

Confidence Intervals and Sets

Throughout we adopt the normal-error model, and wish to say some things about the construction of confidence intervals [and sets] for the parameters $\beta_0, \dots, \beta_{p-1}$.

1. Confidence intervals for one parameter

Suppose we want a confidence intervals for the i th parameter β_i , where $1 \leq i \leq p$ is fixed. Recall that

$$\frac{\hat{\beta}_i - \beta_i}{\sqrt{\sigma^2 [(X'X)^{-1}]_{i,i}}} \sim N(0, 1),$$

and

$$S^2 := \frac{\text{RSS}}{n-p} = \frac{1}{n-p} \sum_{i=1}^n \left(y_i - (X\hat{\beta})_i \right)^2 \sim \frac{\sigma^2 \chi_{n-p}^2}{n-p},$$

and is independent of $\hat{\beta}$. Therefore,

$$\frac{\hat{\beta}_i - \beta_i}{\sqrt{S^2 [(X'X)^{-1}]_{i,i}}} \sim \frac{N(0, 1)}{\sqrt{\chi_{n-p}^2/(n-p)}} = t_{n-p},$$

where the normal and χ^2 random variables are independent. Therefore,

$$\hat{\beta}_i \pm t_{n-p}^{(\alpha/2)} S \sqrt{[(X'X)^{-1}]_{i,i}}$$

is a $(1 - \alpha) \times 100\%$ confidence interval for β_i . This yields a complete analysis of confidence intervals for a univariate parameter; these confidence intervals can also be used for testing, of course.

2. Confidence ellipsoids

The situation is more interesting if we wish to say something about more than one parameter at the same time. For example, suppose we want to know about (β_0, β_1) in the *overly-simplified case that $\sigma = 1$* . The general philosophy of confidence intervals [for univariate parameters] suggests that we look for a random set Ω such that

$$\mathbb{P} \left\{ \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} - \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} \in \Omega \right\} = 1 - \alpha.$$

And we know that

$$\begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} \sim N \left(\begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}, \Sigma \right),$$

where Σ is a 2×2 matrix with

$$\Sigma_{i,j} = \left[(X'X)^{-1} \right]_{i,j}.$$

Here is a possible method: Let

$$\hat{\gamma} := \begin{pmatrix} \hat{\beta}_1 - \beta_1 \\ \hat{\beta}_2 - \beta_2 \end{pmatrix}, \quad \text{so that} \quad \hat{\gamma} \sim N_2(\mathbf{0}, \Sigma).$$

Also, recall that

$$\hat{\gamma}' \Sigma^{-1} \hat{\gamma} \sim \chi_2^2.$$

Therefore, one natural choice for Ω is

$$\Omega := \left\{ \mathbf{x} \in \mathbf{R}^2 : \mathbf{x}' \Sigma^{-1} \mathbf{x} \leq \chi_2^2(\alpha/2) \right\}.$$

What does Ω look like? In order to answer this, let us apply the spectral theorem:

$$\Sigma = \mathbf{P} \mathbf{D} \mathbf{P}' \implies \Sigma^{-1} = \mathbf{P} \mathbf{D}^{-1} \mathbf{P}'.$$

Then we can represent Ω as follows:

$$\begin{aligned} \Omega &= \left\{ \mathbf{x} \in \mathbf{R}^2 : \mathbf{x}' \mathbf{P} \mathbf{D}^{-1} \mathbf{P}' \mathbf{x} \leq \chi_2^2(\alpha/2) \right\} \\ &= \left\{ \mathbf{x} \in \mathbf{R}^2 : (\mathbf{P}' \mathbf{x}) \mathbf{D}^{-1} (\mathbf{P}' \mathbf{x}) \leq \chi_2^2(\alpha/2) \right\} \\ &= \left\{ \mathbf{P} \mathbf{y} \in \mathbf{R}^2 : \mathbf{y}' \mathbf{D}^{-1} \mathbf{y} \leq \chi_2^2(\alpha/2) \right\} \\ &= \left\{ \mathbf{P} \mathbf{y} \in \mathbf{R}^2 : \frac{y_1^2}{\lambda_1} + \frac{y_2^2}{\lambda_2} \leq \chi_2^2(\alpha/2) \right\}. \end{aligned}$$

Consider

$$\mathcal{E} := \left\{ \mathbf{y} \in \mathbf{R}^2 : \frac{y_1^2}{\lambda_1} + \frac{y_2^2}{\lambda_2} \leq \chi_2^2(\alpha/2) \right\}. \quad (1)$$

This is the interior of an ellipsoid, and $\Omega = \mathbf{P}\mathcal{E}$ is the image of the ellipsoid under the “linear orthogonal map” \mathbf{P} . Such sets are called “generalized ellipsoids,” and we have found a $(1 - \alpha) \times 100\%$ confidence [generalized] ellipsoid for $(\beta_1, \beta_2)'$.

The preceding can be generalized to any number of the parameters $\beta_{i_1}, \dots, \beta_{i_k}$, but is hard to work with, as the geometry of Ω can be complicated [particularly if $k \gg 2$]. Therefore, instead we might wish to look for approximate confidence sets that are easier to work with. Before we move on though, let me mention that if you want to know whether or not $H_0 : \beta_1 = \beta_{1,0}, \beta_2 = \beta_{2,0}$, then we can use these confidence bounds fairly easily, since it is not hard to check whether or not $(\beta_{1,0}, \beta_{2,0})'$ is in Ω : You simply compute the scalar quantity

$$(\beta_{1,0}, \beta_{2,0}) \mathbf{\Sigma}^{-1} \begin{pmatrix} \beta_{1,0} \\ \beta_{2,0} \end{pmatrix},$$

and check to see if it is $\leq \chi_2^2(\alpha/2)$! But if you really need to imagine or see the confidence set[s], then this exact method can be unwieldy [particularly in higher dimensions than 2].

3. Bonferonni bounds

Our approximate confidence intervals are based on a fact from general probability theory.

Proposition 1 (Bonferonni’s inequality). *Let E_1, \dots, E_k be k events. Then,*

$$\mathbf{P}(E_1 \cap \dots \cap E_k) \geq 1 - \sum_{j=1}^k \mathbf{P}(E_j^c) = 1 - \sum_{j=1}^k (1 - \mathbf{P}(E_j)).$$

Proof. The event $E_1^c \cup \dots \cup E_k^c$ is the complement of $E_1 \cap \dots \cap E_k$. Therefore,

$$\mathbf{P}(E_1 \cap \dots \cap E_k) = 1 - \mathbf{P}(E_1^c \cup \dots \cup E_k^c),$$

and this is $\geq 1 - \sum_{j=1}^k \mathbf{P}(E_j^c)$ because the probability of a union is at most the sum of the individual probabilities. \square

Here is how we can use Bonferonni’s inequality. Define

$$C_j := \hat{\beta}_j \pm t_{n-p}^{(\alpha/4)} S \sqrt{[(X'X)^{-1}]_{j,j}} \quad (j = 1, 2).$$

We have seen already that

$$P\{\beta_j \in C_j\} = 1 - \frac{\alpha}{2}.$$

[This is why we used $\alpha/4$ in the definition of C_j .] Therefore, Bonferonni's inequality implies that

$$P\{\beta_1 \in C_1, \beta_2 \in C_2\} \geq 1 - \left[\frac{\alpha}{2} + \frac{\alpha}{2} \right] = 1 - \alpha.$$

In other words, (C_1, C_2) forms a “conservative” $(1 - \alpha) \times 100\%$ simultaneous confidence interval for (β_1, β_2) . This method becomes very inaccurate quickly as the number of parameters of interest grows. For instance, if you want Bonferonni confidence sets for $(\beta_1, \beta_2, \beta_3)$, then you need to use individual confidence intervals with confidence level $1 - (\alpha/3)$ each. And for k parameters you need individual confidence level $1 - (\alpha/k)$, which can yield bad performance when k is large.

This method is easy to implement, but usually **very** conservative.¹

4. Scheffé's simultaneous conservative confidence bounds

There is a lovely method, due to Scheffé, that works in a similar fashion to the Bonferonni method; but has also the advantage of being often [far] less conservative! [We are now talking strictly about our linear model.] The starting point of this discussion is a general fact from matrix analysis.

Proposition 2 (The Rayleigh—Ritz inequality). *If $L_{p \times p}$ is positive definite, then for every p -vector \mathbf{b} ,*

$$\mathbf{b}'L^{-1}\mathbf{b} = \max_{\mathbf{h} \neq 0} \left[\frac{(\mathbf{h}'\mathbf{b})^2}{\mathbf{h}'L\mathbf{h}} \right].$$

The preceding is called an inequality because it says that

$$\mathbf{b}'L^{-1}\mathbf{b} \leq \frac{(\mathbf{h}'\mathbf{b})^2}{\mathbf{h}'L\mathbf{h}} \quad \text{for all } \mathbf{h}_{p \times 1},$$

and it also tells us that the inequality is achieved for some $\mathbf{h} \in \mathbb{R}^p$.

¹On the other hand, the Bonferonni method can be applied to a wide range of statistical problems that involve simultaneous confidence intervals [and is not limited to the theory of linear models. So it is well worth your while to understand this important method.

Proof of Rayleigh–Ritz inequality. Recall the Cauchy–Schwarz inequality from your linear algebra course: $(\mathbf{u}'\mathbf{v})^2 \leq \|\mathbf{u}\|^2 \cdot \|\mathbf{v}\|^2$ with equality if and only if $\mathbf{v} = a\mathbf{u}$ for some $a \in \mathbf{R}$. It follows then that

$$\max_{\mathbf{v} \neq \mathbf{0}} \left[\frac{(\mathbf{u}'\mathbf{v})^2}{\mathbf{v}'\mathbf{v}} \right] = \|\mathbf{u}\|^2.$$

Now write $\mathbf{L} := \mathbf{PDP}'$, in spectral form, and change variables:

$$\mathbf{v} := (\mathbf{L}^{1/2})' \mathbf{h},$$

in order to see that

$$\mathbf{u}'\mathbf{v} = \mathbf{u}' (\mathbf{L}^{1/2})' \mathbf{h}, \quad \mathbf{v}'\mathbf{v} = \mathbf{h}'\mathbf{L}\mathbf{h},$$

and so

$$\|\mathbf{u}\|^2 = \max_{\mathbf{h} \neq \mathbf{0}} \left[\frac{(\mathbf{u}' (\mathbf{L}^{1/2})' \mathbf{h})^2}{\mathbf{h}'\mathbf{L}\mathbf{h}} \right].$$

Change variables again: $\mathbf{b} := \mathbf{L}^{1/2}\mathbf{u}$ to see that

$$\|\mathbf{u}\|^2 = \left\| \mathbf{L}^{-1/2}\mathbf{b} \right\|^2 = \mathbf{b}'\mathbf{L}^{-1}\mathbf{b}.$$

This does the job. □

Now let us return to Scheffé's method. Recall that

$$\hat{\boldsymbol{\beta}} - \boldsymbol{\beta} \sim N_p \left(\mathbf{0}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1} \right),$$

so that

$$\frac{1}{\sigma^2} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})' (\mathbf{X}'\mathbf{X}) (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \sim \chi_p^2,$$

and hence

$$\frac{(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})' (\mathbf{X}'\mathbf{X}) (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})}{pS^2} = \frac{(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})' (\mathbf{X}'\mathbf{X}) (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) / p}{S^2} \sim F_{p, n-p}.$$

We apply the Rayleigh–Ritz inequality with $\mathbf{b} := \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}$ and $\mathbf{L} := (\mathbf{X}'\mathbf{X})^{-1}$ in order to find that

$$(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})' (\mathbf{X}'\mathbf{X}) (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) = \max_{\mathbf{h} \neq \mathbf{0}} \left[\frac{\left\{ \mathbf{h}' (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \right\}^2}{\mathbf{h}' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{h}} \right].$$

Therefore,

$$P \left\{ \frac{1}{pS^2} \max_{\mathbf{h} \neq \mathbf{0}} \left[\frac{\left\{ \mathbf{h}' (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \right\}^2}{\mathbf{h}' (\mathbf{X}'\mathbf{X})^{-1} \mathbf{h}} \right] \leq F_{p, n-p}(\alpha) \right\} = 1 - \alpha.$$

Equivalently,

$$\mathbb{P} \left\{ \frac{|\mathbf{h}'\hat{\boldsymbol{\beta}} - \mathbf{h}'\boldsymbol{\beta}|}{\sqrt{\mathbf{h}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{h}}} \leq \sqrt{pS^2 F_{p,n-p}(\alpha)} \text{ for all } \mathbf{h} \in \mathbf{R}^p \right\} = 1 - \alpha.$$

If we restrict attention to a subcollection of \mathbf{h} 's then the probability is even more. In particular, consider only \mathbf{h} 's that are the standard basis vectors of \mathbf{R}^p , in order to deduce from the preceding that

$$\mathbb{P} \left\{ \frac{|\hat{\beta}_j - \beta_j|}{\sqrt{[(\mathbf{X}'\mathbf{X})^{-1}]_{j,j}}} \leq \sqrt{pS^2 F_{p,n-p}(\alpha)} \text{ for all } 1 \leq j \leq p \right\} \geq 1 - \alpha.$$

In other words, we have demonstrated the following.

Theorem 3 (Scheffé). *The following are conservative $(1 - \alpha) \times 100\%$ simultaneous confidence bounds for $(\beta_1, \dots, \beta_p)'$:*

$$\hat{\beta}_j \pm \sqrt{pS^2 [(\mathbf{X}'\mathbf{X})^{-1}]_{j,j} F_{p,n-p}(\alpha)} \quad (1 \leq j \leq p).$$

When the sample size n is very large, the preceding yields an asymptotic simplification. Recall that $\chi_{n-p}^2/(n-p) \rightarrow 1$ as $n \rightarrow \infty$. Therefore, $F_{p,n-p}(\alpha) \approx \chi_p^2(\alpha)/p$ for $n \gg 1$. Therefore, the following are asymptotically-conservative simultaneous confidence bounds:

$$\hat{\beta}_j \pm S \sqrt{[(\mathbf{X}'\mathbf{X})^{-1}]_{j,j} \chi_p^2(\alpha)} \quad (1 \leq j \leq p) \text{ for } n \gg 1.$$

5. Confidence bounds for the regression surface

Given a vector $\mathbf{x} \in \mathbf{R}^p$ of predictor variables, our linear model yields $E\mathbf{y} = \mathbf{x}'\boldsymbol{\beta}$. In other words, we can view our efforts as one about trying to understand the unknown function

$$f(\mathbf{x}) := \mathbf{x}'\boldsymbol{\beta}.$$

And among other things, we have found the following estimator for f :

$$\hat{f}(\mathbf{x}) := \mathbf{x}'\hat{\boldsymbol{\beta}}.$$

Note that if $\mathbf{x} \in \mathbf{R}^p$ is held fixed, then

$$\hat{f}(\mathbf{x}) \sim N \left(f(\mathbf{x}), \sigma^2 \mathbf{x}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x} \right).$$

Therefore, a $(1 - \alpha) \times 100\%$ confidence interval for $f(\mathbf{x})$ [for a fixed predictor variable \mathbf{x}] is

$$\hat{f}(\mathbf{x}) \pm S \sqrt{\mathbf{x}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}} t_{n-p}^{(\alpha/2)}.$$

That is, if we are interested in a confidence interval for $f(\mathbf{x})$ for a fixed \mathbf{x} , then we have

$$\mathbb{P} \left\{ f(\mathbf{x}) = \hat{f}(\mathbf{x}) \pm S \sqrt{\mathbf{x}'(X'X)^{-1}\mathbf{x}} t_{n-p}^{(\alpha/2)} \right\} = 1 - \alpha.$$

On the other hand, we can also apply Scheffé's method and obtain the following simultaneous $(1 - \alpha) \times 100\%$ confidence set:

$$\mathbb{P} \left\{ f(\mathbf{x}) = \hat{f}(\mathbf{x}) \pm S \sqrt{p \mathbf{x}'(X'X)^{-1}\mathbf{x}} F_{p, n-p}(\alpha) \text{ for all } \mathbf{x} \in \mathbf{R}^p \right\} \geq 1 - \alpha.$$

Example 4 (Simple linear regression). Consider the basic regression model

$$Y_i = \alpha + \beta x_i + \varepsilon_i \quad (1 \leq i \leq n).$$

If $\mathbf{w} = (1, w)'$, then a quick computation yields

$$\mathbf{w}'(X'X)^{-1}\mathbf{w} = \frac{[\bar{x}^2 - 2w\bar{x} + w^2]}{ns_x^2},$$

where we recall $s_x^2 := \sum_{i=1}^n (x_i - \bar{x})^2$. Therefore, a simultaneous confidence interval for all of $\alpha + \beta w$'s, as w ranges over \mathbf{R} , is

$$\hat{\alpha} + \hat{\beta}w \pm S \sqrt{\frac{2}{ns_x^2} [\bar{x}^2 - 2w\bar{x} + w^2] F_{2, n-2}(\alpha)}.$$

This expression can be simplified further because:

$$\begin{aligned} \bar{x}^2 - 2w\bar{x} + w^2 &= \bar{x}^2 - (\bar{x})^2 + (\bar{x})^2 - 2w\bar{x} + w^2 \\ &= s_x^2 + (\bar{x} - w)^2. \end{aligned}$$

Therefore, with probability $1 - \alpha$, we have

$$\alpha + \beta w = \hat{\alpha} + \hat{\beta}w \pm S \sqrt{2 \left(\frac{1}{n} + \frac{(\bar{x} - w)^2}{ns_x^2} \right) F_{2, n-2}(\alpha)} \quad \text{for all } w.$$

On the other hand, if we want a confidence interval for $\alpha + \beta w$ for a fixed w , then we can do better using our t -test:

$$\alpha + \beta w = \hat{\alpha} + \hat{\beta}w \pm S \sqrt{\left(\frac{1}{n} + \frac{(\bar{x} - w)^2}{ns_x^2} \right) t_{n-2}^{(\alpha/2)}}.$$

You should check the details of this computation. □

6. Prediction intervals

The difference between confidence and prediction intervals is this: For a confidence interval we try to find an interval around the **parameter** $\mathbf{x}'\boldsymbol{\beta}$. For a prediction interval we do so for the **random variable** $y_0 := \mathbf{x}_0'\boldsymbol{\beta} + \varepsilon_0$, where \mathbf{x}_0 is known and fixed and ε_0 is the “noise,” which is hitherto unobserved (i.e., independent of the vector \mathbf{Y} of observations).

It is not hard to construct a good prediction interval of this type: Note that

$$\hat{y}_0 := \mathbf{x}_0'\hat{\boldsymbol{\beta}}$$

satisfies

$$\hat{y}_0 - y_0 = \mathbf{x}_0'(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) - \varepsilon_0 \sim N\left(0, \sigma^2 \left[\mathbf{x}_0'(X'X)^{-1}\mathbf{x}_0 + 1 \right]\right).$$

Therefore, a prediction interval is

$$\hat{y}_0 \pm S t_{n-p}^{(\alpha/2)} \sqrt{\mathbf{x}_0'(X'X)^{-1}\mathbf{x}_0 + 1}.$$

This means that

$$P \left\{ y_0 \in \hat{y}_0 \pm S t_{n-p}^{(\alpha/2)} \sqrt{\mathbf{x}_0'(X'X)^{-1}\mathbf{x}_0 + 1} \right\} = 1 - \alpha,$$

but note that both y_0 and the prediction interval are now random.