

Distribution

1. Bernoulli

$$X \sim \text{Bernoulli}(p)$$

$$P(X=1) = p, \text{ success}$$

$$\mu = p, \sigma^2 = p(1-p)$$

2. Binomial

$$X \sim \text{Bin}(n, p)$$

n - number of trials

$$\mu = np, \sigma^2 = np(1-p)$$

p - success probability

$$P(X=x) = \frac{n!}{x!(n-x)!} p^x (1-p)^{n-x}$$

◦ Uncertainty in sample proportion $\hat{p} = \frac{x}{n}$

$$\left\{ \begin{array}{l} \mu_{\hat{p}} = p, \text{ Unbiased} \\ \sigma_{\hat{p}} = \end{array} \right.$$

$$\sigma_{\hat{p}} = \frac{\sigma_x}{\sqrt{n}} = \sqrt{\frac{np(1-p)}{n}} = \sqrt{\frac{p(1-p)}{n}}$$

3. Poisson

$$X \sim \text{Poisson}(\lambda)$$

an approximation of Binomial when $n \rightarrow \infty, p \rightarrow 0$

$$\lambda = np, P(X=x) = e^{-\lambda} \cdot \frac{\lambda^x}{x!}$$

$$\mu = \lambda, \sigma_x^2 = \lambda$$

λ : 该度量 X 的概率，单位为单位的 $\frac{1}{\lambda}$;

number of events occur in one unit of time / space

• A rate estimate:

$$\text{If } X \sim \text{Poisson}(\lambda t), \hat{\lambda} = \frac{X}{t}$$

Uncertainty in $\hat{\lambda}$

$$\mu_{\hat{\lambda}} = \lambda, \text{ unbiased}$$

$$\sigma_{\hat{\lambda}} = \frac{\sigma_X}{t} = \frac{\sqrt{\lambda t}}{t} = \sqrt{\frac{\lambda}{t}}$$

4. Hypergeometric $X \sim H(N, R, n)$

N: size of finite population

R: number of items classified as success

n: sample size

X: number of success in the sample

$$P(X=x) = \frac{\binom{R}{x} \left(\frac{nR}{N}\right)^x}{\binom{N}{n}}$$

$$\mu = \frac{nR}{N},$$

$$\sigma^2 = n \left(\frac{R}{N} \right) \left(1 - \frac{R}{N} \right) \left(\frac{N-n}{N-1} \right)$$

5. Geometric $X \sim \text{Geom}(p)$

X: number of trials up to and including the 1st success

$$P(X=x) = p(1-p)^{x-1}$$

$$\mu = \frac{1}{p}, \sigma^2 = \frac{1-p}{p^2}$$

6. Negative Binomial $X \sim NB(r, p)$

X: number of trials up to and including rth success

$r = n$ th success

$$P(X=x) = \binom{x-1}{r-1} p^r (1-p)^{x-r}$$

$$\mu = \frac{r}{p}, \quad \sigma^2 = \frac{r(1-p)}{p^2}$$

$$X \sim NB(r, p) \iff X = Y_1 + \dots + Y_r, \quad Y_i \sim \text{Geom}(p)$$

7. Multinomial

$$X_1, \dots, X_k \sim MN(n, p_1, p_2, \dots, p_k)$$

$$P(X_1, \dots, X_k) = P(X_1=x_1, X_2=x_2, \dots, X_k=x_k) = \frac{n!}{x_1! x_2! x_3! \dots x_k!} p_1^{x_1} p_2^{x_2} \dots p_k^{x_k}$$

$$\text{each } X_i \sim \text{Bin}(n_i, p_i)$$

8. Normal

$$X \sim N(\mu, \sigma^2)$$

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$68\% = \mu \pm \sigma$$

$$95\% = \mu \pm 2\sigma$$

Z -score:

$$99.7\% = \mu \pm 3\sigma$$

$$Z = \frac{x-\mu}{\sigma} \Rightarrow \text{standard normal}$$

Sample mean \bar{X} :

$$\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$$

Linear Function =

$$\textcircled{1} \quad X \sim N(\mu, \sigma^2) \Rightarrow aX+b \sim N(a\mu+b, a^2\sigma^2)$$

$$\textcircled{2} \quad C_1 X_1 + C_2 X_2 \sim N(C_1 \mu_1 + C_2 \mu_2, C_1^2 \sigma_1^2 + C_2^2 \sigma_2^2)$$

$$\textcircled{3} \quad X+Y \sim N(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)$$

$$X-Y \sim N(\mu_X - \mu_Y, \sigma_X^2 + \sigma_Y^2)$$

9. Lognormal

$$Y = e^X, \quad X \sim N(\mu, \sigma^2); \quad \ln Y \sim N(\mu, \sigma^2)$$

$$\mu_Y = e^{\mu + \frac{1}{2}\sigma^2}, \quad \sigma_Y^2 = e^{2\mu + 2\sigma^2} - e^{2\mu + \sigma^2}$$

10. Exponential $X \sim \text{Expon}(\lambda)$

Waiting time

$$\mu = \frac{1}{\lambda}, \quad \sigma^2 = \frac{1}{\lambda^2}$$

$$f(x) = \lambda e^{-\lambda x}$$

$$F(x) = P(X \leq x) = 1 - e^{-\lambda x}$$

- Rate parameter λ from Poisson & $\text{Expon}(\lambda)(T)$

- Memoryless property.

$$P(T > t+s | T > s) = P(T > t)$$

- Uncertainty:

$\mu \pm \frac{1}{\sqrt{n}} \neq \frac{1}{\mu}$, bc $\frac{1}{\mu}$ is not a linear function of μ .

$$\textcircled{4} \quad \hat{\lambda} = \frac{1}{x}, \quad \text{uncertainty} = \frac{\lambda}{n}$$

$$\textcircled{2} \quad \sigma_{\bar{x}} = \frac{1}{\bar{x}\sqrt{n}}$$

$n > 20$ if, estimate $\hat{\theta}$

Point Estimation

① Measure the goodness = mean squared error (MSE)

$\hat{\theta}$: estimator

$\mu_{\hat{\theta}} - \theta$ = bias of estimator

$\sigma_{\hat{\theta}}^2$ = uncertainty = standard error.

$$\text{MSE}_{\hat{\theta}} = (\mu_{\hat{\theta}} - \theta)^2 + \sigma_{\hat{\theta}}^2 = \mu_{(\hat{\theta} - \theta)^2}$$

(bias² + Variance)

② Maximum Likelihood Estimate (MLE)

$n \uparrow$, bias of MLE $\rightarrow 0$; variance of MLE $\rightarrow \text{Min}$

Central Limit Theorem

$$X \sim N(\mu, \sigma)$$

\bar{X} - Sample mean \Rightarrow

S - Sample sum

$$\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$$

$$S_n \sim N(n\mu, n\sigma^2)$$

Confidence Intervals

probability \leftrightarrow random events \leftrightarrow computing method.

confidence \leftrightarrow fixed numerical value

1. Large sample CI for population mean

$n > 30$, $X \sim N(\mu, \sigma^2)$, \bar{X} is approximately normal

$100(1-\alpha)\%$ CI for μ :

• Two-sided:

$$\bar{X} \pm Z_{\frac{\alpha}{2}} \cdot \sigma_{\bar{X}}, \text{ 其中 } \sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}} = \frac{s}{\sqrt{n}}$$

• One-sided:

$$\text{lower CI: } \bar{X} - Z_{\alpha} \cdot \sigma_{\bar{X}}$$

$$\text{upper CI: } \bar{X} + Z_{\alpha} \cdot \sigma_{\bar{X}}$$

2. CI for proportions

由CLT可知，在大数定律下， $\hat{p} \sim N(p, \frac{p(1-p)}{n})$

已知 $X \sim \text{Bin}(n, p)$.

$$\text{令 } \tilde{n} = n+4, \quad \tilde{p} = \frac{x+2}{\tilde{n}}, \quad \text{则 } 100(1-\alpha)\% \text{ CI for } p:$$

• Two-sided

$$\tilde{p} \pm Z_{\frac{\alpha}{2}} \sqrt{\frac{\tilde{p}(1-\tilde{p})}{\tilde{n}}} \quad [0, 1]$$

• One-sided

$$\begin{aligned} \text{lower } \tilde{p} - Z_{\alpha} \sqrt{\frac{\tilde{p}(1-\tilde{p})}{\tilde{n}}} \\ \text{(upper) } \tilde{p} + Z_{\alpha} \sqrt{\frac{\tilde{p}(1-\tilde{p})}{\tilde{n}}} \end{aligned}$$

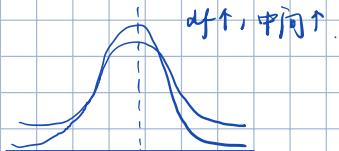
- Traditional method: (要求 success & fail & $n > 10$)

\hat{P} - Large number n 次 trials 中的 成功次数.

$$\text{CI for } P: \hat{P} \pm Z_{\alpha/2} \sqrt{\frac{\hat{P}(1-\hat{P})}{n}}$$

3. Small sample CI for population mean μ

$$n < 30, df = n-1, \left(\frac{\bar{X} - \mu}{s/\sqrt{n}} \right) \sim t_{n-1}$$



要求: ① population 是 normal ② 无 outlier

例: 100(1- α)% CI for μ :

$$\text{Two-sided: } \bar{X} \pm t_{n-1, \alpha/2} \cdot \frac{s}{\sqrt{n}}$$

$$\text{One-sided: lower } \bar{X} - t_{n-1, \alpha} \cdot \frac{s}{\sqrt{n}} \\ (\text{upper}) \quad (+)$$

$$\text{If } \sigma \text{ 已知, 用 } z\text{-distri 而不用 } t = \bar{X} \pm Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

$$\text{Single value } n=1 \text{ 时, } \bar{X} \pm Z_{\alpha/2} \cdot \sigma$$

4. CI for difference between 2 means

$$已知 X \sim N(\mu_X, \sigma_X^2), Y \sim N(\mu_Y, \sigma_Y^2), [n_1, n_2]$$

100(1- α)% CI for $(\mu_X - \mu_Y)$ 是:

$$(\bar{X} - \bar{Y}) \pm Z_{\alpha/2} \cdot \sqrt{\frac{\sigma_X^2}{n_1} + \frac{\sigma_Y^2}{n_2}} \quad (\sigma_X = s_X, \sigma_Y = s_Y)$$

5. CI for difference between 2 proportions.

$X \sim \text{Bin}(n_x, p_x)$, $Y \sim \text{Bin}(n_y, p_y)$

$$\text{令 } \tilde{n}_x = n_x + 2, \tilde{n}_y = n_y + 2, \tilde{p}_x = \frac{x+1}{\tilde{n}_x}, \tilde{p}_y = \frac{y+1}{\tilde{n}_y}$$

則 $100(1-\alpha)\%$ CI for $(p_x - p_y)$ 为：

$$(\tilde{p}_x - \tilde{p}_y) \pm z_{\frac{\alpha}{2}} \cdot \sqrt{\frac{\tilde{p}_x(1-\tilde{p}_x)}{\tilde{n}_x} + \frac{\tilde{p}_y(1-\tilde{p}_y)}{\tilde{n}_y}} \quad [-1, 1]$$

• Traditional method:

$$(\hat{p}_x - \hat{p}_y) \pm z_{\frac{\alpha}{2}} \cdot \sqrt{\frac{\hat{p}_x(1-\hat{p}_x)}{n_x} + \frac{\hat{p}_y(1-\hat{p}_y)}{n_y}}$$

(要 n_x & n_y 大于 10)

6. Small sample CI for difference between 2 means

$100(1-\alpha)\%$ CI for $(\mu_x - \mu_y)$:

$$v = \frac{\left(\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}\right)^2}{\frac{(s_x^2/n_x)^2}{n_x-1} + \frac{(s_y^2/n_y)^2}{n_y-1}}, \text{ round down to nearest integer}$$

$$\bar{X} - \bar{Y} \pm t_{v, \frac{\alpha}{2}} \sqrt{\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}}$$

若 X 与 Y 的 σ^2 相同，则：

$$\bar{X} - \bar{Y} \pm t_{n_x+n_y-2, \frac{\alpha}{2}} \cdot s_p \cdot \sqrt{\frac{1}{n_x} + \frac{1}{n_y}},$$

$$S_p = \sqrt{\frac{(n_x-1)s_x^2 + (n_y-1)s_y^2}{n_x+n_y-2}}, \text{ pooled std.}$$

7. CI with paired data.

100(1- α)% CI for difference of pairs μ_D :

◦ Small Sample ($n < 30$)

$$\bar{D} \pm t_{n-1, \frac{\alpha}{2}} \cdot \frac{s_p}{\sqrt{n}}$$

◦ Large Sample

$$\bar{D} \pm z_{\frac{\alpha}{2}} \cdot \sigma_{\bar{D}} = \bar{D} \pm z_{\frac{\alpha}{2}} \cdot \frac{s_p}{\sqrt{n}}$$

Hypothesis Testing

Steps:

① Define H_0 & H_1 ② Assume H_0 to be true

③ Compute a test statistic

(used to assess the strength of evidence against H_0)

④ Compute P-Value (observed significance level)

(the prob that assuming H_0 is true, the test statistic would have a value whose disagreement with H_0 is as great or greater than the actually observed)

⑤ Conclusion

rule of thumb: $p \leq 0.05$, significant

H_0 is rejected at $\alpha = 0.05$ level

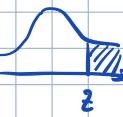
1. Large Sample, Population mean μ .

$n > 30$, Z test,

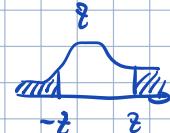
$$Z = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \rightarrow$$

P-value of H_1 : $p = P(Z)$

e.g. $H_1: \mu > \mu_0$, but p-value



$\mu < \mu_0$, $p = P(Z)$



2. Population Proportion

$X \sim \text{Bin}(n, p)$. $\bar{X} \sim np \sim n(p - np) \approx \mathcal{N}(np, np(1-p))$

$$Z = \frac{\hat{p} - p_0}{\sqrt{p_0(1-p_0)/n}} \rightarrow p\text{-value.}$$

3. Small Sample, population mean ($n < 30$)

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}}, \quad df = n-1.$$

• 当 σ 已知, 用 Z-test 而不是 t-test: $Z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$

4. Large sample, difference between 2 means

$$H_0: \mu_x - \mu_y \leq \Delta_0 \dots$$

$$Z = \frac{(\bar{x} - \bar{y}) - \Delta_0}{\sqrt{\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}}}, \quad (\sigma_x = s_x, \sigma_y = s_y)$$

5. Difference between 2 proportions

$$H_0: p_x - p_y \leq 0$$

$$\hat{p}_x = \frac{x}{n_x}, \quad \hat{p}_y = \frac{y}{n_y}, \quad \hat{p} = \frac{x+y}{n_x+n_y}.$$

$$Z = \frac{\hat{p}_x - \hat{p}_y}{\sqrt{\hat{p}(1-\hat{p}) / \left(\frac{1}{n_x} + \frac{1}{n_y}\right)}}$$

6. Small sample, difference between 2 means

$$H_0: \mu_x - \mu_y \leq \Delta_0$$

$$s_p = \sqrt{\frac{[(S_x^2/n_x) + (S_y^2/n_y)]^2}{\left[\frac{(S_x^2/n_x)^2}{(n_x-1)} + \left[\frac{(S_y^2/n_y)^2}{(n_y-1)}\right]\right]}},$$

$$t = \frac{(\bar{x} - \bar{y}) - \Delta_0}{\sqrt{s_p^2/n_x + s_p^2/n_y}}$$

o Equal Variance :

$$s_p = \sqrt{\frac{(n_x-1)s_x^2 + (n_y-1)s_y^2}{n_x+n_y-2}}$$

$$t = \frac{(\bar{x} - \bar{y}) - \Delta_0}{s_p \sqrt{\frac{1}{n_x} + \frac{1}{n_y}}}, \quad df = n_x+n_y-2$$