1 = \dfrac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i + c - (\bar{y} + c))}{\sum_{i=1}^{n}(x_i - \bar{x})^2} = \dfrac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2}$

$i = 0$: $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = (\bar{y} + c) - \left(\dfrac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2}\right)\bar{x}$

(right margin:)

$\beta_1 = \dfrac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2}$

$= \dfrac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i + c - (\bar{y} + c))}{\sum_{i=1}^{n}(x_i - \bar{x})^2}$

$= \dfrac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2} = \hat{\beta}_1$

$\beta_0 = \bar{y} - \beta_1 \bar{x}$  $\hat{\beta}_0 = \dfrac{\sum_{i=1}^{n} y_i + c}{n} - \beta$

$\hat{\beta}_0 = \bar{y} + c - \beta_1 \bar{x}$

$= \beta_0 + c$

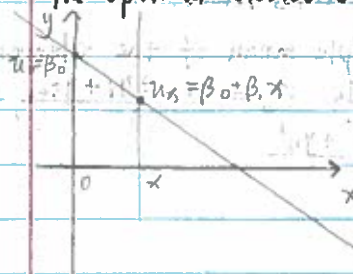**New equations:**

$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x = \beta_0 + c + \hat{\beta}_1 x$

where $\hat{\beta}_0 = \beta_0 + c$

$\hat{\beta}_1 = \beta_1$

(b) The coefficient $\beta_0$ is referred to as the intercept term. It can be interpreted as a summary of the vertical location of the response $Y$, since the effect of changing the constant $c$ of part (a) is directly observable in $\beta_0$. $\beta_0 = u_0$ when $u_x = \beta_0 + \beta_1$. So construct a new intercept $u_x$ at any vertical line $X = x$. The least squares estimate will be $\hat{u}_x = \hat{\beta}_0 + \hat{\beta}_1 x$. What analytical criterion can be used to select $x$ together the optimal choice?



We want to choose a $x$ that minimize the variance. So according to the formula (3.2) in the note:

$\sigma^2_{\hat{u}_x} = \sigma^2\left[\dfrac{1}{n} + \dfrac{(x - \bar{x})^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2}\right]$

So when $x = \bar{x}$, the value of the variance $\sigma^2_{\hat{u}_x}$ becomes the smallest. $\sigma^2_{\hat{u}_x} = \sigma^2\left(\dfrac{1}{n} + \dfrac{(\bar{x} - \bar{x})^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2}\right) = \sigma^2\left(\dfrac{1}{n}\right) = \dfrac{\sigma^2}{n}$  So we want the line of best fit to pass through $(\bar{x}, \bar{y})$.

**Q2.** We are given a 'hat matrix' associated with multiple linear regression with $q$ predictors. $H$ is a linear transformation of an $n$-dimensional response vector $y$ to the $n$-dimensional fitted vector $\hat{y} = Hy$. $\hat{y}$ is obtained by minimizing the SSE, the action taken by $H$ on $y$ is to project it onto the $q$-dimensional subspace $S_q \subset \mathbb{R}^n$ spanned by the $q$ predictors. $\hat{y}$ is the point in $S_q$ closest to $y$.