

3.5.3 Partitioning the Total Sum of Squares

The total sum of squares, which we shall call SST, is

$$SST = \mathbf{Y}^\top \mathbf{Y} = \sum_{i=1}^n y_i^2,$$

where $\mathbf{Y} = (y_1, y_2, \dots, y_n)^\top \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$. It can be easily shown that

$$\frac{SST}{\sigma^2} \sim \chi^2_{(n, \frac{\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta}}{2\sigma^2})}.$$

Recall that the sum of squares of deviations of observed y_i 's from their predicted values is $SSE = \mathbf{Y}^\top (\mathbf{I} - \mathbf{H}) \mathbf{Y} = \mathbf{Y}^\top \mathbf{Y} - \hat{\boldsymbol{\beta}}^\top \mathbf{X}^\top \mathbf{Y}$. The difference

$$\begin{aligned} SSR &= SST - SSE \\ &= \mathbf{Y}^\top \mathbf{Y} - (\mathbf{Y}^\top \mathbf{Y} - \hat{\boldsymbol{\beta}}^\top \mathbf{X}^\top \mathbf{Y}) \\ &= \hat{\boldsymbol{\beta}}^\top \mathbf{X}^\top \mathbf{Y} = \mathbf{Y}^\top \mathbf{H} \mathbf{Y} \end{aligned}$$

represents that portion of SST that is attributed to having fitted the regression, and hence it is called the *sum of squares due to regression*, SSR. It is also often called the *reduction in sum of squares*. Similarly, we can show

$$\frac{SSR}{\sigma^2} \sim \chi^2_{(r(\mathbf{X}), \frac{(\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta})}{2\sigma^2})}.$$

The partitioning of SST can be summarized in a manner that serves as a foundation for developing the traditional analysis of variance (ANOVA) table.

Table 1: Traditional ANOVA table.

Source	df	SS	MS	F-statistics
Regression	$r(\mathbf{X})$	$\hat{\boldsymbol{\beta}}^\top \mathbf{X}^\top \mathbf{Y}$	$\frac{SSR}{r(\mathbf{X})}$	$F(R) = \frac{MSR}{MSE}$
Error	$n - r(\mathbf{X})$	$\mathbf{Y}^\top \mathbf{Y} - \hat{\boldsymbol{\beta}}^\top \mathbf{X}^\top \mathbf{Y}$	$\frac{SSE}{n - r(\mathbf{X})}$	
Total	n	$\mathbf{Y}^\top \mathbf{Y}$		

Note: Since $\mathbf{Y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$ and

$$\begin{aligned} SSR &= \mathbf{Y}^\top \mathbf{H} \mathbf{Y} \quad (\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top), \sim \chi^2_{(r(\mathbf{X}), \frac{\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta}}{2\sigma^2})} \\ SSE &= \mathbf{Y}^\top (\mathbf{I} - \mathbf{H}) \mathbf{Y}, \sim \chi^2_{(n - r(\mathbf{X}), 0)}. \end{aligned}$$

we have $\mathbf{H}(\sigma^2 \mathbf{I})(\mathbf{I} - \mathbf{H}) = \mathbf{0}$,

\Rightarrow SSR is independent of SSE

$$\Rightarrow F(R) \sim F_{[r(\mathbf{X}), n - r(\mathbf{X}), \frac{\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta}}{2\sigma^2}]}$$

Thus under $H_0 : \boldsymbol{\beta}_0 = \mathbf{0}$,

$$F(R) \sim F_{[r(\mathbf{X}), n - r(\mathbf{X}), 0]}.$$

Next, suppose the model had no predictors in it but had simply been $y_i = c_0 + \varepsilon_i \quad i = 1, \dots, n$. Then, the least square estimate of c_0 is

$$\hat{c}_0 = (\mathbf{1}^\top \mathbf{1})^{-1} \mathbf{1}^\top \mathbf{Y} = \frac{1}{n} \mathbf{1}^\top \mathbf{Y} = \frac{1}{n} \sum_{i=1}^n y_i \equiv \bar{y}.$$

The SSM is defined to be the SSR for the linear model without x . Thus,

$$SSM = \hat{c}_0 \mathbf{1}^\top \mathbf{Y} = \frac{1}{n} \mathbf{Y}^\top \mathbf{1} \mathbf{1}^\top \mathbf{Y} = n \bar{y}^2.$$

$$\frac{SSM}{\sigma^2} = \mathbf{Y}^\top \left(\frac{1}{n\sigma^2} \mathbf{1} \mathbf{1}^\top \right) \mathbf{Y}$$

Define \mathbf{A} and $\mathbf{\Sigma}$ as

$$\mathbf{A} \triangleq \frac{1}{n\sigma^2} \mathbf{1} \mathbf{1}^\top$$

and

$$\mathbf{\Sigma} \triangleq \sigma^2 \mathbf{I}$$

Note that

$$\mathbf{A} \mathbf{\Sigma} = \left(\frac{1}{n\sigma^2} \mathbf{1} \mathbf{1}^\top \right) (\sigma^2 \mathbf{I}) = \frac{1}{n} \mathbf{1} \mathbf{1}^\top.$$

Check that

$$(\mathbf{A} \mathbf{\Sigma})^2 = \left(\frac{1}{n} \mathbf{1} \mathbf{1}^\top \right) \left(\frac{1}{n} \mathbf{1} \mathbf{1}^\top \right) = \frac{1}{n^2} \mathbf{1} (\mathbf{1}^\top \mathbf{1}) \mathbf{1}^\top = \frac{1}{n^2} \mathbf{1} (n) \mathbf{1}^\top = \frac{1}{n} \mathbf{1} \mathbf{1}^\top = \mathbf{A} \mathbf{\Sigma},$$

which is idempotent. It follows that

$$\frac{SSM}{\sigma^2} \sim \chi^2_{\left(1, \frac{(\mathbf{1}^\top \mathbf{x} \beta)^2}{2n\sigma^2}\right)}.$$

The *total sum of squares corrected for the mean* is

$$SST_m = SST - SSM = \mathbf{Y}^\top \mathbf{Y} - \mathbf{Y}^\top \mathbf{1} \mathbf{1}^\top \mathbf{Y} = \mathbf{Y} \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \right) \mathbf{Y}.$$

Then, we can show that

$$\frac{SST_m}{\sigma^2} \sim \chi^2_{\left(n-1, \frac{\beta^\top \mathbf{X}^\top \mathbf{X} \beta - \frac{1}{n} (\mathbf{1}^\top \mathbf{X} \beta)^2}{2\sigma^2}\right)}.$$

Similarly, the *sum of squares of regression corrected for the mean* is

$$SSR_m = SSR - SSM = \hat{\beta}^\top \mathbf{X}^\top \mathbf{Y} - \mathbf{Y}^\top \mathbf{1} \mathbf{1}^\top \mathbf{Y} = \hat{\mathbf{b}}^\top \mathbf{Z}^\top \mathbf{Y}.$$

SS explained by the covariables only.

centered and excluded the 1 columns

This is because, we can write

$$\begin{aligned}
 SSR &= \hat{\beta}^\top \mathbf{X}^\top \mathbf{Y} = (\hat{\beta}_0 \quad \hat{\mathbf{b}}) \begin{pmatrix} \mathbf{1}^\top \\ \mathbf{X}_1^\top \end{pmatrix} \mathbf{Y} \\
 &= \hat{\beta}_0 \mathbf{1}^\top \mathbf{Y} + \hat{\mathbf{b}}^\top \mathbf{X}_1^\top \mathbf{Y} \\
 &= \frac{1}{n} (\mathbf{Y}^\top \mathbf{1} - \bar{\mathbf{X}}^\top \hat{\mathbf{b}}) \mathbf{1}^\top \mathbf{Y} + \hat{\mathbf{b}}^\top \mathbf{X}_1^\top \mathbf{Y} \\
 &= \frac{1}{n} \mathbf{Y}^\top \mathbf{1} \mathbf{1}^\top \mathbf{Y} + \hat{\mathbf{b}}^\top (\mathbf{X}_1^\top - \bar{\mathbf{X}} \mathbf{1}^\top) \mathbf{Y} \\
 &= \frac{1}{n} \mathbf{Y}^\top \mathbf{1} \mathbf{1}^\top \mathbf{Y} + \hat{\mathbf{b}}^\top \mathbf{Z}^\top \mathbf{Y} \\
 &\equiv \text{SSM} + \text{SSR}_m.
 \end{aligned}$$

Note that $\text{SSR}_m = \hat{\mathbf{b}}^\top \mathbf{Z}^\top \mathbf{Y} = \hat{\mathbf{b}}^\top (\mathbf{Z}^\top \mathbf{Z}) \hat{\mathbf{b}}$ and

$$\hat{\mathbf{b}} = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{Y} \sim N[\mathbf{b}, (\mathbf{Z}^\top \mathbf{Z})^{-1} \sigma^2],$$

thus

$$\frac{\text{SSR}_m}{\sigma^2} \sim \chi^2_{(k, \frac{\mathbf{b}^\top (\mathbf{Z}^\top \mathbf{Z}) \mathbf{b}}{2\sigma^2})}.$$

Therefore, the analysis of variance corrected for the mean for fitting linear regression can be summarized in the following table:

Table 2 ANOVA table corrected for the mean.

Source	df	SS	MS	F-statistics
Mean	1	$\frac{1}{n} \mathbf{Y}^\top \mathbf{1} \mathbf{1}^\top \mathbf{Y}$	$\frac{\text{SSM}}{1}$	$F(M) = \frac{\text{SSM}}{\text{MSE}}$
Regression	$r(\mathbf{X}) - 1$	$\hat{\mathbf{b}}^\top \mathbf{Z}^\top \mathbf{Y}$	$\frac{\text{SSR}_m}{r(\mathbf{X}) - 1}$	$F(R) = \frac{\text{SSR}_m}{\text{MSE}}$
Error	$n - r(\mathbf{X})$	$\mathbf{Y}^\top \mathbf{Y} - \hat{\beta}^\top \mathbf{X}^\top \mathbf{Y}$	$\frac{\text{SSE}}{n - r(\mathbf{X})}$	
Total	n	$\mathbf{Y}^\top \mathbf{Y}$		

Note:

(1)

$\mathbf{B} \Sigma \mathbf{A}$

Since $\mathbf{1} \mathbf{1}^\top (\mathbf{I} \sigma^2) (\mathbf{I} - \mathbf{H}) = 0$

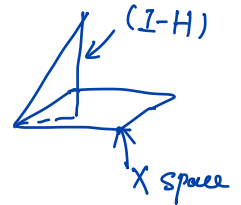
\Rightarrow SSM is independent of SSE

$$\Rightarrow F(M) \sim F'_{[1, n-r(\mathbf{X}), \frac{(\mathbf{1}^\top \mathbf{X} \boldsymbol{\beta})^2}{2n\sigma^2}]}$$

Under $H_0 : \mathbf{1}^\top \mathbf{X} \boldsymbol{\beta} = 0$ ($H_0 : E(\bar{y}) = 0$)

$$F(M) \sim F'_{[1, n-r(\mathbf{X}), 0]}.$$

$\mathbf{1}$ is in the space of \mathbf{X}
then $\mathbf{1} \perp \mathbf{I} - \mathbf{H}$.



(2)

$$\text{Since } \text{SSR}_m = \hat{\mathbf{b}}^\top \mathbf{Z}^\top \mathbf{Y}$$

$$= \mathbf{Y}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{Y}$$

$$\text{and } \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top (\mathbf{I} \sigma^2) (\mathbf{I} - \mathbf{H}) = 0,$$

$$\Rightarrow \text{SSR}_m \text{ is independent of SSE}$$

$$\Rightarrow F(R_m) \sim F'_{[r(\mathbf{X})-1, n-r(\mathbf{X}), \frac{\mathbf{b}^\top (\mathbf{Z}^\top \mathbf{Z}) \mathbf{b}}{2\sigma^2}]} \quad \text{Construct a statistic that contain parameter to be tested.}$$

Under $H_0 : \mathbf{b} = 0$,

$$F(R_m) \sim F_{[r(\mathbf{X})-1, n-r(\mathbf{X}), 0]}.$$

(3)

$$\text{SSM} = \frac{1}{n} \mathbf{Y}^\top \mathbf{1} \mathbf{1}^\top \mathbf{Y}$$

$$\text{SSR}_m = \mathbf{Y}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{Y}$$

$$\text{Since } \mathbf{1} \mathbf{1}^\top (\mathbf{I} \sigma^2) \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top = 0,$$

(this is because $\mathbf{Z} = \mathbf{X}_1 - \mathbf{1} \bar{\mathbf{X}}^\top \Rightarrow \mathbf{1}^\top \mathbf{Z} = n(\bar{\mathbf{X}}^\top - \bar{\mathbf{X}}^\top) = 0$)

Thus, SSM and SSR_m are independent.

3.5.4 General linear hypothesis

The literature of linear models abounds with discussions of different kinds of hypotheses that can be of interest in widely different fields of application. Four hypothesis of particular interest are:

1. $H_0 : \beta = \mathbf{0}$, the hypothesis that all of the elements of β are zero;
2. $H_0 : \beta = \beta_0$, the hypothesis that $\beta_i = \beta_{i0}$ for $i = 1, 2, \dots, k$, that is, that each β_i is equal to some specified value β_{i0} ;
3. $H_0 : \lambda^\top \beta = m$, that some linear combination of the elements of β equals a specified constant;
4. $H_0 : \beta_q = \mathbf{0}$, that some of β_i 's, q of them where $q < k$ are zero.

We will show later that all of the linear hypothesis above and others are special cases of a general procedure even though the calculation of the F -statistics may appear to differ from one hypothesis to another. The general hypothesis we consider is

$$H_0 : K^\top \beta = m, \quad \text{Linear}$$

where β is the $(k+1)$ -dimensional vector of parameters of the model, K^\top is any matrix of size $s \times (k+1)$ and m is a $s \times 1$ vector of specified constants. There is ONLY one restriction on K^\top ; K^\top is assumed to be of full row rank. This means that the linear function of β must be linearly independent. The hypothesis being tested must be made up of linearly independent functions of β and must contain no functions that are linear functions of others therein.

If counting the number of non-zero coeff in β , H_0 will not be linear H_0 .

(a) Testing linear hypothesis

Recall that in parametric linear model with normal error assumption, $Y \sim N(X\beta, \sigma^2 I)$. Then

$$\hat{\beta} = (X^\top X)^{-1} X^\top Y, \quad \hat{\beta} \sim N(\beta, (X^\top X)^{-1} \sigma^2).$$

$\Rightarrow K^\top \hat{\beta} - m \sim N[K^\top \beta - m, K^\top (X^\top X)^{-1} K \sigma^2]$. Next, since

$$\frac{(K^\top (X^\top X)^{-1} K)^{-1}}{\sigma^2} (K^\top (X^\top X)^{-1} K) \sigma^2 = I \quad (\text{symmetric and idempotent}),$$

therefore, $(K^\top (X^\top X)^{-1} K)^{-1}$ is symmetric (why?). Thus, we define

$$Q = (K^\top \hat{\beta} - m)^\top [K^\top (X^\top X)^{-1} K]^{-1} (K^\top \hat{\beta} - m).$$

Then,

$$\frac{Q}{\sigma^2} \sim \chi^2_{s, \frac{1}{2\sigma^2} (K^\top \beta - m)^\top [K^\top (X^\top X)^{-1} K]^{-1} (K^\top \beta - m)}.$$

under $H_0 : K^\top \beta = m \Rightarrow 0$.

unknown.

Next, we intend to show that Q and SSE are independent.

$$\begin{aligned} Q &= (\mathbf{K}^\top \hat{\beta} - \mathbf{m})^\top [\mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K}]^{-1} (\mathbf{K}^\top \hat{\beta} - \mathbf{m}) \\ &= (\mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} - \mathbf{m})^\top [\mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K}]^{-1} (\mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} - \mathbf{m}). \end{aligned}$$

But

$$\mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} - \mathbf{m} = \mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{Y} - \mathbf{X} \mathbf{K} (\mathbf{K}^\top \mathbf{K})^{-1} \mathbf{m}).$$

a shift of \mathbf{Y} . $\therefore \mathbf{b}$.

Therefore,

$$Q = \underbrace{(\mathbf{Y} - \mathbf{X} \mathbf{K} (\mathbf{K}^\top \mathbf{K})^{-1} \mathbf{m})^\top}_{\mathbf{d}^\top} \underbrace{\mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K} (\mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K})^{-1} \mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top}_{\mathbf{B}} \underbrace{(\mathbf{Y} - \mathbf{X} \mathbf{K} (\mathbf{K}^\top \mathbf{K})^{-1} \mathbf{m})}_{\mathbf{d}^\top}.$$

Note that if we let

$$\mathbf{d} = \mathbf{Y} - \mathbf{X} \mathbf{K} (\mathbf{K}^\top \mathbf{K})^{-1} \mathbf{m},$$

$$\begin{aligned} Q &= \mathbf{d}^\top [\mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K} (\mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K})^{-1} \mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top] \mathbf{d} \\ &\equiv \mathbf{d}^\top \mathbf{B} \mathbf{d}, \end{aligned}$$

and

$$\mathbf{d} \sim N(\mathbf{X} \beta - \mathbf{X} \mathbf{K} (\mathbf{K}^\top \mathbf{K})^{-1} \mathbf{m}, \sigma^2 \mathbf{I}).$$

On the other hand,

$$\begin{aligned} SSE &= \mathbf{Y}^\top (\mathbf{I} - \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbf{Y} \\ &= (\mathbf{Y} - \mathbf{X} \mathbf{K} (\mathbf{K}^\top \mathbf{K})^{-1} \mathbf{m})^\top (\mathbf{I} - \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) (\mathbf{Y} - \mathbf{X} \mathbf{K} (\mathbf{K}^\top \mathbf{K})^{-1} \mathbf{m}) \\ &= \mathbf{d}^\top \underbrace{(\mathbf{I} - \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)}_{\mathbf{A}} \mathbf{d} \\ &\equiv \mathbf{d}^\top \mathbf{A} \mathbf{d} \end{aligned}$$

this is because

$$\begin{aligned} \mathbf{X}^\top (\mathbf{I} - \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) &= \mathbf{0}, \\ (\mathbf{I} - \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbf{X} &= \mathbf{0}. \end{aligned}$$

Since

$$\mathbf{A} \mathbf{\Sigma} \mathbf{B} = (\mathbf{I} - \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \sigma^2 \mathbf{I} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K} [\mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K}]^{-1} \mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{0},$$

We have shown Q and SSE are independent. Finally, the test statistics

$$\begin{aligned} F(H) &= \frac{Q/s}{SSE/[N - r(\mathbf{X})]} \\ &= \frac{Q}{s \hat{\sigma}^2} \\ &\sim F'_{s, N-r(\mathbf{X}), \frac{1}{2\sigma^2} (\mathbf{K}^\top \beta - \mathbf{m})^\top [\mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K}]^{-1} (\mathbf{K}^\top \beta - \mathbf{m})}, \end{aligned}$$

where $\hat{\sigma}^2 = \frac{SSE}{N-r(\mathbf{X})}$ is the unbiased estimator of σ^2 .

Under $H_0 : \mathbf{K}^\top \boldsymbol{\beta} = \mathbf{m}$,

$$F(\mathbf{H}) \sim F'_{(s, N-r(\mathbf{X}))}.$$

This F -statistic is thus to be used to test the general linear hypothesis. This F -test can be shown to stem from the likelihood ratio test.

(b) Estimation of $\boldsymbol{\beta}$ under the null hypothesis

A natural question to ask when considering the null hypothesis $H_0 : \mathbf{K}^\top \boldsymbol{\beta} = \mathbf{m}$ is “What is the estimator of $\boldsymbol{\beta}$ under the null hypothesis?” This might be especially pertinent following non-rejection of the hypothesis by the preceding F test. The desired estimator, denoted by $\tilde{\boldsymbol{\beta}}$, is readily obtainable using constrained least squares. Thus, when H_0 is true, $\tilde{\boldsymbol{\beta}}$ is derived so as to minimize the least squares objective function subject to the constraint $\mathbf{K}^\top \boldsymbol{\beta} = \mathbf{m}$. With $2\boldsymbol{\theta}$ as a vector of Lagrange multipliers, we minimize

$$(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + 2\boldsymbol{\theta}^\top (\mathbf{K}^\top \boldsymbol{\beta} - \mathbf{m})$$

with respect to $\boldsymbol{\beta}$ and $\boldsymbol{\theta}$. Note that $2\boldsymbol{\theta}$ is a vector of Lagrange multipliers. After the minimization, it can be easily obtained that

$$\begin{aligned} \tilde{\boldsymbol{\beta}} &= (\mathbf{X}^\top \mathbf{X})^{-1} [\mathbf{X}^\top \mathbf{Y} - \mathbf{K}(\mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K})^{-1} (\mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} - \mathbf{m})] \\ &= \hat{\boldsymbol{\beta}} - \underbrace{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K}(\mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K})^{-1} (\mathbf{K}^\top \hat{\boldsymbol{\beta}} - \mathbf{m})}_{\Delta}, \end{aligned}$$

where $\hat{\boldsymbol{\beta}}$ is the ordinary LS estimate. Note that

$$\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K}(\mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K})^{-1} (\mathbf{K}^\top \hat{\boldsymbol{\beta}} - \mathbf{m}). \quad (*)$$

We have estimated $\boldsymbol{\beta}$ under H_0 . We next show that the corresponding residual sum of squares is $SSE + Q$ where Q is the numerator sum of squares of the F -statistic used in testing the general linear hypothesis.

Without the null hypothesis (the **full model**),

$$SSE = (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}).$$

Under the null hypothesis (the **reduced model**), the sum of squares of residual is

$$\begin{aligned} SSE_{H_0} &= (\mathbf{Y} - \mathbf{X}\tilde{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\tilde{\boldsymbol{\beta}}) \\ &= [\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{X}\tilde{\boldsymbol{\beta}}]^\top [\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{X}\tilde{\boldsymbol{\beta}}] \\ &= [\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{X}(\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}})]^\top [\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{X}(\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}})] \\ &= (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) + (\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}})^\top \mathbf{X}^\top \mathbf{X} (\hat{\boldsymbol{\beta}} - \tilde{\boldsymbol{\beta}}), \end{aligned}$$

cross term.

where the last equation is due to $(\tilde{\beta} - \beta_0)^\top \mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\hat{\beta}) = 0$. From (*),

$$\begin{aligned}
 SSE_{H_0} &= SSE + (\mathbf{K}^\top \hat{\beta} - \mathbf{m})^\top [\mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K}]^{-1} \\
 &\quad \times \mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K} [\mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K}]^{-1} (\mathbf{K}^\top \hat{\beta} - \mathbf{m}) \\
 &= SSE + \underline{(\mathbf{K}^\top \hat{\beta} - \mathbf{m})^\top [\mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K}]^{-1} (\mathbf{K}^\top \hat{\beta} - \mathbf{m})} \\
 &= SSE + Q \quad \text{the non-centric } F \text{ is location parameter.} \\
 &\geq SSE.
 \end{aligned}$$

(c) Four Common Hypothesis.

In this section, we illustrate the expressions for $F(H)$ and $\tilde{\beta}$ for four commonly occurring hypothesis.

1. $H_0 : \beta = \mathbf{0}$, (therefore, $\mathbf{K}^\top = \mathbf{I}$, $s = k + 1$, $\mathbf{m} = \mathbf{0}$)

(1a)

$$\begin{aligned}
 Q &= \mathbf{Y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} ((\mathbf{X}^\top \mathbf{X})^{-1})^{-1} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \\
 &= \mathbf{Y}^\top \mathbf{X} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \\
 &= SSR.
 \end{aligned}$$

(1b)

$$F(\mathbf{H}) = \frac{SSR/r(\mathbf{X})}{SSE/[n - r(\mathbf{X})]}.$$

(1c) Under the null hypothesis,

$$F(\mathbf{H}) \sim F_{(r(\mathbf{X}), n - r(\mathbf{X}))}.$$

(1d)

$$\begin{aligned}
 \tilde{\beta} &= \hat{\beta} - (\mathbf{X}^\top \mathbf{X})^{-1} ((\mathbf{X}^\top \mathbf{X})^{-1})^{-1} (\hat{\beta}) \\
 &= \hat{\beta} - \hat{\beta} \\
 &= \mathbf{0}.
 \end{aligned}$$

2. $H_0 : \beta = \beta_c$ (therefore $\mathbf{K}^\top = \mathbf{I}$, $s = k + 1$, $\mathbf{m} = \beta_c$)

(2a)

$$\begin{aligned}
 Q &= (\mathbf{Y} - \mathbf{X}\beta_c)^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\beta_c) \\
 &= \mathbf{d}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{d}.
 \end{aligned}$$

(2b)

$$F(\mathbf{H}) = \frac{Q}{(k + 1)\hat{\sigma}^2}.$$

(2c) Under the null hypothesis,

$$F(\mathbf{H}) \sim F_{(k+1, n-(k+1))}.$$

$$(2d) \quad \tilde{\beta} = \hat{\beta} - (\hat{\beta} - \beta_c) = \beta_c.$$

3. $H_0 : \lambda^\top \beta = m$ (therefore $\mathbf{K}^\top = \lambda^\top, s = 1, \mathbf{m} = m$)

(3a)

$$\begin{aligned} Q &= (\lambda^\top \hat{\beta} - m)^\top [\lambda^\top (\mathbf{X}^\top \mathbf{X})^{-1} \lambda]^{-1} (\lambda^\top \hat{\beta} - m) \\ &= (\lambda^\top \hat{\beta} - m)^2 / \lambda^\top (\mathbf{X}^\top \mathbf{X})^{-1} \lambda. \end{aligned}$$

(3b)

$$F(H) = \frac{Q}{\hat{\sigma}^2}.$$

(3c) Under the null hypothesis,

$$F(H) \sim F_{(1, n-r(\mathbf{X}))}.$$

Note that

$$\sqrt{F(H)} = \frac{\sqrt{Q}}{\hat{\sigma}} \sim t_{n-r(\mathbf{X})}.$$

(3d)

$$\begin{aligned} \tilde{\beta} &= \hat{\beta} - (\mathbf{X}^\top \mathbf{X})^{-1} \lambda [\lambda^\top (\mathbf{X}^\top \mathbf{X})^{-1} \lambda]^{-1} (\lambda^\top \hat{\beta} - m) \\ &= \hat{\beta} - \frac{(\lambda^\top \hat{\beta} - m)}{\lambda^\top (\mathbf{X}^\top \mathbf{X})^{-1} \lambda} (\mathbf{X}^\top \mathbf{X})^{-1} \lambda. \end{aligned}$$

Note that

$$\lambda^\top \hat{\beta} - m \sim N(\lambda^\top \beta - m, \lambda^\top (\mathbf{X}^\top \mathbf{X})^{-1} \lambda \sigma^2).$$

~~3~~ $H_0 : \beta_q = 0$ where $\beta = \begin{bmatrix} \beta_q \\ \beta_p \end{bmatrix}$

In this case, $\mathbf{K}^\top = [\mathbf{I}_q \quad \mathbf{0}]$, $p = k + 1 - q$, $\mathbf{m} = \mathbf{0}$, $s = q$.

Partition $\hat{\beta}$ and $(\mathbf{X}^\top \mathbf{X})^{-1}$ such that

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_q \\ \hat{\beta}_p \end{bmatrix} \text{ and } (\mathbf{X}^\top \mathbf{X})^{-1} = \begin{bmatrix} \mathbf{T}_{qq} & \mathbf{T}_{qp} \\ \mathbf{T}_{pq} & \mathbf{T}_{pp} \end{bmatrix}.$$

Note that $\mathbf{K}^\top \hat{\beta} = \hat{\beta}_q$ and $[\mathbf{K}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{K}]^{-1} = \mathbf{T}_{qq}^{-1}$.

$$(4a) \quad Q = \hat{\beta}_q^\top \mathbf{T}_{qq}^{-1} \hat{\beta}_q.$$

$$(4b) \quad F(H) = \frac{Q}{\hat{\sigma}^2}.$$

(4c) Under H_0 , $F(H) \sim F_{(q, n-r(\mathbf{X}))}$.

$$\tilde{\beta} = \hat{\beta}_p - \mathbf{T}_{pq} \mathbf{T}_{qq}^{-1} \hat{\beta}_q.$$

$$(4d) \quad \tilde{\beta} = \begin{bmatrix} \mathbf{0}_q \\ \hat{\beta}_p - \mathbf{T}_{pq} \mathbf{T}_{qq}^{-1} \hat{\beta}_q \end{bmatrix}.$$

1. LS is not robust since it is for mean regress excluding the heavy-tail error distribution

38

3.6* Quantile regression and the least absolute deviation (see Koenker 2008).

Much of applied statistics may be viewed as an elaboration of the linear regression model and associated estimation methods of least squares. As a competitive alternative to the least squares regression, quantile regression (Koenker and Bassett, 1978, Econometrica) is a widely used statistical tool amongst researchers and practitioners for the modeling and inference of conditional quantile function. By allowing varying nature across different points of conditional distribution, the quantile regression model is able to accommodate heterogeneous effect in the population and is robust to the outliers of the response.

3.6.1 Quantile.

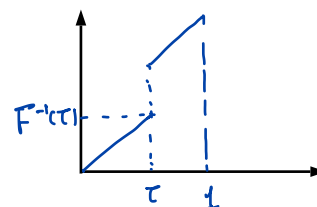
population level definition (sample level definition)?

Any real-valued random variable Y may be characterized by its (right-continuous) distribution function

$$F(y) = P(Y \leq y),$$

whereas for any $0 < \tau < 1$,

$$F^{-1}(\tau) = \inf\{y : F(y) \geq \tau\}$$



is called the τ th quantile of Y . The median $F^{-1}(1/2)$ plays the central role.

say, median regress

3.6.2 Linear quantile regression.

Given a covariate vector $X \in R^{p+1}$, the τ th conditional quantile function of $Y \in R$ is modeled as

$$Q_{Y|X}(\tau) = X^\top \beta_\tau, \quad (1)$$

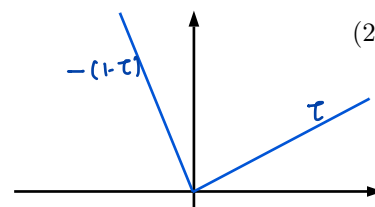
for certain specific $\tau \in (0, 1)$ of interest, and β_τ is $(p+1)$ -vector usually including an intercept. For model (1) with complete data, a classical estimate of β_τ , denoted as $\hat{\beta}_\tau^c$, is obtained by minimizing

$X^\top \hat{\beta}_\tau$ is the τ -th quantile of Y .

$$R(\beta_\tau) \equiv \sum_{i=1}^n \rho_\tau(y_i - x_i^\top \beta_\tau)$$

over β_τ , where (y_i, x_i) are iid copies of (Y, X) , $i = 1, \dots, n$,

$$\rho_\tau(u) = u(\tau - I(u < 0))$$



is the quantile check loss function and $I(\cdot)$ is the indicator function. The beauty of this method rests in its simplicity and ease of inference. Recall that in least square estimator has a closed-form solution by solving the so-called normal equation. In quantile regression, we proceed similarly, but we need to exercise some caution about the differentiation step. The objective function $R(\beta_\tau)$ is piecewise linear and continuous. It is differentiable except at the points at which one or more residuals, $y_i - x_i^\top \beta_\tau$, are zero. At such points, $R(\beta_\tau)$ has directional derivatives in all directions, depending, however, on the direction of evaluation. The

$$\begin{aligned} \min_{\hat{x}} \mathbb{E}_x \rho_\tau(X - \hat{x}) &= \min_{\hat{x}} \left\{ \int_{-\infty}^{\hat{x}} (1-\tau)(x - \hat{x}) dF(x) + \int_{\hat{x}}^{\infty} \tau(x - \hat{x}) dF(x) \right\} \\ \frac{\partial}{\partial \hat{x}} \downarrow &= (1-\tau) \int_{-\infty}^{\hat{x}} (-1) dF(x) \quad \frac{\partial}{\partial \hat{x}} \downarrow = +\tau \int_{\hat{x}}^{\infty} (-1) dF(x) \\ &= -(1-\tau) \int_{-\infty}^{\hat{x}} dF(x) - \tau \int_{\hat{x}}^{\infty} dF(x) = 0 \\ &\Rightarrow (1-\tau) F(\hat{x}) - \tau(1 - F(\hat{x})) = 0 \\ &\Rightarrow F(\hat{x}) = \tau \Rightarrow \hat{x} = \inf\{x : F(x) = \tau\}. \end{aligned}$$

directional derivative of R in direction w is given by

$$\begin{aligned}
 \nabla R(\beta_\tau, w) &\equiv \frac{d}{dt} R(\beta_\tau + tw)|_{t=0} \\
 &= \frac{d}{dt} \sum_{i=1}^n [y_i - x_i^\top (\beta_\tau + tw)] \{\tau - I[y_i - x_i^\top (\beta_\tau + tw) < 0]\}|_{t=0} \\
 &= - \sum_{i=1}^n \psi_\tau(y_i - x_i^\top \beta_\tau, -x_i^\top w) x_i^\top w,
 \end{aligned} \tag{3}$$

where

$$\psi_\tau(u, v) = \begin{cases} \tau - I(u < 0) & \text{if } u \neq 0; \\ \tau - I(v < 0) & \text{if } u = 0. \end{cases}$$

If, at a point $\hat{\beta}_\tau$, the directional derivatives are all nonnegative (that is, $\nabla R(\beta_\tau) \geq 0$ for all $w \in \mathbb{R}^p$ with $\|w\| = 1$), then $\hat{\beta}_\tau$ minimizes $R(\beta_\tau)$. This is a natural generalization of simply setting $\nabla R(\beta_\tau) = 0$ when $R(\cdot)$ is smooth. It simply requires that the function is increasing as we move away from the point $\hat{\beta}_\tau$ regardless of the direction in which we decide to move.

The computation of $\hat{\beta}_\tau^c$ is straightforward with the help of linear programming. More details of the computation of quantile regression can be found in Koenker (2008). There is a vast literature on the estimation and inference for one or several percentile levels for model (1); see Koenker (2003, 2008), He (1997, 2003), among many others.

* 3.6.3 Asymptotic properties of quantile regression.

The conditional distribution function of Y_i s will be written as

$$P(Y_i < y | x_i) = F_{Y_i}(y | x_i) = F_i(y)$$

and so

$$Q_{Y_i}(\tau | x_i) = F_{Y_i}^{-1}(\tau | x_i) \equiv \xi_i(\tau).$$

We will employ the following regularity conditions to explore the asymptotic behavior of the estimator $\hat{\beta}_n(\tau)$, the minimizer of

$$\sum_{i=1}^n \rho_\tau(y_i - x_i^\top \beta_\tau).$$

- **Condition (A1).** The distribution functions $\{F_i\}$ are absolutely continuous, with continuous densities $f_i(\xi)$ uniformly bounded away from 0 and ∞ at the points $\xi_i(\tau), i = 1, 2, \dots$
- **Condition (A2).** There exists positive definite matrices D_0 and $D_1(\tau)$ such at

$$(i) \lim_{n \rightarrow \infty} n^{-1} \sum x_i x_i^\top = D_0;$$

$$\begin{aligned}
 \hat{\beta}_{LAD} &= \min_{\beta} \frac{1}{n} \sum_{i=1}^n |y_i - x_i^\top \beta| \\
 Q_{Y|X}(\tau=0.5) &= X\beta \Leftrightarrow Y = X\beta + \varepsilon, \text{ med}(\varepsilon|X) = 0. \\
 \text{without assumption that } \varepsilon \perp\!\!\!\perp X \\
 \text{sgn}(\cdot) &= \mathcal{S}(\cdot).
 \end{aligned}$$

- (ii) $\lim_{n \rightarrow \infty} n^{-1} \sum f_i(\xi_i(\tau)) x_i x_i^\top = D_1(\tau);$
- (iii) $\max_{i=1, \dots, n} \|x_i\| / \sqrt{n} \rightarrow 0.$

Conditions A2(i) and A2(iii) are familiar throughout the literature on M-estimators for regression models; some variant of them is necessary to ensure that a **Lindeberg condition is satisfied**. Condition A2(ii) is really a matter of notational convenience and could be deduced from Condition A2(i) and a slightly strengthened version of Condition A1.

Theorem 4.1 *Under Conditions A1 and A2,*

$$\sqrt{n}\{\hat{\beta}_n(\tau) - \beta_0(\tau)\} \rightarrow N(0, \tau(1-\tau)D_1^{-1}D_0D_1^{-1})$$

in distribution as $n \rightarrow \infty$. In particular, in the iid error model (that is, $Y_i = x_i^\top \beta + u_i$, $i = 1, \dots, n$ with iid errors u_i have a common distribution function F),

$$\sqrt{n}\{\hat{\beta}_n(\tau) - \beta_0(\tau)\} \rightarrow N(0, \omega^2 D_0^{-1}),$$

in distribution as $n \rightarrow \infty$, where $\omega^2 = \tau(1-\tau)/f_i^2(\xi_i(\tau))$.

Remark: The asymptotic properties above were established under fixed design, which can be extended to handle random design without further difficulties.

3.6.4 The least absolute deviation.

In the iid error model (random design) $Y_i = x_i^\top \beta + u_i$, $i = 1, \dots, n$, where u_i have a common distribution function F , the LAD is a special case of quantile regression when $\tau = 0.5$. Its estimate $\hat{\beta}_{LAD}$ is defined to be the minimizer of

$$\min_{\beta} \sum_{i=1}^n |Y_i - X_i^\top \beta|.$$

Hence, the following theorem can be established in a straightforward fashion under regularity conditions.

Theorem 4.2 *Under Conditions A1 and A2,*

$$\sqrt{n}\{\hat{\beta}_{LAD} - \beta_0\} \rightarrow N(0, \omega^2 V^{-1})$$

in distribution as $n \rightarrow \infty$, where $\omega^2 = 1/4f^2(0)$ and $V = E(XX^\top)$.