

Some ANCOVA theory

Let X be a random variable

Let X_i be iid $\sim f_x$ (pdf or pmf: probability density or mass function) with mean μ and variance σ^2 .

Think of X_i as repeated observations from the same population or the same statistic calculated for repeated experiments.

$$\text{var}(X) \equiv \sigma_x^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n - 1}$$

Let Y be another random variable. Let f_{xy} be the joint pdf (or pmf) of X and Y . f_x and f_y are called marginal pdf's. We now have an additional characteristic of the joint pdf: the covariance of X and Y .

$$\text{cov}(X, Y) = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{n - 2}$$

If X is independent of Y then $f_{xy} = f_x f_y$ and $\text{cov}(X, Y) = 0$.

(We won't use it today, but $\text{cor}(X, Y) = \text{cov}(X, Y) / (\sigma_x \sigma_y)$.)

Without any normality requirement, it is easy to show that

$$E(aX + bY) = aE(X) + bE(Y), \text{ but not } E(XY) = E(X)E(Y)$$

and

$$\text{cov}(aX, bY) = ab \text{cov}(X, Y) \text{ and } \text{cov}(aX + bY, cZ) = ac \text{cov}(X, Z) + bc \text{cov}(Y, Z)$$

and

$$\text{var}(aX + bY + c) = a^2 \text{var}(X) + b^2 \text{var}(Y) + 2ab \text{cov}(X, Y).$$

What is $\text{var}(aX + bY + cZ)$?

$$\begin{aligned}\text{var}(aX + bY + cZ) &= \text{var}((aX + bY) + cZ) \\ &= \text{var}(aX + bY) + c^2\text{var}(Z) + 2c \text{cov}(aX + bY, Z) \\ &= a^2\text{var}(X) + b^2\text{var}(Y) + 2ab \text{cov}(X, Y) + c^2\text{var}(Z) \\ &\quad + 2c(a\text{cov}(X, Z) + b\text{cov}(Y, Z)) \\ &= a^2\text{var}(X) + b^2\text{var}(Y) + c^2\text{var}(Z) + 2ab \text{cov}(X, Y) \\ &\quad + 2ac \text{cov}(X, Z) + 2bc \text{cov}(Y, Z)\end{aligned}$$

In simple linear regression, where $\hat{\beta}$ is a vector of length 2

$$\text{var}(\hat{\beta}) = \sigma^2[X'X]^{-1}$$

where X here is a matrix with the first column all 1's and the second column equal to the n explanatory variables, x_i .

Using standard matrix properties, $[X'X]$ has diagonal elements n and $\sum x_i^2$, and off-diagonal elements equal to $\sum x_i$. Also

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = g \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

where g is $\frac{1}{ad-bc}$, so $\text{var}(\hat{\beta}_0)$ and $\text{var}(\hat{\beta}_1)$ can be explicitly written out. In practice we need to use S.E. $(\hat{\beta}_j)$ where σ^2 is replaced by an estimator, $\hat{\sigma}^2 = SS_{res}/df$. Note the effect of “centering” the explanatory variable. Also note that it is CI's and p-values that need normality.

F-test of two nested models:

The numerator is an estimator of σ^2 under the null hypothesis that the extra components of the “larger” model are useless, and the denominator is always an estimator of σ^2 . The numerator is $(SS_{small} - SS_{big})/dfn$ with dfn equal to dfn , the difference in the number of parameters between the two models. The denominator is SS/dfd for the “larger” model where $dfd = n - pl$ and n is the number of subjects in the regression and pl is the number of parameters in the “larger” model. Under the null distribution, this F statistic has a central F distribution with dfn and dfd degrees of freedom.

Scheffe multiple (infinite) comparison procedure for contrasts:

$$C = \sum_{i=1}^r c_i Y_i$$

$$\text{var}(C) \equiv \sigma_C^2 = \sum_{i=1}^r c_i^2 \text{var}(Y_i) + 2 \sum_{i < j} c_i c_j \text{cov}(Y_i Y_j)$$

We need a value of m such that the experiment-wise error rate of any number of confidence intervals of the form $C \pm m \sigma_C$ is bounded by α . Scheffe found m to be $\sqrt{(r-1)F_{(\alpha, r-1, dfd)}}$, where dfd is the df of σ_C .

Summarizing adjusted means in a model with a single covariate fixed at value x_0 and T treatments and different slopes and intercepts for each treatment:

Pick a few meaningful values of x_0 such as Q1, Q2, Q3.

Let $\widehat{\mu}_j$ represent $E(Y|X = x_0, T = t_j)$. The model says that adjusted mean $\widehat{\mu}_j = \widehat{\alpha}_j + x_0 \widehat{\beta}_j$. So $\text{var}(\widehat{\mu}_j) \equiv \sigma_{\widehat{\mu}_j}^2 = \text{var}(\widehat{\alpha}_j) + x_0^2 \text{var}(\widehat{\beta}_j) + 2x_0 \text{cov}(\widehat{\alpha}_j, \widehat{\beta}_j)$.

A confidence region can be written as $\widehat{\mu}_j \pm t_{(1-\alpha/2, df)} \sigma_{\widehat{\mu}_j}$.

Multiple testing that $E(Y|X = x_0, T = t_1)$ differs from $E(Y|X = x, T = t_2)$ in a model with a single covariate fixed at value x_0 and T treatments and different slopes and intercepts for each treatment:

Let $\widehat{\mu}_j$ represent $E(Y|X = x_0, T = t_j)$. The model says $\widehat{\mu}_j = \widehat{\alpha}_j + x_0 \widehat{\beta}_j$. So $\text{var}(\widehat{\mu}_j) = \text{var}(\widehat{\alpha}_j) + x_0^2 \text{var}(\widehat{\beta}_j) + 2x_0 \text{cov}(\widehat{\alpha}_j, \widehat{\beta}_j)$. For any pair of levels of treatment, the β 's are uncorrelated. So $\text{var}(\widehat{\mu}_1 - \widehat{\mu}_2) = \text{var}(\widehat{\mu}_1) + \text{var}(\widehat{\mu}_2)$.

To find the region of x 's where there is a "significant difference" in adjusted outcomes between a pair of treatments, make an infinite number of CIs (for all x 's, or all in a reasonable range) using Scheffe's method, and see which ones exclude zero for the difference.