Lecture 17

§1 ANOVA 中的 statistics

1. Definition: total sum of square (TSS / SST) (总平方和统计量)

TSS 可被表示为
$$TSS = y^T y - \frac{1}{n} y^T J y$$
$$= y^T (I - \frac{1}{n} J) y$$

其中 J 为所有元素均为 1 的矩阵

证明:
$$TSS = \sum_{i=1}^{n} (y_i - \bar{y})^2$$
$$= \sum_{i=1}^{n} y_i^2 - n\bar{y}^2$$
$$= y^T y - (\frac{1}{n} J y)^T (\frac{1}{n} J y) \quad (\frac{1}{n} J y = [\bar{y}, \bar{y}, \cdots, \bar{y}]^T = \bar{y} \cdot 1)$$
$$= y^T y - \frac{1}{n} y J y \quad (\frac{1}{n} J \text{ 为 idempotent: } \frac{1}{n} J \cdot \frac{1}{n} J = \frac{1}{n} J)$$

2. Property: $I - \frac{1}{n} J$ 的性质

① $I - \frac{1}{n} J$ 为 symmetric matrix, 即 $(I - \frac{1}{n} J)^T = I - \frac{1}{n} J$

证明:
$$(I - \frac{1}{n} J)^T = I^T - \frac{1}{n} J^T = I - \frac{1}{n} J$$

注: 若 matrix A 为 symmetric, 则 $y^T A y$ 为 quadratic form, 因此 TSS 为 quadratic form

② $I - \frac{1}{n} J$ 为 idempotent matrix, 即 $(I - \frac{1}{n} J)(I - \frac{1}{n} J) = I - \frac{1}{n} J$

证明:
$$(I - \frac{1}{n} J)(I - \frac{1}{n} J) = I - \frac{2}{n} J + \frac{1}{n} J \frac{1}{n} J = I - \frac{2}{n} J + \frac{1}{n} J = I - \frac{1}{n} J$$

③ $I - \frac{1}{n} J$ 的 rank 为 $n-1$ (为 SST 的 d.f.)

证明:
$$rank(I - \frac{1}{n} J) = rank(I) - rank(\frac{1}{n} J) = n - 1$$

3. Definition: sum of square errors (SSE / RSS) (误差/残差平方和)

SSE 可被表示为
$$SSE = \hat{e}^T \hat{e}$$
$$= y^T (I - H) y$$

其中 J 为所有元素均为 1 的矩阵

证明:
$$SSE = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$
$$= \hat{e}^T \hat{e}$$
$$= [(I-H)y]^T [(I-H)y]$$
$$= y^T (I-H)^T (I-H) y$$
$$= y^T (I-H) y \quad (I-H \text{ 为 symmetric 且 idempotent})$$

注: SSE 为 quadratic form, $\text{rank}(I-H) = \text{rank}(I) - \text{rank}(H) = n - \sum_{i=1}^{n} h_{ii} = n-2$.
因此 $d.f.(SSE) = n-2$

4. <u>Definition</u>: sum of squared due to the regression model (SSR/SSReg) (回归平方和)

SSR 可被表示为
$$SSR = y^T H y - \frac{1}{n} y^T J y$$
$$= y^T (H - \frac{1}{n} J) y$$

证明:
$$SSR = \sum_{i=1}^{n} (\hat{y}_i - \bar{y}_i)^2$$
$$= \sum_{i=1}^{n} (\hat{y}_i^2 - 2\hat{y}_i \bar{y}_i + \bar{y}_i^2)$$
$$= \sum_{i=1}^{n} \hat{y}_i^2 - n\bar{y} \quad (\text{由于} \sum_{i=1}^{n} \hat{y}_i = \sum_{i=1}^{n} y_i, \text{有} \sum_{i=1}^{n}(-2\hat{y}_i\bar{y}_i + \bar{y}_i^2) = -2n\bar{y}_i^2 + n y_i^2 = -ny_i^2)$$
$$= \hat{y}^T \hat{y} - \frac{1}{n} y^T J y$$
$$= (Hy)^T (Hy) - \frac{1}{n} y^T J y$$
$$= y^T H y - \frac{1}{n} y^T J y \quad (H \text{ 为 symmetric 且 idempotent})$$

注: SSR 为 quadratic form, $\text{rank}(H - \frac{1}{n}J) = \text{rank}(H) - \text{rank}(\frac{1}{n}J) = 2-1 = 1$
因此 $d.f.(SSR) = 1$

5. <u>Lemma</u>: $TSS = SSE + SSR$ (Cochran's theorem)

若 observations $(X, y)$ 服从 linear model
$$y = X\beta + \varepsilon$$
其中 $\varepsilon \sim N(0, \sigma^2 I)$
定义 $\hat{y} = X\hat{\beta}$ (fitted value), 则有
$$TSS = SSE + SSR$$
即 $y^T (I - \frac{1}{n}J) y = y^T (I-H) y + y^T (H - \frac{1}{n}J) y$

证明: (第一种证法)
$$TSS = y^T (I - \frac{1}{n}J) y$$
$$= y^T (I - H + H - \frac{1}{n}J) y$$
$$= y^T (I-H) y + y^T (H - \frac{1}{n}J) y$$
$$= SSE + SSR$$

证明: (第二种证法)
由于 $TSS = y^T y - \frac{1}{n} y^T J y$, 先处理 $y^T y$:
$$y^T y = (y - X\hat{\beta} + X\hat{\beta})^T (y - X\hat{\beta} + X\hat{\beta})$$
$$= (y - X\hat{\beta})^T (y - X\hat{\beta}) + (y - X\hat{\beta})^T X\hat{\beta} + (X\hat{\beta})^T (y - X\hat{\beta}) + (X\hat{\beta})^T X\hat{\beta}$$
$$= \hat{e}^T \hat{e} + \hat{e}^T X\hat{\beta} + (X\hat{\beta})^T \hat{e} + \hat{\beta} X^T X \hat{\beta}$$
其中,
① $\hat{e}^T X\hat{\beta} = (X\hat{\beta})^T \hat{e} = 0$  由于 $\hat{e}^T X = (y - Hy)^T X = y^T X - y^T X = 0$

$$\textcircled{2} \quad \hat{\beta} X^T X \hat{\beta} = \hat{y}^T H \hat{y} \quad \text{由于} \quad X\hat{\beta} = \hat{y} = H\hat{y}$$

因此，
$$\hat{y}^T \hat{y} = \hat{e}^T \hat{e} + \hat{y}^T H \hat{y}$$
$$TSS = \underbrace{\hat{e}^T \hat{e}}_{SSE} + \underbrace{\hat{y}^T H \hat{y} - \frac{1}{n} \hat{y}^T J \hat{y}}_{SSR}$$
$$= SSE + SSR$$

b. **Property**：SSE / RSS 的期望（用SSE估计 $\sigma^2$）

对于 $SSE = \hat{e}^T \hat{e} = y^T(I-H)y$，有
$$E(SSE) = (n-2)\sigma^2$$

**证明**：
$$E(SSE) = E(\hat{e}^T \hat{e})$$
$$= E(\hat{y}^T(I-H)\hat{y})$$
$$= E(trace(\hat{y}^T(I-H)\hat{y})) \quad (\text{由于} \ \hat{y}^T(I-H)\hat{y} \ \text{为 scalar})$$
$$= E(trace((I-H)\hat{y}\hat{y}^T)) \quad (\text{由于} \ trace(AB) = trace(BA))$$
$$= trace[(I-H)E(yy^T)]$$
$$= trace[(I-H)(\sigma^2 I + X\beta\beta^T X^T)] \quad (E(\hat{y}\hat{y}^T) = Var(\hat{y}) + E(\hat{y})E(\hat{y}^T))$$
$$= trace[(I-H)\sigma^2 + X\beta\beta^T X^T - X(X^T X)^{-1}X^T X\beta\beta^T X^T]$$
$$= trace(I-H)\sigma^2$$
$$= (n-2)\sigma^2$$

注：由此证明了 $S^2 = \dfrac{RSS}{n-2}$ 为 $\sigma^2$ 的无偏估计量

7. **R中的 ANOVA**

The anova command is one way in R to produce an ANOVA table, in addition to analysing it. For example, for the 654-point SLR problem in Assignment 2, question 1:

```
a2 = read.table("data.txt",sep="_",header=T) # Load the data set
fev <- a2$fev; age <- a2$age
mod1 = lm(fev~age)
anova(mod1)
```

```
## Analysis of Variance Table
##
## Response: fev
##            Df     Sum      Sq    Mean Sq    F value     Pr(>F)
## age         1   280.92   280.919   872.18   < 2.2e-16   ***
## Residuals 652   210.00     0.322
## —
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```