3a) Regression line. Let $a, b \in \mathbb{R}^n$, $m_a = avg(a) = \frac{1^T a}{n}$, $m_b = avg(b) = \frac{1^T b}{n}$

$$S_a = std(a) = \frac{1}{\sqrt{n}} \|a - m_a 1\|, \quad S_b = std(b) = \frac{1}{\sqrt{n}} \|b - m_b 1\|$$

We assume the vectors are not constant ($S_a \neq 0$ and $S_b \neq 0$) and write the correlation coefficient as $\quad \rho = \frac{1}{n} \frac{(a - m_a 1)^T (b - m_b 1)}{S_a S_b}$

We considered the problem of fitting a straight line to the points $(a_k, b_k)$ by minimizing $J = \frac{1}{n} \sum_{k=1}^{n} (c_1 + c_2 a_k - b_k)^2 = \frac{1}{n} \|c_1 1 + c_2 a - b\|^2$

Show that the optimal coefficients are $c_2 = \rho S_b / S_a$ and $c_1 = m_b - m_a c_2$. Show that for those values of $c_1$ and $c_2$, we have $J = (1 - \rho^2) S_b^2$.

Let $a_0 = a - m_a 1$ and $b_0 = b - m_b 1$

thus $S_a = \frac{1}{\sqrt{n}} \|a_0\|$ and $S_b = \frac{1}{\sqrt{n}} \|b_0\|$ and $\rho = \frac{1}{n} \cdot \frac{a_0^T b_0}{S_a S_b}$

$$J = \frac{1}{n} \|c_1 + c_2 a - b\|^2 = \frac{1}{n} \|c_1 + c_2 a_0 - b_0\|^2$$

$$\frac{\partial J}{\partial c_1} = \frac{2}{n} 1^T (c_1 1 + c_2 a_0 - b_0) = 2(c_1 + m_a c_2 - m_b) = 0$$

Thus $\underline{c_1 = m_b - m_a c_2}$ 🔒

Then $J(c_2) = \frac{1}{n} \|c_2 a_0 - b_0\|^2 = S_a^2 c_2^2 + S_b^2 - 2 c_2 \rho S_a S_b$

$J'(c_2) = 0$ so $2 S_a^2 c_2 - 2 \rho S_a S_b = 0$

$$c_2 = \frac{\rho S_a S_b}{S_a^2 c_2} \implies \text{thus } \underline{c_2 = \rho S_b / S_a}$$ 🔒

Plug: $J(c_2) = \frac{1}{n} \|c_2 a_0 - b_0\|^2 = c_2^2 \frac{\|a_0\|^2}{n} + \frac{\|b_0\|^2}{n} - 2 c_2 \frac{a_0^T b_0}{n}$

$$J = S_a^2 c_2^2 + S_b^2 - 2\rho S_a S_b c_2$$

$$= S_a^2 \left(\rho \frac{S_b}{S_a}\right)^2 + S_b^2 - 2\rho S_a S_b \left(\rho \frac{S_b}{S_a}\right)$$

$$= \rho^2 S_b^2 + S_b^2 - 2\rho^2 S_b^2$$

$$= S_b^2 - \rho^2 S_b^2$$

Thus $\quad \underline{J = (1 - \rho^2) S_b^2}$ 🔒

3b) orthogonal distance regression

$\forall p \in (a_k, b_k)$, the vertical deviation from the straight line defined by $y = c_1 + c_2 x$ is given by $e_k = |c_1 + c_2 a_k - b_k|$

The orthogonal distance of $(a_k, b_k)$ to the line: $d_k = \dfrac{|c_1 + c_2 a_k - b_k|}{\sqrt{1 + c_2^2}}$

We can find the straight line that minimizes the sum of the squared orthogonal distance $\sum J = \dfrac{1}{n} \sum_{k=1}^{n} d_k^2 = \dfrac{\| c_1 \mathbf{1} + c_2 a - b \|^2}{n(1 + c_2^2)}$

i) Show that the optimal value of $c_1$ is $c_1 = m_b - m_a c_2$ as for the least squares fit.

Let $r(c_1) = c_1 \mathbf{1} + c_2 a - b$

$\dfrac{d}{dc_1} \| r(c_1) \|^2 = 2 \cdot \mathbf{1}^T r(c_1) = 2(n c_1 + c_2 \mathbf{1}^T a - \mathbf{1}^T b) = 0$

$c_1 = \dfrac{\mathbf{1}^T b}{n} - \dfrac{c_2 \mathbf{1}^T a}{n}$

$\underline{\underline{c_1 = m_b - m_a c_2}}$

ii) $J = \dfrac{s_a^2 c_2^2 + s_b^2 - 2\rho s_a s_b c_2}{1 + c_2^2}$. Set $\dfrac{dJ}{c_2} = 0$. Then $\rho c_2^2 + \left( \dfrac{s_a}{s_b} - \dfrac{s_b}{s_a} \right) c_2 - \rho = 0$

If $\rho = 0$ and $s_a = s_b$, any value of $c_2$ is optimal. If $\rho = 0$ and $s_a \neq s_b$ the quadratic eq. has a unique solution $c_2 = 0$. If $\rho \neq 0$, the quadratic eq. has 2 positive and a negative root. Show that the solution that minimizes $J$ is the root $c_2$ with the same sign as $\rho$.

$\rho c_2^2 + \left( \dfrac{s_a}{s_b} - \dfrac{s_b}{s_a} \right) c_2 - \rho = 0$

If $\rho \neq 0$, the quadratic eq. has a positive and a negative root

Let $r_1$ and $r_2$ be roots. Thus $r_1 r_2 < 0$.

$J'(c_2) = \dfrac{(2 s_a^2 c_2 - 2\rho s_a s_b)(1 + c_2^2) - (s_a^2 c_2^2 + s_b^2 - 2\rho s_a s_b c_2)(2 c_2)}{(1 + c_2^2)^2}$

$J'(0) = \dfrac{-2\rho s_a s_b}{1} = -2\rho(s_a s_b)$

If $\rho > 0$ then $J'(0) < 0$: The function is decreasing at 0. So positive $c_2$ is a minimum and negative $c_2$ is a max.

If $\rho < 0$ then $J'(0) > 0$: the function is increasing at 0. So, negative $c_2$ is a minimum and positive $c_2$ is a maximum.

Thus $\rho > 0$, $a$ and $b$ increase so the best fit slope should be positive
$\rho < 0$, $a$ and $b$ decrease so the best fit slope should be negative

Thus the solution that minimizes $J$ is the root $c_2$ with the same sign as $\rho$.