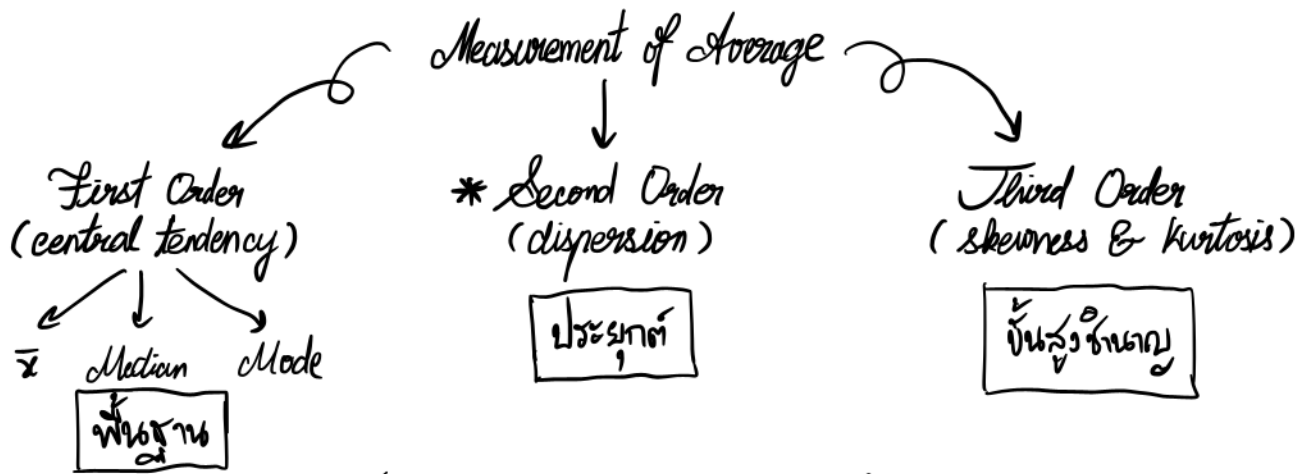
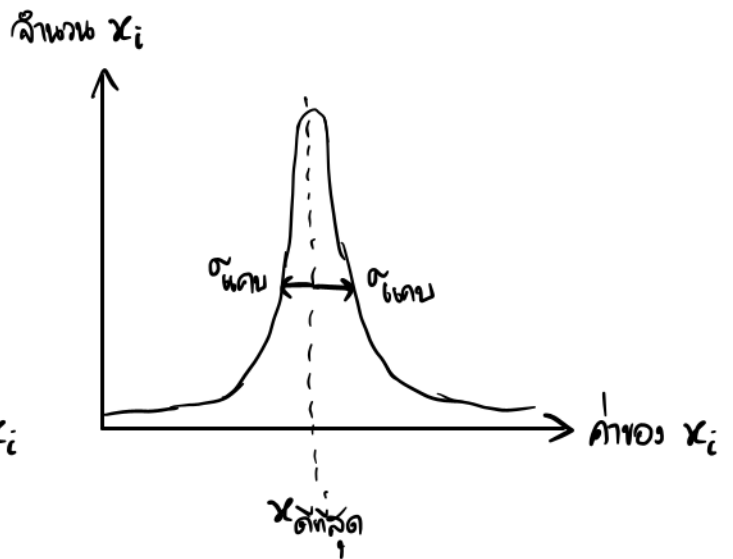
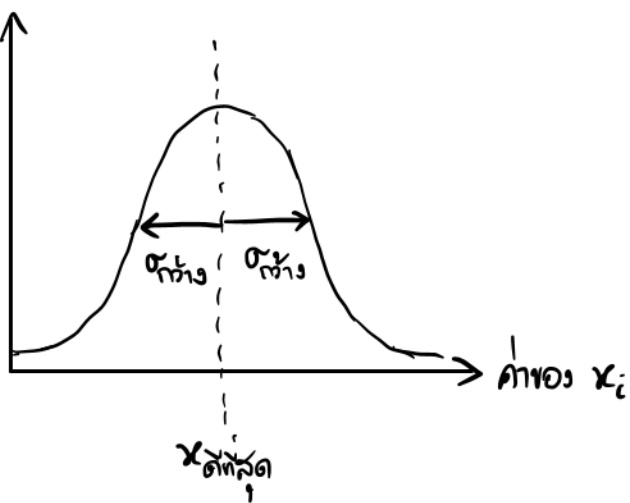


# Data Analysis (DA101) : ข้อมูลตัวแปรเดียว และหลายตัวแปร



การกระจาย/ส่วนเบี่ยงเบนเป็นการวัดข้อมูลทางสถิติลำดับที่สอง (ให้บอกความเป็นกลุ่มของข้อมูล)

จำนวน  $x_i$  ←  $x_i$  เป็น SAMPLE ข้อมูล



สมมติเหตุการณ์การวัดค่า โดยทั่วไปนั้นเครื่องมือที่ใช้วัดค่า มักมีความคลาดเคลื่อนในตัวเสมอ  
 ไม่มีเครื่องมือใดบนโลกที่เที่ยงตรง (วัดค่าได้ตรงกับความจริงที่สุด เรียกว่า Accuracy) และ  
 วัดละเอียด (วัดได้ค่าเดิมซ้ำ ๆ เรียกว่า Precision) โดยทั้ง 2 ค่า มีผลกับส่วนเบี่ยงเบน  
 โดยตรง กราฟยิ่งแคบ (ด้านขวา) ยิ่งดี ค่าเกาะกลุ่มกัน ทวีปไปมักใช้ส่วนเบี่ยงเบนมาตรฐาน  
 (STANDARD DEVIATION) แต่ก็มีค่าอื่น ๆ ที่นำมาใช้แทนได้ แต่ไม่ค่อยนิยมมากนัก

หมายเหตุ ในบางหน้าอาจมีการนำเคล็ดลับมาใช้ หากไม่เข้าใจก็อ่านไปแล้ว เชื่อไปก่อนว่าจริงก็ได้  
 แต่หากต้องการรู้เนื้อหา สามารถอ่านเองได้ ไม่ต้องพึ่งคนอื่น

โดยก่อนเริ่มพูดถึงค่าสถิติมาเรานั่น ต้องมาพูดถึงหัวใจของค่าต่างๆ การสถิติก่อน นักสถิติมองว่า การเก็บหรือวัดข้อมูล ครั้งเดียว ไม่มีความน่าเชื่อถือ เพราะเรามีอาจทราบได้ว่าค่า นั้นสามารถปองจัดรวมเป็นจริงได้หรือไม่ นักสถิติจึงเลือกเก็บข้อมูลหลายๆ "ตัวอย่าง" จาก กลุ่ม "ประชากร" และไม่นับรวม แต่แล้วก็ยังไม่ทราบวิธีนำข้อมูลเหล่านั้นมาใช้ประโยชน์ ได้อย่างเต็มที่ ก็เพราะมันเยอะไปหมด ต้องคิดหาวิธีเอาตัวแทนข้อมูลเหล่านั้นออกมา เพื่อให้คนที่เกี่ยวข้องอ่านข้อมูลได้มากขึ้น เข้าใจได้อย่างลึกซึ้งถึง สิ่งที่ข้อมูลต้องการแสดงออกมา หรือสิ่งที่คาดว่าจะเป็นได้จากข้อมูลนั้นๆ จึงคิด "ค่าคาดหวัง" (EXPECTED VALUE) ออกมา โดยให้แนวคิดของความน่าจะเป็นในการพบข้อมูลนั้น กระจายออกมาในทุกๆ ข้อมูล (PROBABILITY DISTRIBUTION) นิยามคุณเป็น factor ของค่าเหล่านั้นๆ แล้วนำมารวม กัน จะได้ตัวแทนข้อมูลออกมา เขียนในรูปคณิตศาสตร์ว่า

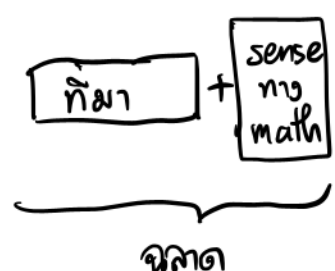
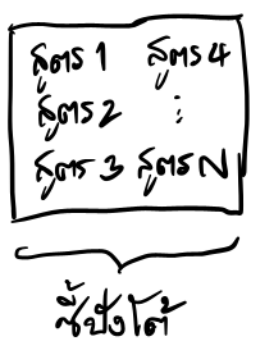
$$E[X] = \sum_{i=1}^N p_i X_i$$

↑ EXPECTED VALUE     
 ↑ PROBABILITY OF  $X_i$      
 DISCRETE VARIABLE  $i$

โดย  $\sum p_i = 1$  เสมอ  
 หมายความว่า  $E[X] = \text{constant}$   
 คือ  $E[YE[X]] = E[X]E[Y]$

ให้เข้าใจเสมอว่า 1 = 100%

ซึ่งค่าค่านี้เป็นที่มาของค่าทางสถิติหลายๆ ค่ามาก ซึ่งในที่นี่จะยกมาบางค่าเท่าที่อยากชวนให้ใคร่รู้สัก นานต้องการแปลโลก การวิเคราะห์ข้อมูลทางสถิติ ในแบบต่างๆ (แปลกๆ) ขอให้ผู้อ่านทำการศึกษาค้นคว้าด้วยตนเอง เพราะมันจะทำให้เราอยากอ่าน และเข้าไปเรื่อยๆ และขอให้ตั้งคำถามไว้เสมอว่า ทำไม? เพราะวิชาสถิติหลายคนมองว่าเป็นวิชาท่องสูตร แล้วทำโจทย์ แต่จริงๆ แล้ว ทุกอย่างเชื่อมโยงกันมา และนำมาใช้ประโยชน์ได้ทั้งนั้น



1. STANDARD DEVIATION (S.D. หรือ  $\sigma$ )

$$\sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{N}} \rightarrow \sigma^2 = \text{ความแปรปรวน (VARIANCE) มาตรฐานนั่นเองในรูป var(x)}$$

2. MEAN DEVIATION (M.D.)

$$M.D. = \frac{\sum |x_i - \bar{x}|}{N}$$

3. INTERQUARTILE RANGE (I.Q.R.)

$$I.Q.R. = Q_3 - Q_1$$

4. INTERQUARTILE DEVIATION (Q.D.)

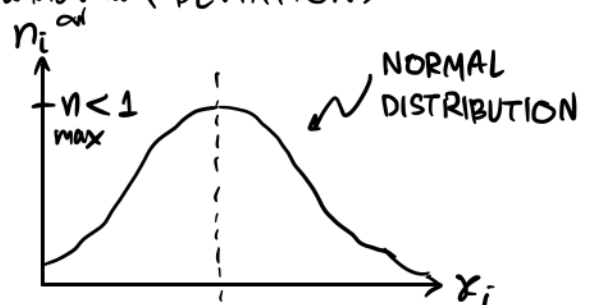
$$Q.D. = \frac{Q_3 - Q_1}{Q_3 + Q_1}$$

การประยุกต์ใช้ความแปรปรวนในข้อมูลทางสถิติ

① NORMALIZATION เมื่อต้น : ค่า z

โดยทั่วไปข้อมูลจะกระจายไปเรื่อยๆ เมื่อนำมาพล็อตเป็น histogram แล้วจะมีแบบที่ดูไม่คงที่ (ผลรวมของพื้นที่ใต้กราฟเป็นผลรวมของข้อมูล  $(\sum x_i)$ ) แต่เมื่อผ่านกระบวนการผลรวมของข้อมูลจะเป็นแบบเดียวกับ 100% หรือ 1.00 (ทั้งหมด) โดยเทียบอัตราส่วนความเบี่ยงเบนจากค่าเฉลี่ย (ERROR) กับส่วนเบี่ยงเบนมาตรฐาน (DEVIATION)

ผลต่าง  
จะได้ค่า  $z_i = \frac{x_i - \bar{x}}{\sigma}$  ซึ่งมีรูปร่างหน้าตา  $\rightarrow$



PDF (PROBABILITY DISTRIBUTION FUNCTION)

โดยเขียนสมการได้หรือจะตั้งชื่อว่า NORMAL ได้เป็ฟ  $P(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}x^2)$  (\*)

"95% CONFIDENCE INTERVAL"

ผลในทางที่นิยมต่อมาน่าจะข้อมูลส่วนใหญ่ (95% - 97.5% ของข้อมูล) ซึ่งอยู่บริเวณกึ่งกลาง (ในที่นี้  $\bar{x}$ ) จะรายงานค่า  $A \pm B$  ซึ่งหมายถึงข้อมูลส่วนใหญ่มีค่าตั้งแต่  $A-B$  ไปถึง  $A+B$  ซึ่งค่า  $A$  คือ  $\bar{x}$  (กึ่งกลาง) และ  $B$  คือขอบเขต คำนวณ

จากการหาพื้นที่ใต้กราฟ  $P(x)$  คือ

$$A = \int_{-y}^y P(x) dx = 0.95$$

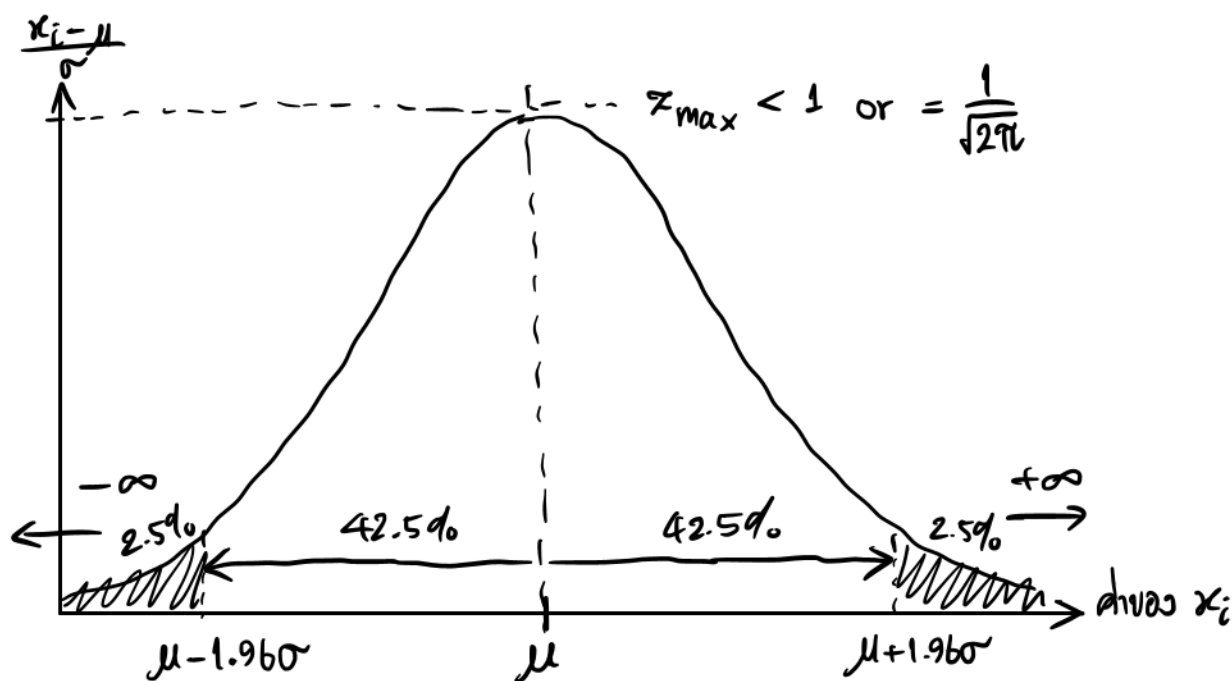
↖ เขตคูลิสเล็กน้อย

เมื่อทำการวิจัยทางคณิตศาสตร์ จะได้  $y = 1.96$  ซึ่งฟังก์ชัน  $P(x)$  เป็นการถ่วงน้ำหนักสำหรับ  $\mu$  (เขียนแทน  $\bar{x}$ ) = 0 และ  $\sigma = 1$  สำหรับกรณีทั่วไปจะใช้ฟังก์ชัน

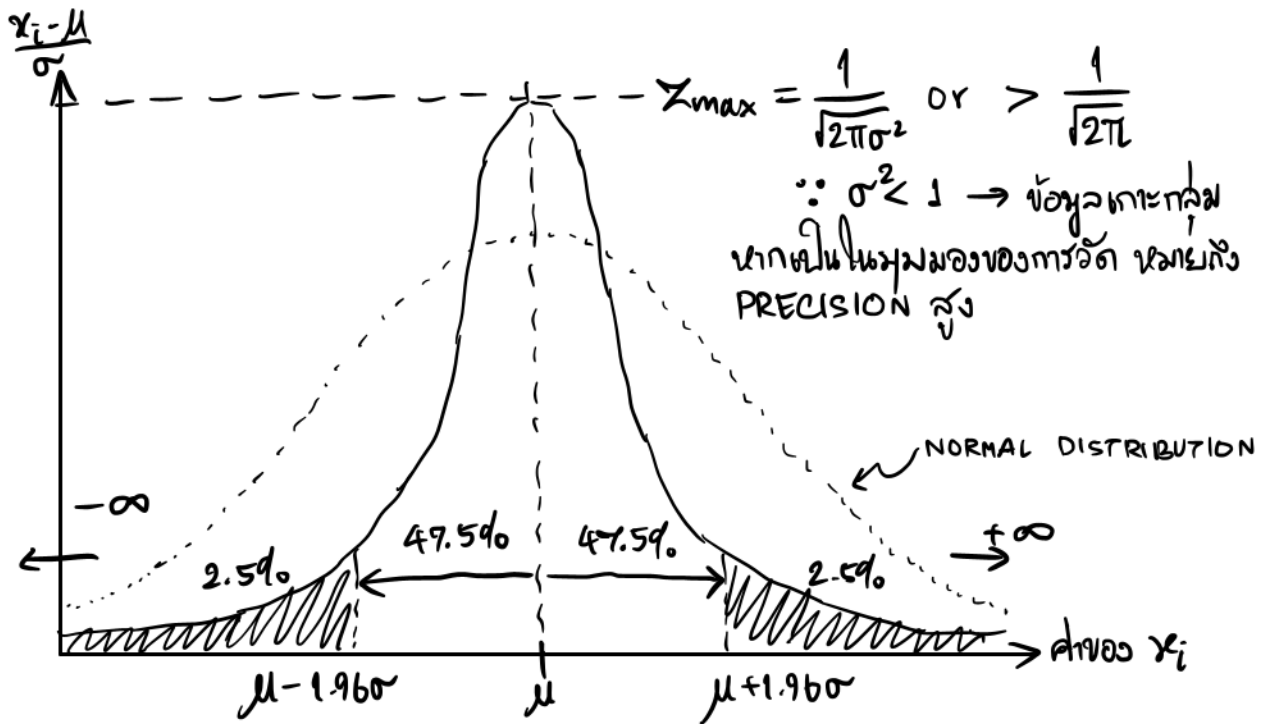
$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) (*) \text{ เพิ่มโตม}$$

เรียกว่า SYMMETRICAL DISTRIBUTION (GENERAL FORM)

ดังนั้นเมื่อคิด factor  $\sigma$  เข้ามาเกี่ยวข้อง จะได้  $A \pm B = \mu \pm 1.96\sigma$



(กราฟที่กระจายแบบ NORMAL)



$\because \sigma^2 \downarrow \rightarrow$  ข้อมูลเกาะกลุ่ม  
หากเป็นในมุมมองของการวัด หมายถึง  
PRECISION สูง

(กรณีไม่กระจายแบบ NORMAL / สมมาตรทั่วไป)

กลุ่มข้อมูลจำพวก 5% ในความเป็นจริงมักมองว่าเป็นประชากรส่วนน้อย, ค่า ERROR, ค่า NOISE, ความผิดพลาดในการวัด, ฯลฯ ซึ่งไม่น่ามาคิด

## ② VARIANCE ( $\sigma^2$ , $\text{Var}(x)$ , $\text{Cov}(x, x)$ )

ทั่วไปเราอาจใช้สูตร  $\text{STDEV}^2$  มาคิด แต่พอจำจริงแล้ว สูตรพวกนี้มาจากไหน?

หาก  $X$  เป็นตัวแปรสุ่มแบบไม่ต่อเนื่อง (random discrete variable) ความแปรปรวนนิยามเป็นผลรวมของผลคูณของกำลังสองของผลต่างกับโอกาสที่จะเจอตัวแปรที่มีค่านั้นในทาง

สถิติ หรือเขียนในรูปทั่วไปว่า  $\text{Var}(x) = \sum_{i=1}^N p_i (x_i - \mu)^2$  ซึ่ง  $p_i = \frac{n_i}{N}$  ← จำนวน  $x_i$

หาก  $x_i$  เป็นตัวแปรเต็มจะได้  $p_i = \frac{1}{N}$  ซึ่งจะได้รูปพิเศษ (SPECIAL CASE) ที่เราเขียน

กันในวิชาสถิติ ม.ปลาย - มหาลัย ว่า  $\text{Var}(x) = \frac{\sum (x_i - \mu)^2}{N}$  ( $\geq 0$ ) เสมอ

เราอาจมอง VARIANCE ในรูปแบบดั้งเดิมในทางของค่าคาดหวัง (EXPECTED VALUE)

ว่า  $\text{Var}(X) = E[(X_i - \mu)^2] = \sum p_i (x_i - \mu)^2$  ก็จะได้รูปเดิม

เพราะหากเราใช้ค่าคาดหวังมาหาที่มาของค่าเฉลี่ย :  $\mu = E[X] = \sum p_i x_i$  ;  $p_i = \frac{1}{N}$

$\therefore \mu = \bar{x} = \frac{\sum x_i}{N}$  ได้เช่นกัน

ในอีกมุมมอง VARIANCE, จาก  $\mu = E[X]$  บวก  $\text{Var}(X) = E[(X_i - \mu)^2]$

ได้ว่า  $\text{Var}(X) = E[(X - E[X])^2]$

$$= E[X^2 - 2XE[X] + (E[X])^2] ; E[c\xi] = cE(\xi)$$

$$= E[X^2] - E[X]^2$$

$$= \frac{\sum X^2}{N} - \mu^2 = \frac{\sum X^2}{N} - \left(\frac{\sum X}{N}\right)^2$$

ซึ่งเป็นสูตรที่เราเรียนกันมาโดยไม่ทราบที่มาที่ไป

ค่าอื่นๆ ที่มาจาก  $E[X]$

- WEIGHTED MEAN:  $\bar{x} = E[X]$

$$= \sum p_i x_i ; p_i = \frac{w_i}{\sum w_i}$$

$$= \sum \frac{w_i}{\sum w_i} x_i$$

$$= \frac{\sum w_i x_i}{\sum w_i}$$

- ROOT MEAN SQUARED:  $x_{\text{RMS}} = \sqrt{E[X^2]}$  ← MATCHING  $\square^2$  AND  $\sqrt{\square}$

$$= \sqrt{\sum p_i x_i^2}$$

$$= \sqrt{\frac{\sum X^2}{N}}$$

- CONTINUOUS CASE:  $E[X] = \lim_{\Delta x \rightarrow 0} \sum_{i=1}^N p_i x_i$  ← PDF แทนเพราะจะได้มาเหมือนกัน กับ  $\sum p_i = 1 : \int_{\mathbb{R}} f(x) dx = 1$

$$E[X] = \int_{\mathbb{R}} x f(x) dx$$

### สมบัติของ VARIANCE

1.  $\text{Var}(X) \geq 0$

2.  $\text{Var}(X) = 0$  if  $x = \text{constant}$

3.  $\text{Var}(X+a) = \text{Var}(X)$

4.  $\text{Var}(aX) = a^2 \text{Var}(X)$

5.  $\text{Var}(aX+bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X, Y)$  } ทั่วไป?

6.  $\text{Var}(aX-bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) - 2ab \text{Cov}(X, Y)$  }

๕. ผลรวมของ  $\text{Var}(X_i)$  ที่  $X_i$  และ  $X_j$  ซึ่งไม่ใช่ค่าเดียวกัน และไม่มีความสัมพันธ์กันเลย (เขียนในรูปคณิตศาสตร์ได้ว่า  $\text{Cov}(X_i, X_j) = 0, \forall (i \neq j)$ ) ได้ว่า

$$\text{Var}\left(\sum_{i=1}^N X_i\right) = \sum_{i=1}^N \text{Var}(X_i)$$

หรือในรูปที่อ่านง่ายขึ้นเล็กน้อยว่า  $\sigma_{\text{SUM}}^2 = \sigma_1^2 + \sigma_2^2 + \sigma_3^2 + \dots + \sigma_N^2$

### ③ COVARIANCE ( $\text{Cov}(X), \sigma(X, Y), \sigma_{XY}$ )

COVARIANCE คือ VARIANCE ในรูปแบบของข้อมูล 2 ตัวที่มีความแปรปรวน "ร่วมกัน" ทั่วไป แตกต่างหรือไม่เท่าไร หากอธิบายง่าย ๆ คือ ข้อมูล 2 ตัวนี้ขยับกันไปจากค่าเฉลี่ยของตัวเองเท่าไร เขียนในรูปคณิตศาสตร์ได้ว่า

$$\begin{aligned} \text{Cov}(X, Y) &= E[(X - E[X])(Y - E[Y])] \\ &= E[XY - XE[Y] - YE[X] + E[X]E[Y]] \\ &= E[XY] - E[X]E[Y] - E[X]E[Y] + E[X]E[Y] \\ &= E[XY] - E[X]E[Y] \end{aligned}$$

หรือหากมอง  $X, Y$  เป็น DISCRETE VARIABLES :

$$\text{Cov}(X, Y) = \sum_{i=1}^N p_i (x_i - E[X])(y_i - E[Y])$$

ซึ่ง  $p_i$  ในที่นี้ คือโอกาสที่จะเจอคู่อันดับ  $(x_i, y_i)$  / อาจพูดได้ว่า  $p_{XY}(x, y) = P(X \cap Y)$

ย้อนกลับไปปรากฏการณ์ของเบย์ (Bayes' Theorem) !

ให้  $A$  และ  $B$  เป็นเหตุการณ์

① หาก  $A$  และ  $B$  ไม่เกี่ยวข้อง (independent) :

$$P(A \cap B) = P(A) \cdot P(B)$$

② หาก  $A$  และ  $B$  เกี่ยวข้องกัน (dependent) :

$$P(A \cap B) = P(A|B) \cdot P(B)$$

$P(A|B)$  หมายถึงโอกาสที่  $A$  จะเกิด ก็ต่อเมื่อมีเหตุการณ์  $B$  เกิดขึ้นแล้ว

เช่น พร้อมมาทำงานสาย ขณะที่ ฝนตก

A

B

ข้อ ② เป็นรูปทั่วไปของความน่าจะเป็นของความสัมพันธ์ของข้อ 2 เหตุการณ์  
แบบได้ใจจริงขึ้นมาขนาดความสัมพันธ์:

$$P(A \cap B) = P(B \cap A)$$

$$P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$$

$$\therefore P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} \rightsquigarrow \text{Bayes' Theorem}$$

ทบทวน COVARIANCE!

หากเรามอง Cov. ของตัวแปรเดิม,  $\text{Cov}(X, X)$  จะได้

$$\begin{aligned} \text{Cov}(X, X) &= E[X \cdot X] - E[X]E[X] \\ &= E[X^2] - E[X]^2 \\ &= \text{Var}(X) \end{aligned}$$

๖๖๖๖๖๖  $X$  ๖๖๖๖  $Y$  ไม่สัมพันธ์กัน  $\rightarrow \text{Cov}(X, Y) = 0$   
ถ้า... ๖๖๖๖

สมบัติพิเศษบางประการ

ให้  $X, Y, W, V = \text{random variables} \in \mathbb{R}$  ๖๖๖๖  $a, b, c, d = \text{constants} \in \mathbb{R}$

1.  $\text{Cov}(X, a) = 0$

2.  $\text{Cov}(X, X) = \text{Var}(X)$

3.  $\text{Cov}(X, Y) = \text{Cov}(Y, X)$

4.  $\text{Cov}(aX, bY) = ab\text{Cov}(X, Y)$

5.  $\text{Cov}(X+a, Y+b) = \text{Cov}(X, Y)$

6.  $\text{Cov}(aX+bY, cW+dV) = ac\text{Cov}(X, W) + ad\text{Cov}(X, V) + bc\text{Cov}(Y, W) + bd\text{Cov}(Y, V)$  } ทำไม?

๗.  $\sigma_{XY}^2 \leq \sigma_X^2 \sigma_Y^2$

↑  
คล้ายๆ กันอะไร?

(ลองหาอ่าน TRIANGLE INEQUALITY)  
(หาอ่านเรื่อง Vector:  
Cauchy-Schwarz inequality)



มาลองดูแคลคูลัสพิเศษของ COVARIANCE ว่าสามารถทำอะไรได้บ้าง?

$$\text{Cov}(X, Y) = E[XY] - E[X]E[Y]$$

$$\rightarrow \text{Cov}(X, Y) = 0 \rightarrow E[XY] = E[X]E[Y]$$

$$\sum \sum (p_{ij} x_i y_j) = (\sum p_i x_i) (\sum p_i y_i)$$

For discrete  $x_i, y_i$ ;  $p_i = \frac{1}{N}$

$$\sum \sum p_x p_y x_i y_i = (\sum p_x x_i) (\sum p_y y_i)$$

$$\therefore \sum (p_x x_i) \sum (p_y y_i) = \sum (p_x x_i) \sum (p_y y_i) \blacksquare$$

หากนี้ ตัวแปร  $x$  ไม่ขึ้นกับ  $y$  ( $y$  อาจเป็น constant),

$$N \cdot c \sum x_i = N \cdot c \sum x_i \rightarrow \text{เป็นจริง} \blacksquare$$

④ NORMALIZATION พื้นฐาน: ค่าสัมประสิทธิ์ความแปรปรวน

COEFFICIENT OF VARIATION ( $C_V$ ) หรือ RELATIVE S.D. คือการหาค่า

ส่วนเบี่ยงเบนมาตรฐานมาคิดเทียบกับค่าเฉลี่ย หมายถึงว่าค่าที่วัด/สำรวจได้ เบี่ยงเบนไปเป็นกี่เท่า / กี่% ของค่าที่ควรเป็น เช่น  $2.52 \text{ cm} \pm 5\%$  ← 0.05 เท่าของ 2.52 นั่นคือ 0.13 cm เขียนในรูปคณิตศาสตร์ได้ว่า

$$C_V = \frac{\sigma}{\mu}$$

บางครั้งหากทราบงาน S.D. = 950 อาจจะไม่จ้องเยอะ แต่พอมารู้ทีหลังว่าค่าที่วัดได้คือ 2,500,680 ซึ่งโผล่ออกไปที เพราะเพียง  $\pm 0.04\%$  เท่านั้น (น้อยมาก) แต่ถ้าค่าที่วัดได้เป็น 1,000 แล้ว S.D. เป็น 950 คงต้องตกใจเป็นอย่างมากเพราะค่าเพียงเพียงเกือบเท่าตัวของค่าจริง (แล้ว) แต่ไม่เหมาะแก่ค่าเล็กๆ เพราะเมื่อ  $\mu \rightarrow 0$  จะได้  $C_V \rightarrow \pm \infty$  ซึ่งไม่สวยเท่าไรนัก

⑤ NORMALIZATION เบื้องต้น : ค่า STANDARD ERROR /  $\sigma_{\text{mean}}$

STANDARD ERROR (SE,  $\sigma_{\bar{x}}$ ) มีความแปรปรวนในการใช้งานมากกว่า CV เพราะถูกหลักการการคิด sample (ตัวอย่าง) ของข้อมูลหลายๆ ชุดมากกว่า โดยคิดจากค่าคาดหวังของความแปรปรวน POPULATION

$\sigma_{\bar{x}}^2 = E[E[\sigma^2]]$ $\sigma_{\bar{x}}^2 = \frac{1}{n} \sum \frac{\sigma^2}{N}$ $\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n}$ $\therefore \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} \rightarrow \sigma_{\text{sample}} \text{ ใช้ } N-1$	โดย $\sigma \approx \sigma_x$ ← SAMPLE	หรือ $\sigma^2 = \sigma_1^2 + \sigma_2^2 + \dots + \sigma_N^2$ $\sigma^2 = N\sigma_x^2$ $\therefore \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{N}} = \frac{\sigma}{\sqrt{n}}$
--	--	--

ซึ่งค่านี้นิยมใช้กับการวัดค่าด้วยเครื่องมือต่างๆ หลายครั้งหาค่าที่คลั่งที่สุด (ค่าอะไร? อย่งไร?) สามารถใช้  $\sigma_{\bar{x}}$  เป็นค่า uncertainty ได้ชัดเจนกว่าค่าจากเครื่องมือแล้ว เหมาะสม (if  $\sigma_{\bar{x}} > \delta x$ )  $\delta x$  คือ UNCERTAINTY OF MEASUREMENT ซึ่งหาค่า  $\sigma_{\bar{x}}$  มาใช้รายงาน 95% confidence interval ( $\mu \pm 1.96\sigma_{\bar{x}}$ ) ได้ด้วย

⑥ NORMALIZATION ประยุกต์ : ค่าสัมประสิทธิ์สหสัมพันธ์

ค่าสัมประสิทธิ์สหสัมพันธ์ (PEARSON CORRELATION COEFFICIENT,  $r, \rho$ )

คือค่าที่บอกความสัมพันธ์ของข้อมูล 2 ชุดว่าสัมพันธ์/เกี่ยวเนื่องกันหรือไม่ อย่งไร

ซึ่งช่วงของค่าบ่งชี้ได้ดังนี้:

$$r = \begin{cases} +x & ; \text{สัมพันธ์แปรผันตาม } 0 < r \leq 1 \\ 0 & ; \text{ไม่สัมพันธ์กันเลย, } r = 0 \\ -x & ; \text{สัมพันธ์แปรผกผัน } -1 \leq r < 0 \end{cases}$$

โดยค่า  $r$  คำนวณได้จากค่าความสัมพันธ์ของ 2 ตัวแปร (COVARIANCE) มาทำการ

NORMALIZATION กับ  $\sigma_x \sigma_y$  (เหมือนหาความสัมพันธ์ INEQUALITY ใน Cov.)

ที่ว่า  $\sigma_{xy}^2 \leq \sigma_x^2 \sigma_y^2 \rightarrow \left| \frac{\sigma_{xy}}{\sigma_x \sigma_y} \right| \leq 1 \rightarrow -1 \leq \frac{\sigma_{xy}}{\sigma_x \sigma_y} \leq 1$

จะได้ว่า

$$r_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} \quad \leftarrow \text{(*)}$$

หรือ

$$= \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y}$$

หรือ

$$= \frac{E[XY] - E[X]E[Y]}{\sqrt{E[X^2] - E[X]^2} \sqrt{E[Y^2] - E[Y]^2}}$$

หรือ

$$= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

หรือ

$$= \frac{1}{n-1} \sum z_x z_y \quad \leftarrow \text{มาจากไหน?}$$

ซึ่งทั้งหมดที่กล่าวมาล้วนมีค่าคล้ายคลึงกับ COSINE'S SIMILARITY:

$$\cos \theta = \frac{\vec{A} \cdot \vec{B}}{AB} = \frac{\sum A_i B_i}{\sqrt{\sum A_i^2} \sqrt{\sum B_i^2}} ; \quad \begin{matrix} A_i = NE(X) \\ B_i = NE(Y) \end{matrix}$$

$$\rightarrow A = |\vec{A}| = \sqrt{\vec{A} \cdot \vec{A}} = \sqrt{\sum A_i^2}$$

$$\rightarrow \vec{A} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} \rightarrow \vec{A} \cdot \vec{A} = \sum A_i^2, \quad \vec{B} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

⊕ มีศาสตร์ที่น่าสนใจ น้อยคนจะรู้จักเขา (ไม่ค่อยได้ยินใครพูดกับเรื่องนี้สักเท่าไร) คือ

เรื่อง DIRECTIONAL STATISTICS (CIRCULAR STATISTICS)



และทศนิยมต่อตำแหน่งสุดท้ายตรงกัน

สำหรับความคลาดเคลื่อน ในตัวเลขที่มีสำคัญเพียง 1 จำนวนเท่านั้น เช่น

หน่วย ก.	สูง	179.8	± 0.4	cm	หรือ
		1.798	± 0.004	m	หรือ
		1798	± 4	mm	

อย่า รายงานค่าเป็น 179.82 ± 0.4 cm หรือ

179.8 ± 0.43 cm \*

แต่ในบางครั้งอาจเจอ 179.82 ± 0.45 cm แต่ไม่นิยมเป็นอย่างมาก

\* เป็นตัวเลขที่มีสำคัญต่ำที่สุด (1) ให้ปัดเลขก่อน จึงรายงาน เช่น ±0.45 → ±0.5

2. ค่าระหว่งการคำนวณ (ก่อนรายงาน) ให้เก็บไว้จนถึงตัวก่อนสุดท้าย หรือ เก็บไว้

5-6 หลัก หรือ เก็บไว้ให้มากที่สุด  $\leftarrow$  ควรทำ เพราะอาจเกิดความคลาดเคลื่อนจากการคำนวณมาบ้าง (พลาดง่าย)

การวัด 8x ใช้อย่างไร?

1. เขาคงเครื่องมือ หากอ่านเพียงครั้งเดียว เช่น ไม่บรรทัดทานความยาว ซึ่งไม่บรรทัด

มีสเกลละเอียดสุด 1 mm แต่ตามนุษย์อ่านได้ละเอียดที่สุด (แล้วยังแม่นยำอยู่)

คือ ครึ่งสเกล ในที่นี้ คือ  $\frac{1}{2}$  mm หรือ 0.05 cm ดังนั้นเราสามารถรายงานค่า

ได้เป็น 2.65 cm แต่จะเป็น 2.653 cm ไม่ได้ (ต้องมีสเกลเวอร์บขึ้น)

อุปกรณ์ที่มี VERNIER SCALE เช่น VERNIER CALIPER, MICROMETER, ฯลฯ

2. อ่านต่อมวี่เดิม ซ้ำ หลายครั้ง และนำค่า N ครั้ง จะได้ค่าที่ไม่เหมือนกัน

ออกมา N ค่า :  $x_1, x_2, x_3, \dots, x_N \rightarrow X$  คือค่าที่วัดได้  
ส่วน  $x$  คือค่าจริง (พระเจ้าเท่านั้นที่รู้)

วิธีที่ขมขื่นที่สุดและง่ายที่สุดคือ พิจารณาผลต่างจากค่าจริงกำลังสองอย่างง่าย

ความเป็นจริงทางคณิตศาสตร์ นั่นคือค่า  $(x_i - x)^2$  เพราะ  $\sum (x_i - x)^2 > 0$  เสมอ

แต่  $\sum (x_i - x) = 0$

LEAST SQUARES

ค่า  $x$  จะเป็นอย่างไร ใช้แคลคูลัสในการหาด้วยหลักการ MINIMIZE นั่นคือ

$\frac{d}{dx} f(x) = 0$  จะได้จุดวิกฤตออกมา

$$\frac{d}{dx} \sum_{i=1}^N (x_i - x)^2 = 0$$

$$-2 \sum_{i=1}^N (x_i - x) = 0$$

$$\sum_{i=1}^N x_i = \sum_{i=1}^N x = Nx$$

$$\therefore x = \frac{\sum x_i}{N} \rightarrow \text{นั่นคือค่าเฉลี่ยเลขคณิต} (\bar{x})$$

ซึ่ง  $x$  เป็นค่า  $x$  ดีสุด (THE BEST ESTIMATE ของ  $x$ ) ดังนั้น  $\frac{\sum x_i}{N}$  น่าจะ  
เป็น  $x$  ที่พระเจ้ารู้แน่ชัด แต่ทำเราวัด  $N$  บางอย่าง  $(\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N x_i)$

จะเป็นค่า  $x$  ที่แท้จริงในอุดมคติ แปลว่า ยิ่งวัดค่าเยอะยิ่งแม่นยำขึ้น (ควรจะ)

ดังนั้นความเอนางของค่า  $x_i$  จาก  $x$  กำลังสองเฉลี่ย จะเป็น  $\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 = \sigma^2$  (VARIANCE)

(แต่บางครั้งมักใช้  $N-1$  เพื่อความถูกต้อง (แต่ปริมาณมากๆ เริ่มไม่ค่อยต่าง)

ถ้ามีการแจกแจงทางสถิติ : NORMAL DIST., VARIANCE,  $\sigma_x^2$  ได้กล่าวไปข้างต้นแล้ว

ตัวอย่าง วัดความยาว  $l$  ของแท่งไม้ 5 ครั้ง ได้ผลเป็น

24.25, 24.26, 24.22, 24.28, 24.24 cm

จงรายงานผลความยาวแท่งไม้นี้ (ตอบ  $24.25 \pm 0.01$  cm)

หรือ  $24.247 \pm 0.009$  cm

$i$	$x_i$	$(x_i - \bar{x})^2$
		คิดในตารางนี้

ตามหนังสือ

การรวมความคลาดเคลื่อน (Propagation of uncertainties)

สมมติวัดค่า  $x$  ได้  $x = \bar{x} \pm \delta x$  และ  $y$  ได้  $y = \bar{y} \pm \delta y$

1. ผลบวก :  $x + y = (\bar{x} + \bar{y}) \pm (\delta x + \delta y)$

$x + 2y = x + y + y = (\bar{x} + 2\bar{y}) \pm (\delta x + 2\delta y)$

2. ผลต่าง :  $x - y = (\bar{x} - \bar{y}) \pm (\delta x + \delta y)$

บวก

3. ผลคูณ และ ผลหาร

$x = \bar{x} \pm \delta x = \bar{x} \left(1 \pm \frac{\delta x}{\bar{x}}\right)$

$y = \bar{y} \pm \delta y = \bar{y} \left(1 \pm \frac{\delta y}{\bar{y}}\right)$

$xy = \bar{x} \cdot \bar{y} \left(1 \pm \left(\frac{\delta x}{\bar{x}} + \frac{\delta y}{\bar{y}}\right) \pm \frac{\delta x \delta y}{\bar{x} \bar{y}}\right)$

สังเกต  $\left(\frac{\delta x \delta y}{\bar{x} \bar{y}}\right) \ll \frac{\delta x}{\bar{x}}$

$\ll \frac{\delta y}{\bar{y}}$

$\therefore xy = \bar{x} \cdot \bar{y} \pm \left(\frac{\delta x}{\bar{x}} + \frac{\delta y}{\bar{y}}\right) \cdot \bar{x} \bar{y}$

ในการคูณและหาร  $\frac{\delta A}{A}$  มาบวกกัน สิ่งนี้คือ RELATIVE UNCERTAINTY

4. ทฤษฎีการเชิงอนุพันธ์ : ใช้ปริมาณของพิกัด  $\sigma$  แทน  $\sigma$  ถ้า  $\Sigma$  คือ Sigma

5. กรณีทั่วไป เช่น  $z = f(x)$  ใช้การประมาณการเบี่ยงเบนคือ  $\frac{\delta z}{\delta x} = \frac{d}{dx} z$

ดังนั้น  $\delta z = \delta x \frac{d}{dx} f(x)$

จะได้ว่า  $z = f(\bar{x}) \pm \delta x f'(x)|_{x=\bar{x}}$

6. กรณีทั่วไป 2 มิติ  $z = f(x, y) \rightarrow \delta z = \frac{\partial f}{\partial x} \delta x + \frac{\partial f}{\partial y} \delta y$

$\rightarrow (\delta z)^2 = \left(\frac{\partial f}{\partial x}\right)^2 (\delta x)^2 + \left(\frac{\partial f}{\partial y}\right)^2 (\delta y)^2 + 2\left(\frac{\partial f}{\partial x}\right)\left(\frac{\partial f}{\partial y}\right) \delta x \delta y$

$= 0 \because \sum_x \sum_y (\dots) = 0$  / ค่าบวกหักลบ

หากเป็นการกระจายสุ่มพหุคูณ  $\delta A$  เป็น  $\sigma_{mean}$  ได้ว่า

$$e_{\sigma_{mean}} = \sqrt{\left(\frac{\partial f}{\partial x}\right)^2 \sigma_x^2 + \left(\frac{\partial f}{\partial y}\right)^2 \sigma_y^2}$$

$$\therefore e_f = f(\bar{x}, \bar{y}) \pm e_{\sigma_{mean}}$$

### 7. ผลรวมทางสถิติ

จากสมบัติผลรวมของ  $\text{Var}(X_i)$  ได้ว่า  $\sigma^2 = \sigma_1^2 + \sigma_2^2 + \dots + \sigma_N^2$

ใช้แทนได้ทั้งหมดว่า  $\sigma = \sqrt{\sum \sigma_{mean}^2}$

เช่น 1.  $x+y = (\bar{x} + \bar{y}) \pm \sqrt{(\delta x)^2 + (\delta y)^2}$  ซึ่งดีกว่า

3.  $xy = \bar{x}\bar{y} \pm \bar{x}\bar{y} \sqrt{\left(\frac{\delta x}{\bar{x}}\right)^2 + \left(\frac{\delta y}{\bar{y}}\right)^2}$

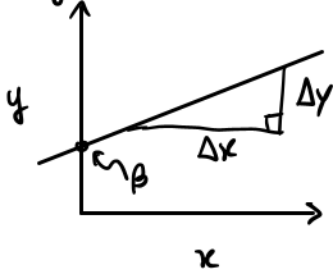
### การเขียนกราฟจากข้อมูลการทดลอง

กราฟ มีหน้าที่ทำให้เราเห็นถึงความสัมพันธ์และพัฒนารอบข้อมูล แต่ละแกน - นิยมที่จะนำข้อมูลมาเขียนในรูปของเส้นตรงเนื่องจากตีความง่าย คือความชัน และมองรูปร่างว่าอะไรคืออะไร

เราพิจารณากราฟความสัมพันธ์ของ 2 ปริมาณ  $y$  กับ  $x$  ในความสัมพันธ์  $y(x)$

เช่น  $y = ax + b$ ,  $xy = c$ ,  $y = ax^n$ ,  $y = ae^{ax}$ ,  $y = \tan(x^2)$ , ฯลฯ

1.  $y = ax + b$

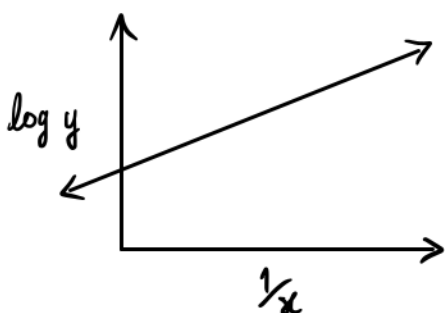


$$\text{ความชัน} = \frac{\Delta y}{\Delta x}$$

$$y\text{-intercept} = \beta = b$$

2. รูปอื่น ๆ เช่น  $\log y = a \cdot \frac{1}{x} + b$

ให้ใช้แกนตั้งเป็น  $\log y$  และแกนนอนเป็น  $\frac{1}{x}$



จะมีรูปเป็นเส้นตรง

โดยทั่วไปข้อมูลที่เรากำลังทำกราฟนี้ไว้

เส้นต่อเนื่อง (continuous) แต่เป็น

จุดข้อมูลที่ไม่ต่อเนื่อง (discrete)

ดังนั้นเราควรวางจุดที่วัดได้ให้อยู่ใน

รูปของเส้นตรงให้ได้อย่างที่ถูกต้อง

ไม่ใช่ใช้วิธีอื่น หรือ เช็กไปมาเทียบ

นั่นมีใช้เส้นตรง



# ข้อมูลที่มีความสัมพันธ์แบบเส้นตรงและการถดถอย (LINEAR REGRESSION)

สมมติ ทรงกระบอก ให้หาความสัมพันธ์ของปริมาตร  $V$  เทียบกับ  $h$  ว่ามีความสัมพันธ์

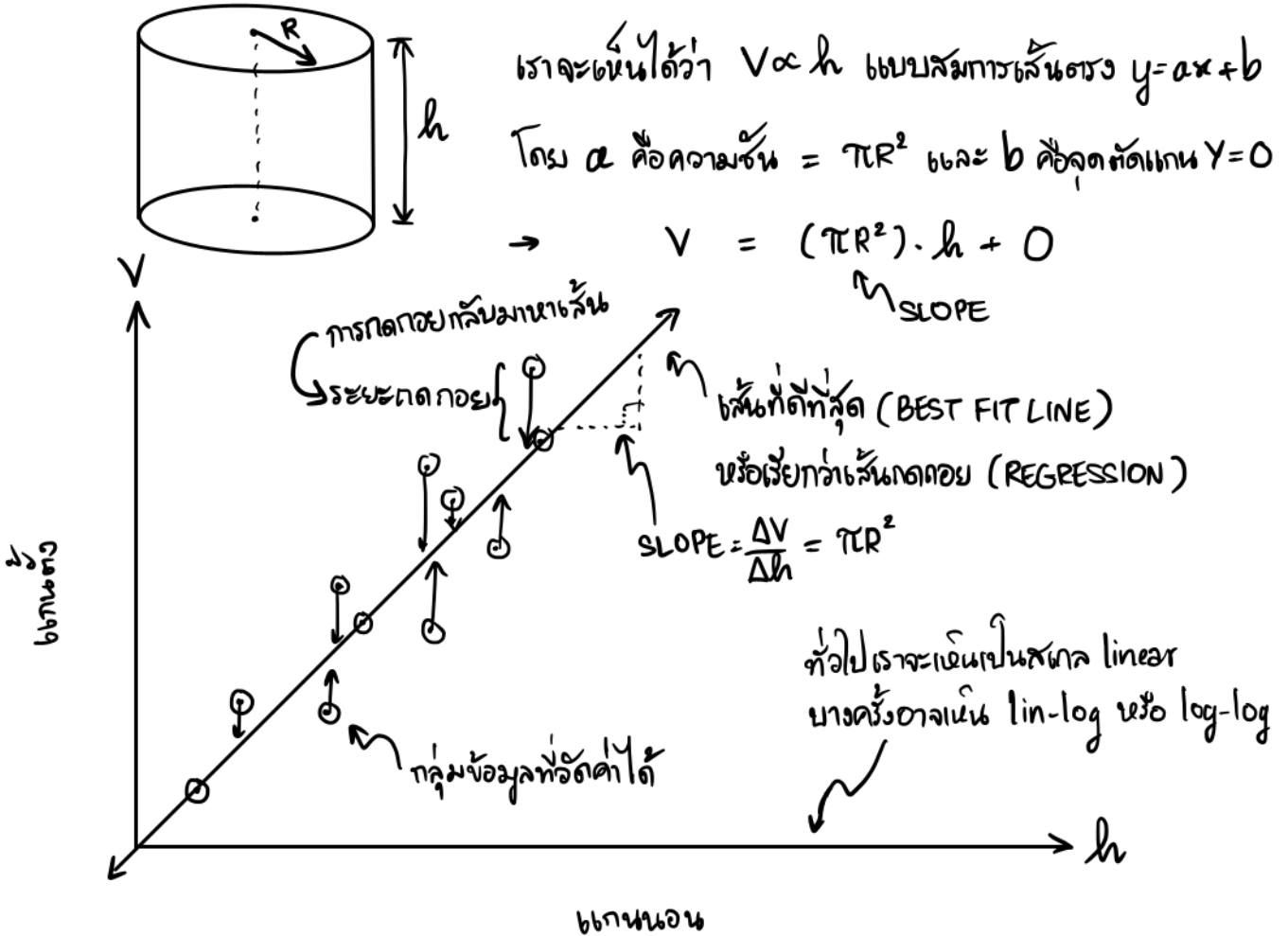
กันอย่างไร เริ่มจากใช้ความสัมพันธ์  $V = \pi R^2 h$

เราจะเห็นว่าได้ว่า  $V \propto h$  แบบสมการเส้นตรง  $y = ax + b$

โดย  $a$  คือความชัน =  $\pi R^2$  และ  $b$  คือจุดตัดแกน  $Y=0$

$$\rightarrow V = (\pi R^2) \cdot h + 0$$

↖ SLOPE



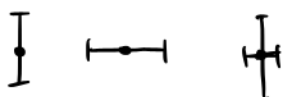
การลากเส้น BEST FIT ด้วยมือนั้น อาศัยความชำนาญ และความสามารถในการประมาณค่าที่ดีที่สุดออกมาได้ ซึ่งเป็นพื้นฐานที่ตองของนักวิทยาศาสตร์, วิศวกร, นักวิเคราะห์, ฯลฯ

โดยเราต้องการลากให้ผลรวมระยะถดถอยมีค่าน้อยที่สุด (MINIMIZATION) หรือให้เส้นตรงมากที่สุด แต่ในเชิงเวลาจริง เราก็ใช้เครื่องคำนวณหาเส้นนั้นออกมา เพราะมันแม่นยำกว่ามือเรามาก ๆ แต่ใครคิดสูตรพวกนั้นให้เครื่องจักร?

V	h
...	...

ตารางพล็อต

\* เพิ่มเติม หากข้อมูลมีความคลาดเคลื่อนมาก จะแสดงเป็น ERROR BARS



สมมติตัวแปรเก็บค่าได้เป็นคู่อันดับ  $(x_i, y_i)$  ซึ่งอาจสร้างสมการเส้นตรง  
ได้เป็น  $\bar{y} = ax_i + b$

จากทฤษฎี LEAST SQUARES ซึ่งนำมาหา BEST ESTIMATE เราต้องการให้ค่า

$$S = \sum_{i=1}^n (y_i - \bar{y})^2 \text{ น้อยที่สุด,}$$

$$S = \sum_{i=1}^n (y_i - ax_i - b)^2$$

หา  $b$  :  $\frac{\partial S}{\partial b} = 0$

$$0 = -2 \sum_{i=1}^n (y_i - ax_i - b)$$

$$\sum_{i=1}^n b = \sum_{i=1}^n y_i - a \sum_{i=1}^n x_i$$

$$b = \frac{1}{n} (\sum y_i - a \sum x_i)$$

$$\therefore b = \bar{y} - a\bar{x} \quad *$$

หา  $a$  :  $\frac{\partial S}{\partial a} = 0$

$$0 = -2 \sum_{i=1}^n x_i (y_i - ax_i - b)$$

$$0 = \sum_{i=1}^n (x_i y_i - ax_i^2 - bx_i)$$

$$0 = \sum_{i=1}^n (x_i y_i - (\bar{y} - a\bar{x})x_i - ax_i^2)$$

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

$$\therefore a = \frac{\sum x_i (y_i - \bar{y})}{\sum x_i (x_i - \bar{x})} \quad *$$