

Regression Analysis

Multiple Linear Regression

Nicoleta Serban, Ph.D.

Professor

School of Industrial and Systems Engineering

Assumptions and Diagnostics



1

About This Lesson



2

Multiple Linear Regression: Model

Data: $\{(x_{1,1}, \dots, x_{1,p}), y_1\}, \dots, \{(x_{n,1}, \dots, x_{n,p}), y_n\}$

Model: $y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p} + \varepsilon_i, i = 1, \dots, n$

Assumptions:

- *Linearity/Mean Zero Assumption:* The relationship between the response variable and each predicting variable is linear. (For each $j, j = 1, \dots, p, y_i$ and x_{ij} are linearly related, $i = 1, \dots, n$.) $E(\varepsilon_i) = 0$
- *Constant Variance Assumption:* $\text{Var}(\varepsilon_i) = \sigma^2$
- *Independence Assumption:* $\{\varepsilon_1, \dots, \varepsilon_n\}$ are independent random variables
- Assumption that $\varepsilon_i \sim \text{Normal}$ for confidence/prediction intervals, hypothesis testing



3

Properties of the Errors & Residuals

Properties of (true) errors:

- $E(\varepsilon_i) = 0$
- $\text{Var}(\varepsilon_i) = \sigma^2$
- $E(\boldsymbol{\varepsilon}) = \mathbf{0}$
- $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}$

Properties of the (estimated) residuals:

- $\hat{\boldsymbol{\varepsilon}} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$
- $E(\hat{\boldsymbol{\varepsilon}}) = \mathbf{0}$ (or $E(\hat{\varepsilon}_i) = 0$)
- $\mathbf{V}(\hat{\boldsymbol{\varepsilon}}) = \sigma^2(\mathbf{I} - \mathbf{H})$ (or $\text{Var}(\hat{\varepsilon}_i) = \sigma^2(1 - h_{i,i})$)
 - Where \mathbf{H} is the hat matrix, and $h_{i,i}$ is the i -th element on its diagonal



4

Properties of the Errors & Residuals

Properties of (true) errors:

- $E(\varepsilon_i) = 0$
- $\text{Var}(\varepsilon_i) = \sigma^2$
- $E(\boldsymbol{\varepsilon}) = \mathbf{0}$
- $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}$

Properties of the (estimated) residuals:

- $\hat{\boldsymbol{\varepsilon}} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$
- $E(\hat{\boldsymbol{\varepsilon}}) = \mathbf{0}$ (or $E(\hat{\varepsilon}_i) = 0$)
- $\mathbf{V}(\hat{\boldsymbol{\varepsilon}}) = \sigma^2(\mathbf{I} - \mathbf{H})$ (or $\text{Var}(\hat{\varepsilon}_i) = \sigma^2(1 - h_{i,i})$)
 - Where \mathbf{H} is the hat matrix, and $h_{i,i}$ is the i -th element on its diagonal

- While the true errors have constant variance, the estimated residuals do not.

- To use the estimated residuals for assessing the model assumptions, we need to standardize:

$$r_i = \hat{\varepsilon}_i / (\hat{\sigma} \sqrt{1 - h_{i,i}})$$

Residuals Analysis

Standardized Residual Values: $r_i = \hat{\varepsilon}_i / (\hat{\sigma} \sqrt{1 - h_{i,i}})$

Graphical assessment of MLR assumptions:

- Plot standardized residuals r_i against each predictor
 - *Linearity*
- Plot standardized residuals r_i against fitted values
 - *Constant Variance & Independence*
- QQ normal plot & histogram
 - *Normality*

Residuals Analysis

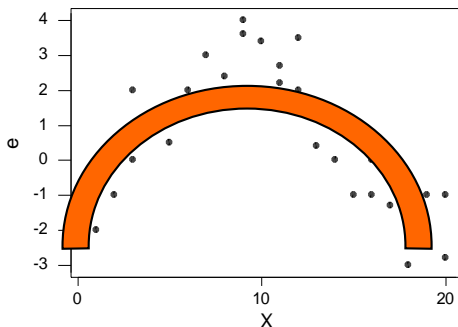
Standardized Residual Values: $r_i = \hat{\varepsilon}_i / (\hat{\sigma} \sqrt{1 - h_{i,i}})$

Graphical assessment of MLR assumptions:

- Plot standardized residuals r_i against each predictor
 - *Linearity*
 - Plot standardized residuals r_i against fitted values
 - *Constant Variance & Independence*
 - QQ normal plot & histogram
 - *Normality*
- We evaluate the normality assumption using the residuals, not the response variable.
 - We do not check the predicting variables for normality.
 - However, if the distribution of a predicting variable is strongly skewed, it is possible that the linearity assumption with respect to that variable will not hold.

Residual Analysis: Linearity Assumption

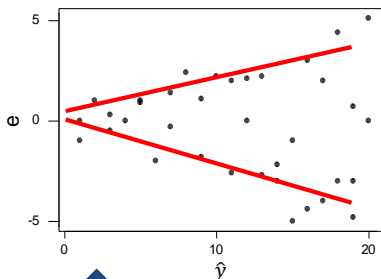
Linearity: Plot the residuals against each predicting variable.



This shows that there may be a non-linear relationship between X and Y .

Residual Analysis: Constant Variance Assumption

Constant Variance: Plot the residuals against fitted values.



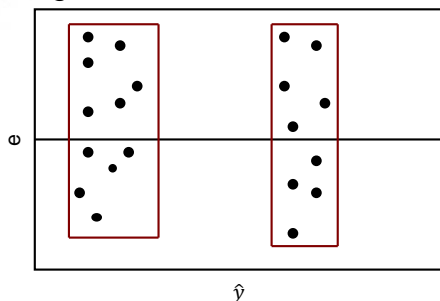
The residuals show larger variance as the predicting variable increases.



Here, it is an example for which σ^2 is not constant.

Residual Analysis: Independence Assumption

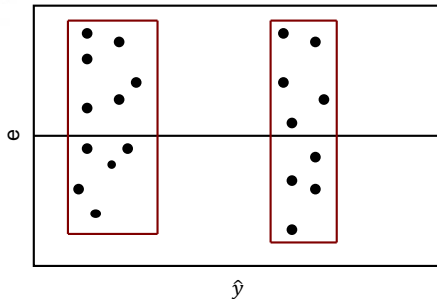
Independence (*uncorrelated errors*): Plot the residuals against fitted values.



- There are clusters of residuals.
- The independence assumption does not hold.

Residual Analysis: Independence Assumption

Independence (*uncorrelated errors*): Plot the residuals against fitted values.



- There are clusters of residuals.
- The independence assumption does not hold.

- Using residual analysis, we are actually checking for uncorrelated errors, not independence.
- Independence is a more complicated matter. If the data are from a randomized experiment, then independence holds, but most data are from observational studies.
- We commonly correct for selection bias in observational studies using controlling variables.

Checking the Assumption of Normality

One way to check this assumption in a regression is using a **Normal Probability (Q-Q) Plot**

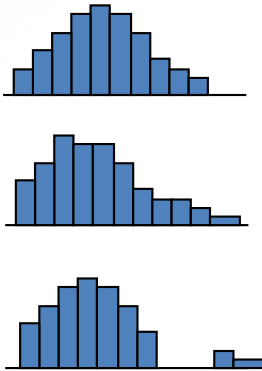
| | |
|---------|---|
| y-axis: | e_i |
| x-axis: | $\Phi^{-1}\left(\frac{r_i - 3/8}{n + 1/4}\right)$ |

r_i = rank of e_i (between 1, n)
 Φ = CDF of Normal Distribution

- Let the R statistical software do this for you!
- A straight line in a normal probability plot implies that the assumption is valid
- **Curvature (especially at the ends)** shows non-normality

Residual Analysis: Normality Assumption

A complementary approach for checking for the normality assumption is by plotting the histogram of the residuals.



Normality Assumption:

The residuals should have an approximately symmetric, unimodal distribution, with no gaps in the data.

Predicting Variable Transformation

- If the model fit is inadequate, it does not mean that a regression is not useful.
- One problem might be that one or more predicting variable X might not have a linear relationship with the response variable Y .
- To model the nonlinear relationship, we transform X by some nonlinear function such as

$$f(x) = x^a$$

or

$$f(x) = \log(x)$$

Normality Transformation

Problem: Constant variance or/and normality assumption

Solution: Transform the response variable from y to \hat{y}^* via

$$\hat{y}^* = y^\lambda$$

where the value of λ depends on how $\text{Var}(y)$ changes as x changes.

$$\sigma_y(x) \propto \text{const} \quad \lambda = 1 \quad (\text{don't transform})$$

$$\sigma_y(x) \propto \sqrt{\mu_x} \quad \lambda = 1/2 \quad \hat{y}^* = \sqrt{y}$$

$$\sigma_y(x) \propto \mu_x \quad \lambda = 0 \quad \hat{y}^* = \ln(y)$$

$$\sigma_y(x) \propto \mu_x^2 \quad \lambda = -1 \quad \hat{y}^* = \frac{1}{y}$$



15

Outliers in Regression

A data point far from the majority of the data (in y and/or any x) may be called an *outlier*, especially if it does not follow the general trend of the rest of the data.

- Data points that are far from the means of the X s or near the edge of the observation space are called *leverage points*.
- A data point that is far from the means of y and/or an x is called an *influential point* if it influences the fit of the regression.
- Excluding a leverage point may or may not the regression fit significantly, thus a leverage point may or may not be an influential point.

The upshot: Sometimes there are good reasons to exclude subsets of data (e.g., errors in data entry or experimental errors). Sometimes an outlier belongs in the data. Outliers should always be examined.



16

Checking for Outliers

Cook's Distance:
$$D_i = \frac{(\hat{Y}_{(i)} - \hat{Y})^T (\hat{Y}_{(i)} - \hat{Y})}{(k + 1) \hat{\sigma}^2}$$

where $\hat{Y}_{(i)}$ are the fitted values from the model fitted without the i^{th} observation (i.e., excluding the i^{th} observation from the data) and \hat{Y} are the fitted values from the model fitted with the i^{th} observation (i.e., including all observations).

Cook's Distance measures how much the estimated parameter values in the regression model change when the i^{th} observation is removed.

Rule of Thumb: $D_i > 4/n$, $D_i > 1$, OR any "large" D_i should be investigated.



17

Checking for Outliers

Cook's Distance:
$$D_i = \frac{(\hat{Y}_{(i)} - \hat{Y})^T (\hat{Y}_{(i)} - \hat{Y})}{(k + 1) \hat{\sigma}^2}$$

where $\hat{Y}_{(i)}$ are the fitted values from the model fitted without the i^{th} observation (i.e., excluding the i^{th} observation from the data) and \hat{Y} are the fitted values from the model fitted with the i^{th} observation (i.e., including all observations).

Cook's Distance measures how much the estimated parameter values in the regression model change when the i^{th} observation is removed.

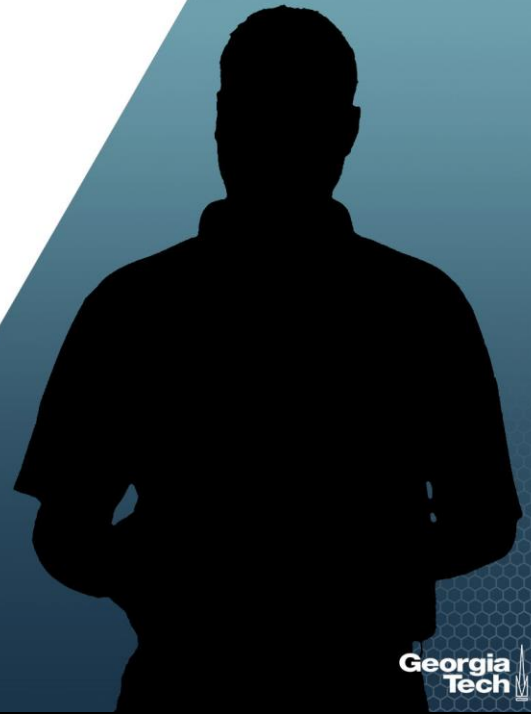
Rule of Thumb: $D_i > 4/n$, $D_i > 1$, OR any "large" D_i should be investigated.

- Outliers: are those few observations with much larger Cook's distance than the rest of observations;
- If a large number of outliers, then they probably point to a heavy tailed distribution rather than truly extreme values.



18

Summary



Georgia
Tech