S — MAS (1)

- population
  - parameter : characteristic of population
  - sample
  - statistics : characteristic of sample
  - variable : characteristic of elements
  - data : value of variable

- Phương pháps collect data
  - retrospective study : các data có từ quá khứ
  - observational study : data từ quan sát, đo đạc
  - experment study : data từ thực nghiệm
  - simulation study : using models → data
  - survey : ⌈ sample
    ⌊ population : CENSUS

- Type of data
  - qualitative (định tính) : gender, color, major, place, size (~những thứ để để phân loại)
  - quatitative —— discrete (đếm rời rạc)
    (định lượng) ↘ continuous (liên tục)

- Sampling method
  - representative : lấy sao cho đại diện đc population
  - replacement / without replacement :
    chọn xg bỏ ra (k° lấy nx) / chọn xg bỏ lại (có thể lấy tiếp)
  - non random (not representative)
  - random sampling :
    + simple random
    + stratified : bốc đều, từ các nhóm (class), mỗi class simple random

HONG HA

## I. Basic probability formulas

- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

- $P(A \mid B) = \dfrac{P(A \cap B)}{P(B)}$

- $P(A \mid B) = \dfrac{P(B \mid A) \cdot P(A)}{P(B)}$

- If A, B independent: $P(A \cap B) = P(A) \cdot P(B)$

## II. Discrete random variables

- $\mathcal{M} = E(x) = \sum x_i \cdot P(x = x_i)$

- $\sigma^2 = V(x) = \sum (x_i - \mathcal{M})^2 \cdot P(x = x_i)$

$$= \sum x_i^2 . P(x=x_i) - \mathcal{M}^2$$

- $E(ax + by) = a.E(x) + b.E(y)$
- $V(ax + by) = a^2 . V(x) + b^2 . V(y)$
- Probability mass function: $f(x_i) = P(x=x_i)$
- Cumulative distribution function: $F(x_i) = P(x \leq x_i)$
- Some special distribution:
  1. Discrete uniform distribution
     - $P(x=X_i) = \dfrac{1}{n}$
     - $\mathcal{M} = \dfrac{a+b}{2}$
     - $\sigma^2 = \dfrac{(b-a+1)^2 - 1}{12}$
  2. Binomial distribution
     - $P(x=k) = nCk . p^k . (1-p)^{n-k}$
     - $\mathcal{M} = n.p$
     - $\sigma^2 = n.p . (1-p)$
  3. Poisson distribution
     - $P(x=k) = \dfrac{e^{-\lambda.T}}{k!} (\lambda.T)^k$
     - $\mathcal{M} = \lambda.T$
     - $\sigma^2 = \lambda.T$
  4. Hypergeometric distribution
     - $P(x=k) = \dfrac{KCk . (N-K)C(n-k)}{NCn}$
     - $\mathcal{M} = n.p$
     - $\sigma^2 = n.p.(1-p). \dfrac{N-n}{N-1}$
  5. Geometric distribution
     - $P(x=k) = (1-p)^{k-1} . p$
     - $\mathcal{M} = \dfrac{1}{p}$
     - $\sigma^2 = \dfrac{1-p}{p^2}$

6. Negative binomial distribution
- $P(x=k) = (k-1)C(r-1) \cdot p^r \cdot (1-p)^{k-r}$
- $\mathcal{M} = \dfrac{r}{p}$
- $\sigma^2 = \dfrac{r \cdot (1-p)}{p^2}$

## III. Continuous random variable

- Probability density function $f(x)$: $P(a<x<b) = \displaystyle\int_a^b f(x)\, d_x$

- Cumulative distribution function $F(x)$:
    - $F(x_i) = P(x \le x_i)$
    - $F(x_i)' = f(x_i)$

- $\mathcal{M} = E(x) = \displaystyle\int_{-\infty}^{+\infty} x \cdot f(x)\, d_x$

- $E(x^n) = \displaystyle\int_{-\infty}^{+\infty} x^n \cdot f(x)\, d_x$

- $\sigma^2 = V(x) = \displaystyle\int_{-\infty}^{+\infty} x^2 \cdot f(x)\, d_x - \mathcal{M}^2$

- Some special distribution:
    1. Continuous uniform distribution
        - $f(x) = \dfrac{1}{b-a}$ , $a \le x \le b$

          $= 0$ , elsewhere
        - $\mathcal{M} = \dfrac{a+b}{2}$
        - $\sigma^2 = \dfrac{(b-a)^2}{12}$

    2. Normal distribution $N(\mathcal{M}, \sigma^2)$
        - $z = \dfrac{x - \mathcal{M}}{\sigma}$
        - $f(z) = \dfrac{1}{\sqrt{2\Pi}} \cdot e^{\frac{x^2}{2}}$
        - $\phi(x) = p(z<x_i)$

- ○ $\phi(-x) = 1 - \phi(x)$

3. Normal distribution approximate binomial and poisson distribution

    a. Binomial ($np > 5$ and $n(1-p) > 5$)

        ■ $z = \dfrac{x - n.p}{\sqrt{n.p.(1-p)}}$

        ■ $P(X_{BINORM} \leq a) = P(X_{NORMAL} \leq a+0.5)$

        ■ $P(X_{BINORM} \geqq a) = P(X_{NORMAL} \geqq a-0.5)$

    b. Poisson

        ■ $z = \dfrac{x - \lambda}{\sqrt{\lambda}}$

        ■ $P(X_{POISSON} \leq a) = P(X_{NORMAL} \leq a+0.5)$

        ■ $P(X_{POISSON} \geqq a) = P(X_{NORMAL} \geqq a-0.5)$

4. Exponential distribution

- ○ $f(x) = \lambda \cdot e^{-\lambda.T}$ , $x > 0$

- ○ $= 0$ , elsewhere

- ○ $P(x \geqq a) = e^{-\lambda.a}$ ,$(a > 0)$

- ○ $\mathcal{M} = \dfrac{1}{\lambda}$

- ○ $\sigma^2 = \dfrac{1}{\lambda^2}$

**IV. Descriptive statistic** (Take a sample of size n from population N)

- ● Sample mean: $\overline{x} = \dfrac{\sum x_i}{n}$

- ● Sample median: $L = \dfrac{n+1}{2}$ so Median $= \dfrac{x_{ceil(L)} + x_{floor(L)}}{2}$

- ● Mode: Số phần tử xuất hiện nhiều nhất

- ● Range: max - min

- ● Sample variance: $s^2 = \dfrac{\sum (\overline{x} - x_i)^2}{n-1}$

- ● Quatiles:

$\circ$ $L_1 = \dfrac{n+1}{4}$ so $Q_1 = \dfrac{x_{ceil(L_1)} + x_{floor(L_1)}}{2}$

$\circ$ $L_2 = \dfrac{n+1}{2}$ so $Q_2 = \dfrac{x_{ceil(L_2)} + x_{floor(L_2)}}{2}$

$\circ$ $L_3 = \dfrac{3.(n+1)}{4}$ so $Q_3 = \dfrac{x_{ceil(L_3)} + x_{floor(L_3)}}{2}$

## V. Sampling distribution

- Population mean $\mathcal{M}$, variance $\sigma^2$. Sample size n. *(Normal distribution or n > 30)*:

  $\circ$ Phân phối của $\overline{X}$ có dạng: $N(\mathcal{M}, \dfrac{\sigma^2}{n})$

  $\circ$ Phân phối của $\overline{X}_1 - \overline{X}_2$ có dạng: $N(\mathcal{M}_1 - \mathcal{M}_2, \dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2})$

- For proportion of population p, sample size n. *(np $\geqq$ 5 or n.(1-p) $\geqq$ 5)*:

  $\circ$ Phân phối của $\widehat{P}$ có dạng: $N(P, \dfrac{P.(1-P)}{n})$

  $\circ$ Phân phối của $\widehat{P_1} - \widehat{P_2}$ có dạng: $N(P_1 - P_2, \dfrac{P_1.(1-P_1)}{n_1} + \dfrac{P_2.(1-P_2)}{n_2})$

## VI. Statistical intervals - Test claims for one sample

- $(l, u) = (\overline{X} - E, \overline{X} + E)$
- width = 2E
- P-value = 2 . $P(Z > |Z_0|)$

1. Population variance known

   $\circ$ $E = z_{\alpha/2} \cdot \dfrac{\sigma}{\sqrt{n}}$

   $\circ$ $z_0 = \dfrac{\overline{X} - \mathcal{M}}{\sigma / \sqrt{n}}$

2. Population variance unknown

   $\circ$ n > 30:

     ▪ $E = z_{\alpha/2} \cdot \dfrac{S}{\sqrt{n}}$

- $z_0 = \dfrac{\overline{X} - \mathcal{M}}{S/\sqrt{n}}$

  - $n \le 30$:

    - $E = t_{\alpha/2,\,n-1} \cdot \dfrac{S}{\sqrt{n}}$

    - $t_0 = \dfrac{\overline{X} - \mathcal{M}}{S/\sqrt{n}}$

- For propotion:

  - $(l,\,u) = (\widehat{P} - E,\ \widehat{P} + E)$

  - $E = z_{\alpha/2} \cdot \sqrt{\dfrac{P.(1-P)}{n}}$

  - $z_0 = \dfrac{\widehat{P} - P}{\sqrt{\dfrac{P.(1-P)}{n}}}$

  - Nếu đề không cho $\widehat{P}$, mặc định $\widehat{P} = 0.5$

- Nếu là one-side thì tương tự nhưng thay $\alpha/2$ thành $\alpha$


**VII. Test claims for 2 samples** (2 population independent, normal distribution or both $n_1, n_2 > 30$)

- $(l,\,u) = (\overline{X_1} - \overline{X_2} - E,\ \overline{X_1} - \overline{X_2} + E)$

1. Population variance known

   - $E = z_{\alpha/2} \cdot \sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}$

   - $z_0 = \dfrac{\overline{X_1} - \overline{X_2} - \Delta_0}{\sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_2^2}{n_2}}}$

2. Population variance unknown

   - Assume $\sigma_1^2 = \sigma_2^2$

     - Degree of freedom: $df = n_1 + n_1 + 2$

     - $S_p^2 = \dfrac{(n_1 - 1).S_1^2 + (n_2 - 1).S_2^2}{n_1 + n_2 - 2}$

- $E = t_{\alpha/2, df} \cdot \sqrt{\dfrac{S_p^2}{n_1} + \dfrac{S_p^2}{n_2}}$

- $t_0 = \dfrac{\overline{X_1} - \overline{X_2} - \Delta_0}{\sqrt{\dfrac{S_p^2}{n_1} + \dfrac{S_p^2}{n_2}}}$

  - Not assume $\sigma_1^2 = \sigma_2^2$

    - Degree of freedom: $df = \dfrac{\left(\dfrac{S_1^2}{n_1} + \dfrac{S_2^2}{n_2}\right)^2}{\dfrac{S_1^4}{n_1^2 \cdot (n_1 - 1)} + \dfrac{S_2^4}{n_2^2 \cdot (n_2 - 1)}}$

    - $E = t_{\alpha/2, df} \cdot \sqrt{\dfrac{S_1^2}{n_1} + \dfrac{S_2^2}{n_2}}$

    - $t_0 = \dfrac{\overline{X_1} - \overline{X_2} - \Delta_0}{\sqrt{\dfrac{S_1^2}{n_1} + \dfrac{S_2^2}{n_2}}}$

- For propotion:

  - $(l, u) = (\widehat{P_1} - \widehat{P_2} - E \ , \ \widehat{P_1} - \widehat{P_2} + E)$

  - $E = z_{\alpha/2} \cdot \sqrt{\dfrac{\widehat{P_1} \cdot (1 - \widehat{P_1})}{n_1} + \dfrac{\widehat{P_2} \cdot (1 - \widehat{P_2})}{n_2}}$

  - $\widehat{P} = \dfrac{x_1 + x_2}{n_1 + n_2}$ (trong đó $x_i = n \cdot \widehat{P_i}$)

  - $z_0 = \dfrac{\widehat{P_1} - \widehat{P_2} - \Delta_0}{\sqrt{\widehat{P} \cdot (1 - \widehat{P}) \cdot \left(\dfrac{1}{n_1} + \dfrac{1}{n_2}\right)}}$

**VIII. Linear Regression**

- $S_{XY} = \sum \left(x_i - \overline{x}\right)\left(y_i - \overline{y}\right) = \sum x_i y_i - n \cdot \overline{x} \cdot \overline{y}$

- $S_{XX} = \sum \left(x_i - \overline{x}\right)^2 = \sum x_i^2 - n \cdot \overline{x}^2$

- $S_{YY} = \sum\left(y_i - \bar{y}\right)^2 = \sum y_i^2 - n \cdot \bar{y}^2$

- Slope: $\widehat{\beta}_1 = \dfrac{S_{XY}}{S_{XX}} = \dfrac{\sum x_i y_i - n \cdot \bar{x} \cdot \bar{y}}{\sum x_i^2 - n \cdot \bar{x}^2}$

- Intercept: $\widehat{\beta}_0 = \bar{y} - \widehat{\beta}_1 \cdot \bar{x}$

- Error sum of square: $SS_E = \sum\left(y_i - \widehat{y}_i\right)^2$

- Regression sum of square: $SS_R = \sum\left(\widehat{y}_i - \bar{y}\right)^2$

- Total sum of square: $SS_T = \sum\left(y_i - \bar{y}\right)^2$

- $SS_E + SS_R = SS_T$

- Standard error: $\widehat{\sigma} = \sqrt{\dfrac{SS_E}{n-2}}$

- Coefficient of correlation: $R = \sqrt{\dfrac{SS_R}{SS_T}} = \dfrac{S_{XY}}{\sqrt{S_{XX} \cdot S_{YY}}}$

- Test claims about the slope *(df = n-2)*:

  - $se(\widehat{\beta}_1) = \sqrt{\dfrac{\widehat{\sigma}^2}{S_{XX}}}$

  - $t_0 = \dfrac{\widehat{\beta}_1 - \beta_{1,0}}{se(\widehat{\beta}_1)}$

- Test claims about the intercept *(df = n-2)*:

  - $se(\widehat{\beta}_0) = \sqrt{\widehat{\sigma}^2 \cdot \left(\dfrac{1}{n} + \dfrac{\bar{x}^2}{S_{XX}}\right)}$

  - $t_0 = \dfrac{\widehat{\beta}_0 - \beta_{0,0}}{se(\widehat{\beta}_0)}$

- Test claims about the coefficient of correlation *(df = n-2)*: $t_0 = \dfrac{R - 0}{\sqrt{\dfrac{1 - R^2}{n-2}}}$