

27411 Chemometrics and Biological Data Analysis

Lecture 2: MLR - Multiple Linear Regression (eNote3)

Per Bruun Brockhoff

DTU Compute, Statistics and Data Analysis
Building 324, Room 220
Danish Technical University
2800 Lyngby – Danmark
e-mail: perbb@dtu.dk

DTU Compute
Department of Applied Mathematics and Computer Science

Regression Analysis- overview

- Description (Estimation)
- Inference (Hypothesis testing, confidence bands)
- Prediction (of future values)
- Validation/model diagnostics/data check.

DTU Compute
Department of Applied Mathematics and Computer Science

Overview

- 1 Regression Analysis - overview
- 2 Simple Linear Regression
- 3 Multiple Linear Regression (MLR)
- 4 Assumptions of MLR
- 5 Regression Analysis
 - Cross products matrix
 - Inference
 - Centering and Scaling
 - Interpretation
 - Validation
- 6 Multi-collinearity!!!

DTU Compute
Department of Applied Mathematics and Computer Science

Simple Linear Regression

- Model: $y_i = b_0 + b_1x_i + \varepsilon_i$
- "Guesses" Y given knowledge of X : $\hat{y}_i = \hat{b}_0 + \hat{b}_1x_i$
 - Error in guessing: $\sum_i (y_i - \hat{y}_i)^2$
- Guessing Y WITHOUT any X -knowledge: $y_i^{guess} = \bar{y}$
 - Error in guessing: $\sum_i (y_i - \bar{y})^2$ ("usual" Y -variance)

DTU Compute
Department of Applied Mathematics and Computer Science

Explained variance

$$R^2 = \frac{\sum_i (y_i - \bar{y})^2 - \sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2} = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2}$$

- R^2 = squared correlation between y_i s and \hat{y}_i s
- R^2 = squared correlation between y_i s and x_i s
 - (\hat{b}_0, \hat{b}_1) is chosen to maximize R^2
 - (\hat{b}_0, \hat{b}_1) is chosen to minimize squared residual variation.
 - Least Squares (LS) estimates.

Multiple Linear Regression (MLR)

$$y_i = b_0 + b_1x_{i1} + b_2x_{i2} + \dots + b_px_{ip} + \varepsilon_i, \quad i = 1, \dots, n$$

E.g. $n = 6$ objects and $p = 3$ x -variables:

$$y_1 = b_0 + b_1x_{11} + b_2x_{12} + b_3x_{13} + \varepsilon_1$$

$$y_2 = b_0 + b_1x_{21} + b_2x_{22} + b_3x_{23} + \varepsilon_2$$

$$y_3 = b_0 + b_1x_{31} + b_2x_{32} + b_3x_{33} + \varepsilon_3$$

$$y_4 = b_0 + b_1x_{41} + b_2x_{42} + b_3x_{43} + \varepsilon_4$$

$$y_5 = b_0 + b_1x_{51} + b_2x_{52} + b_3x_{53} + \varepsilon_5$$

$$y_6 = b_0 + b_1x_{61} + b_2x_{62} + b_3x_{63} + \varepsilon_6$$

Multiple Linear Regression (MLR)

Or equivalently:

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \\ y_6 \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} \\ 1 & x_{21} & x_{22} & x_{23} \\ 1 & x_{31} & x_{32} & x_{33} \\ 1 & x_{41} & x_{42} & x_{43} \\ 1 & x_{51} & x_{52} & x_{53} \\ 1 & x_{61} & x_{62} & x_{63} \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ b_3 \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \varepsilon_4 \\ \varepsilon_5 \\ \varepsilon_6 \end{bmatrix}$$

Or equivalently:

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \boldsymbol{\varepsilon}$$

Assumptions of MLR

- Model structure is OK:

$$E(y_i | x_{i1}, x_{i2}, \dots, x_{ip}) = b_0 + b_1x_{i1} + b_2x_{i2} + \dots + b_px_{ip}$$

- Errors $\varepsilon_1, \dots, \varepsilon_n$
- are independent

- have homogeneous variance: $\text{Var}(\varepsilon_i) = \sigma^2$

- The $n \times p$ -matrix \mathbf{X} is of full rank:

- No x -variable can be perfectly expressed by a linear combination of the others. (No *perfect redundancy*)
- A consequence: The number of objects must exceed the number of variables: $n > p$

Regression Analysis

- Description: $\hat{\mathbf{b}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ (Unbiased: $E(\hat{\mathbf{b}}) = \mathbf{b}$)
- Inference: $\text{Var}(\hat{\mathbf{b}}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}$ (Minimal)
- Prediction: $\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{b}} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$

$$\hat{y}_i = \hat{b}_0 + \hat{b}_1 x_{i1} + \hat{b}_2 x_{i2} + \dots + \hat{b}_p x_{ip}$$

- Validation: "Look at" $e_i = y_i - \hat{y}_i$

Matrix of cross products ($n = 6, p = 3$)

$$\mathbf{X}'\mathbf{X} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ x_{11} & x_{21} & x_{31} & x_{41} & x_{51} & x_{61} \\ x_{12} & x_{22} & x_{32} & x_{42} & x_{52} & x_{62} \\ x_{13} & x_{23} & x_{33} & x_{43} & x_{53} & x_{63} \end{bmatrix} \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} \\ 1 & x_{21} & x_{22} & x_{23} \\ 1 & x_{31} & x_{32} & x_{33} \\ 1 & x_{41} & x_{42} & x_{43} \\ 1 & x_{51} & x_{52} & x_{53} \\ 1 & x_{61} & x_{62} & x_{63} \end{bmatrix}$$

$$= \begin{bmatrix} 6 & \sum_i x_{i1} & \sum_i x_{i2} & \sum_i x_{i3} \\ \sum_i x_{i1} & \sum_i x_{i1}^2 & \sum_i x_{i1}x_{i2} & \sum_i x_{i1}x_{i3} \\ \sum_i x_{i2} & \sum_i x_{i2}x_{i1} & \sum_i x_{i2}^2 & \sum_i x_{i2}x_{i3} \\ \sum_i x_{i3} & \sum_i x_{i3}x_{i1} & \sum_i x_{i3}x_{i2} & \sum_i x_{i3}^2 \end{bmatrix}$$

Matrix of cross products

IF the variables are mean-centered AND we omit the first column (of ones):

$$\begin{aligned} (\mathbf{X}'\mathbf{X})_{jk} &= \sum_i (x_{ij} - \bar{x}_{.j})(x_{ik} - \bar{x}_{.k}) \\ &= (n-1)\text{Corr}(\mathbf{x}_j, \mathbf{x}_k)\sqrt{\text{Var}(\mathbf{x}_j)}\sqrt{\text{Var}(\mathbf{x}_k)} \end{aligned}$$

Matrix of cross products

IF the variables are mean-centered AND standardized: (still omitting the first column):

$$\begin{aligned} (\mathbf{X}'\mathbf{X})_{jk} &= \sum_i \frac{(x_{ij} - \bar{x}_{.j})(x_{ik} - \bar{x}_{.k})}{\sqrt{\text{Var}(\mathbf{x}_j)}\sqrt{\text{Var}(\mathbf{x}_k)}} \\ &= (n-1)\text{Corr}(\mathbf{x}_j, \mathbf{x}_k) \end{aligned}$$

Inference

- $\hat{\sigma}^2 = \sum_i (y_i - \hat{y}_i)^2 / (n - p - 1)$
- Testing the model: $H_0 : b_1 = b_2 = \dots = b_p$

$$F = \frac{\sum_i (\hat{y}_i - \bar{y})^2}{\sum_i (y_i - \hat{y}_i)^2 / (n - p - 1)}, \sim F(p, n - p - 1)$$

- Testing individual parameters: $H_0 : b_k = 0$

$$t = \frac{\hat{b}_k}{SE(\hat{b}_k)}, \sim t(n - p - 1)$$

$$SE((\hat{b}_k))^2 = k\text{th diagonal element of } \hat{\sigma}^2 (\mathbf{X}'\mathbf{X})^{-1}$$

DTU Compute
Department of Applied Mathematics and Computer Science

Centering and scaling in MLR

- X-Centering changes \hat{b}_0 BUT leaves everything else UNCHANGED
- X-scaling changes $\hat{b}_1, \hat{b}_2, \dots, \hat{b}_p$ by the scaling factors
 - Variance and inference is unchanged!
- MLR is *invariant* to centering and scaling

DTU Compute
Department of Applied Mathematics and Computer Science

Interpretation of parameters

- Each parameter expresses the ADDITIONAL effect of the variable AFTER the effects of all other variables have been removed
- Effects of other variables are *corrected for*
- Expresses the expected change in y given a change of the individual variable leaving all other variables fixed.

DTU Compute
Department of Applied Mathematics and Computer Science

Interpretation of parameters

$$y_i = b_0 + b_1 x_{i1} + b_2 x_{i2} + \varepsilon_i$$

\hat{b}_1 can be found by:

- Regress \mathbf{y} on \mathbf{x}_2 , compute residuals: e_i^y
- Regress \mathbf{x}_1 on \mathbf{x}_2 , compute residuals: e_i^x
- Regress e^y on e^x .

DTU Compute
Department of Applied Mathematics and Computer Science

Validation: Assumptions

- Model structure is satisfactory
- Homogeneous variance: $\text{Var}(\varepsilon_i) = \sigma^2$
- Independent errors
- (Data is OK)

Residuals

- Residuals: $e_i = y_i - \hat{y}_i$
- Studentized Residuals I: $e_i^* = \frac{e_i}{\hat{\sigma}\sqrt{1-h_{ii}}}$
- Studentized Residuals II: $e_i^* = \frac{e_i}{\hat{\sigma}(i)\sqrt{1-h_{ii}}}$
- Leverage of observation i : $h_{ii} = (X(X^tX)^{-1}X^t)_{ii}$
 - Express the (potential) influence of the observation
 - Measures the relative X-position of the observation
- $\hat{\sigma}(i)$: Estimated in a model fit WITHOUT using the i 'th observation.

Validation: Tools

- Residual investigations
- Plots and tests
- Investigating models that allows for EXTENDED structure
 - Interactions
 - Non-linearities
- Cross-Validation and Jackknifing
- Checking for outliers and influential observations

Influential/outlying observations

- Residuals vs. leverage
- Another method: DFBETA: measure the change of beta when leaving out an observation.
- Outliers: Extreme y-values. Look at all plots!

Additional structure

- Potentially: plot residuals vs. each of the X-variables
- Include e.g. x_i^2 (and maybe also x_i^3) in the model and check whether these terms are significant!
- Include interaction terms: Allowing e.g. the effect of ELEVATION to depend on COUNTY:
 - Include the product of COUNTY and ELEVATION in the model

Resampling

- Cross-Validation (CV)
- Jackknifing (Leave-on-out CV)
- Bootstrapping
- A good generic approach:
 - Split the data into a TRAINING and a TEST set.
 - Use Cross-validation on the TRAINING data
 - Check the model performance on the TEST-set
 - MAYBE: REPEAT all this many times (Repeated Double Cross Validation)

Multi-collinearity!!!

- IF too much redundant X-information: model fit becomes UNSTABLE!
 - Poor prediction performance!
- Classic approach: Variable selection!
- Better approach: "Biased regression methods"
 - Principal Component Regression (PCR)
 - Partial Least Squares Regression (PLSR)
 - Ridge and Lasso Regression (RR)
 - Many others exist