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

# Lecture 4
# Process Modeling

Chris A. Mack

*Adjunct Associate Professor*

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

© Chris Mack, 2016    Data to Decisions    1

---

# Process Modeling

http://www.itl.nist.gov/div898/handbook/pmd/section1/pmd1.htm

- Process modeling is the concise description of the total variation in one quantity $y$ (called the response variable) by partitioning it into
  - A deterministic component given by a mathematical function of one or more other quantities, $x_1, x_2, ...$ and possibly unknown coefficients $\beta_0, \beta_1, ...$
  - A random component $\varepsilon$ that follows a particular probability distribution

$$y = f(x, \beta) + \varepsilon \qquad x = (x_1, x_2, ...); \quad \beta = (\beta_0, \beta_1, ...)$$

© Chris Mack, 2016    Data to Decisions    2

---

# Process Modeling

$$y = f(x, \beta) + \varepsilon$$

- Generally, we require $E[\varepsilon] = 0$.
  - Thus $f(x, \beta)$ describes the *average* response, $E[y]$, if the experiment is repeated many times, not the actual response for a given trial
- Our three tasks in modeling:
  - Find the equation from $f(x, \beta)$ that meets our goals
  - Find the values of the coefficients $\beta$ that are "best" in some sense
  - Characterize the nature of $\varepsilon$ (distribution of errors)

© Chris Mack, 2016    Data to Decisions    3

---

# Process Modeling

$$y = f(x, \beta) + \varepsilon$$

- The perfect model has
  - The correct set of input variables $x_1, x_2, ...$
  - The correct model form $f(x, \beta)$
  - The correct values for the coefficients $\beta_0, \beta_1, ...$
  - The correct probability distribution for $\varepsilon$, including parameters such as its standard deviation $\sigma$
- Picking the right model (form and predictor variables) is called modeling building
- Finding the best estimate of the parameter values and the properties of the random variable $\varepsilon$ is called regression

© Chris Mack, 2016    Data to Decisions    4

---

# Model Generalizability

$$y = f(x, \beta) + \varepsilon$$

- The three aspects of our model (equation, coefficients, and errors) can have different levels of generalizability
  - We often want to know the levels of generalizability
- Example: model of thermal stress on polymer
  - The equation form applies to all materials (under certain conditions)
  - The parameters change for different materials
  - The errors are a function of measurement and experimental methods, independent of materials

© Chris Mack, 2016    Data to Decisions    5

---

# Some Terminology

$$y = f(x, \beta) + \varepsilon$$

- $y$ = response variable, response, dependent variable
- $x$ = predictor variable, explanatory variable, independent variable, predictor, regressor
- Our "model" is both the function $f(x, \beta)$ and the assumed distribution of $\varepsilon$

© Chris Mack, 2016    Data to Decisions    6

## Regression

- Regression involves three things:
  - Data (a response variable as a function of one or more predictor variables)
  - Model (fixed form and predictor variables, but with unknown parameters)
  - Method (a statistical regression technique appropriate for the model and the data to find the "best" values of the model parameters)
- High quality regression requires high quality in all three items

© Chris Mack, 2016        Data to Decisions        7

## The Model

- Statistical Relationship: $y_i = f(x_i, \beta) + \varepsilon_i = \hat{y}_i + \varepsilon_i$
- Functional Relationship: $\hat{y} = f(x, \beta)$ or $E[Y|X] = f(X, \beta)$

  $X, Y$ = random variables (probability terminology)
  $\hat{y}$ = predicted (mean) response
  $y_i$ = measured response for i[th] data point
  $x_i$ = value of explanatory variable for i[th] data point
  $\beta_k$ = true model parameters (can never be known)
  $b_k$ = best fit model parameters for this data set (sample); our estimate for $\beta_k$.
  $\varepsilon_i$ = true value of i[th] residual (from true model, not known)
  $e_i$ = actual i[th] residual for the current model

© Chris Mack, 2016        Data to Decisions        8

## Example Model

- Straight line model:
  - $f(x, \beta) = \beta_0 + \beta_1 x$
  - $\hat{y}_i = \beta_0 + \beta_1 x_i$
  - $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$
- Regression produces the "best" estimate of the model given the data $(x_i, y_i)$:
  - $y_i = b_0 + b_1 x_i + e_i$

© Chris Mack, 2016        Data to Decisions        9

## Models for Linear Regression

- We use linear regression for linear-parameter models: $\hat{y}$ is directly proportional to each unknown model coefficient
  - $\hat{y} = \sum_k \beta_k f_k(x)$ for bivariate data
  - Example: $\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 \ln(x)$
- Multivariate data: two or more explanatory variables (we'll call them $x_1$, $x_2$, etc.)

© Chris Mack, 2016        Data to Decisions        10

## Nonlinear Regression

- We call our regression nonlinear if it is nonlinear in the coefficients
  - Linear regression: $\hat{y} = \beta_0 + \beta_1 \ln(x) + \beta_2 x^3$
  - Nonlinear regression: $\hat{y} = \beta_0 e^{\beta_1 x}$
- Linear regression is relatively easy
  - Numerically stable with unique solution given a reasonable definition of "best" fit
- Nonlinear regression is relatively hard

© Chris Mack, 2016        Data to Decisions        11

## Lecture 4: What have we learned?

- What are the three tasks in process modeling?
- Explain the relationship between model building and regression
- What are the two major outputs of regression?
- Define "linear regression"

© Chris Mack, 2016        Data to Decisions        12