

CHE384, From Data to Decisions: Measurement, Uncertainty, Analysis, and Modeling

Lecture 37

Independence of Residuals

Chris A. Mack
Adjunct Associate Professor

<http://www.lithoguru.com/scientist/statistics/>

© Chris Mack, 2016 Data to Decisions 1

Assumptions in OLS Regression

1. ε is a random variable that does not depend on x (i.e., the model is perfect, it properly accounts for the role of x in predicting y)
2. $E[\varepsilon_i] = 0$ (the population mean of the true residual is zero); this will always be true for a model with an intercept
3. All ε_i are independent of each other (uncorrelated for the population, though not for a sample)
4. All ε_i have the same probability density function (pdf), and thus the same variance (called homoscedasticity)
5. $\varepsilon \sim N(0, \sigma_\varepsilon)$ (the residuals, and thus the y_i , are normally distributed)
6. The values of each x_i are known exactly

© Chris Mack, 2016 Data to Decisions 2

Residuals Are Not Independent

- Residuals are variation unexplained by the fitted model
- Assumption 3 (all residuals are independent of each other) is never true for a sample
 - Residuals depend on the fitted regression function, which depends on the same data that the residuals come from
 - With p fitting parameters, n residuals only have $n - p$ degrees of freedom
- For n sufficiently large compared to p , we can ignore this dependence

© Chris Mack, 2016 Data to Decisions 3

Residuals in Sequence

- Other factors can prevent residual independence
 - Model error (usually a missing predictor variable)
 - Time dependence: sample aging, measurement drift
 - Spatial dependence: where the measurement was taken
- A **run sequence plot** shows the residuals in time sequence or other natural order to look for **systematic variation**
 - If time order corresponds to changing predictor values, the drift/aging may be hidden in the functional relationship
 - **Randomization** in experimental design prevents this
- A **lag plot** can make systematic variation more visible

© Chris Mack, 2016 Data to Decisions 4

Run Sequence Plots

NBS measurements for a standard weight (Data Set 1)

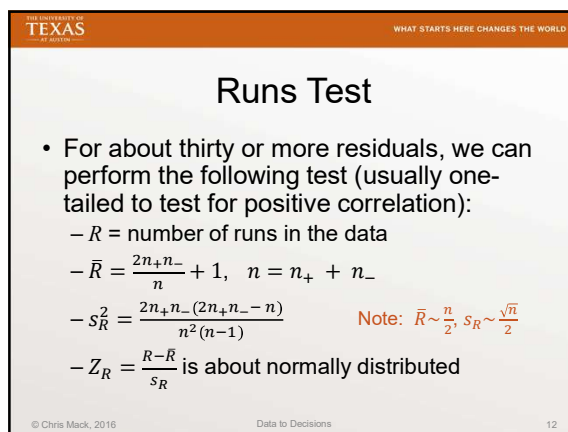
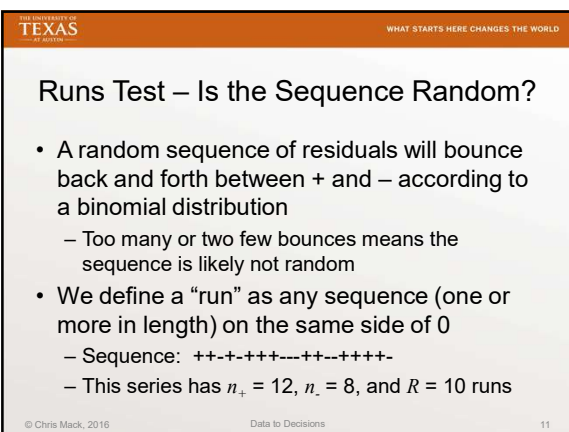
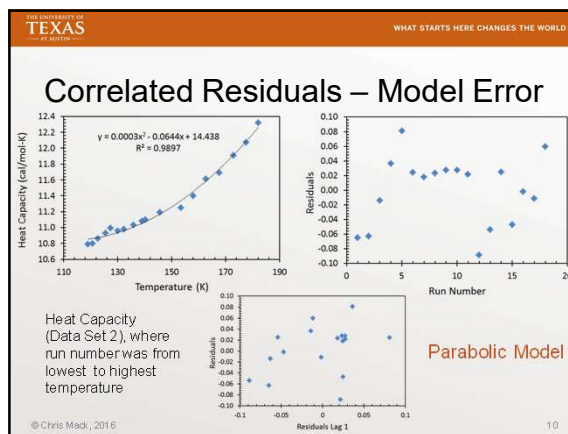
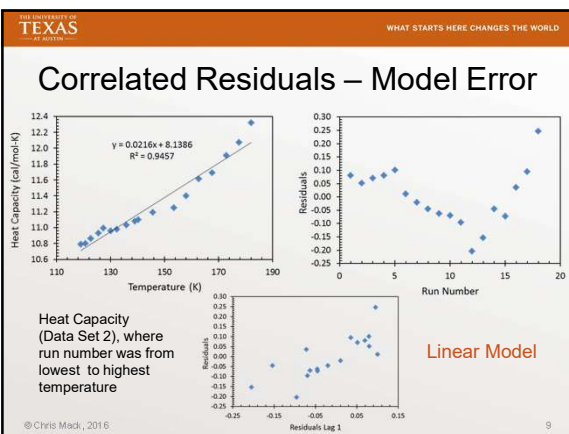
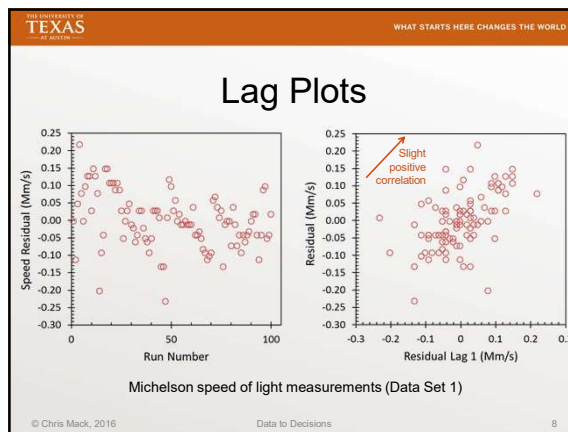
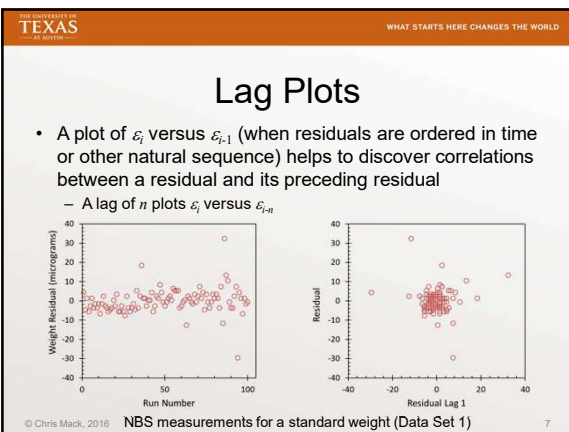
Intensity measurements for a chaotic laser (Data Set 1)

© Chris Mack, 2016 Data to Decisions 5

Run Sequence Plots

Here, run number is the patient in sequence of measurement (Data Set 2).
(The small number of data points makes any conclusions tentative.)

© Chris Mack, 2016 Data to Decisions 6



THE UNIVERSITY OF
TEXAS
AT ARLINGTON

WHAT STARTS HERE CHANGES THE WORLD

Lecture 37: What have we learned?

- What can cause correlated (non-independent) residuals?
- Be able to generate and interpret a lag plot if the data sequence (order) is known
- Why is randomization of data order important in experimental design?
- What is a runs test and how is it performed?

© Chris Mack, 2016

Data to Decisions

13