SOME CALCULATIONS —

FOR A GIVEN DATASET, A COLLECTION OF
$(x,y)$ COORDINATE PAIRS $\{x_i, y_i\}_{i=1}^{n}$ :

| $i =$ | $x_i =$ | $y_i =$ |
|-------|---------|---------|
| 1 | $x_1$ | $y_1$ |
| 2 | $x_2$ | $y_2$ |
| $\vdots$ | $\vdots$ | $\vdots$ |
| $n$ | $x_n$ | $y_n$ |

LET $\bar{x}$, $s_x^2$ BE THE MEAN AND VARIANCE FOR $x$

$$\bar{x} = \frac{\sum x_i}{n}, \quad s_x^2 = \frac{\sum (x_i - \bar{x})^2}{n}, \quad \text{SIMILARLY FOR } \bar{y}, s_y^2.$$

THEN, THE DATASET IN STANDARD UNITS (SU) IS

| $i$ | $x_i$ (SU) | $y_i$ (SU) |
|-----|-----------|-----------|
| 1 | $\dfrac{x_1 - \bar{x}}{s_x}$ | $\dfrac{y_1 - \bar{y}}{s_y}$ |
| 2 | $\dfrac{x_2 - \bar{x}}{s_x}$ | $\dfrac{y_2 - \bar{y}}{s_y}$ |
| $\vdots$ | | |
| $n$ | $\dfrac{x_n - \bar{x}}{s_x}$ | $\dfrac{y_n - \bar{y}}{s_y}$ |

AND WE DENOTE

$\dfrac{x_i - \bar{x}}{s_x}$ BY $x_i$(SU),

SIMILARLY FOR $y_i$(SU)

$\left( \text{AND, OF COURSE } s_x = \sqrt{s_x^2} \text{ AND } s_y = \sqrt{s_y^2} \right)$

WE DEFINE (COVARIANCE) OF X AND y TO BE

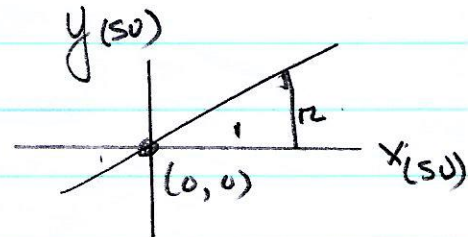$$\sum \frac{(x_i - \bar{x})(y_i - \bar{y})}{n}$$

AND WE LET

$$r = \frac{1}{s_x s_y} \cdot \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y})$$

↗ ,

CALLED
CORRELATION
COEFFICIENT

$$= \frac{1}{n} \sum \underbrace{\frac{(x_i - \bar{x})}{s_x}}_{x_{i \,(su)}} \cdot \underbrace{\frac{(y_i - \bar{y})}{s_y}}_{y_{i \,(su)}}$$

THE [LINEAR] "REGRESSION LINE TO THE MEAN"
IS DEFINED TO BE

$$y_{(su)} = r \; x_{(su)}$$



WHICH CONTAINS (0, 0) WITH SLOPE r

(i.e. WHEN x AND y ARE EXPRESSED IN
STANDARD UNITS, THE REGRESSION LINE
● CONTAINS (0,0) AND HAS SLOPE r.
● r MEASURES "STRENGTH OF LINEAR RELATIONSHIP"
● r IS NOT SLOPE OF REGRESSION LINE
FOR X AND y IN ORIGINAL UNITS

THE REGRESSION LINE IN ORIGINAL UNITS IS

$$\frac{y - \bar{y}}{s_y} = r \frac{x - \bar{x}}{s_x}$$

OR

$$y = r\left(\frac{x - \bar{x}}{s_x}\right) s_y + \bar{y}$$

$$= \boxed{r \frac{s_y}{s_x}} x + \boxed{\bar{y} - r \frac{s_y}{s_x} \bar{x}}$$

THINK
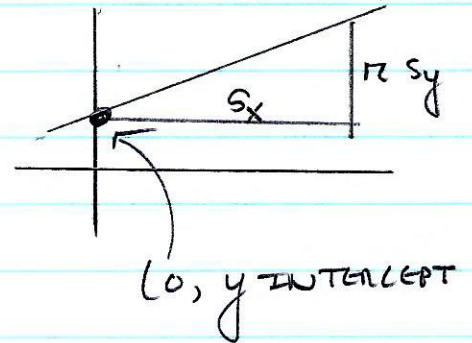$y = mx + b$

SLOPE        INTERCEPT

WHICH CONTAINS $(\bar{x}, \bar{y})$



$(0, y \text{ INTERCEPT})$

SO, GIVEN $x$, WE "ESTIMATE $y$" BY COMPUTING $y$ IN ABOVE EQUATION

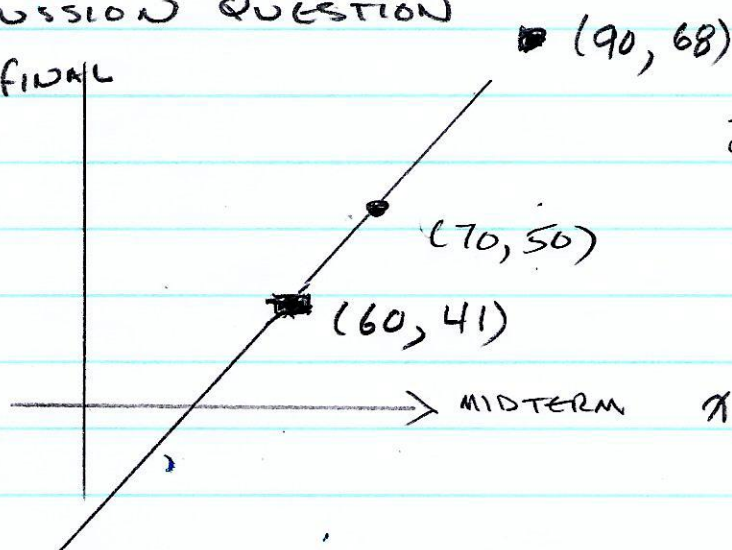THE MANNER IN WHICH VALUES FOR $y$ ARE ACTUALLY DISTRIBUTED ABOUT "ESTIMATE FOR $y$ GIVEN $x$" IS ANOTHER STORY.

THE APPROPRIATENESS OF THINKING ABOUT RELATIONSHIP BETWEEN $x$ AND $y$ AS "LINEAR" IS ANOTHER QUESTION.

DISCUSSION QUESTION



y FINAL

$(90, 68)$

$(70, 50)$

$(60, 41)$

MIDTERM   $x$

$\bar{x} = 70$  $\bar{y} = 50$

$S_x = 10$  $S_y = 12$

$r = 0.75$

ESTIMATE FINAL EXAM FOR MIDTERM OF 90

$$\hat{y} = r \frac{S_y}{S_x} x + \left( \bar{y} - r \frac{S_y}{S_x} \bar{x} \right)$$

$$= 0.75 \, \frac{12}{10} x + 50 - 0.75 \, \frac{12}{10} \cdot 70$$

$$= .9 x + (50 - .9(70))$$

$$= .9 x - 13$$

WHEN $x = 90$,  $\hat{y} = .9(90) - 13 = 68$

WHEN $\hat{x} = 60$,  $\hat{y} = .9(60) - 13 = 41$