

6/23/11

Lecture #17

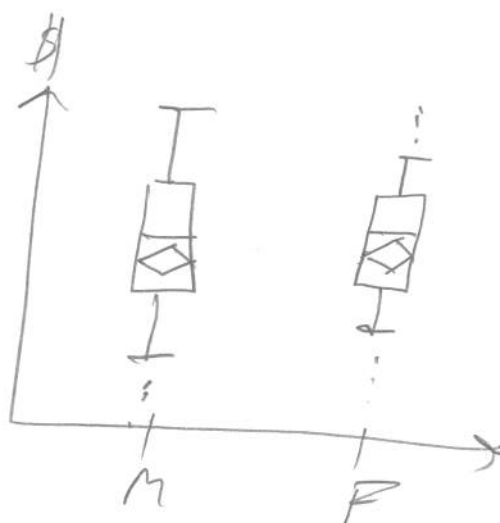
Plan

- 2-prop &
2-samp hyp.
testing

Yesterday we talked about bivariate analysis of categorical data. Here's some data:

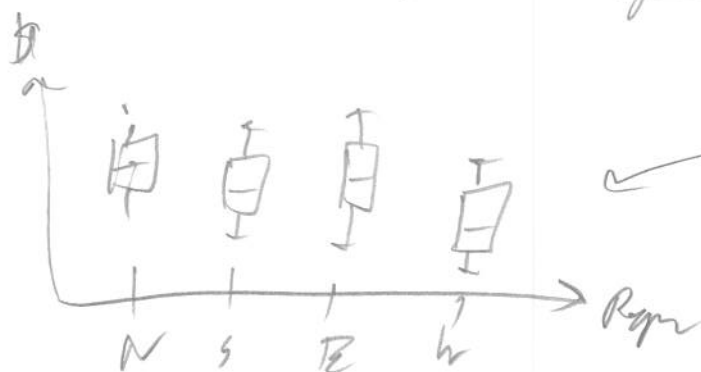


(proportions) ~~data~~



Hygiene data

In this class, the categorical variable will only have two states
 \Rightarrow two "groups" which are different



comparisons
 such as those we
 call ANOVA
 (start 102, 112)

The natural thing to do is to test!

- 1) If the proportions between the two groups are the same (i.e. $p_1 = p_2$), 2-prop z-test
 - 2) If the means between the two groups are the same (i.e. $\mu_1 = \mu_2$), 2-sample t-test
- Or... to build CIs for the difference in the proportions or means.

Same w/ 2-prop tests and CIs 2-param

$$\hat{p}_1 \sim N(p_1, \sqrt{\frac{p_1(1-p_1)}{n_1}}), \quad \hat{p}_2 \sim N(p_2, \sqrt{\frac{p_2(1-p_2)}{n_2}})$$

We can observe the difference:

$$D = \hat{p}_1 - \hat{p}_2$$

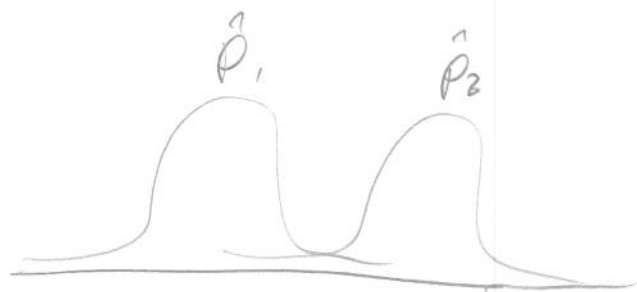
error of diff

of course for
there to be
true result
 $n\hat{p}_1 > 10, n(1-\hat{p}_1) > 10$
 $n\hat{p}_2 > 10, n(1-\hat{p}_2) > 10$

The difference we get in an experiment is $d = \hat{p}_1 - \hat{p}_2$ a draw from

$$D = \hat{p}_1 - \hat{p}_2$$

What is O's estimator? We're substituting the numbers:



$$SE[\hat{p}_1 - \hat{p}_2] = \sqrt{Var[\hat{p}_1 - \hat{p}_2]}$$

Assume both groups are ind.

$$\hat{p}_1 - \hat{p}_2 \sim N(p_1 - p_2, \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}})$$

$$= \sqrt{Var \hat{p}_1 + Var \hat{p}_2}$$

$$= \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$$

Now we do an experiment. we get d , we realize:

$$Z = \frac{D - E[D]}{SE[D]}, \quad z = \frac{d - (p_1 - p_2)}{\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}} \quad (\text{exact})$$

Let's see an example: p#37

35 of 60 customers in the East prefer red fabric ~~32~~ 32 of 72 customers in the West prefer red fabric. Is there any real difference in red fabric preferences or difference due to chance? Test at $\alpha = 5\%$.

Let p_1 be the % of customers in the E; let p_2 be the % of customers in the West

$H_0: p_1 - p_2 = 0$ i.e. no difference

$H_a: p_1 - p_2 \neq 0$ i.e. there is a difference (two-tailed)

$$\hat{p}_1 = \frac{33}{60} = .5833, \quad \hat{p}_2 = \frac{32}{72} = .4444$$

LP

Let's test it... ~~and~~ compute Z statistic first:

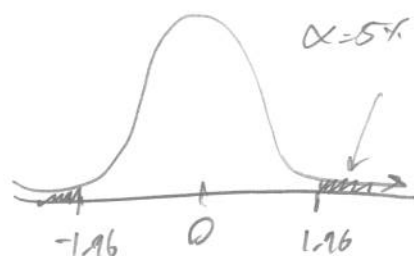
$$Z = \frac{d - (0)}{\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}} \approx \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}} = \frac{.5833 - .4444}{\sqrt{\frac{.5833(1-.5833)}{60} + \frac{.4444(1-.4444)}{72}}}$$

don't know

p_1, p_2
just diff!

$$= \frac{.1389}{0.08645} = 1.607 < Z_{0.025} = 1.96$$

\Rightarrow fail to reject



$$p_{val} = P(|Z| \geq 1.607) = 2 \cdot .054 = .108$$

$$> \alpha = 0.05$$

\Rightarrow fail to reject

Compute a ^{95%} CI for the difference:

$$\begin{aligned} CI_{\hat{p}_1 - \hat{p}_2, 95\%} &= [d \pm Z_{0.025} \cdot SE(D)] \\ &\approx [\hat{p}_1 - \hat{p}_2 \pm Z_{0.025} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}] \\ &= [.1389 \pm 1.96 \cdot 0.08645] \\ &= [-0.031, 0.308] \end{aligned}$$

Since

$0 \in CI_{\hat{p}_1 - \hat{p}_2, 95\%} \Rightarrow$ fail to reject ^{hypothesis} $p_1 - p_2 = 0$ at $\alpha = 5\%$

You will need to study left-tailed and right-tailed on your own

Now for the more interesting case (para to road graph)
 we are looking for difference in means. θ -params!

$$\bar{X}_1 \sim N(\mu_1, (\frac{\sigma_1}{\sqrt{n_1}})^2), \quad \bar{X}_2 \sim N(\mu_2, (\frac{\sigma_2}{\sqrt{n_2}})^2)$$

$$SE[\bar{X}_1 - \bar{X}_2] = \sqrt{Var[\bar{X}_1 - \bar{X}_2]} = \sqrt{Var[\bar{X}_1] + Var[\bar{X}_2]} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

$$D = \bar{X}_1 - \bar{X}_2 \sim N(\mu_1 - \mu_2, (\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}})^2) = N(\mu_1 - \mu_2, (\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}})^2)$$

$d = \bar{X}_1 - \bar{X}_2$ is a draw from $D = \bar{X}_1 - \bar{X}_2$

if $\sigma_1 = \sigma_2$

Further complication: What if we don't know σ_1, σ_2 ?

We only have $s_1, s_2 \Rightarrow T$ distribution \Rightarrow (Need both pop's var)
 but what is ~~the~~ SE and what is the df?

$$T = \frac{d - E[D]}{SE[D]}$$

We consider two situations:

① σ_1, σ_2 different: if $\frac{s_1^2}{s_2^2} > 3$ or $\frac{s_2^2}{s_1^2} > 3$

$$SE[\bar{X}_1 - \bar{X}_2] \approx \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$$df = \min\{n_1, n_2\}$$

do not use formula ~~on p. 445~~
 it's too complicated

② $\sigma = \sigma_1 = \sigma_2$ same:

$$\frac{s_1^2}{s_2^2} = 3$$

$$SE[\bar{X}_1 - \bar{X}_2] \approx s_{pooled} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \quad \text{where } s_{pooled} = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}}, \quad df = n_1 + n_2 - 2$$

Example p929. Website testing between two designs, new, d.d.
 300 visitors: 169 old site $\bar{x}_{old} = \$253$, sold = \$130, 131 new site $\bar{x}_{new} = \$328$,
 Let μ_1 be the mean for the new website
 Let μ_2 be the mean for the old website.
 Histogram passes Q-Q plot for normality

Is the new site better for revenue generation? $\alpha = 5\%$

$H_0: \mu_1 - \mu_2 = 0$, $H_a: \mu_1 - \mu_2 > 0$ (right tailed)

Which 2-samp t-test do we do?

$$\frac{s_1^2}{s_2^2} = \frac{161^2}{130^2} = 1.53 < 3 \Rightarrow \text{pooled var...}$$

$$df = n_1 + n_2 - 2 = 290 \Rightarrow \text{use } z$$

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} = \frac{328 - 253}{\sqrt{\frac{(131-1)161^2 + (169-1)130^2}{131+169-2} \left(\frac{1}{131} + \frac{1}{169} \right)}}$$

$$= \frac{75}{199.39 \cdot 0.116} = \frac{75}{16.803} = 4.46 > t_{0.05, 290} \approx z_{0.05} = 1.96$$

\Rightarrow reject H_0

$$p\text{-val} = P(T_{290} \geq 4.46) \approx P(Z \geq 4.46) \approx 0 < \alpha = 5\%$$

\Rightarrow reject H_0

The new website is probably better than the old website

Airbag restraints: Experiment the new airbags crash cars
detection + inflation 70ms, New design to do it faster.

Crash test = 4 cars old system $\bar{X} = 72ms$, $s = 5ms$

Crash test = 5 cars new system $\bar{X} = 67ms$, $s = 2.8ms$, $\alpha = 1\%$

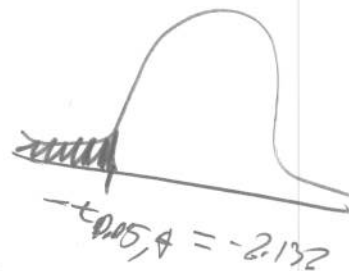
Let μ_1 be ^{the} mean of new system

Let μ_2 be the mean of old system

$$H_0: \mu_1 - \mu_2 \geq 0$$

$$H_a: \mu_1 - \mu_2 < 0 \quad (\text{left-tailed})$$

$$\frac{s_2^2}{s_1^2} = \frac{5ms^2}{2.8ms^2} = 3.2 > 3 \Rightarrow \text{approx}$$



$$t = \frac{d - (0)}{SE(d)} = \frac{\bar{X} - \bar{X}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} = \frac{67 - 72}{\sqrt{\frac{2.8^2}{5} + \frac{5^2}{4}}} = \frac{-5}{2.79} = -1.78$$

$$df = \min\{n_1, n_2\} = \min\{4, 5\} = 4$$

$$> -2.132$$

\Rightarrow fail to reject

$$p\text{-val} = P(T_4 \leq -1.78) = 0.0748 > \alpha = 5\% \Rightarrow \text{fail to reject}$$

Ignore 'partial comparisons' pP36-441

One with categorical vs. interval bivariate analysis!
red ch 10 for ka! also interval vs. since

$$X, Y \text{ independent} \quad P(X=x, Y=y) = P(X=x) P(Y=y)$$

$$\text{or conv} \quad f_{XY}(x, y) = f_X(x) f_Y(y)$$

$$E(XY) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy f_{XY}(x, y) dx dy = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy f_X(x) f_Y(y) dx dy$$

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty}$$