

# Unit 2: Linear models: Inference

Eugene Katsevich

September 22, 2021

We now understand the least squares estimator  $\hat{\beta}$  from geometric and algebraic points of view. In Unit 2, we will switch to a probabilistic perspective to derive inferential statements for linear models, in the form of hypothesis tests and confidence intervals. In order to facilitate this, we will assume that the error terms are normally distributed:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \text{where } \boldsymbol{\epsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n). \quad (1)$$

## 1 Building blocks for linear model inference

First we put in place some building blocks: The multivariate normal distribution (Section 1.1), the distributions of linear regression estimates and residuals (Section 1.2), and estimation of the noise variance  $\sigma^2$  (Section 1.3).

### 1.1 The multivariate normal distribution

Recall that a random vector  $\mathbf{w} \in \mathbb{R}^d$  has a multivariate normal distribution with mean  $\boldsymbol{\mu}$  and covariate matrix  $\boldsymbol{\Sigma}$  if it has probability density

$$p(\mathbf{w}) = \frac{1}{\sqrt{(2\pi)^d \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{w} - \boldsymbol{\mu})\right).$$

These random vectors have lots of special properties, including:

- (Linear transformation) If  $\mathbf{w} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then  $\mathbf{A}\mathbf{w} + \mathbf{b} \sim N(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$ .
- (Independence) If  $\begin{pmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \end{pmatrix} \sim N\left(\begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{12}^T & \boldsymbol{\Sigma}_{22} \end{pmatrix}\right)$ , then  $\mathbf{w}_1 \perp\!\!\!\perp \mathbf{w}_2$  if and only if  $\boldsymbol{\Sigma}_{12} = \mathbf{0}$ .

An important distribution related to the multivariate normal is the  $\chi_d^2$  (chi-squared with  $d$  degrees of freedom) distribution, defined as

$$\chi_d^2 \equiv \sum_{j=1}^d w_j^2 \quad \text{for } w_1, \dots, w_d \stackrel{\text{i.i.d.}}{\sim} N(0, 1).$$

### 1.2 The distributions of linear regression estimates and residuals

The most important distributional result in linear regression is that

$$\hat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1}). \quad (2)$$

Indeed, by the linear transformation property of the multivariate normal distribution,

$$\mathbf{y} \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n) \implies \hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \sim N((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta, (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I}_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}) \\ = N(\beta, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}).$$

Next, let's consider the joint distribution of  $\hat{\mu} = \mathbf{X}\hat{\beta}$  and  $\hat{\epsilon} = \mathbf{y} - \mathbf{X}\hat{\beta}$ . We have

$$\begin{pmatrix} \hat{\mu} \\ \hat{\epsilon} \end{pmatrix} = \begin{pmatrix} \mathbf{H}\mathbf{y} \\ (\mathbf{I} - \mathbf{H})\mathbf{y} \end{pmatrix} = \begin{pmatrix} \mathbf{H} \\ \mathbf{I} - \mathbf{H} \end{pmatrix} \mathbf{y} \sim N\left(\begin{pmatrix} \mathbf{H} \\ \mathbf{I} - \mathbf{H} \end{pmatrix} \mathbf{X}\beta, \begin{pmatrix} \mathbf{H} \\ \mathbf{I} - \mathbf{H} \end{pmatrix} \cdot \sigma^2 \mathbf{I} \begin{pmatrix} \mathbf{H} & \mathbf{I} - \mathbf{H} \end{pmatrix}\right) \\ = N\left(\begin{pmatrix} \mathbf{X}\beta \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \sigma^2 \mathbf{H} & \mathbf{0} \\ \mathbf{0} & \sigma^2 (\mathbf{I} - \mathbf{H}) \end{pmatrix}\right). \quad (3)$$

In other words,

$$\hat{\mu} \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{H}) \quad \text{and} \quad \hat{\epsilon} \sim N(\mathbf{0}, \sigma^2 (\mathbf{I} - \mathbf{H})), \quad \text{with} \quad \hat{\mu} \perp \hat{\epsilon}. \quad (4)$$

Since  $\hat{\beta}$  is a deterministic function of  $\hat{\mu}$  (in particular,  $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \hat{\mu}$ ), it also follows that

$$\hat{\beta} \perp \hat{\epsilon}. \quad (5)$$

### 1.3 Estimation of the noise variance $\sigma^2$

We can't quite do inference for  $\beta$  based on the distributional result (2) because the noise variance  $\sigma^2$  is unknown to us. Intuitively, since  $\sigma^2 = \mathbb{E}[\epsilon_i^2]$ , we can get an estimate of  $\sigma^2$  by looking at the quantity  $\|\hat{\epsilon}\|^2$ . To get the distribution of this quantity, we need the following lemma:

**Lemma 1.1.** *Let  $\mathbf{w} \sim N(\mathbf{0}, \mathbf{P})$  for some projection matrix  $\mathbf{P}$ . Then,  $\|\mathbf{w}\|^2 \sim \chi_d^2$ , where  $d = \text{trace}(\mathbf{P})$  is the dimension of the subspace onto which  $\mathbf{P}$  projects.*

*Proof.* Let  $\mathbf{P} = \mathbf{U}\mathbf{D}\mathbf{U}^T$  be an eigenvalue decomposition of  $\mathbf{P}$ , where  $\mathbf{U}$  is orthogonal and  $\mathbf{D}$  is a diagonal matrix with  $D_{ii} \in \{0, 1\}$ . We have  $\mathbf{w} \stackrel{d}{=} \mathbf{U}\mathbf{D}\mathbf{z}$  for  $\mathbf{z} \sim N(\mathbf{0}, \mathbf{I}_n)$ . Therefore,

$$\|\mathbf{w}\|^2 = \|\mathbf{D}\mathbf{z}\|^2 = \sum_{i:D_{ii}=1} z_i^2 \sim \chi_d^2, \quad \text{where } d = |\{i : D_{ii} = 1\}| = \text{trace}(\mathbf{D}) = \text{trace}(\mathbf{P}).$$

□

Recall that  $\mathbf{I} - \mathbf{H}$  is a projection onto the  $(n - p)$ -dimensional space  $C(\mathbf{X})^\perp$ , so by Lemma 1.1 and equation (4) we have

$$\|\hat{\epsilon}\|^2 \sim \sigma^2 \chi_{n-p}^2. \quad (6)$$

From this result, it follows that  $\mathbb{E}[\|\hat{\epsilon}\|^2] = n - p$ , so

$$\hat{\sigma}^2 \equiv \frac{1}{n - p} \|\hat{\epsilon}\|^2 \quad (7)$$

is an unbiased estimate for  $\sigma^2$ . Why does the denominator need to be  $n - p$  rather than  $n$  for the estimator above to be unbiased? The reason for this is that the residuals  $\hat{\epsilon}$  are the projection of the true noise vector  $\epsilon$  onto the lower-dimensional subspace  $C(\mathbf{X})^\perp$ . To see this, note that

$$\hat{\epsilon} = (\mathbf{I} - \mathbf{H})\mathbf{y} = (\mathbf{I} - \mathbf{H})(\mathbf{X}\beta + \epsilon) = (\mathbf{I} - \mathbf{H})\epsilon. \quad (8)$$

## 2 Hypothesis testing

Typically two types of null hypotheses are tested in a regression setting: Those involving one-dimensional parameters and those involving multi-dimensional parameters. For example, consider the null hypotheses  $H_0 : \beta_j = 0$  and  $H_0 : \beta_S = \mathbf{0}$  for  $S \subseteq \{0, 1, \dots, p-1\}$ , respectively. We discuss tests of these two kinds of hypothesis in Sections 2.1 and 2.2, and then discuss the power of these tests in Section 2.3.

### 2.1 Testing a one-dimensional parameter

***t*-test for a single coefficient.** The most common question to ask in a linear regression context is: Is the  $j$ th predictor associated with the response, when controlling for the other predictors? In the language of hypothesis testing, this corresponds to the null hypothesis

$$H_0 : \beta_j = 0. \quad (9)$$

According to (2), we have  $\hat{\beta}_j \sim N(0, \sigma^2/s_j^2)$ , where, as we learned in Unit 1,

$$s_j^2 \equiv [(\mathbf{X}^T \mathbf{X})_{jj}^{-1}]^{-1} = \|\mathbf{x}_{*j}^\perp\|^2. \quad (10)$$

Therefore,

$$\frac{\hat{\beta}_j}{\sigma/s_j} \sim N(0, 1), \quad (11)$$

and we are tempted to define a level  $\alpha$  test of the null hypothesis (9) based on this normal distribution. While this is infeasible since we don't know  $\sigma^2$ , we can substitute in the unbiased estimate (7) derived in Section 1.3. Then,

$$\text{SE}_j \equiv \frac{\hat{\sigma}}{s_j} \quad \text{is the standard error of } \hat{\beta}_j, \quad (12)$$

which is an approximation to the standard deviation of  $\hat{\beta}_j$ . Dividing  $\hat{\beta}_j$  by its standard error gives us the *t*-statistic

$$t_j \equiv \frac{\hat{\beta}_j}{\text{SE}_j} = \frac{\hat{\beta}_j}{\sqrt{\frac{1}{n-p} \|\hat{\boldsymbol{\epsilon}}\|^2 / s_j}}. \quad (13)$$

This statistic is *pivotal*, in the sense that it has the same distribution for any  $\boldsymbol{\beta}$  such that  $\beta_j = 0$ . Indeed, we can rewrite it as

$$t_j = \frac{\frac{\hat{\beta}_j}{\sigma/s_j}}{\sqrt{\frac{\sigma^{-2} \|\hat{\boldsymbol{\epsilon}}\|^2}{n-p}}}. \quad (14)$$

Recalling the independence of  $\hat{\boldsymbol{\beta}}$  and  $\hat{\boldsymbol{\epsilon}}$  (5), the scaled chi square distribution of  $\|\hat{\boldsymbol{\epsilon}}\|^2$  (6), the standard normal distribution of  $\frac{\hat{\beta}_j}{\sigma/s_j}$  (11), we find that

$$\text{under } H_0 : \beta_j = 0, \quad t_j \sim \frac{N(0, 1)}{\sqrt{\frac{1}{n-p} \chi_{n-p}^2}}, \quad \text{with numerator and denominator independent.} \quad (15)$$

The latter distribution is called the *t distribution with  $n - p$  degrees of freedom* and denoted  $t_{n-p}$ . This paves the way for the two-sided *t*-test:

$$\phi_t(\mathbf{X}, \mathbf{y}) = \mathbb{1}(|t_j| > t_{n-p}(1 - \alpha/2)), \quad (16)$$

where  $t_{n-p}(1 - \alpha/2)$  denotes the  $1 - \alpha/2$  quantile of  $t_{n-p}$ . Note that, by the law of large numbers,

$$\frac{1}{n-p} \chi_{n-p}^2 \xrightarrow{P} 1 \quad \text{as } n-p \rightarrow \infty, \quad (17)$$

so for large  $n-p$  we have  $t_j \sim t_{n-p} \approx N(0, 1)$ . Hence, the  $t$ -test is approximately equal to the following  $z$ -test:

$$\phi_t(\mathbf{X}, \mathbf{y}) \approx \phi_z(\mathbf{X}, \mathbf{y}) \equiv \mathbb{1}(|t_j| > z(1 - \alpha/2)), \quad (18)$$

where  $z(1 - \alpha/2)$  is the  $1 - \alpha/2$  quantile of  $N(0, 1)$ . The  $t$ -test can also be defined in a one-sided fashion, if power against one-sided alternatives is desired.

**Example: One-sample model.** Consider the intercept-only linear regression model  $y = \beta_0 + \epsilon$ , and let's apply the  $t$ -test derived above to test the null hypothesis  $H_0 : \beta_0 = 0$ . We have  $\hat{\beta}_0 = \bar{y}$ . Furthermore, we have

$$\text{SE}_0^2 = \frac{\hat{\sigma}^2}{n}, \quad \text{where } \hat{\sigma}^2 = \frac{1}{n-1} \|\mathbf{y} - \bar{y}\mathbf{1}_n\|^2. \quad (19)$$

Hence, we obtain the  $t$  statistic

$$t = \frac{\hat{\beta}_0}{\text{SE}_0} = \frac{\sqrt{n}\bar{y}}{\sqrt{\frac{1}{n-1} \|\mathbf{y} - \bar{y}\mathbf{1}_n\|^2}}. \quad (20)$$

According to the theory above, this test statistic has a null distribution of  $t_{n-1}$ .

**Example: Two-sample model.** Suppose we have  $x_1 \in \{0, 1\}$ , in which case the linear regression  $y = \beta_0 + \beta_1 x_1 + \epsilon$  becomes a two-sample model. We can rewrite this model as

$$y_i \sim \begin{cases} N(\beta_0, \sigma^2) & \text{for } x_i = 0; \\ N(\beta_0 + \beta_1, \sigma^2) & \text{for } x_i = 1. \end{cases} \quad (21)$$

It is often of interest to test the null hypothesis  $H_0 : \beta_1 = 0$ , i.e. that the two groups have equal means. Let's define

$$\bar{y}_0 \equiv \frac{1}{n_0} \sum_{i:x_i=0} y_i, \quad \bar{y}_1 \equiv \frac{1}{n_1} \sum_{i:x_i=1} y_i, \quad \text{where } n_0 = |\{i : x_i = 0\}| \text{ and } n_1 = |\{i : x_i = 1\}|. \quad (22)$$

Then, we have seen before that  $\hat{\beta}_0 = \bar{y}_0$  and  $\hat{\beta}_1 = \bar{y}_1 - \bar{y}_0$ . We can compute that

$$s_1^2 \equiv \|\mathbf{x}_{*1}^\perp\|^2 = \|\mathbf{x}_{*1} - \frac{n_1}{n} \mathbf{1}\|^2 = n_1 \frac{n_0^2}{n^2} + n_0 \frac{n_1^2}{n^2} = \frac{n_0 n_1}{n} = \frac{1}{\frac{1}{n_0} + \frac{1}{n_1}} \quad (23)$$

and

$$\hat{\sigma}^2 = \frac{1}{n-2} \left( \sum_{i:x_i=0} (y_i - \bar{y}_0)^2 + \sum_{i:x_i=1} (y_i - \bar{y}_1)^2 \right). \quad (24)$$

Therefore, we arrive at a  $t$ -statistic of

$$t = \frac{\sqrt{\frac{1}{\frac{1}{n_0} + \frac{1}{n_1}}} (\bar{y}_1 - \bar{y}_0)}{\sqrt{\frac{1}{n-2} \left( \sum_{i:x_i=0} (y_i - \bar{y}_0)^2 + \sum_{i:x_i=1} (y_i - \bar{y}_1)^2 \right)}}. \quad (25)$$

Under the null hypothesis, this statistic has a distribution of  $t_{n-2}$ .

***t*-test for a contrast among coefficients.** Given a vector  $\mathbf{c} \in \mathbb{R}^p$ , the quantity  $\mathbf{c}^T \boldsymbol{\beta}$  is sometimes called a *contrast*. For example, suppose  $\mathbf{c} = (1, -1, 0, \dots, 0)$ . Then,  $\mathbf{c}^T \boldsymbol{\beta} = \beta_1 - \beta_2$  is the difference in effects of the first and second predictors. We are sometimes interested in testing whether such a contrast is equal to zero, i.e.  $H_0 : \mathbf{c}^T \boldsymbol{\beta} = 0$ . While this hypothesis can involve two or more of the predictors, the parameter  $\mathbf{c}^T \boldsymbol{\beta}$  is still one-dimensional and therefore we can still apply a *t*-test. Going back to the distribution  $\hat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1})$ , we find that

$$\mathbf{c}^T \hat{\boldsymbol{\beta}} \sim N(\mathbf{c}^T \boldsymbol{\beta}, \sigma^2 \mathbf{c}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c}).$$

Therefore, under the null hypothesis that  $\mathbf{c}^T \boldsymbol{\beta} = 0$ , we can derive that

$$\frac{\mathbf{c}^T \hat{\boldsymbol{\beta}}}{\hat{\sigma} \sqrt{\mathbf{c}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{c}}} \sim t_{n-p}, \quad (26)$$

giving us another *t*-test. Note that the *t*-tests described above can be recovered from this more general formulation by setting  $\mathbf{c} = \mathbf{e}_j$ , the indicator vector with  $j$ th coordinate equal to 1 and all others equal to zero.

## 2.2 Testing a multi-dimensional parameter

***F*-test for a group of coefficients.** Now we move on to the case of testing a multi-dimensional parameter:  $H_0 : \boldsymbol{\beta}_S = \mathbf{0}$  for some  $S \subseteq \{0, 1, \dots, p-1\}$ . In other words, we would like to test

$$H_0 : \mathbf{y} = \mathbf{X}_{*,S} \boldsymbol{\beta}_{-S} + \boldsymbol{\epsilon} \quad \text{versus} \quad H_1 : \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\epsilon}. \quad (27)$$

To test this hypothesis, let us fit least squares coefficients  $\hat{\boldsymbol{\beta}}_{-S}$  and  $\hat{\boldsymbol{\beta}}$  for the partial model as well as the full model. If the partial model fits well, then the residuals  $\mathbf{y} - \mathbf{X}_{*,S} \hat{\boldsymbol{\beta}}_{-S}$  from this model will not be much larger than the residuals  $\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}$  from the full model. To quantify this intuition, let us recall our analysis of variance decomposition from Unit 1:

$$\|\mathbf{y} - \mathbf{X}_{*,S} \hat{\boldsymbol{\beta}}_{-S}\|^2 = \|\mathbf{X} \hat{\boldsymbol{\beta}} - \mathbf{X}_{*,S} \hat{\boldsymbol{\beta}}_{-S}\|^2 + \|\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}\|^2. \quad (28)$$

Let's consider the ratio

$$\frac{\|\mathbf{y} - \mathbf{X}_{*,S} \hat{\boldsymbol{\beta}}_{-S}\|^2 - \|\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}\|^2}{\|\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}\|^2} = \frac{\|\mathbf{X} \hat{\boldsymbol{\beta}} - \mathbf{X}_{*,S} \hat{\boldsymbol{\beta}}_{-S}\|^2}{\|\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}\|^2}, \quad (29)$$

which is the relative increase in the residual sum of squares when going from the full model to the partial model. Let us rewrite this ratio in terms of projection matrices. Let  $\mathbf{H}$  be the projection matrix for the full model, and let  $\mathbf{H}_{-S}$  be the projection matrix for the partial model. Note that  $\mathbf{H} - \mathbf{H}_{-S}$  is the projection matrix onto the  $|S|$ -dimensional space  $C(\mathbf{X}) \cap C(\mathbf{X}_{-S})^T$ . We have

$$\frac{\|\mathbf{X} \hat{\boldsymbol{\beta}} - \mathbf{X}_{*,S} \hat{\boldsymbol{\beta}}_{-S}\|^2}{\|\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}\|^2} = \frac{\|(\mathbf{H} - \mathbf{H}_{-S}) \mathbf{y}\|^2}{\|(\mathbf{I} - \mathbf{H}) \mathbf{y}\|^2}. \quad (30)$$

Under the null hypothesis, we have

$$\frac{\|(\mathbf{H} - \mathbf{H}_{-S}) \mathbf{y}\|^2}{\|(\mathbf{I} - \mathbf{H}) \mathbf{y}\|^2} = \frac{\|(\mathbf{H} - \mathbf{H}_{-S}) \boldsymbol{\epsilon}\|^2}{\|(\mathbf{I} - \mathbf{H}) \boldsymbol{\epsilon}\|^2}. \quad (31)$$

Since the projection matrices in the numerator and denominator project onto orthogonal subspaces, we have  $(\mathbf{H} - \mathbf{H}_{-S}) \boldsymbol{\epsilon} \perp (\mathbf{I} - \mathbf{H}) \boldsymbol{\epsilon}$ , with  $\|(\mathbf{H} - \mathbf{H}_{-S}) \boldsymbol{\epsilon}\|^2 \sim \sigma^2 \chi_{|S|}^2$  and  $\|(\mathbf{I} - \mathbf{H}) \boldsymbol{\epsilon}\|^2 \sim \sigma^2 \chi_{n-p}^2$ .

Renormalizing numerator and denominator to have expectation 1 under the null, we arrive at the  $F$ -statistic

$$F \equiv \frac{(\|\mathbf{y} - \mathbf{X}_{*,-S}\hat{\boldsymbol{\beta}}_{-S}\|^2 - \|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2)/|S|}{\|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2/(n-p)}. \quad (32)$$

We have derived that under the null hypothesis,

$$F \sim \frac{\chi^2_{|S|}/|S|}{\chi^2_{n-p}/(n-p)}, \quad \text{with numerator and denominator independent.} \quad (33)$$

This distribution is called the  $F$ -distribution with  $|S|$  and  $n-p$  degrees of freedom, and denoted  $F_{|S|,n-p}$ . Denoting by  $F_{|S|,n-p}(1-\alpha)$  the  $1-\alpha$  quantile of this distribution, we arrive at the  $F$ -test

$$\phi_F(\mathbf{X}, \mathbf{y}) \equiv \mathbb{1}(F > F_{|S|,n-p}(1-\alpha)). \quad (34)$$

**Example: Testing for any significant coefficients except the intercept.** Suppose  $\mathbf{x}_{*,0} = \mathbf{1}_n$  is an intercept term. Then, consider the null hypothesis  $H_0 : \beta_1 = \dots = \beta_{p-1} = 0$ . In other words, the null hypothesis is the intercept-only model and the alternative hypothesis is the regression model with an intercept and  $p-1$  additional predictors. In this case,  $S = \{1, \dots, p-1\}$  and  $-S = \{0\}$ . The corresponding  $F$  statistic is

$$F \equiv \frac{(\|\mathbf{y} - \bar{y}\mathbf{1}\|^2 - \|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2)/(p-1)}{\|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2/(n-p)}, \quad (35)$$

with null distribution  $F_{p-1,n-p}$ .

**Example: Testing for equality of group means in  $C$ -groups model.** As a further special case, consider the  $C$ -groups model from Unit 1. Recall the ANOVA decomposition

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\bar{y}_{c(i)} - \bar{y})^2 + \sum_{i=1}^n (y_i - \bar{y}_{c(i)})^2 = \text{SSB} + \text{SSW}. \quad (36)$$

The  $F$ -statistic in this case becomes

$$F = \frac{\sum_{i=1}^n (\bar{y}_{c(i)} - \bar{y})^2 / (C-1)}{\sum_{i=1}^n (y_i - \bar{y}_{c(i)})^2 / (n-C)} = \frac{\text{SSB} / (C-1)}{\text{SSW} / (n-C)}, \quad (37)$$

with null distribution  $F_{C-1,n-C}$ .

### 2.3 Power

So far we've been focused on finding the null distributions of various test statistics in order to construct tests with Type-I error control. Now let's shift our attention to examining the power of these tests.

**The power of a  $t$ -test.** Consider the  $t$ -test of the null hypothesis  $H_0 : \beta_j = 0$ . Suppose that, in reality,  $\beta_j \neq 0$ . What is the probability the  $t$ -test will reject the null hypothesis? To answer this question, recall that  $\hat{\beta}_j \sim N(\beta_j, \sigma^2/s_j^2)$ . Therefore,

$$t = \frac{\hat{\beta}_j}{\text{SE}_j} = \frac{\beta_j}{\text{SE}_j} + \frac{\hat{\beta}_j - \beta_j}{\text{SE}_j} \sim N\left(\frac{\beta_j s_j}{\sigma}, 1\right). \quad (38)$$

Here we have made the approximation  $\text{SE}_j \approx \frac{\sigma}{s_j}$ , which is pretty good when  $n - p$  is large. Therefore, the power of the two-sided  $t$ -test is

$$\mathbb{E}[\phi_t] = \mathbb{P}[\phi_t = 1] \approx \mathbb{P}[|t| > z_{1-\alpha/2}] \approx \mathbb{P}\left[\left|N\left(\frac{\beta_j s_j}{\sigma}, 1\right)\right| > z_{1-\alpha/2}\right]. \quad (39)$$

Therefore, the quantity  $\frac{\beta_j s_j}{\sigma}$  determines the power of the  $t$ -test. To understand  $s_j$  a little better, let's assume that the rows  $\mathbf{x}_{i*}$  of the model matrix are drawn i.i.d. from some distribution  $(x_0, \dots, x_{p-1})$ . Then we have roughly

$$\mathbf{x}_{*j}^\perp \approx \mathbf{x}_{*j} - \mathbb{E}[\mathbf{x}_{*j} | \mathbf{X}_{*,j}], \quad (40)$$

so  $x_{ij}^\perp \approx x_{ij} - \mathbb{E}[x_{ij} | \mathbf{x}_{i,-j}]$ . Hence,

$$s_j^2 \equiv \|\mathbf{x}_{*j}^\perp\|^2 \approx n\mathbb{E}[(x_j - \mathbb{E}[x_j | \mathbf{x}_{-j}])^2] = n\mathbb{E}[\text{Var}[x_j | \mathbf{x}_{-j}]]. \quad (41)$$

Hence, we can rewrite the alternative distribution (38) as

$$t \sim N\left(\frac{\beta_j \cdot \sqrt{n} \cdot \sqrt{\mathbb{E}[\text{Var}[x_j | \mathbf{x}_{-j}] ]}}{\sigma}, 1\right). \quad (42)$$

We can see clearly now how the power of the  $t$ -test varies with the effect size  $\beta_j$ , the sample size  $n$ , the degree of collinearity  $\mathbb{E}[\text{Var}[x_j | \mathbf{x}_{-j}]]$ , and the noise standard deviation  $\sigma$ .

**The power of an  $F$ -test.** Now let's turn our attention to computing the power of the  $F$ -test. We have

$$F = \frac{\|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{X}_{*,S}\hat{\boldsymbol{\beta}}_{-S}\|^2/|S|}{\|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2/|n-p|} = \frac{\|(\mathbf{H} - \mathbf{H}_{-S})\mathbf{y}\|^2/|S|}{\|(\mathbf{I} - \mathbf{H})\mathbf{y}\|^2/|n-p|} \approx \frac{\|(\mathbf{H} - \mathbf{H}_{-S})\mathbf{y}\|^2/|S|}{\sigma^2}. \quad (43)$$

To calculate the distribution of the numerator, we need to introduce the notion of a non-central chi-squared random variable.

**Definition 2.1.** For some vector  $\boldsymbol{\mu} \in \mathbb{R}^d$ , suppose  $\mathbf{z} \sim N(\boldsymbol{\mu}, \mathbf{I}_d)$ . Then, we define the distribution of  $\|\mathbf{z}\|^2$  as the non-central chi-square random variable with  $d$  degrees of freedom and noncentrality parameter  $\|\boldsymbol{\mu}\|^2$  and denote this distribution by  $\chi_d^2(\|\boldsymbol{\mu}\|^2)$ .

It can be shown that if  $\mathbf{P}$  is a projection matrix and  $\mathbf{y} = \boldsymbol{\mu} + \boldsymbol{\epsilon}$ , then  $\frac{1}{\sigma^2}\|\mathbf{P}\mathbf{y}\|^2 \sim \chi_{\text{tr}(\mathbf{P})}^2(\frac{1}{\sigma^2}\|\mathbf{P}\boldsymbol{\mu}\|^2)$ . It therefore follows that

$$F \approx \frac{\|(\mathbf{H} - \mathbf{H}_{-S})\mathbf{y}\|^2/|S|}{\sigma^2} \sim \frac{1}{|S|}\chi_{|S|}^2(\|(\mathbf{H} - \mathbf{H}_{-S})\mathbf{X}\boldsymbol{\beta}\|^2) = \frac{1}{|S|}\chi_{|S|}^2\left(\frac{1}{\sigma^2}\|\mathbf{X}_{*,S}^\perp\boldsymbol{\beta}_S\|^2\right). \quad (44)$$

Assuming as before that the rows of  $\mathbf{X}$  are samples from a joint distribution, we can write

$$\|\mathbf{X}_{*,S}^\perp\boldsymbol{\beta}_S\|^2 \approx n\boldsymbol{\beta}_S^T \mathbb{E}[\text{Var}[\mathbf{x}_S | \mathbf{x}_{-S}]]\boldsymbol{\beta}_S. \quad (45)$$

Therefore,

$$F \sim \frac{1}{|S|}\chi_{|S|}^2\left(\frac{n\boldsymbol{\beta}_S^T \mathbb{E}[\text{Var}[\mathbf{x}_S | \mathbf{x}_{-S}]]\boldsymbol{\beta}_S}{\sigma^2}\right), \quad (46)$$

which is similar in spirit to equation (42).

**Power when predictors are added to the model.** As we know, the outcome of a regression is a function of the predictors that are used. What happens to the  $t$ -test  $p$ -value for  $H_0 : \beta_j = 0$  when a predictor is added to the model? To keep things simple, let's consider the

$$\text{true underlying model: } y = \beta_0 x_0 + \beta_1 x_1 + \epsilon. \quad (47)$$

Let's consider the power of testing  $H_0 : \beta_0 = 0$  in the regression models

$$\text{model 0: } y = \beta_0 x_0 + \epsilon \quad \text{versus} \quad \text{model 1: } y = \beta_0 x_0 + \beta_1 x_1 + \epsilon. \quad (48)$$

There are four cases based on  $\text{cor}[\mathbf{x}_{*0}, \mathbf{x}_{*1}]$  and the value of  $\beta_1$  in the true model:

1.  $\text{cor}[\mathbf{x}_{*0}, \mathbf{x}_{*1}] \neq 0$  and  $\beta_1 \neq 0$ . In this case, in model 0 we have omitted an important variable that is correlated with  $\mathbf{x}_{*0}$ . Therefore, the meaning of  $\beta_0$  differs between model 0 and model 1, so it may not be meaningful to compare the  $p$ -values arising from these two models.
2.  $\text{cor}[\mathbf{x}_{*0}, \mathbf{x}_{*1}] \neq 0$  and  $\beta_1 = 0$ . In this case, we are adding a null predictor that is correlated with  $\mathbf{x}_{*0}$ . Recall that the power of the  $t$ -test hinges on the quantity  $\frac{\beta_j \cdot \sqrt{n} \cdot \sqrt{\mathbb{E}[\text{Var}[x_j | \mathbf{x}_{-j}] ]}}{\sigma}$ . Adding the predictor  $x_1$  has the effect of reducing the conditional predictor variance  $\mathbb{E}[\text{Var}[x_j | \mathbf{x}_{-j}] ]$ , therefore reducing the power. This is a case of *predictor competition*.
3.  $\text{cor}[\mathbf{x}_{*0}, \mathbf{x}_{*1}] = 0$  and  $\beta_1 \neq 0$ . In this case, we are adding a non-null predictor that is orthogonal to  $\mathbf{x}_{*0}$ . While the conditional predictor variance  $\mathbb{E}[\text{Var}[x_j | \mathbf{x}_{-j}] ]$  remains the same due to orthogonality, the residual variance  $\sigma^2$  is reduced when going from model 0 to model 1. Therefore, in this case adding  $x_1$  to the model increases the power for testing  $H_0 : \beta_0 = 0$ . This is a case of *predictor collaboration*.
4.  $\text{cor}[\mathbf{x}_{*0}, \mathbf{x}_{*1}] = 0$  and  $\beta_1 = 0$ . In this case, we are adding an orthogonal null variable, which does not change the conditional predictor variance or the residual variance, and therefore keeps the power of the test the same.

In conclusion, adding a predictor can either increase or decrease the power of a  $t$ -test. Similar reasoning can be applied to the  $F$ -test.

### 3 Confidence and prediction intervals

In addition to hypothesis testing, we often want to construct confidence intervals for the coefficients.

**Confidence interval for a coefficient.** Under  $H_0 : \beta_j = 0$ , we showed that  $\frac{\hat{\beta}_j}{\hat{\sigma}/s_j} \sim t_{n-p}$ . The same argument shows that for arbitrary  $\beta_j$ , we have

$$\frac{\hat{\beta}_j - \beta_j}{\hat{\sigma}/s_j} \sim t_{n-p}. \quad (49)$$

We can use this relationship to construct a confidence interval for  $\beta_j$  as follows:

$$\begin{aligned} 1 - \alpha &= \mathbb{P}[|t_{n-p}| \leq t_{n-p}(1 - \alpha/2)] = \mathbb{P}\left[\left|\frac{\hat{\beta}_j - \beta_j}{\hat{\sigma}/s_j}\right| \leq t_{n-p}(1 - \alpha/2)\right] \\ &= \mathbb{P}\left[\beta_j \in \left[\hat{\beta}_j - \frac{\hat{\sigma}}{s_j} t_{n-p}(1 - \alpha/2), \hat{\beta}_j + \frac{\hat{\sigma}}{s_j} t_{n-p}(1 - \alpha/2)\right]\right] \\ &\equiv \mathbb{P}[\beta_j \in I_j]. \end{aligned} \quad (50)$$

The confidence interval  $I_j$  defined above therefore has  $1 - \alpha$  coverage.



**Confidence interval for  $\mathbb{E}[y|\mathbf{x}_0]$ .** Suppose now that we have a new predictor vector  $\mathbf{x}_0 \in \mathbb{R}^p$ . The mean of the response for this predictor vector is  $\mathbb{E}[y|\mathbf{x}_0] = \mathbf{x}_0^T \boldsymbol{\beta}$ . Plugging in  $\mathbf{x}_0$  for  $\mathbf{c}$  in the relation (26), we obtain

$$\frac{\mathbf{x}_0^T \hat{\boldsymbol{\beta}} - \mathbf{x}_0^T \boldsymbol{\beta}}{\hat{\sigma} \sqrt{\mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0}} \sim t_{n-p}.$$

From this we can derive that

$$\mathbf{x}_0^T \hat{\boldsymbol{\beta}} \pm \hat{\sigma} \sqrt{\mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0} \cdot t_{n-p}(1 - \alpha/2) \quad (51)$$

is a  $1 - \alpha$  confidence interval for  $\mathbf{x}_0^T \boldsymbol{\beta}$ .

**Prediction interval for  $y|\mathbf{x}_0$ .** Instead of creating a confidence interval for a point on the regression line, we may want to create a confidence interval for a new draw  $y_0$  of  $y$  for  $\mathbf{x} = \mathbf{x}_0$ , i.e. a *prediction interval*. Note that

$$y_0 - \mathbf{x}_0^T \hat{\boldsymbol{\beta}} = \mathbf{x}_0^T \boldsymbol{\beta} + \epsilon_0 - \mathbf{x}_0^T \hat{\boldsymbol{\beta}} = \epsilon_0 + \mathbf{x}_0^T (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) \sim N(0, \sigma^2 + \sigma^2 \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0). \quad (52)$$

Therefore, we have

$$\frac{y_0 - \mathbf{x}_0^T \hat{\boldsymbol{\beta}}}{\hat{\sigma} \sqrt{1 + \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0}} \sim t_{n-p}, \quad (53)$$

which leads to the  $1 - \alpha$  prediction interval

$$\mathbf{x}_0^T \hat{\boldsymbol{\beta}} \pm \hat{\sigma} \sqrt{1 + \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0} \cdot t_{n-p}(1 - \alpha/2). \quad (54)$$

## 4 Practical considerations

**Practical versus statistical significance.** You can have a statistically significant effect that is not practically significant. The hypothesis testing framework is most useful in the case when the signal to noise ratio is relatively small. Otherwise, constructing a confidence interval for the effect size is a more meaningful approach.

**Correlation versus causation, and Simpson's paradox.** We need to be very careful when interpreting linear regression coefficients, which can be sensitive to the choice of other predictors to include. You can get misleading conclusions if you omit important variables from the regression. A special case of this is *Simpson's paradox*, where an important discrete variable is omitted. Consider the example in Figure 1.

**Dealing with correlated predictors.** It depends on the goal. If we're trying to tease apart effects of correlated predictors, then we have no choice but to proceed as usual despite lower power. Otherwise, we can test predictors in groups via the  $F$ -test to get higher power at the cost of lower "resolution."

### Kidney stone treatment [\[ edit \]](#)

Another example comes from a real-life medical study<sup>[15]</sup> comparing the success rates of two treatments for kidney stones.<sup>[16]</sup> The table below shows the success rates and numbers of treatments for treatments involving both small and large kidney stones, where Treatment A includes open surgical procedures and Treatment B includes closed surgical procedures. The numbers in parentheses indicate the number of success cases over the total size of the group.

Treatment Stone size	Treatment A	Treatment B
Small stones	Group 1 93% (81/87)	Group 2 87% (234/270)
Large stones	Group 3 73% (192/263)	Group 4 69% (55/80)
Both	78% (273/350)	83% (289/350)

The paradoxical conclusion is that treatment A is more effective when used on small stones, and also when used on large stones, yet treatment B appears to be more effective when considering both sizes at the same time. In this example, the "lurking" variable (or [confounding variable](#)) causing the paradox is the size of the stones, which was not previously known to researchers to be important until its effects were included.

Which treatment is considered better is determined by which success ratio (successes/total) is larger. The reversal of the inequality between the two ratios when considering the combined data, which creates Simpson's paradox, happens because two effects occur together:

1. The sizes of the groups, which are combined when the lurking variable is ignored, are very different. Doctors tend to give cases with large stones the better treatment A, and the cases with small stones the inferior treatment B. Therefore, the totals are dominated by groups 3 and 2, and not by the two much smaller groups 1 and 4.
2. The lurking variable, stone size, has a large effect on the ratios; i.e., the success rate is more strongly influenced by the severity of the case than by the choice of treatment. Therefore, the group of patients with large stones using treatment A (group 3) does worse than the group with small stones, even if the latter used the inferior treatment B (group 2).

Based on these effects, the paradoxical result is seen to arise because the effect of the size of the stones overwhelms the benefits of the better treatment (A). In short, the less effective treatment B appeared to be more effective because it was applied more frequently to the small stones cases, which were easier to treat.<sup>[16]</sup>

Figure 1: An example of Simpson's paradox (source: Wikipedia).

**Model selection.** We need to ask ourselves: Why do we want to do model selection? It can either be for prediction purposes or for inferential purposes. If it is for prediction purposes, then we can apply cross-validation to select a model and we don't need to think very hard about statistical significance. If it is for inference, then we need to be more careful. There are various classical model selection criteria (e.g. AIC, BIC), but it is not entirely clear what statistical guarantee we are getting for the resulting models. A simpler approach is to apply a *t*-test for each variable in the model, apply a multiple testing correction to the resulting *p*-values, and report the set of significant variables and the associated guarantee. Re-fitting the linear regression after model selection leads us into some dicey inferential territory due to selection bias. This is the subject of ongoing research and the jury is still out on the best way of doing this.

## 5 R demo

Let's put into practice what we've learned in Units 1 and 2.

```
houses_data = read_tsv("../data/Houses.dat")
houses_data

## # A tibble: 100 x 7
##   case taxes  beds baths  new price  size
##   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1     1     1  3104     4     2     0  280.  2048
```

```
## 2      2 1173      2      1      0 146.    912
## 3      3 3076      4      2      0 238.   1654
## 4      4 1608      3      2      0 200.   2068
## 5      5 1454      3      3      0 160.   1477
## 6      6 2997      3      2      1 500.   3153
## 7      7 4054      3      2      0 266.   1355
## 8      8 3002      3      2      1 290.   2075
## 9      9 6627      5      4      0 587.   3990
## 10     10 320      3      2      0 70.    1160
## # ... with 90 more rows
```

```
# explore the variables
```

```
# should we model beds/baths as categorical or continuous?
```

```
# running a regression and interpreting the summary
```

```
# hypothesis tests, confidence intervals
```

```
# interactions
```

```
# confidence bands
```