

# 1.017/1.010 Class 23

## Analyzing Regression Results

---

### Analyzing and Interpreting Regression Results

Least-squares estimation methods provide a way to fit linear regression models (e.g. polynomial curves) to data. Once a model is obtained it is useful to be able to quantify:

1. The significance of the regression
2. The accuracy of the parameter estimates and predictions

The significance of the regression can be analyzed with an **ANOVA** approach. Estimation and prediction accuracy are related to the **means and variances** of the regression parameters.

### Regression ANOVA

The regression term is not significant (it does not explain any of the  $y$  variability) if the following hypothesis is true:

$$H_0: E[y(x)] = h(x)A = a_1$$

That is, the mean of  $y$  is a constant that does not depend on the independent variable  $x$ .

This hypothesis can be tested with a statistic based on the following sums-of-squares:

$$SST = \sum_{i=1}^n (y_i - m_y)^2 \quad ; \quad m_y = \frac{1}{n} \sum_{i=1}^n y_i$$

$$\begin{aligned} SSE &= [Y - H\hat{A}]'[Y - H\hat{A}] \\ &= \sum_{i=1}^n \left[ y_i - (\hat{a}_1 + \hat{a}_2 x_i + \hat{a}_3 x_i^2) \right]^2 \end{aligned}$$

$$SSR = SST - SSE$$

**$SST$**  measures the  $y$  variability if the regression model **is not** used.

**$SSE$**  measures the  $y$  variability if the regression model **is** used.

**$SSR$**  measures the  $y$  variability explained by the regression model.

The statistic used to test significance of the regression is the ratio of the mean sums of squares for regression and error:

$$MSR = \frac{SSR}{m-1}$$

$$MSE = \frac{SSE}{n-m}$$

$$F_R(MSR, MSE) = \frac{MSR}{MSE}$$

$E[MSR]$  depends on the magnitudes of the regression coefficients  $a_2, \dots, a_m$  while  $E[MSE]$  does not. Therefore, their ratio is sensitive to the magnitude of these coefficients.

When  $H_0$  is true  $F_R$  follows an **F distribution** with degree of freedom parameters  $\nu_R = m-1$  and  $\nu_E = n-m$ . The rejection region and  $p$  values are derived from this distribution. If  $F_R$  is large and  $p$  is small,  $H_0$  is rejected and the regression is significant.

### ANOVA Table for Linear Regression:

Source	SS	df	MS	$F$	$p$
Regression	$SSR$	$\nu_R = m - 1$	$MSR = \frac{SSR}{\nu_R}$	$F_R = \frac{MSR}{MSE}$	$p = 1 - F_{F, \nu_R, \nu_E}(F)$
Error	$SSE$	$\nu_E = n - m$	$MSE = \frac{SSE}{\nu_E}$		
Total	$SST$	$\nu_T = n - 1$			

The  $R$ -squared coefficient is:

$$R^2 = \frac{SSR}{SST}$$

$R^2$  is often used to describe the quality of a regression fit.  $R^2 = 1$  is a perfect fit.

The internal MATLAB function `regress` provides the  $R^2$ ,  $F_R$ , and  $p$  values obtained from the regression ANOVA.

### Properties of Regression Parameters and Predictions

The estimates of parameters  $a_1, a_2, \dots, a_m$  obtained in a regression analysis have the general form:

$$\hat{A} = [H'H]^{-1} H'Y = WY$$

$$\hat{a}_i = \sum_{j=1}^n W_{ij} y_j ; \quad i = 1 \dots m$$

So the estimates are **linear combinations** of the measurements  $[y_1 y_2 \dots y_n]$ , with each measurement weighted by a coefficient  $W_{ij}$  that depends only on the known  $x$  values  $[x_1 x_2 \dots x_n]$ . In this respect, regression parameter estimates are similar to the sample mean, which is also a linear combination of measurements.

Each regression **parameter estimate** is a random variable with its own CDF. Its mean and variance may be found from the estimation and measurement equations and the assumed statistical properties of the random residuals  $e_i \dots E[e_i] = 0, Var[e_i] = \sigma_e^2$ , which are assumed to be independent :

$$E[\hat{a}_i] = a_i ; \quad i = 1 \dots m$$

$$Var[\hat{a}_i] = \sigma_e^2 \{ [H'H]^{-1} \}_{ii} \approx s_e^2 \{ [H'H]^{-1} \}_{ii} ; \quad i = 1 \dots m$$

The unknown residual error variance  $\sigma_e^2$  can be approximated by:

$$\sigma_e^2 \approx s_e^2 = MSE = \frac{1}{n-m} \sum_{i=1}^n [y_i - (\hat{a}_1 + \hat{a}_2 x_i + \hat{a}_3 x_i^2)]^2$$

The least-squares regression parameters are **unbiased** and **consistent** .

The **prediction** derived from the regression parameters is also a random variable that is a linear combination of the measurements. Example for quadratic regression model discussed in class:

$$\hat{y}(x) = h(x)\hat{A} = \hat{a}_1 + \hat{a}_2 x + \hat{a}_3 x^2 \quad ; \quad h(x) = \begin{bmatrix} 1 & x & x^2 \end{bmatrix}$$

Mean and variance of this prediction at any  $x$  are:

$$E[\hat{y}(x)] = E[h(x)\hat{A}] = h(x)A = E[y(x)] = a_1 + a_2 x + a_3 x^2$$

$$Var[\hat{y}(x)] = h(x)\sigma_e^2[H'H]^{-1}h'(x) \approx h(x)s_e^2[H'H]^{-1}h'(x)$$

These results also apply for other  $h(x)$ .

## Regression Parameter Confidence Intervals

When the sample size  $n$  is **large** the regression parameters are approximately normally distributed and the CDF of each estimate is completely defined by its mean and variance:

$$F_{\hat{a}_i}(\hat{a}_i) \sim N(E[\hat{a}_i], Var[\hat{a}_i]) = N(a_i, \sigma_e^2 \{[H'H]^{-1}\}_{ii})$$

The procedure for deriving large sample confidence intervals and for testing hypotheses is the same as for the sample mean.

The  $1-\alpha$  **two-sided large sample confidence interval** is:

$$\hat{a}_i - z_U SD[\hat{a}_i] \leq a_i \leq \hat{a}_i + z_L SD[\hat{a}_i]$$

$$\hat{a}_i - z_U s_e \{[H'H]^{-1}\}_{ii}^{1/2} \leq a_i \leq \hat{a}_i + z_L s_e \{[H'H]^{-1}\}_{ii}^{1/2}$$

where  $z_L$  and  $z_U$  are obtained from the unit normal distribution ( $z_L = -1.96$  and  $z_U = +1.96$  for  $\alpha = 0.05$ ):

$$z_L = F_z^{-1}\left(\frac{\alpha}{2}\right) \quad z_U = F_z^{-1}\left(1 - \frac{\alpha}{2}\right)$$

When the sample size  $n$  is **small** and the residual errors are **normally distributed** the regression parameters are **t distributed** with  $\nu = n - m$  degrees of freedom. The two-sided confidence intervals are computed as above, with  $F_z$  replaced by  $F_{t,\nu}$ .

The regression coefficient confidence intervals are evaluated by the internal MATLAB function `regress`.

## Regression Prediction Confidence Intervals

When the sample size  $n$  is **large** the regression prediction is approximately normally distributed with a CDF completely defined by its mean and variance:

$$F_{\hat{y}}[\hat{y}(x)] \sim N(E[\hat{y}(x)], \text{Var}[\hat{y}(x)]) = N[h(x)A, h(x)\sigma_e^2[H'H]^{-1}h'(x)]$$

The  $1-\alpha$  **two-sided large sample confidence interval** is:

$$\hat{y}(x) - z_U SD[\hat{y}(x)] \leq y(x) \leq \hat{y}(x) + z_U SD[\hat{y}(x)]$$

$$\hat{y}(x) - z_U [h(x)s_e^2 [H'H]^{-1}h'(x)]^{1/2} \leq y(x) \leq \hat{y}(x) + z_U [h(x)s_e^2 [H'H]^{-1}h'(x)]^{1/2}$$

where  $z_L, z_U$ , and the prediction standard deviation are obtained from the equations given earlier and  $\sigma_e^2$  is approximated by  $s_e^2$ .

When the sample size  $n$  is **small** and the residual errors are **normally distributed** the regression prediction is **t distributed** with  $\nu = n - 2$  degrees of freedom. The two-sided confidence interval is computed as in the large sample case, with  $F_z$  replaced by  $F_{t,\nu}$ .

The regression prediction confidence interval depends on  $x$  and widens for  $x$  far from the values  $[x_1, x_2 \dots x_n]$  corresponding to measurements. This interval is evaluated by the internal MATLAB function `regress`.



Copyright 2003 Massachusetts Institute of Technology  
Last modified Oct. 8, 2003