# ANalysis Of VAriance (ANOVA)

$y = X\beta + \epsilon, \ \epsilon \sim N(\mathbf{0}, \ \sigma^2 I)$

Let $X_1 = \mathbf{1}, \ X_m = X,$ and $X_{m+1} = I$.

Suppose $X_2, \ldots, X_{m-1}$ are design matrices satisfying

$\mathcal{C}(X_1) \subset \mathcal{C}(X_2) \subset \cdots \subset \mathcal{C}(X_{m-1}) \subset \mathcal{C}(X_m).$

Let $r_j = \text{rank}(X_j) \ \forall \ j = 1, \ldots, m+1.$

Let $P_j = P_{X_j} \quad \forall\, j = 1, \ldots, m+1$. Then

$$
\begin{aligned}
\sum_{i=1}^{n} (y_i - \bar{y}.)^2 &= y'(I - P_1)y = y'(P_{m+1} - P_1)y \\
&= y' \left( \sum_{j=2}^{m+1} P_j - \sum_{j=1}^{m} P_j \right) y \\
&= y'(P_{m+1} - P_m + P_m - P_{m-1} + \cdots + P_2 - P_1)y \\
&= y'(P_{m+1} - P_m)y + \ldots + y'(P_2 - P_1)y \\
&= \sum_{j=1}^{m} y'(P_{j+1} - P_j)y.
\end{aligned}
$$

The sums of squares in the equation

$$\boldsymbol{y}'(\boldsymbol{I} - \boldsymbol{P}_1)\boldsymbol{y} = \sum_{j=1}^{m} \boldsymbol{y}'(\boldsymbol{P}_{j+1} - \boldsymbol{P}_j)\boldsymbol{y}$$

are often arranged in an ANOVA table.

| Sum of Squares | Sum of Squares |
|---|---|
| $\mathbf{y}'(\mathbf{P}_2 - \mathbf{P}_1)\mathbf{y}$ | $SS(2 \mid 1)$ |
| $\mathbf{y}'(\mathbf{P}_3 - \mathbf{P}_2)\mathbf{y}$ | $SS(3 \mid 2)$ |
| $\vdots$ | $\vdots$ |
| $\mathbf{y}'(\mathbf{P}_m - \mathbf{P}_{m-1})\mathbf{y}$ | $SS(m \mid m-1)$ |
| $\mathbf{y}'(\mathbf{P}_{m+1} - \mathbf{P}_m)\mathbf{y}$ | $SSE = \mathbf{y}'(\mathbf{I} - \mathbf{P}_X)\mathbf{y}$ |
| $\mathbf{y}'(I - \mathbf{P}_1)\mathbf{y}$ | $SSTo = \sum_{i=1}^{n}(y_i - \bar{y}.)^2$ |

Note that $\forall\, j = 1,\, \ldots,\, m$

$$
\begin{aligned}
(\boldsymbol{P}_{j+1} - \boldsymbol{P}_j)(\boldsymbol{P}_{j+1} - \boldsymbol{P}_j) &= \boldsymbol{P}_{j+1}\boldsymbol{P}_{j+1} - \boldsymbol{P}_{j+1}\boldsymbol{P}_j - \boldsymbol{P}_j\boldsymbol{P}_{j+1} + \boldsymbol{P}_j\boldsymbol{P}_j \\
&= \boldsymbol{P}_{j+1} - \boldsymbol{P}_j - \boldsymbol{P}_j + \boldsymbol{P}_j \\
&= \boldsymbol{P}_{j+1} - \boldsymbol{P}_j.
\end{aligned}
$$

Also, $\forall\, j < \ell$

$$
\begin{aligned}
(\boldsymbol{P}_{j+1} - \boldsymbol{P}_j)(\boldsymbol{P}_{\ell+1} - \boldsymbol{P}_\ell) &= \boldsymbol{P}_{j+1}\boldsymbol{P}_{\ell+1} - \boldsymbol{P}_{j+1}\boldsymbol{P}_\ell - \boldsymbol{P}_j\boldsymbol{P}_{\ell+1} + \boldsymbol{P}_j\boldsymbol{P}_\ell \\
&= \boldsymbol{P}_{j+1} - \boldsymbol{P}_{j+1} - \boldsymbol{P}_j + \boldsymbol{P}_j \\
&= \boldsymbol{0}.
\end{aligned}
$$

Using these facts and previous facts about distributions of quadratic forms, it can be shown that

$$\frac{\boldsymbol{y}'(\boldsymbol{P}_{j+1} - \boldsymbol{P}_j)\boldsymbol{y}}{\sigma^2} \sim \chi^2_{r_{j+1}-r_j}(\boldsymbol{\beta}'\boldsymbol{X}'(\boldsymbol{P}_{j+1} - \boldsymbol{P}_j)\boldsymbol{X}\boldsymbol{\beta}/\sigma^2)$$

for all $j = 1, \ldots, m$ and that these $m$ $\chi^2$ random variables are mutually independent.

| Sum of Squares | Degrees of Freedom | DF |
|---|---|---|
| $\boldsymbol{y}'(\boldsymbol{P}_2 - \boldsymbol{P}_1)\boldsymbol{y}$ | $\mathrm{rank}(\boldsymbol{X}_2) - \mathrm{rank}(\boldsymbol{X}_1)$ | $r_2 - 1$ |
| $\boldsymbol{y}'(\boldsymbol{P}_3 - \boldsymbol{P}_2)\boldsymbol{y}$ | $\mathrm{rank}(\boldsymbol{X}_3) - \mathrm{rank}(\boldsymbol{X}_2)$ | $r_3 - r_2$ |
| $\vdots$ | $\vdots$ | $\vdots$ |
| $\boldsymbol{y}'(\boldsymbol{P}_m - \boldsymbol{P}_{m-1})\boldsymbol{y}$ | $\mathrm{rank}(\boldsymbol{X}_m) - \mathrm{rank}(\boldsymbol{X}_{m-1})$ | $r - r_{m-1}$ |
| $\boldsymbol{y}'(\boldsymbol{P}_{m+1} - \boldsymbol{P}_m)\boldsymbol{y}$ | $\mathrm{rank}(\boldsymbol{X}_{m+1}) - \mathrm{rank}(\boldsymbol{X}_m)$ | $n - r$ |
| $\boldsymbol{y}'(\boldsymbol{I} - \boldsymbol{P}_1)\boldsymbol{y}$ | $\mathrm{rank}(\boldsymbol{X}_{m+1}) - \mathrm{rank}(\boldsymbol{X}_1)$ | $n - 1$ |

For $j = 1, \ldots, m-1$ we have

$$F_j = \frac{y'(P_{j+1} - P_j)y/(r_{j+1} - r_j)}{y'(I - P_X)y/(n-r)}$$

$$\sim F_{r_{j+1}-r_j, n-r}(\beta'X'(P_{j+1} - P_j)X\beta/\sigma^2).$$

For $j = 1, \ldots, m-1$, define

$$MS(j+1 \mid j) = \frac{SS(j+1 \mid j)}{r_{j+1} - r_j} = \frac{y'(P_{j+1} - P_j)y}{r_{j+1} - r_j}.$$

# ANOVA Table

| Sum of Squares | Degrees of Freedom | Mean Square |
|---|---|---|
| $SS(2 \mid 1)$ | $r_2 - 1$ | $MS(2\mid1)$ |
| $SS(3 \mid 2)$ | $r_3 - r_2$ | $MS(3\mid2)$ |
| $\vdots$ | $\vdots$ | $\vdots$ |
| $SS(m \mid m-1)$ | $r - r_{m-1}$ | $MS(m\mid m-1)$ |
| $SSE$ | $n - r$ | $MSE$ |
| $SSTO$ | $n - 1$ | |

Note that

$$
\begin{aligned}
SS(j+1 \mid j) &= \mathbf{y}'(\mathbf{P}_{j+1} - \mathbf{P}_j)\mathbf{y} \\
&= \mathbf{y}'(\mathbf{P}_{j+1} - \mathbf{P}_j + \mathbf{I} - \mathbf{I})\mathbf{y} \\
&= \mathbf{y}'(\mathbf{I} - \mathbf{P}_j - \mathbf{I} + \mathbf{P}_{j+1})\mathbf{y} \\
&= \mathbf{y}'(\mathbf{I} - \mathbf{P}_j)\mathbf{y} - \mathbf{y}'(\mathbf{I} - \mathbf{P}_{j+1})\mathbf{y} \\
&= SSE_{\text{REDUCED}} - SSE_{\text{FULL}}
\end{aligned}
$$

Thus, $SS(j + 1 \mid j)$ is the amount the error sum of square decreases when $y$ is projected onto $\mathcal{C}(X_{j+1})$ instead of $\mathcal{C}(X_j)$.

$SS(j + 1 \mid j), j = 1, \ldots, m - 1$ are called *"Sequential Sums of Squares."*

SAS calls these "Type I Sums of Squares. "

The statistic

$$F_j = \frac{MS(j+1 \mid j)}{MSE}$$

can be used to test

$$H_0 : \mathrm{E}(\boldsymbol{y}) \in \mathcal{C}(\boldsymbol{X}_j) \text{ vs. } \mathrm{H_A} : \mathrm{E}(\boldsymbol{y}) \in \mathcal{C}(\boldsymbol{X}_{j+1}) \setminus \mathcal{C}(\boldsymbol{X}_j).$$

The noncentrality parameter is

$$
\begin{aligned}
\boldsymbol{\beta}'\boldsymbol{X}'(\boldsymbol{P}_{j+1} - \boldsymbol{P}_j)\boldsymbol{X}\boldsymbol{\beta}/\sigma^2 &= \frac{\boldsymbol{\beta}'\boldsymbol{X}'(\boldsymbol{P}_{j+1} - \boldsymbol{P}_j)'(\boldsymbol{P}_{j+1} - \boldsymbol{P}_j)\boldsymbol{X}\boldsymbol{\beta}}{\sigma^2} \\
&= ||\,(\boldsymbol{P}_{j+1} - \boldsymbol{P}_j)\boldsymbol{X}\boldsymbol{\beta}\,||^2\,/\sigma^2 \\
&= ||\,\boldsymbol{P}_{j+1}\mathrm{E}(\boldsymbol{y}) - \boldsymbol{P}_j\mathrm{E}(\boldsymbol{y})\,||^2\,/\sigma^2.
\end{aligned}
$$

If $H_0$ is true, $\boldsymbol{P}_{j+1}\mathrm{E}(\boldsymbol{y}) = \boldsymbol{P}_j\mathrm{E}(\boldsymbol{y}) = \mathrm{E}(\boldsymbol{y})$.

Thus, the $NCP = 0$ under $H_0$.

## Example: Multiple Regression

$$
\begin{aligned}
X_1 &= \mathbf{1} \\
X_2 &= [\mathbf{1},\ x_1] \\
X_3 &= [\mathbf{1},\ x_1,\ x_2] \\
&\vdots \\
X_m &= [\mathbf{1},\ x_1,\ \ldots,\ x_{m-1}]
\end{aligned}
$$

$SS(j + 1 \mid j)$ is the decrease in SSE that results when the explanatory variable $x_j$ is added to a model containing an intercept and explanatory variables $x_1, \ldots, x_{j-1}$.

# Example: Test for Linear Trend and Test for Lack of Linear Fit

$$X_1 = \mathbf{1}, \quad X_2 = \begin{bmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1 \\ 1 & 2 \\ 1 & 2 \\ 1 & 2 \\ 1 & 3 \\ 1 & 3 \\ 1 & 3 \end{bmatrix}, \quad X_3 = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix}$$

# $F$ Test for a Linear Trend

Let $\mu_i$ = mean yield for a plot that received $i$ units of fertilizer $(i = 1, 2, 3)$.

$$\frac{MS(2|1)}{MSE}$$

can be used to test

$$H_0 : \mu_1 = \mu_2 = \mu_3 \Longleftrightarrow \mu_i = \beta_0 \ \forall \ i = 1, 2, 3 \text{ for some } \beta_0 \in I\!R$$

versus

$$H_A : \mu_i = \beta_0 + \beta_1(i) \quad i = 1, 2, 3 \text{ for some } \beta_0 \in I\!R, \ \beta_1 \in I\!R \setminus \{0\}.$$

# $F$ Test for Lack of Linear Fit

The statistic

$$\frac{MS(3|2)}{MSE}$$

can be used to test

$$H_0 : \mu_i = \beta_0 + \beta_1(i) \quad i = 1, 2, 3 \text{ for some } \beta_0, \beta_1 \in I\!R$$

versus

$$H_A : \text{There does not exist } \beta_0, \beta_1 \in I\!R \text{ such that}$$

$$\mu_i = \beta_0 + \beta_1(i) \qquad \forall \, i = 1, 2, 3.$$

The lack of fit test is a reduced vs. full model $F$ test.

Thus, we can also obtain this test by testing

$$H_0 : \boldsymbol{C\beta} = \boldsymbol{d} \quad \text{vs.} \quad H_A : \boldsymbol{C\beta} = \boldsymbol{d}$$

for appropriate $\boldsymbol{C}$ and $\boldsymbol{d}$.

$$\boldsymbol{\beta} = \left[ \begin{array}{c} \mu_1 \\ \mu_2 \\ \mu_3 \end{array} \right] \qquad \boldsymbol{C} = ? \qquad \boldsymbol{d} = ?$$