

Descriptive Statistics Summary

One Variable

$$\text{mean } \bar{x} = \frac{\sum x_i}{n}$$

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

median
quartiles

$$\text{variance} = s^2$$

$$\text{outliers} = 1.5 \times \text{IQR}$$

Two variables

x = explanatory variable y = response variable

Regression line

$$y = a + bx$$

$$b = \frac{rs_y}{s_x}$$

$$\text{minimizes } D^2 = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - (a + bx_i))^2$$

$$a = \bar{y} - b\bar{x}$$

a is vertical intercept

b is slope

Correlation, r

Coefficient of determination, r^2

$$r \approx +1$$

$$r \approx -1$$

$$r^2 = \frac{\text{Variance in predicted } \hat{y} \text{ values}}{\text{Variance in observed } y \text{ values}} = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2}$$

$$r \approx 0$$

Residuals $y_i - \hat{y}_i$

all close to 0 if line fits well

Sample statistics vs. population parameters

\bar{x}, s

μ, σ

Random Variables

Discrete $\mu = \sum p_i x_i$ $\sigma = \sqrt{\sum p_i (x_i - \mu)^2}$ example: Binomial, $B(n,p)$

Continuous $\mu = \int_{-\infty}^{\infty} x f(x) dx$ $\sigma = \sqrt{\int_{-\infty}^{\infty} f(x) (x - \mu)^2 dx}$ example: Normal, $N(\mu, \sigma)$

Inference Summary

Sampling Distribution (approximately normal)

Means \bar{x}

$$\mu_{\bar{x}} = \mu, \quad \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

Proportions \hat{p}

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

Confidence Interval: Means

Confidence Interval: Proportions

one sample

$$\sigma \text{ known: } \left(\bar{x} - z \frac{\sigma}{\sqrt{n}}, \bar{x} + z \frac{\sigma}{\sqrt{n}} \right)$$

one sample

$$\hat{p} = \frac{X}{n}$$

one sample

$$\sigma \text{ not known: } \left(\bar{x} - t \frac{s}{\sqrt{n}}, \bar{x} + t \frac{s}{\sqrt{n}} \right)$$

$$\left(\hat{p} - z \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \hat{p} + z \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right)$$

two samples

σ not known:

$$\left(\bar{x}_1 - \bar{x}_2 - t \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}, \bar{x}_1 - \bar{x}_2 + t \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \right)$$

two samples

$$\hat{p}_1 = \frac{X_1}{n_1} \quad \hat{p}_2 = \frac{X_2}{n_2}$$

$$SE_{\hat{p}} = \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

(matched pairs, find CI for differences using one sample methods)

$$(\hat{p}_1 - \hat{p}_2 - z SE_{\hat{p}}, \hat{p}_1 - \hat{p}_2 + z SE_{\hat{p}})$$

Hypothesis Tests: Means	Hypothesis Tests: Proportions
<p><u>one sample</u> σ known: $H_0 : \text{mean} = \mu$</p> $z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \text{ standard normal}$ <p><u>one sample</u> σ not known: $H_0 : \text{mean} = \mu$</p> $t = \frac{\bar{x} - \mu}{s/\sqrt{n}} \text{ t distrib, } df = n - 1$ <p><u>two samples</u> σ not known: $H_0 : \mu_1 = \mu_2$</p> $t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \text{ t distrib, } df = \min(n_1 - 1, n_2 - 1)$ <p><u>many samples</u> largest $s < 2 \cdot$ smallest s</p> $H_0 : \mu_1 = \mu_2 = \dots = \mu_k$ <p>ANOVA $F(k - 1, N - k) \quad F = \frac{MSG}{MSE}$</p> <p>$k$ groups, N data altogether</p>	<p><u>one sample</u> $H_0 : \text{prop} = p_0$</p> $z = \frac{\hat{p} - p_0}{\sqrt{p_0(1 - p_0)/n}} \text{ standard normal}$ <p><u>two samples</u> $H_0 : p_1 = p_2$</p> $z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1 - \hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} \quad \hat{p} = \frac{X_1 + X_2}{n_1 + n_2}$ <p>std norm pooled sample</p> <p><u>many samples</u> $H_0 : \text{Independence}$</p> <p>Chi Square (all counts ≥ 5)</p> $\chi^2 = \sum \frac{(\text{observed} - \text{expected})^2}{\text{expected}}$