

Linear Least Squares Curve Fitting (Regression) Simulations -II

The Excel file **data_records_2021** contains 5 sets of y measurements (**y1, y2,..., y5**) obtained for $-10 \leq x \leq +10$.

Use 5-fold cross validation in conjunction with the residual sum of squares (RSS) (see eqn 1 below) and R2 (coefficient of determination, eqn 2 below) metrics to estimate a polynomial model that best explains the given data. Show plots and statistics as appropriate to justify your selection. Also show the coefficients of the chosen polynomial model.

$$\text{RSS: } \sum_{i=1}^m (y_i - \hat{y}_i)^2 \quad (1)$$

$$\text{R2: } 1 - \frac{\text{RSS}}{\sum_{i=1}^m (y_i - \hat{\mu})^2} = 1 - \frac{\overline{\text{RSS}}}{\widehat{\text{Var}}(y)} \quad (2)$$

In Eqn (2) $\hat{\mu}$ is the estimated mean and $\widehat{\text{Var}}(y)$ is the estimated variance of a data record, $\{y_i, i = 1, 2, \dots, m\}$.

(Note: RSS and R2 can be used to evaluate performance on both, training data and test data)

Be alert to possible numerical errors arising out of ill-conditioning of $\mathbf{A}^T \mathbf{A}$ which may affect solution accuracy of the normal equations, especially for larger model orders. As before, use only the basic tools available for matrix multiplication, transpose, inverse, solution of equations, etc. (this way you may learn a lot) and not use any built-in functions that your programming environment may offer to directly solve the problem (this way you may not learn a lot).