

Distributions

- Sampling dist. for mean of sample of size n from population with mean μ , standard deviation σ

$$\bar{X} \sim \text{Norm}(\mu, \sigma/\sqrt{n}).$$

- Binomial: For $X \sim \text{Binom}(n, p)$,

$$\mu_X = np, \quad \sigma_X = \sqrt{np(1-p)}.$$

When $np \geq 10$ and $n(1-p) \geq 10$, then X has approx. dist. $\text{Norm}(np, \sqrt{np(1-p)})$.

- Sampling dist. for sample proportion has

$$\mu_{\hat{p}} = p, \quad \sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}}.$$

When $N > 20n$, $np \geq 10$, and $n(1-p) \geq 10$, then \hat{p} has approx. dist. $\text{Norm}(p, \sqrt{p(1-p)/n})$.

Probability

- conditional probability $\Pr(B | A) = \frac{\Pr(A \text{ and } B)}{\Pr(A)}$
- Bayes' rule $\Pr(B | A) = \frac{\Pr(A | B)\Pr(B)}{\Pr(A)}$
- total probability $\Pr(A) = \Pr(A | B)\Pr(B) + \Pr(A | B^c)\Pr(B^c)$

Inference Procedures

- Level C Confidence Intervals (general):

$$(\text{estimate}) \pm (\text{critical value})(\text{approx. std. error})$$

- 1-sample proportion:

- CIs for p , $\text{SE} = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$

- z-test (when \hat{p} approx. normal)

$$\text{test stat. } (\mathbf{H}_0: p = p_0): z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$$

- 1-sample t : test statistic when $\mathbf{H}_0: \mu = \mu_0$

$$t = \frac{\bar{x} - \mu_0}{\text{SE}}, \quad \text{SE} = \frac{s}{\sqrt{n}}, \quad df = n - 1$$

- 2-sample t : test statistic when $\mathbf{H}_0: \mu_1 - \mu_2 = 0$

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\text{SE}}, \quad \text{SE} = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

As conservative estimate, take $df = \min(n_1, n_2) - 1$

- Chi-square test statistic:

$$\chi^2 = \sum \frac{[(\text{observed count}) - (\text{expected count})]^2}{\text{expected count}}$$

contingency table: $df = (\# \text{rows} - 1)(\# \text{columns} - 1)$

goodness-of-fit: $df = (\#\text{groups}) - 1 - (\#\text{est. params})$

- Model utility test:

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} = \frac{b_1}{\text{SE}_{b_1}}, \quad \text{with } df = n - 2$$

- F-test in ANOVA: $F = \frac{\text{MSG}}{\text{MSE}}$, where

$$df_{\text{numer}} = (\#\text{groups}) - 1, \quad \text{and}$$

$$df_{\text{denom}} = (\text{sample size}) - (\#\text{groups})$$

Miscellaneous

- Determine sample size, 1-proportion settings:

To have margin of error no larger than a desired size M , take,

$$n \geq \left(\frac{z^*}{M}\right)^2 \tilde{p}(1-\tilde{p}),$$

where \tilde{p} is an estimate of p (take $\tilde{p} = 0.5$ if no estimate is available)

Combinations of Random Variables

If X, Y are random variables, a, b are numbers, then

- $E(aX) = aE(X)$
- $E(X \pm Y) = E(X) \pm E(Y)$
- $\text{Var}(aX) = a^2 \cdot \text{Var}(X), \quad \text{or} \quad \text{SD}(aX) = |a| \cdot \text{SD}(X)$
- Moreover, if X, Y are independent,

$$\sigma_{X \pm Y}^2 = \sigma_X^2 + \sigma_Y^2, \quad \text{or} \quad \sigma_{X \pm Y} = \sqrt{\sigma_X^2 + \sigma_Y^2}.$$

Least Squares Regression

The coefficients (from data) are given by

$$b_1 = r \frac{s_y}{s_x}, \quad b_0 = \bar{y} - b_1 \bar{x}$$