

STA610 Lab01 ANOVA and REANOVA

Hun Kang

2024-08-30

- Write down your answers in any blank sheet and submit your work in paper during the lab.
- Your work will not be graded. As long as you submit, you will get a full credit.
- For those who missed the lab today, you can submit it via email to me for half credit.

1. ANOVA as a linear regression

Consider an one-way ANOVA model where an i th unit of a j th group is modeled as

$$y_{ij} = \theta_j + \epsilon_{ij}, \quad i \in [n_j], j \in [m] \quad (n_j = n \forall j)$$

$\epsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$ $j = 1 \dots m$

Another way to look at this model is as follows. First, we reparameterize the group mean parameters θ_j as

$$\theta_j = \begin{cases} \alpha & j = 1 \\ \alpha + \beta_{j-1} & j \geq 2 \end{cases}$$

Then you can check that the above model can be written as

$$\mathbf{y} = \alpha \mathbf{1}_{mn} + \mathbf{X}\beta + \epsilon$$

where $\mathbf{1}_{mn}$ is a vector of 1 of length mn , $\mathbf{y} = [y_{11}, \dots, y_{nm}]^T$, similarly for ϵ and $\beta = [\beta_1, \dots, \beta_{m-1}]^T$.

Q1-1: Write down the form of a design matrix \mathbf{X} .

The hypothesis of testing no difference in group means can be written as

$$H_0 : \theta_1 = \dots = \theta_p \quad \text{vs} \quad H_1 : \theta_i \neq \theta_j \quad \exists (i, j)$$

Q1-2: Re-express the above H_0 using β .

This tells us that ANOVA can be seen as a linear regression, and its hypothesis testing is equivalent to testing a submodel of linear regression. We do not proceed from here, but the general idea is as follows. Note that we wrote ANOVA decomposition as $SST = SSA + SSW$ where

$$SST = \sum_{j=1}^m \sum_{i=1}^n (y_{ij} - \bar{y})^2$$

length of error using $c(\mathbf{I}_n)$

$$SSA = \sum_{j=1}^m \sum_{i=1}^n (\bar{y}_j - \bar{y})^2$$

$SST - SSW$ by using $c(\mathbf{I}_m \mathbf{X})$

$$SSW = \sum_{j=1}^m \sum_{i=1}^n (y_{ij} - \bar{y}_j)^2$$

length of error using $c(\mathbf{I}_m \mathbf{X})$

For a matrix \mathbf{X} , we write $c(\mathbf{X})$ the vector space its columns expand, i.e., column space. Also, we write P_X an orthogonal projection matrix onto $c(\mathbf{X})$. For convenience, let $Z = \text{cbind}(\mathbf{1}_{mn}, \mathbf{X})$. The dimension of $c(Z)$ is m .

$$y_{ij} = \theta_j + \varepsilon_{ij}, \quad \varepsilon_{ij} \sim N(0, \sigma^2) \quad i\text{-th unit of } j\text{-th group}$$

3 groups, 2 units

ANOVA

$$\begin{bmatrix} y_{11} \\ y_{12} \\ y_{21} \\ y_{22} \\ y_{31} \\ y_{32} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \alpha \\ \beta_1 \\ \beta_2 \end{bmatrix} + \begin{bmatrix} \varepsilon_{11} \\ \varepsilon_{12} \\ \varepsilon_{21} \\ \varepsilon_{22} \\ \varepsilon_{31} \\ \varepsilon_{32} \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \alpha \\ \beta_1 \\ \beta_2 \end{bmatrix} + \varepsilon$$

$$= \alpha \mathbf{1} + X\beta + \varepsilon$$

$$= Z \begin{bmatrix} \alpha \\ \beta \end{bmatrix} + \varepsilon$$

$$\theta_1 = \alpha$$

$$\theta_2 = \alpha + \beta_1$$

$$\theta_3 = \alpha + \beta_2$$

\vdots

$$\theta_m = \alpha + \beta_{m-1}$$

$$H_0: \theta_j = 0$$

$$\left(H_0: \theta_j \text{ equal} \right. \\ \left. \text{vs } H_1: \theta_j \text{ different } \right)$$

$$M_1: y = \alpha \mathbf{1} + \varepsilon$$

$$M_2: y = \alpha \mathbf{1} + \underbrace{X\beta}_{\text{significant or not}} + \varepsilon$$

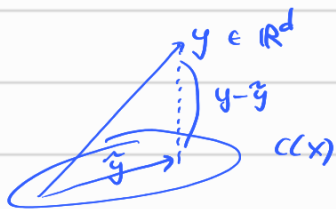
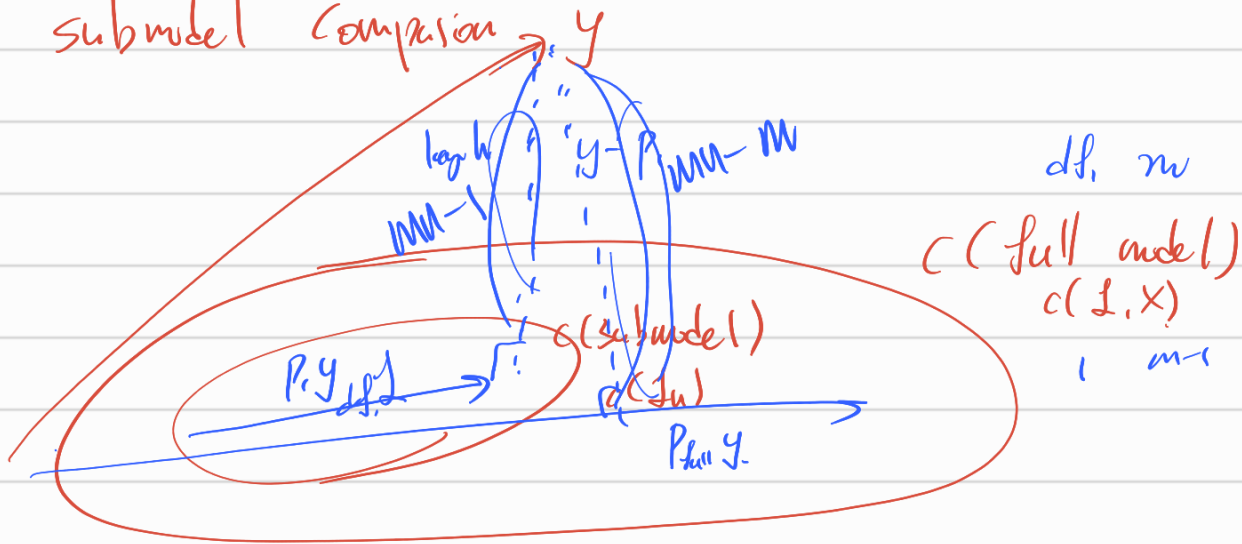
M_1 submodel of M_2 .

$$c(\underline{I_n}) \quad \text{vs} \quad c(\underline{[I_n, X]})$$

i) project y onto $c(\cdot)$: $\tilde{y} = Py$.

ii) Compare the residuals $\|y - P_y\|^2$ vs $\|y - P_X y\|^2$

submodel comparison



$$\hat{y} \in C(X) \quad \text{s.t.} \quad (y - \hat{y})^T \begin{pmatrix} \cdot \\ \cdot \end{pmatrix} = 0 \quad \text{element of } C(X)$$

$$(y - X\beta)^T X\beta = 0$$

$$y \in \mathbb{R}^d, \quad \beta \neq 0$$

$$(X^T y - X^T X \beta)^T \beta = 0$$

$X\beta$: any element in $C(X)$

$$= 0$$

$$X^T y = X^T X \beta$$

$$X = \begin{pmatrix} | & | & | \\ \cdot & \cdot & \cdot \\ | & | & | \end{pmatrix} \quad \text{full rank}$$

$$\beta^* = (X^T X)^{-1} X^T y$$

$$X\beta^* = \underbrace{X(X^T X)^{-1} X^T}_{P \text{ projection matrix}} y \quad \text{closest to } y \text{ among } C(X)$$

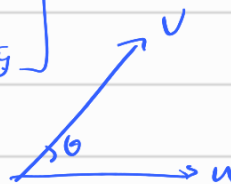
$$X = I_n$$

$$P_{Lu} = I_n (I_n^T I_n)^{-1} I_n^T y = \frac{1}{n} (I_n^T I_n) y$$

$$= \frac{1}{n} \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix} y = \begin{bmatrix} \bar{y} \\ \vdots \\ \bar{y} \end{bmatrix}$$

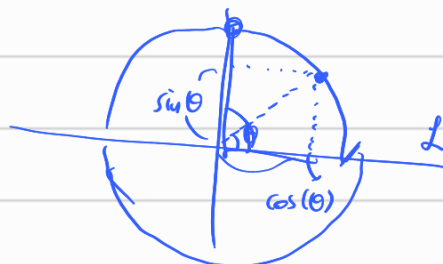
$$u, v \in \mathbb{R}^d$$

$$u^T v = \sum_{i=1}^d u_i v_i = \langle u, v \rangle$$



$$= \|u\| \|v\| \cos \theta$$

$$= 0 \quad (\theta = 90^\circ)$$



$$(I-P_1)y^T(I-P_1)y = y^T(I-P_1)^T(I-P_1)y = y^T(I-P_1)y$$

1. One can see that $SST = \|(I - P_1)y\|^2 = y^T(I - P_1)y$, $SSA = y^T(P_Z - P_1)y$ and $SSW = y^T(I - P_Z)y$.
2. $(I - P_1)y$ is a vector y projected onto the column space orthogonal to $c(1_{mn})$, whose dimension is $mn - 1$.
3. $(I - P_Z)y$ is a vector y projected onto the column space orthogonal to $c(Z)$, whose dimension is $mn - m = m(n - 1)$.
4. $(P_Z - P_1)y$ is a vector y projected onto the subspace of $c(Z)$ that are orthogonal to $c(1_{mn})$, whose dimension is $m - 1$.

Using these, the F-statistics of testing H_0 is

$$F(y) = \frac{SSA/(m-1)}{SSW/m(n-1)}$$

SSA: norm of y projected onto $c(Z) \cap c(1_n)^\perp$

norm of y projected onto $c(1_n)$

Note that $SSA = y^T(P_Z - P_1)y = y^T(I - P_1)y - y^T(I - P_Z)y$. A different interpretation of $y^T(I - P_Z)y$ is to see it as a residual sum of squares of a model with a design matrix Z . In this aspect, large SSA means that by including X , we see a large decrease in the residuals, i.e., the full model with Z fits the data better than the intercept only model.

```
library(tidyverse)
URL <- "https://campus.murraystate.edu/academic/faculty/cmecklin/STA565/wheat.txt"
wheat <- read.table(URL, header=TRUE)
```

```
str(wheat)
```

```
## 'data.frame': 30 obs. of 3 variables:
## $ variety : chr "A" "A" "A" "A" ...
## $ location: int 1 2 3 4 5 6 1 2 3 4 ...
## $ yield : num 35.3 31 32.7 36.8 37.2 33.1 33.7 32.2 31.4 32.7 ...
```

```
lm1 = lm(yield ~ 1, data = wheat)
lm2 = lm(yield ~ 1 + as.factor(location), data = wheat)
anova(lm1, lm2)
```

test of sub add!

```
## Analysis of Variance Table
##
## Model 1: yield ~ 1
## Model 2: yield ~ 1 + as.factor(location)
## Res.Df RSS Df Sum of Sq F Pr(>F)
## 1 29 151.34
## 2 24 103.60 5 47.742 2.212 0.08631 .
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Q1-3: Conduct the f-test of treatment effect as learned in the class and check that the test statistics is the same.

anova(lm2)

2. REANOVA and covariance structure

Consider a REANOVA model

$$y_{ij} = \mu + \underbrace{a_j}_{\text{random}} + \epsilon_{ij}, \quad i \in [n_j], j \in [m] \quad (n_j = n^{\vee} j)$$

this is to model Covariance matrix.

$$a_j \stackrel{iid}{\sim} N(0, \tau^2)$$

$$\epsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2), \quad E[a_j \epsilon_{ij}] = 0$$

indep
 $\Rightarrow \text{Cov}(y_{1j}, y_{2j}) = \text{Cov}(\mu + a_j + \epsilon_{1j}, \mu + a_j + \epsilon_{2j}) = \tau^2$

where μ , σ^2 and τ^2 are some unknown fixed parameter.

In the class, we saw that $E(y_{ij}) = \mu$, $V(y_{ij}) = \sigma^2 + \tau^2$ and $\text{Cov}(y_{1j}, y_{2j}) = \tau^2$. Also, since y_{ij} is a sum of Gaussian random variables, it itself also follows normal distribution. From this, we can write a joint distribution of $\mathbf{y}_j = [y_{1j}, \dots, y_{n_j}]^T$ for a group j as

$$\mathbf{y}_j \sim N(\mu \mathbf{1}_{n_j}, \Sigma_j)$$

$\mathbf{y}_j = \begin{pmatrix} \vdots \end{pmatrix} \sim \text{MVN} \left(\begin{pmatrix} \vdots \end{pmatrix}, \begin{pmatrix} \sigma^2 + \tau^2 & \tau^2 & \dots \\ \tau^2 & \sigma^2 + \tau^2 & \dots \\ \vdots & \vdots & \ddots \end{pmatrix} \right)$

Q2-1: Write down the covariance matrix Σ_j as follows:

$$(\Sigma_j)_{kl} = \begin{cases} ? & (k = l) \\ ? & (k \neq l) \end{cases}$$

*$\text{Cov}(y_{ij}, y_{ij}) = \sigma^2 + \tau^2$
 τ^2*

Combining all \mathbf{y}_j , we can rewrite the above REANOVA model in a matrix-vector form as

$$\mathbf{y} \sim N(\mu \mathbf{1}_{mn}, \Sigma)$$

$\text{Cov}(y_{i1}, y_{i2}) = \sigma^2 + \tau^2$

Q2-2: What is $\text{Cov}(y_{i1}, y_{i2})$? Using this, how can we write Σ ?

In contrast, if we change the model as

$$y_{ij} = \mu + \underbrace{\alpha_j}_{\text{constant}} + \epsilon_{ij}, \quad i \in [n_j], j \in [m] \quad (n_j = n^{\vee} j)$$

$$\epsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$$

$V(y_{ij}) = \sigma^2$ $\text{Cov}(y_{1j}, y_{2j}) = 0$

where α_j is also some unknown fixed parameter, then we have

$$\mathbf{y}_j \sim N((\mu + \alpha_j) \mathbf{1}_{n_j}, \sigma^2 I_{n_j})$$

$\begin{bmatrix} \sigma^2 & & 0 \\ & \ddots & \\ 0 & & \sigma^2 \end{bmatrix}$

Q2-3: Write down the joint model for \mathbf{y}

The key takeaway is that, by adding another source of randomness a_j for group-wise variation, marginally we are modelling the covariance structure of \mathbf{y} as Σ .

