# Multivariate Statistics

## Two sample T-test for Difference in means

$$t_{data} = \frac{\overline{x_1} - \overline{x_2}}{\sqrt{(s_1^2/n_1) + (s_2^2/n_2)}}$$

$df = n_1 - 1 + n_2 - 1 = n_1 + n_2 - 2$

will be minimum

**Q** Page - 189   ie $n_1 - 1 = \frac{25 \cdot 29 - 1}{3006 \cdot 64}$

$n_2 - 1 = 894 - 1$
$= 802$

**S-1** State hypotheses

$H_0 = \mu_1 = \mu_2$ (no-difference)

$H_1 : \mu_1 \neq \mu_2$ (means are different)

This is two-tail test.

**S-y** conclusion

$P = 0.508 > 0.05$

do not reject the null hypothesis.

**S-2** compute test statistics

$$t = \frac{\overline{x_1} - \overline{x_2}}{\sqrt{(s_1^2/n_1) + (s_2^2/n_2)}} = \frac{1.5714 - 1.5361}{\sqrt{\frac{(1.3126)^2}{2529} + \frac{(1.3251)^2}{894}}}$$

$= 0.65995$

**S-3** p-value (use 2 - table because sample is 8 b

so t-test en

$1 - 0.7421$
$= 0.2579$

P-Value $= 2 \cdot P(t > t_{data}) = P(t > 0.65995)$
$= 2 \times 0.257 = 0.514$
$= 2 \times (0.254) = 0.508$  0.508

from z-table now 0.6 & colum 0.06 = 0.6 + 0.06 = 0.66
valu
is 0.7454, own right tail = 1 - 0.7454 = 0.2546

# Two Sample Z-Test for Difference in Proportions

$$Z_{data} = \frac{P_1 - P_2}{\sqrt{P_{pooled} \cdot (1 - P_{pooled})\left(\left(\frac{1}{n_1}\right) + \left(\frac{1}{n_2}\right)\right)}}$$

where $P_{pooled} = \frac{x_1 + x_2}{n_1 + n_2}$

**Ans**

$x_1 = 707$ of $n_1 = 2529$ customers on the training set belongs to no the voice mail plan while $n_2 = 215$ of $n_2 = 804$ customers in the tests set.

So $P_1 = \frac{x_1}{n_1} = \frac{707}{2529} = 0.2796$

$P_2 = \frac{x_2}{n_2} = \frac{215}{804} = 0.2674$

$P_{pooled} = \frac{x_1 + x_2}{n_1 + n_2} = \frac{707 + 215}{2529 + 804} = 0.2766$

**Hypothesis**

$H_0 : \pi_1 = \pi_2 \quad vs \quad H_1 : \pi_1 \neq \pi_2$

$Z_{data} = \frac{P_1 - P_2}{\sqrt{P_{pooled} \cdot (1 - P_{pooled})\left(\left(\frac{1}{n_1}\right) + \left(\frac{1}{n_2}\right)\right)}}$

$= \frac{0.2796 - 0.2674}{\sqrt{0.2766 \cdot (0.7234)\left((1/2529) + (1/804)\right)}}$

$= 0.6736$

-3

p-value is

$$p\text{-value} = 2 \cdot P(z > 0.6736)$$

$$= 0.5006$$

S-4   conclution

$$= 0.2517 \times 2 = 0.502$$

$$1 - 0.7485 = 0.2515$$

$$p\text{-value} = 0.5006 > 0.05$$

so do not reject $H_0$.

## TEST FOR HOMOGENEITY OF PROPORTIONS (chi-square distribution)

observed frequency (Marital Status)

| Dataset | Married | Single | Other | Total |
|---|---|---|---|---|
| Training set | 410 | 340 | 250 | 1000 |
| Test set | 95 | 85 | 70 | 250 |
| Total | 505 | 425 | 320 | 1250 |

### S-1   Hypothesis set

$H_0$ :  $P_{married, training} = P_{married, test}$

$P_{single, training} = P_{single, test}$

$P_{other, training} = P_{other, test}$

$H_a$ : At least one of the claims in $H_0$ is wrong.

S-2    expected frequency

married, trainy = $\dfrac{(1000)}{1250}$ = 4

$$\text{expected frequency} = \dfrac{(\text{row total})\,(\text{column total})}{\text{grand total}}$$

## Expected frequencies

| Dataset | Married | Single | Other | Total |
|---|---|---|---|---|
| Training set | 404 | 340 | 256 | 1000 |
| Test set | 101 | 85 | 64 | 250 |
| Total | 505 | 425 | 320 | 1250 |

**S-3** calculate test statistic $x^2_{data}$

| cell | Observed frq. | Expected frq. | $\dfrac{(Obs - Exp)^2}{Exp}$ |
|---|---|---|---|
| Married, trainy | 410 | 404 | $\dfrac{(410-404)^2}{404} = 0.09$ |
| Married, test | 95 | 101 | 0.36 |
| Single, training | 340 | 340 | 0 |
| Single, test | 85 | 85 | 0 |
| Other, trainy | 250 | 256 | 0.14 |
| Other, test | 70 | 64 | 0.56 |

$$\text{sum} = x^2_{data} = 1.15$$

p-value $= P(x^2 > \chi^2_{deta}) = P(\chi^2 > 1.15) = 0.5627$

① test of independent : $df = (r-1)(c-1) = 2$

ie $df = (no. of rows -1)(no. of clcum)$

$= (2-1)(3-1) = 1 \times 2 = 2$

② goodnen of fit

$p-value \approx 0.5627$

$df = k-1$ (k=no. of categories)

$= 6-1$

P-5 **Conclusion**

p-value $< \alpha$ then reject Ho

but Her $0.5627 \not< 0.05$ then observed frequencies represent proportion that are significantly different for the trainos & test deta.

Backspace

Pause
Break

Insert

Lenovo

Home

End

NmLk

/

*

PgUp

7
Home

8
↑

9
PgUp

+

Delete

PgDn

4
←

5

6
→

Enter

shift

1
End

2
↓

3
PgDn

Enter

↑

Ctrl

←

↓

→

0

Ins

.
Del