# MATH 4720 / MSSC 5720

## Instructor: Mehdi Maadooliat

**Review for Final Exam**

**Department of Mathematical and Statistical Sciences**

(1) A study of potential age discrimination considers promotions among middle managers in a large company. The data are as follows:

| | Age | | | |
|---|---|---|---|---|
| | Under 30 | 30-39 | 40-49 | 50 and over |
| Promoted | 9 | 29 | 32 | 10 |
| Not promoted | 41 | 41 | 48 | 40 |

Is there a statistically significant relationship between age and promotion at $\alpha = 0.05$ ?

① $H_o$: Age and Promotion are indep

$H_a$: " " " " are NOT indep.

| | Age | | | | Total |
|---|---|---|---|---|---|
| | Under 30 | 30-39 | 40-49 | 50 and over | |
| Promoted | 9 | 29 | 32 | 10 | 80 |
| NOT ~ | 41 | 41 | 48 | 40 | 170 |
| Total | 50 | 70 | 80 | 50 | 250 |

Observed ($O_{ij}$'s)

Assumption: All $E_{ij} > 5$ ✓

Expected Counts ($E_{ij}$'s)

| | Under 30 | 30-39 | 40-49 | 50 and over | |
|---|---|---|---|---|---|
| Promoted | $\frac{50 \times 80}{250} = 16$ | $\frac{70 \times 80}{250} = 22.4$ | $\frac{80 \times 80}{250} = 25.6$ | $\frac{50 \times 80}{250} = 16$ | 80 |
| Not ~ | $\frac{50 \times 170}{250} = 34$ | $\frac{70 \times 170}{250} = 47.6$ | $\frac{80 \times 170}{250} = 54.4$ | $\frac{50 \times 170}{250} = 34$ | 170 |
| Total | 50 | 70 | 80 | 50 | 250 |

2

$$\underline{T.S} \quad \chi^2 = \sum_{i,j} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$$\chi^2 = \frac{(9-16)^2}{16} + \frac{(29-22.4)^2}{22.4} + \frac{(32-25.6)^2}{25.6} + \frac{(10-16)^2}{16}$$

$$+ \frac{(41-34)^2}{34} + \frac{(41-47.6)^2}{47.6} + \frac{(48-54.4)^2}{54.4} + \frac{(40-34)^2}{34}$$

$$\Longrightarrow \chi^2 = 13.025$$

$$\Longrightarrow \chi^2 = 13.025$$

$$\text{Rej } H_0 \text{ if } \quad \chi^2 > \underbrace{\chi^2_{0.05}}_{7.815} \left( df = (4-1)(2-1) \right)$$

$$\Longrightarrow 13.025 > 7.815 \Longrightarrow \text{ there is a significant relationship between age and promotion.}$$

(2) A linear regression model is to be tested from a data on 12 male patients with congestive heart failures. The dependent variable is the cardiac index Y and the independent variable is the weight X. The Minitab output of the regression and the residual analysis is given below.

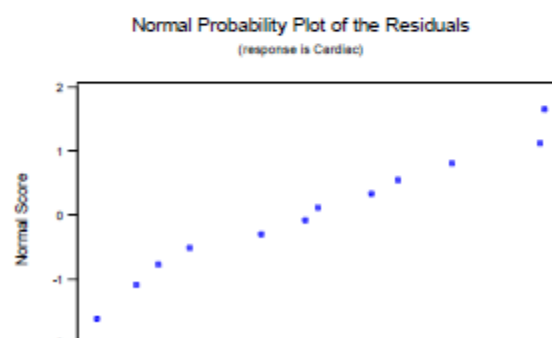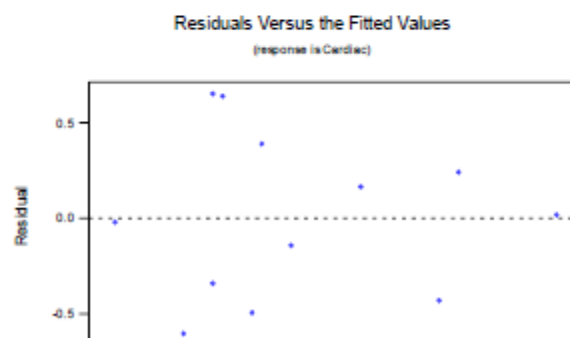## Regression Analysis

```
The regression equation is
Cardiac Index = 0.909 + 0.0126 Weight(kg)

Predictor          Coef          StDev              T            P
Constant         0.9090        0.7837           1.16        0.273
Weight(kg)       0.01256       0.01023          1.23        0.247

S = 0.4499        R-Sq = 13.1%        R-Sq(adj) = 4.4%
```

Residuals Versus the Fitted Values
(response is Cardiac)

Normal Probability Plot of the Residuals
(response is Cardiac)



(b) Intercept : the expected response (y) for $x = 0$ ⟹ $0.909$

A 95% CI of $\beta_0$ is $\hat{\beta}_0 \pm t_{a/2}se(\hat{\beta}_0)$, i.e., $0.909 \pm 2.229 \times 0.7837$, i.e. $(-.0838, 2.656)$.

Slop : the amount of change in cardiac index for a unit change in weight.

: $0.01256$

A 95% CI of $\beta_1$ is $\hat{\beta}_1 \pm t_{a/2}se(\hat{\beta}_1)$, i.e., $0.01256 \pm 2.229 \times 0.01023$, i.e. $(-0.01024, 0.03536)$.

(a) What is the least square line of best fit describing the relationship between Y and the X.?
(b) Interpret the intercept and the slope of the line, and find the 95% confidence intervals for
the true intercept and the true slop parameters.
(c) Test the hypothesis that the cardiac index depends on the weight. Use $\alpha = 0.10$. You must
also write the value of the test statistics and the p-value.
(d) What are the assumptions for the above regression analysis? From the Minitab output,
explain if all these assumptions are satisfied. If not, then what would you suggest to
modify the regression model.
(e) How much variability in the cardiac index is due the patients' weights?
(f) Based on the information presented here, predict the cardiac index for a patient of weight
of 80 kg? Is this a reliable prediction? Explain.
(g) Assume that the 95% confidence interval for the mean cardiac index of patients of weight
80 kg is (-0.1, 3.8), and the prediction interval is (-1.1, 5.5). Interpret these intervals.

(c)

$H_o: \beta_1 = 0$

$H_a: \beta_1 \neq 0$

TS. $t = \dfrac{\hat{\beta_1}}{se(\hat{\beta_1})} = \dfrac{0.01256}{0.01023} \simeq 1.23$

p-value $= 2 P(t(df=n-2) > 1.23) = 2(0.1234)$

$= 0.247 \implies$ Fail to rej $H_o$

Not enough evidence to support the
liner relationship between cardiac index & weight

5

(d) the error ($\varepsilon_i$'s) are iid normal (QQ-plot of residuals confirm that)

$Var(\varepsilon_i) = constant$ (residual plot confirms that.)

(f) $y = 0.909 + 0.0126 (80)$

it's not a reliable prediction since we reject $H_o: \beta_1 = 0$ (fail to)

(g) — We are 95% confident that on average cardiac index of patients with 80 kg weight is $(-0.1, 3.8)$

— we predict with 95% chance that cardiac index a patient with 80 kg weight is $(-1.1, 5.5)$

(3) Researchers in the development of new treatments for cancer patients often evaluate the effectiveness of new treatments by reporting the proportion of patients who survive for a specific period of time after completion of the treatment. In a recent study on 600 patients for comparing an old treatment with a new treatment, the following data reports the number of patients surviving at least 5 years after the treatment.

| | # of Patients | # survived |
|---|---|---|
| Old Treatment | 450 | 200 |
| New Treatment | 150 | 74 |

(a) Find a 95% confidence interval of proportion of patients who will survive under the old treatment. Also find a 95% confidence interval of proportion of patients who will survive under the new treatment. You must check all the necessary assumptions.

(b) Is there a sufficient evidence that the new treatment is better than the old treatment at $\alpha = 0.05$ ? You must check all the necessary assumptions and state the null and the alternative hypotheses.

@ 95% CI of survival under old treatment.

$$\hat{\pi}_0 \pm z_{\alpha/2} \sqrt{\frac{\hat{\pi}_0(1-\hat{\pi}_0)}{n_0}} \quad \text{where}$$

$$\begin{cases} \hat{\pi}_0 = \frac{200}{450} = 0.4\overline{4}4 \\ n_0 = 450 \\ z_{\alpha/2} = 1.96 \end{cases}$$

$(0.399, 0.490)$

Assumption

$$n_0\hat{\pi}_0 = 450(0.4\overline{4}) > 5$$
$$n_0(1-\hat{\pi}_0) = 450(1-0.4\overline{4}) > 5 \quad \checkmark$$

$$0.4\overline{4} \pm 1.96 \sqrt{\frac{0.4\overline{4}(1-0.4\overline{4})}{450}} = 0.4\overline{4}4 \pm 0.046$$

95% CI for survival under new treatment

$$\hat{\pi}_n \pm z_{\alpha/2}\sqrt{\frac{\hat{\pi}_n(1-\hat{\pi}_n)}{n_n}} \quad \text{wehere}$$

$$\begin{cases} \hat{\pi}_n = \dfrac{74}{150} = 0.49\overline{3} \\[2mm] n_n = 150 \\[2mm] z_{\alpha/2} = 1.96 \end{cases}$$

Assumption:

$$n\hat{\pi}_n = 150\,(0.49\overline{3}) > 5$$
$$n(1-\hat{\pi}_n) = 150\,(1-0.49\overline{3}) > 5$$ ✓

95% CI for $\pi_n$ is $\quad 0.49\overline{3} \pm \overset{1.96}{\sqrt{\dfrac{0.49\overline{3}(1-0.49\overline{3})}{150}}} = 0.49\overline{3} \pm 0.080$

$$\boxed{(0.413,\ 0.573\ )}$$

ⓑ $H_o$: $\pi_o = \pi_n$

$H_a$: $\pi_o < \pi_n$

Assumptions: $n_o \hat{\pi}_o > 5$, $n_o(1-\hat{\pi}_o) > 5$ ⎫ ✓

$n_n \hat{\pi}_n > 5$, $n_n(1-\hat{\pi}_n) > 5$ ⎬

T.S  $Z = \dfrac{\hat{\pi}_o - \hat{\pi}_n}{\sqrt{\dfrac{\hat{\pi}_o(1-\hat{\pi}_o)}{n_o} + \dfrac{\hat{\pi}_n(1-\hat{\pi}_n)}{n_n}}} = \dfrac{0.44\overline{4} - 0.49\overline{3}}{\sqrt{\dfrac{0.44\overline{4}(1-0.44\overline{4})}{450} + \dfrac{0.49\overline{3}(1-0.49\overline{3})}{150}}} = -1.039$

Reject $H_o$ if $\underset{-1.039}{Z} < \underset{-1.645}{-Z_\alpha}$ $\implies$ Fail to reject $\left(H_o: \pi_o = \pi_n\right)$

NOT suff evidence that the new treatment
is better

(4) A local doctor suspects that there is a seasonal trend in the occurrence of the common cold. The doctor collected the following information from a sample of 800 cases of patients with the common cold over the past few years. The data is given below

| Season | Frequency |
|--------|-----------|
| Winter | 258 |
| Spring | 215 |
| Summer | 170 |
| Fall | 157 |

(a) State the null hypothesis to test the seasonal trend on the occurrence of the common cold.

(b) Is there sufficient evidence that there is a seasonal trend in the occurrence of the common cold at $\alpha = 0.01$ ?

(4) (a) $H_o: \pi_1 = \pi_2 = \pi_3 = \pi_4 = \frac{1}{4}$

$$\chi^2 = \sum_{i=1}^{4} \frac{(O_i - E_i)^2}{E_i}$$

$H_a: \pi_i \neq \frac{1}{4}$ for some $i = 1, 2, 3, 4$

(b) $N = 258 + 215 + 170 + 157 = 800$

$\hat{\pi}_1 = \frac{258}{800} = 0.323$  $\hat{\pi}_2 = \frac{215}{800} = 0.269$  $\hat{\pi}_3 = \frac{170}{800} = 0.212$  $\hat{\pi}_4 = \frac{157}{800} = 0.196$

$$\Rightarrow \chi^2 = \frac{(258-200)^2}{200} + \frac{(215-200)^2}{200} + \frac{(170-200)^2}{200} + \frac{(157-200)^2}{200}$$

$$= 31.69$$

Reject $H_0$ if $\underset{31.69}{\underbrace{\chi^2}} > \underset{11.34}{\underline{\chi^2_{0.01} (df = 4-1)}}$

Suff Evid that, there is a seasonal trend in the occurrence of the common cold.

(5) Four methods of training individuals to use a special piece of equipment have been developed. To compare the methods, 24 employees with similar experience were selected and six assigned to each training method. After the training, the employees were asked to perform a number of tasks and a score of efficiency assigned to each employee. The ANOVA and Means for the results are shown below.

| Source | df | MS | | Method | A | B | C | D |
|--------|----|----|----|--------|---|---|---|---|
| Method | 3 | 142.041 | | Mean | 30.333 | 24.167 | 33.667 | 23.667 |
| Error | 20 | 24.615 | | St.Dev. | 5.0 | 6.1 | 4.0 | 4.5 |

(a) Is there sufficient evidence that there is a difference in the four training methods at $\alpha = 0.05$.? State also the null and the alternative hypotheses.
(b) Using a multiple comparisons procedures (LSD or Tukey's), determine the pairs of methods that are significantly different at $\alpha = 0.05$. **Show your work.**
(c) What assumptions did you make for the above analyses?

(a) $H_o : \mu_1 = \mu_2 = \mu_3 = \mu_4$

$H_a : \mu_i \neq \mu_j$ for some $(i,j)$

Reject $H_o \Rightarrow$ suff. evid. that there is a diff in the four training methods.

T.S $\quad F = \dfrac{MS_{Trt}}{MS_{Err}} = \dfrac{SS_{Trt}/df_{Trt}}{SS_{Err}/df_{Err}} = \dfrac{142.041}{24.615} = 5.77$

Rej $H_o$ if $\underset{\underset{5.77}{\smile}}{F} > F_{0.05}(df_1 = 4-1, df_2 = 24-4)$

$F_{0.05}(3, 20) = 3.09$

| Method | A | B | C | D |
|---|---|---|---|---|
| Mean | 30.333 | 24.167 | 33.667 | 23.667 |
| St.Dev. | 5.0 | 6.1 | 4.0 | 4.5 |

(6) Fisher's LSD: $|\bar{y}_i - \bar{y}_j| > t_{\alpha/2} \sqrt{MSE(\frac{1}{n_i} + \frac{1}{n_j})}$

$t_{\alpha/2}(df=20) = 2.086$

$MSE = 24.615$

$\frac{1}{n_i} + \frac{1}{n_j} = \frac{1}{6} + \frac{1}{6} = \frac{1}{3}$

$t_{\alpha/2}\sqrt{MSE(\frac{1}{n_i} + \frac{1}{n_j})} = 2.086\sqrt{24.615(\frac{1}{3})} = 5.975$

Amethod

A vs B : $|30.333 - 24.167| = 6.166 \quad > 5.975$

A vs C : $|30.333 - 33.667| = 3.334 \quad < 5.975$

A vs D : $|30.333 - 23.667| = 6.666 \quad > 5.975$

Methods B and D are different from method A

B vs C : $|24.167 - 33.667| = 9.5 \quad > 5.975$

B vs D : $|24.167 - 23.667| = 0.5 \quad < 5.975$

B and C are different

C vs D : $|33.667 - 23.667| = 10 \quad > 5.975 \Rightarrow$ C and D are different

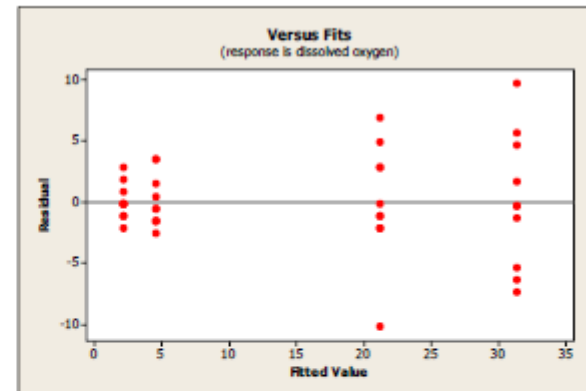(C) — equallity of the variances , (Homogeneity of variances):

$$\sigma_1 = \sigma_2 = \sigma_3 = \sigma_4$$

— Data for each method is following a normal distribution.

(6) An ANOVA is to be performed to compare four groups. The means and the standard deviations are as follows:

|          | Mean | St. Dev |
|----------|------|---------|
| Group 1  | 2.2  | 1.476   |
| Group 2  | 4.6  | 2.119   |
| Group 3  | 21.2 | 4.733   |
| Group 4  | 31.4 | 5.522   |

1. What assumption of ANOVA is not satisfied?

2. What you need to do so that a valid ANOVA can be performed to compare the four groups (You may refer to the graph on the right)?



Versus Fits
(response is dissolved oxygen)

(a) Homogeneity of variances

(b) transform the data.