

# 1 Inferences about regression parameters

For our linear regression model  $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ , we have not made any assumptions about the data to calculate the parameter estimates  $\hat{\beta}_0$  and  $\hat{\beta}_1$ , since we are simply applying the method of least squares. To conduct inference about regression parameters, however, we need to make the following assumptions about the errors  $\varepsilon_i$ :

1. **Linearity**: the above model is actually correct, we are not neglecting any terms,
2. **Independence**: the errors  $\varepsilon_i$  are independent,
3. **Normality**: the errors  $\varepsilon_i$  follow a normal distribution, and
4. **Homogeneity of variance**:  $\text{Var}(\varepsilon_i) = \sigma_\varepsilon^2$  for all observations.

## 2 Inference for the regression coefficients $\beta_0$ and $\beta_1$

With the assumptions above, we can show that the least squares estimators are unbiased and have the following variances:

$$\text{Var}(\hat{\beta}_0) = \sigma_\varepsilon^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right) \text{ and } \text{Var}(\hat{\beta}_1) = \frac{\sigma_\varepsilon^2}{S_{xx}}.$$

We can estimate  $\sigma_\varepsilon^2$  by  $\hat{\sigma}_\varepsilon^2 = s_\varepsilon^2 = \frac{1}{n-2} \sum (y_i - \hat{y}_i)^2$ . Thus we can construct confidence intervals for  $\beta_0$  and  $\beta_1$  with:

$$\hat{\beta}_0 \pm t_{\alpha/2} s_\varepsilon \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}} \text{ and } \hat{\beta}_1 \pm t_{\alpha/2} \frac{s_\varepsilon}{\sqrt{S_{xx}}}.$$

In a similar way, we can conduct tests of hypotheses as illustrated in the text on page 591.

## 3 Inference for $E(y_{n+1})$ and $y_{n+1}$ (Confidence intervals and prediction intervals)

Two common uses of regression analyses are to predict the mean value of  $y$  at a given  $x_{n+1}$  value ( $E(y_{n+1})$ ), and to predict the value of an individual observation at a given  $x_{n+1}$  value ( $y_{n+1}$ ). In both cases, the estimated value is obtained by evaluating the least-squares prediction equation at  $x_{n+1}$ . The difference is that an interval estimate for an individual value should be much wider than for a mean. Thus we have:

$$\hat{y}_{n+1} = \hat{\beta}_0 + \hat{\beta}_1 x_{n+1} = \bar{y} + \hat{\beta}_1 (x_{n+1} - \bar{x}),$$

as the estimate for either  $E(y_{n+1})$  or  $y_{n+1}$ . For the confidence interval for  $E(y_{n+1})$  we have

$$\hat{y}_{n+1} \pm t_{\alpha/2} s_\varepsilon \sqrt{\frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{S_{xx}}}.$$

We can compute these intervals for individual values of  $x_{n+1}$ , or we can compute them for many values of  $x_{n+1}$  to create confidence bands. Interval estimates for an individual observation at a given  $x_{n+1}$  value are called prediction intervals, and are given by:

$$\hat{y}_{n+1} \pm t_{\alpha/2} s_\varepsilon \sqrt{1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{S_{xx}}}.$$

The extra '1' under the square root sign arises from the extra variability for an individual observation. As with confidence interval calculations, often several prediction intervals are calculated to create prediction bands. Check the SAS code for this lecture to see how to generate confidence bands and prediction bands.

## 4 The ANOVA table

A summary of the regression results that is used for many linear model analyses is called the analysis of variance (ANOVA) table. It is based on the idea that the variability of the  $y$  values about their mean can be partitioned into two sources: one part describing variability in  $y$  'explained' by the regression, and a second part due to error in predicting  $y$  from the regression model. This relationship can be seen by noting that:

$y_i - \bar{y} = (\hat{y}_i - \bar{y}) + (y_i - \hat{y}_i)$  (also see Figure 11-10 of page 582 of the Ott text). When each side is squared and summed (and the crossproduct term vanishes) we get the result:

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2,$$

which is also expressed as  $SSY = SS(\text{Regression}) + SS(\text{Residual})$ . The term  $SS(\text{Regression})$  is called the sum of squares due to regression. These sums of squares are presented in the ANOVA table, and then  $SS(\text{Regression})$  and  $SS(\text{Residual})$  are divided by their degrees of freedom to obtain mean squares, denoted  $MS(\text{Regression})$  and  $MS(\text{Residual})$ , respectively.  $MS(\text{Regression})$  and  $MS(\text{Residual})$  are statistically independent, and if the null hypothesis

of  $H_0 : \beta_1 = 0$  is true, then they both estimate the population variance  $\sigma^2$ . Thus we can take the ratio  $F = \text{MS}(\text{Regression}) / \text{MS}(\text{Residual})$  and use it as a test of  $H_0 : \beta_1 = 0$ .  $F$  follows an F distribution with numerator degrees of freedom =  $df_1=1$  and denominator degrees of freedom =  $df_2= n-2$ . This test is equivalent to the t-test of  $H_0 : \beta_1 = 0$  presented earlier, and in fact  $F_{\alpha,1,n-2} = t_{\alpha/2,n-2}^2$ .