

Random Variables

- Generally, when $f(x)$ is the pmf for (discrete) r.v. X :

$$E(X^j) = \sum_x x^j f(x)$$

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

Sample statistics

For iid sample X_1, \dots, X_n from population with mean μ , sd σ ,

$$Y = \sum_i X_i : E(Y) = n\mu, \quad SD(Y) = \sigma \sqrt{n} \quad (\text{multiply})$$

$$\bar{X} = Y/n : E(\bar{X}) = \mu, \quad SD(\bar{X}) = \sigma / \sqrt{n} \quad (\text{divide})$$

- When the population is normal, or n large, then

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

- When the X_i are Bernoulli (i.e., $\text{Binom}(1, p)$), then $Y = X_1 + X_2 + \dots + X_n \sim \text{Binom}(n, p)$. Moreover, when $np \geq 10$ and $n(1-p) \geq 10$, then Y is approx. normal:

$$Y \sim \text{Norm}(np, \sqrt{np(1-p)}) \quad \text{and}$$

$$\hat{p} = \frac{Y}{n} \sim \text{Norm}(p, \sqrt{p(1-p)/n})$$

Inference Procedures

- Level C Confidence Intervals (general):

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

- 1-sample proportion:

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

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

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

- 1-sample t : test statistic when $H_0: \mu = \mu_0$

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

- 2-sample t : test statistic when $H_0: \mu_1 - \mu_2 = 0$

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

Use t -distribution with conservative estimate for degrees of freedom $df = \min(n_1, n_2) - 1$

- 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:

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

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

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

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

Probability

- Conditional probability: $P(B | A) = \frac{P(A \text{ and } B)}{P(A)}$

- Bayes' rule: $P(B | A) = \frac{P(A | B) P(B)}{P(A)}$

- Total probability: $P(A) = P(A | B) P(B) + P(A | B^c) P(B^c)$

Miscellaneous

- Sample standard deviation $s = \sqrt{\frac{1}{n-1} \sum_i (x_i - \bar{x})^2}$

Combinations of Random Variables

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

$$E(aX) = a E(X)$$

$$E(X \pm Y) = E(X) \pm E(Y)$$

$$\text{Var}(aX) = a^2 \cdot \text{Var}(X), \quad \text{or} \quad SD(aX) = |a| \cdot SD(X)$$

- 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

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