Restorative
Dentistry
& Endodontics

**RDE**

Check for updates

# Statistical notes for clinical researchers: analysis of covariance (ANCOVA)

Hae-Young Kim ⓘ *

Department of Health Policy and Management, College of Health Science, and Department of Public Health Science, Graduate School, Korea University, Seoul, Korea

**\*Correspondence to**
**Hae-Young Kim, DDS, PhD**
Professor, Department of Health Policy and Management, Korea University College of Health Science, and Department of Public Health Science, Korea University Graduate School, 145 Anam-ro, Seongbuk-gu, Seoul 02841, Korea.
E-mail: kimhaey@korea.ac.kr

**ORCID iDs**
Hae-Young Kim ⓘ
https://orcid.org/0000-0003-2043-2575

Previously, we discussed analysis of variance (ANOVA) and simple linear regression, which commonly share continuous dependent variables. While ANOVA uses categorical variables as independent variables, regression uses mainly continuous variables for them. However, we may want to include both kinds of variables in analysis. A statistical model with continuous dependent variables and both types of independent variables is called a general linear model (GLM). In this section, we discuss analysis of covariance (ANCOVA) as a type of GLM models. An ANCOVA is similar to an ANOVA model, but it includes a continuous variable as well as categorical variables as independent variables, being a mixture model of ANOVA and regression models.

## RAISING A QUESTION ON IGNORING COVARIATES

An example data is composed of 3 variables, treatment effect, treatment methods (Tx; 2 groups), and age in **Table 1**. Our interest is on comparison of treatment effects by 2 Tx, experimental and control groups. We may consider independent $t$-test, ignoring age variable. Distribution of treatment effects of 2 groups is depicted in **Figure 1A**. We could obtain the $p$ value, 0.048, by applying independent student's $t$-test comparing treatment groups and conclude the treatment effect of treatment group is superior to that of control group.

Now let's look into the relationship between treatment effect and age. We can notice a trend that higher age is related to higher treatment effect in **Figure 1B**. The positive correlation between effect and age is quantitatively measured by a Pearson correlation coefficient, 0.805 ($p < 0.001$). The positive correlation is further analyzed by regression analysis. We get regression equations for pooled sample of both groups as well as for each group as following:

For pooled sample: Treatment effect = 24.2 + 0.74 × Age + Error   (1)

For experimental group: Treatment effect = 47.5 + 0.33 × Age + Error   (2)

For control group: Treatment effect = 6.21 + 1.03 × Age + Error   (3)

Meanwhile, the mean age of subjects in the experimental group is 44.83 years, which is higher than that of the control group, 43.58 years. We may be suspicious that the difference of effect between two groups is partly attributed to the age difference between two groups, because age

Generated by KAMJE PRESS

**Table 1.** Data of treatment groups, treatment effect (Effect), and age (Age)

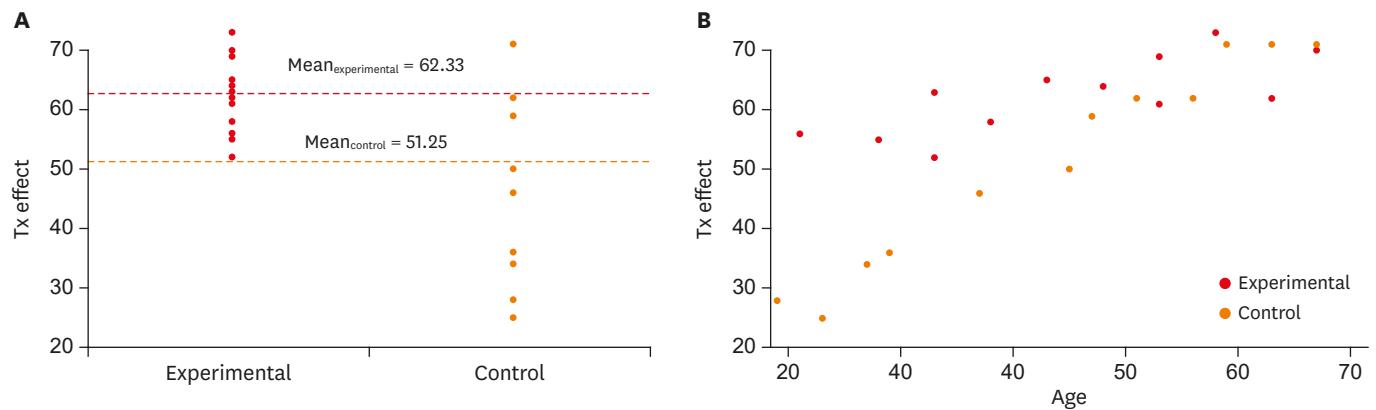| Group | | | | | | | | | | | | | Mean | SD |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Experimental group | | | | | | | | | | | | | | |
| Effect | 56 | 55 | 63 | 52 | 58 | 65 | 64 | 61 | 69 | 73 | 62 | 70 | 62.33 | 6.39 |
| Age | 21 | 28 | 33 | 33 | 38 | 43 | 48 | 53 | 53 | 58 | 63 | 67 | 44.83 | 14.52 |
| Control group | | | | | | | | | | | | | | |
| Effect | 28 | 25 | 71 | 62 | 50 | 46 | 34 | 59 | 36 | 71 | 62 | 71 | 51.25 | 17.19 |
| Age | 19 | 23 | 67 | 56 | 45 | 37 | 27 | 47 | 29 | 59 | 51 | 63 | 43.58 | 16.36 |



**Figure 1.** Scatter plot of the example data in **Table 1**. (A) Distribution of treatment effects of 2 groups; (B) Scatter plot of treatment effects and age by treatment groups.

is positively correlated with treatment effect. How can we resolve this issue? There is a clear need to consider the covariate, age, into the model to control its possible influence.

## ANCOVA MODEL: COMPARING MEANS CONSIDERING COVARIATES

To compare 2 means, we can apply ANOVA as well, which is applicable in comparing 2 or more group means. The result shows significant difference between two groups ($p = 0.048$), which is exactly the same with that from the independent $t$-test in **Figure 2C**. Still, the possible covariate, age, is ignored. The model including 2 groups explains the variation of effect as much as corrected model sum of squares of 737.042 among total sum of squares of 4,435.958 in **Figure 2C**. **Figure 2B** displays the proportion of errors as 0.83, which is proportion of Error sum of squares of 3,698.917 among total sum of squares of 4,435.958. The proportion of errors represents the portion of variation that the model cannot explain. Also, we find the proportion of explained variance, R-squared, is 0.166, which represents that only 16.6% of variance in the response variable is explained by this model.

The ANOVA model can be performed using GLM procedure. The result is expressed as a GLM equation in **Figure 2A**, as:

$$\text{Treatment effect} = 51.25 + 11.08 \times \text{Tx} + \text{error} \quad (4)$$

where Tx = 0 for control group and Tx = 1 for experimental group.

We can obtain the mean effects of experimental and control groups as 62.33 (= 51.25 + 11.08) and 51.25, exactly the same as which appears above.
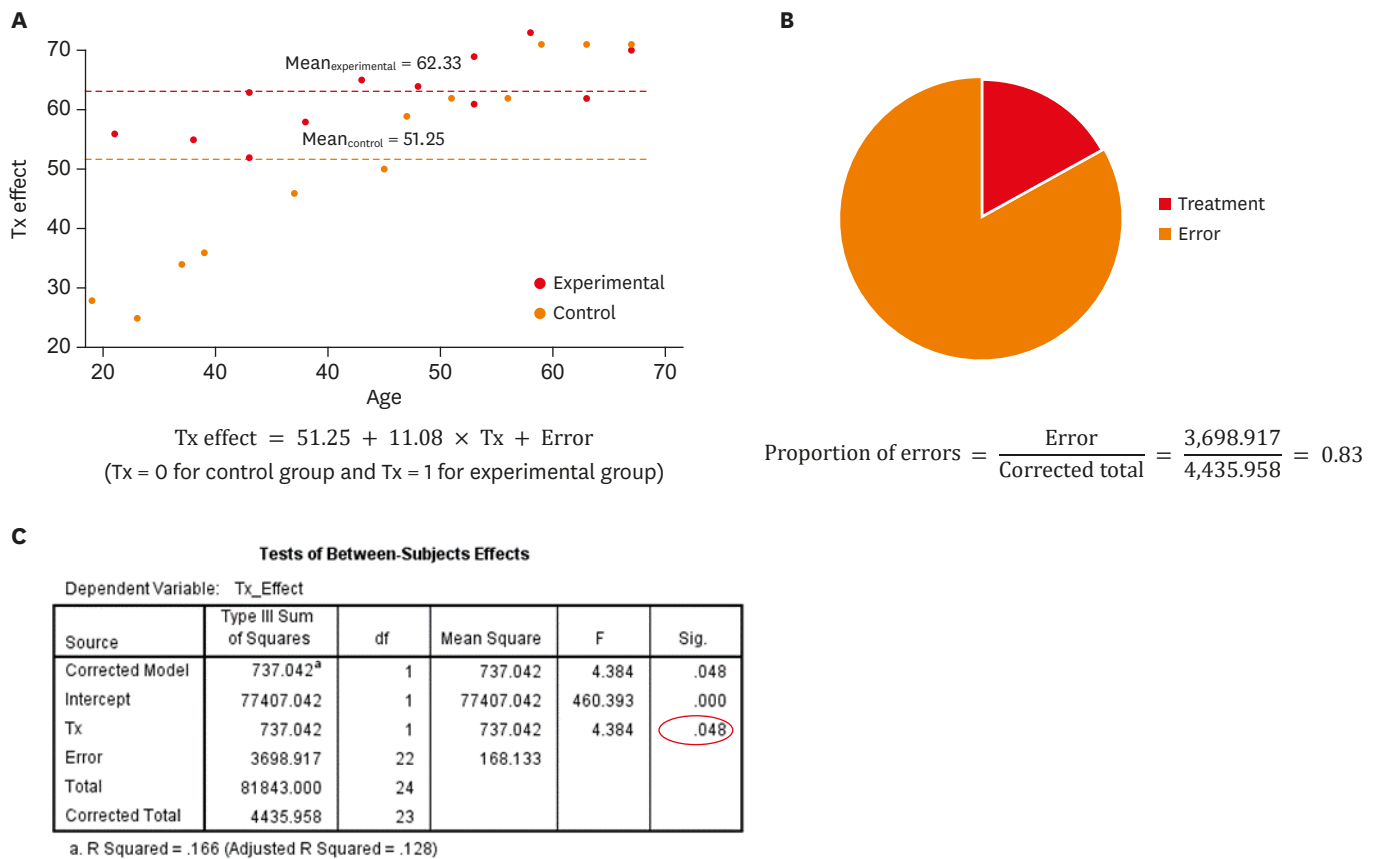
Tx effect $= 51.25 + 11.08 \times$ Tx $+$ Error
(Tx $= 0$ for control group and Tx $= 1$ for experimental group)

Proportion of errors $= \dfrac{\text{Error}}{\text{Corrected total}} = \dfrac{3{,}698.917}{4{,}435.958} = 0.83$

**Tests of Between-Subjects Effects**

Dependent Variable: Tx_Effect

| Source | Type III Sum of Squares | df | Mean Square | F | Sig. |
|---|---|---|---|---|---|
| Corrected Model | 737.042$^a$ | 1 | 737.042 | 4.384 | .048 |
| Intercept | 77407.042 | 1 | 77407.042 | 460.393 | .000 |
| Tx | 737.042 | 1 | 737.042 | 4.384 | .048 |
| Error | 3698.917 | 22 | 168.133 | | |
| Total | 81843.000 | 24 | | | |
| Corrected Total | 4435.958 | 23 | | | |

a. R Squared = .166 (Adjusted R Squared = .128)

**Figure 2.** An analysis of variance (ANOVA) model considering 2 treatment methods (Tx). (A) Depiction of data and model; (B) Proportion of errors; (C) ANOVA table.

To solve the question whether different age levels influence the degree of group difference in treatment effect level, we include age into the model. We insert the covariate, age, into the previous ANOVA model, constructing an ANCOVA model. The result is shown in **Figure 3**. The resulting ANCOVA equation is:

$$\text{Tx effect} = 19.71 + 10.18 \times \text{Tx} + 0.72 \times \text{Age} + \text{Error} \quad (5)$$

where Tx = 0 for control group and Tx = 1 for experimental group.

The difference of effect between 2 groups has changed slightly from 11.08 in Equation 4 to 10.18 in Equation 5. The size of intercept has reduced greatly by around 31 because it has been adjusted by age. One unit increase of age is related to an increase of 0.72 unit in treatment effect. The proportion of errors has decreased greatly from 0.83 to 0.21 in **Figure 3B**. The reason is because age explained a big portion of variability in the response variable (gray colored segment).

As appeared in **Figure 3C**, the proportion explained by the ANCOVA model has improved up to 78.8%, mainly due to the contribution of age variable. The *p* values of Tx and age are 0.01 and < 0.001, respectively, which represent a highly significant result. The inclusion of covariate which is highly correlated with response can remove a considerable portion of errors, reducing the proportion of errors. In contrast, the explanation ability and significance of factors increase in the ANCOVA model.
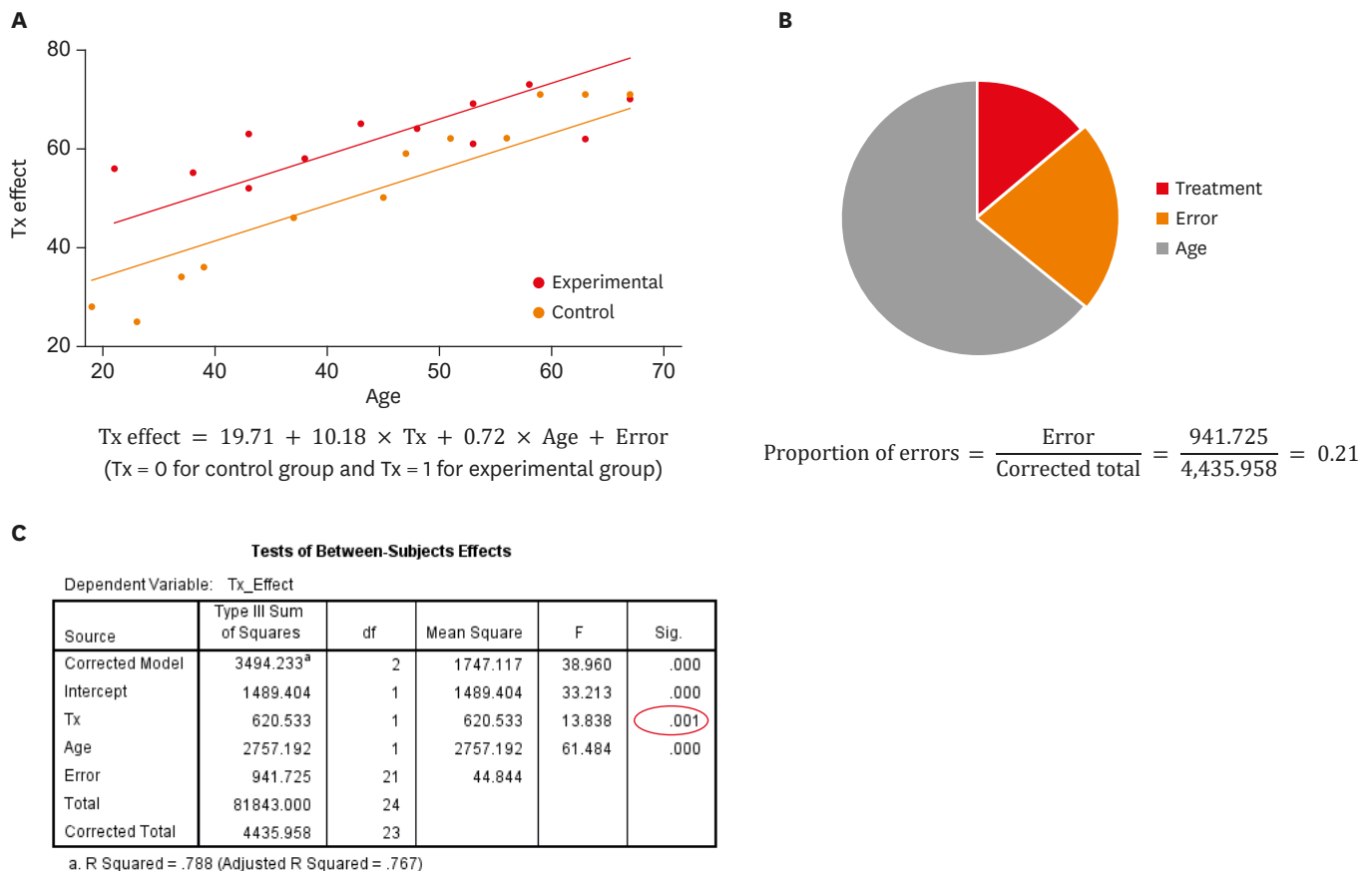
Tx effect = 19.71 + 10.18 × Tx + 0.72 × Age + Error
(Tx = 0 for control group and Tx = 1 for experimental group)

$$\text{Proportion of errors} = \frac{\text{Error}}{\text{Corrected total}} = \frac{941.725}{4,435.958} = 0.21$$

**Tests of Between-Subjects Effects**

Dependent Variable:  Tx_Effect

| Source | Type III Sum of Squares | df | Mean Square | F | Sig. |
|---|---|---|---|---|---|
| Corrected Model | 3494.233$^a$ | 2 | 1747.117 | 38.960 | .000 |
| Intercept | 1489.404 | 1 | 1489.404 | 33.213 | .000 |
| Tx | 620.533 | 1 | 620.533 | 13.838 | .001 |
| Age | 2757.192 | 1 | 2757.192 | 61.484 | .000 |
| Error | 941.725 | 21 | 44.844 | | |
| Total | 81843.000 | 24 | | | |
| Corrected Total | 4435.958 | 23 | | | |

a. R Squared = .788 (Adjusted R Squared = .767)

**Figure 3.** An analysis of covariance (ANCOVA) model considering 2 treatment methods (Tx) and age as main effects. (A) Depiction of data and model; (B) Proportion of errors; (C) analysis of variance (ANOVA) table.
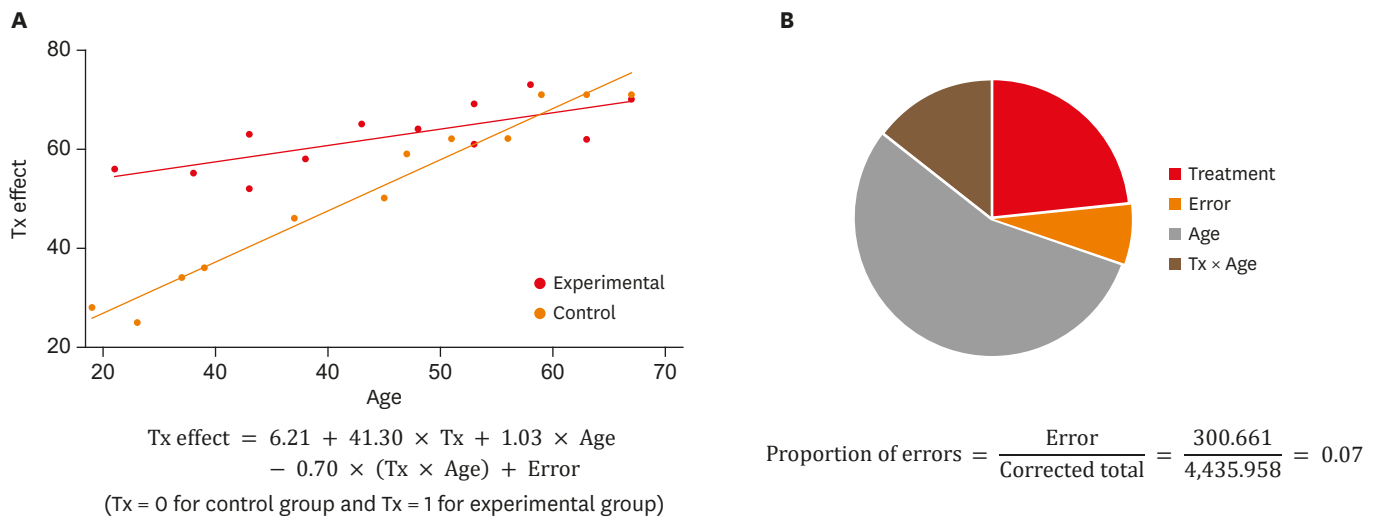
It is noticeable that the slope of age is the same as 0.72 for both treatment groups in **Figure 3A**, which is restricted by the assumption of the ANCOVA model. The slope is similar to that of pooled sample, 0.74 as appeared in Equation 1. However, the slopes of 2 groups may actually be different because the slopes of 2 groups seem substantially different from one another, 0.33 in Equation 2 and 1.03 in Equation 3.

## ANCOVA MODEL WITH INTERACTION

An ANCOVA model with interaction term is often called 'a moderated regression,' specifically [1]. Now we consider including an interaction term between group and age into the previous ANCOVA model, to assess if there is a significant difference in slopes of 2 groups. In **Figure 3C**, we find that the interaction term, Tx × Age, is statistically significant ($p < 0.001$), which supports the need of interaction term. By applying the model, the proportion of errors has decreased dramatically to 7%, as a considerable portion of variance is explained by the interaction term (**Figure 4B**). Also, 93.2% of total variance is explained by the model (R-squared = 0.932, **Figure 4C**).

The construction of ANCOVA model with interaction results in the model as follows:

Tx effect = 6.21 + 41.30 × Tx + 1.03 × Age − 0.70 × (Tx × Age) + Error   (6)

**A**



$$\text{Tx effect} = 6.21 + 41.30 \times \text{Tx} + 1.03 \times \text{Age} - 0.70 \times (\text{Tx} \times \text{Age}) + \text{Error}$$

(Tx = 0 for control group and Tx = 1 for experimental group)

**B**



$$\text{Proportion of errors} = \frac{\text{Error}}{\text{Corrected total}} = \frac{300.661}{4,435.958} = 0.07$$

**C**

**Tests of Between-Subjects Effects**

Dependent Variable: Tx_Effect

| Source | Type III Sum of Squares | df | Mean Square | F | Sig. |
|---|---|---|---|---|---|
| Corrected Model | 4135.297[a] | 3 | 1378.432 | 91.693 | .000 |
| Intercept | 1720.055 | 1 | 1720.055 | 114.418 | .000 |
| Tx | 1016.591 | 1 | 1016.591 | 67.624 | .000 |
| Age | 2413.784 | 1 | 2413.784 | 160.565 | .000 |
| Tx * Age | 641.064 | 1 | 641.064 | 42.644 | .000 |
| Error | 300.661 | 20 | 15.033 | | |
| Total | 81843.000 | 24 | | | |
| Corrected Total | 4435.958 | 23 | | | |

a. R Squared = .932 (Adjusted R Squared = .922)

**Figure 4.** Information on a model considering 2 treatment methods (Tx), age, and their interaction effect (Tx × Age). (A) Depiction of data and model; (B) Proportion of errors; (C) analysis of variance (ANOVA) table.

where Tx = 0 for control group and Tx = 1 for experimental group.

Immediately, Equation 6 can create two separate models for both control and experimental groups. The resulting Equation 7 and Equation 8 is exactly the same with the results obtained by simple regression, Equation 2 and Equation 3, respectively.

For control group (Tx = 0): Tx effect = 6.21 + 1.03 × Age + Error  (7)

For experimental group (Tx = 1): Tx effect = (6.21 + 41.30) + (1.03 − 0.70) × Age + Error = 47.51 + 0.33 × Age  (8)

## REFERENCES

1. Leppink J. Analysis of covariance (ANCOVA) *vs.* moderated regression (MODREG): why the interaction matters. Health Prof Educ 2018;4:225-232.
   **CROSSREF**

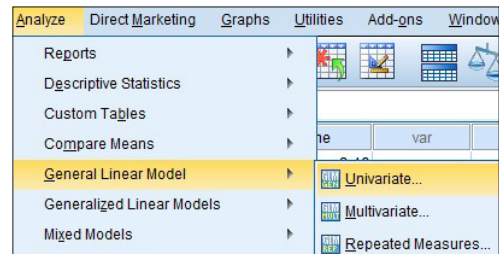**Appendix 1.** Procedure of analysis for analysis of covariance (ANCOVA) using IBM SPSS

The procedure of ANCOVA using IBM SPSS Statistics for Windows Version 23.0 (IBM Corp., Armonk, NY, USA) is as follows.
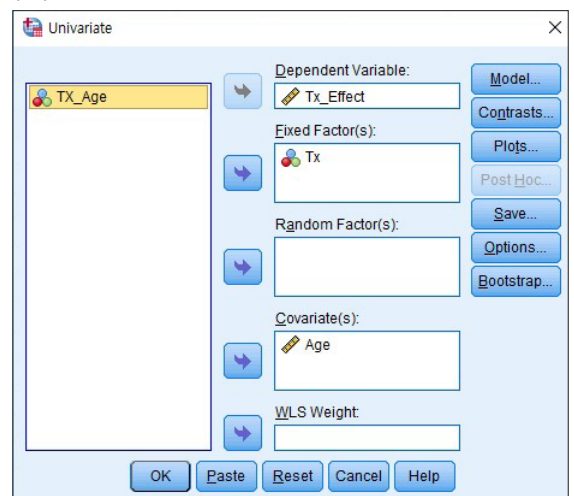
1. ANCOVA model (without interaction)

(A) Data

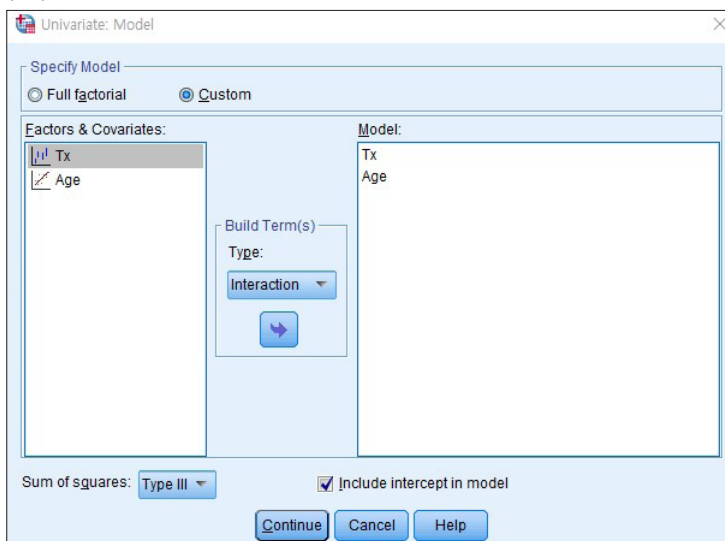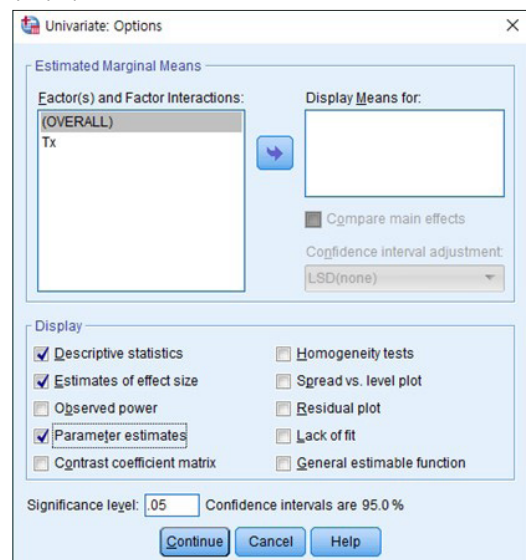| Tx_Effect | Age | Tx |
|---|---|---|
| 56 | 21 | 0 |
| 28 | 19 | 1 |
| 55 | 28 | 0 |
| 25 | 23 | 1 |
| 71 | 67 | 1 |
| 63 | 33 | 0 |
| 52 | 33 | 0 |
| 62 | 56 | 1 |
| 50 | 45 | 1 |
| 58 | 38 | 0 |
| 46 | 37 | 1 |
| 34 | 27 | 1 |
| 65 | 43 | 0 |
| 59 | 47 | 1 |
| 64 | 48 | 0 |
| 61 | 53 | 0 |
| 36 | 29 | 1 |
| 69 | 53 | 0 |
| 73 | 58 | 0 |
| 62 | 63 | 0 |
| 71 | 59 | 1 |
| 62 | 51 | 1 |
| 70 | 67 | 0 |
| 71 | 63 | 1 |

(B) Analyze-Regression-Linear



(B-1) Variables



(B-2) Model



(B-3) Options



(continued to the next page)

**Appendix 1.** (Continued) Procedure of analysis for analysis of covariance (ANCOVA) using IBM SPSS

(C-1) Descriptive statistic

### Descriptive Statistics

Dependent Variable: Tx_Effect

| Tx | Mean | Std. Deviation | N |
|---|---|---|---|
| 0 | 62.33 | 6.387 | 12 |
| 1 | 51.25 | 17.189 | 12 |
| Total | 56.79 | 13.888 | 24 |

(C-2) Analysis of variance (ANOVA) table

### Tests of Between-Subjects Effects

Dependent Variable: Tx_Effect

| Source | Type III Sum of Squares | df | Mean Square | F | Sig. | Partial Eta Squared |
|---|---|---|---|---|---|---|
| Corrected Model | 3494.233ᵃ | 2 | 1747.117 | 38.960 | .000 | .788 |
| Intercept | 1489.404 | 1 | 1489.404 | 33.213 | .000 | .613 |
| Tx | 620.533 | 1 | 620.533 | 13.838 | .001 | .397 |
| Age | 2757.192 | 1 | 2757.192 | 61.484 | .000 | .745 |
| Error | 941.725 | 21 | 44.844 | | | |
| Total | 81843.000 | 24 | | | | |
| Corrected Total | 4435.958 | 23 | | | | |

a. R Squared = .788 (Adjusted R Squared = .767)

(C-3) Estimated marginal means

### Estimates

Dependent Variable: Tx_Effect

| Tx | Mean | Std. Error | 95% Confidence Interval Lower Bound | Upper Bound |
|---|---|---|---|---|
| 0 | 61.881ᵃ | 1.934 | 57.859 | 65.903 |
| 1 | 51.702ᵃ | 1.934 | 47.680 | 55.724 |

a. Covariates appearing in the model are evaluated at the following values: Age = 44.21.

### 2. ANCOVA model with interaction

(B-2) Model (interaction term included)



(C-2) ANOVA table

Dependent Variable: Tx_Effect

| Source | Type III Sum of Squares | df | Mean Square | F | Sig. | Partial Eta Squared |
|---|---|---|---|---|---|---|
| Corrected Model | 4135.297ᵃ | 3 | 1378.432 | 91.693 | .000 | .932 |
| Intercept | 1720.055 | 1 | 1720.055 | 114.418 | .000 | .851 |
| Tx | 1016.591 | 1 | 1016.591 | 67.624 | .000 | .772 |
| Age | 2413.784 | 1 | 2413.784 | 160.565 | .000 | .889 |
| Tx * Age | 641.064 | 1 | 641.064 | 42.644 | .000 | .681 |
| Error | 300.661 | 20 | 15.033 | | | |
| Total | 81843.000 | 24 | | | | |
| Corrected Total | 4435.958 | 23 | | | | |

a. R Squared = .932 (Adjusted R Squared = .922)

(C-3) Estimated marginal means

Dependent Variable: Tx_Effect

| Tx | Mean | Std. Error | 95% Confidence Interval Lower Bound | Upper Bound |
|---|---|---|---|---|
| 0 | 62.127ᵃ | 1.120 | 59.790 | 64.464 |
| 1 | 51.896ᵃ | 1.120 | 49.559 | 54.232 |

a. Covariates appearing in the model are evaluated at the following values: Age = 44.21.