



4. Process Modeling

The goal for this chapter is to present the background and specific analysis techniques needed to construct a statistical model that describes a particular scientific or engineering process. The types of models discussed in this chapter are limited to those based on an explicit mathematical function. These types of models can be used for prediction of process outputs, for calibration, or for process optimization.

1. [Introduction](#)

1. [Definition](#)
2. [Terminology](#)
3. [Uses](#)
4. [Methods](#)

2. [Assumptions](#)

1. [Assumptions](#)

3. [Design](#)

1. [Definition](#)
2. [Importance](#)
3. [Design Principles](#)
4. [Optimal Designs](#)
5. [Assessment](#)

4. [Analysis](#)

1. [Modeling Steps](#)
2. [Model Selection](#)
3. [Model Fitting](#)
4. [Model Validation](#)
5. [Model Improvement](#)

5. [Interpretation & Use](#)

1. [Prediction](#)
2. [Calibration](#)
3. [Optimization](#)

6. [Case Studies](#)

1. [Load Cell Output](#)
2. [Alaska Pipeline](#)
3. [Ultrasonic Reference Block](#)
4. [Thermal Expansion of Copper](#)

[Detailed Table of Contents: Process Modeling](#)

[References: Process Modeling](#)

[Appendix: Some Useful Functions for Process Modeling](#)



4. Process Modeling - Detailed Table of Contents [4.]

The goal for this chapter is to present the background and specific analysis techniques needed to construct a statistical model that describes a particular scientific or engineering process. The types of models discussed in this chapter are limited to those based on an explicit mathematical function. These types of models can be used for prediction of process outputs, for calibration, or for process optimization.

1. [Introduction to Process Modeling](#) [4.1.]
 1. [What is process modeling?](#) [4.1.1.]
 2. [What terminology do statisticians use to describe process models?](#) [4.1.2.]
 3. [What are process models used for?](#) [4.1.3.]
 1. [Estimation](#) [4.1.3.1.]
 2. [Prediction](#) [4.1.3.2.]
 3. [Calibration](#) [4.1.3.3.]
 4. [Optimization](#) [4.1.3.4.]
 4. [What are some of the different statistical methods for model building?](#) [4.1.4.]
 1. [Linear Least Squares Regression](#) [4.1.4.1.]
 2. [Nonlinear Least Squares Regression](#) [4.1.4.2.]
 3. [Weighted Least Squares Regression](#) [4.1.4.3.]
 4. [LOESS \(aka LOWESS\)](#) [4.1.4.4.]
2. [Underlying Assumptions for Process Modeling](#) [4.2.]
 1. [What are the typical underlying assumptions in process modeling?](#) [4.2.1.]
 1. [The process is a *statistical* process.](#) [4.2.1.1.]
 2. [The means of the random errors are zero.](#) [4.2.1.2.]
 3. [The random errors have a constant standard deviation.](#) [4.2.1.3.]
 4. [The random errors follow a normal distribution.](#) [4.2.1.4.]
 5. [The data are randomly sampled from the process.](#) [4.2.1.5.]
 6. [The explanatory variables are observed without error.](#) [4.2.1.6.]
3. [Data Collection for Process Modeling](#) [4.3.]
 1. [What is design of experiments \(DOE\)?](#) [4.3.1.]
 2. [Why is experimental design important for process modeling?](#) [4.3.2.]
 3. [What are some general design principles for process modeling?](#) [4.3.3.]
 4. [I've heard some people refer to "optimal" designs, shouldn't I use those?](#) [4.3.4.]
 5. [How can I tell if a particular experimental design is good for my application?](#) [4.3.5.]
4. [Data Analysis for Process Modeling](#) [4.4.]
 1. [What are the basic steps for developing an effective process model?](#) [4.4.1.]
 2. [How do I select a function to describe my process?](#) [4.4.2.]
 1. [Incorporating Scientific Knowledge into Function Selection](#) [4.4.2.1.]
 2. [Using the Data to Select an Appropriate Function](#) [4.4.2.2.]
 3. [Using Methods that Do Not Require Function Specification](#) [4.4.2.3.]
 3. [How are estimates of the unknown parameters obtained?](#) [4.4.3.]
 1. [Least Squares](#) [4.4.3.1.]
 2. [Weighted Least Squares](#) [4.4.3.2.]

4. [How can I tell if a model fits my data?](#) [4.4.4.]
 1. [How can I assess the sufficiency of the functional part of the model?](#) [4.4.4.1.]
 2. [How can I detect non-constant variation across the data?](#) [4.4.4.2.]
 3. [How can I tell if there was drift in the measurement process?](#) [4.4.4.3.]
 4. [How can I assess whether the random errors are independent from one to the next?](#) [4.4.4.4.]
 5. [How can I test whether or not the random errors are distributed normally?](#) [4.4.4.5.]
 6. [How can I test whether any significant terms are missing or misspecified in the functional part of the model?](#) [4.4.4.6.]
 7. [How can I test whether all of the terms in the functional part of the model are necessary?](#) [4.4.4.7.]
5. [If my current model does not fit the data well, how can I improve it?](#) [4.4.5.]
 1. [Updating the Function Based on Residual Plots](#) [4.4.5.1.]
 2. [Accounting for Non-Constant Variation Across the Data](#) [4.4.5.2.]
 3. [Accounting for Errors with a Non-Normal Distribution](#) [4.4.5.3.]
5. [Use and Interpretation of Process Models](#) [4.5.]
 1. [What types of predictions can I make using the model?](#) [4.5.1.]
 1. [How do I estimate the average response for a particular set of predictor variable values?](#) [4.5.1.1.]
 2. [How can I predict the value and estimate the uncertainty of a single response?](#) [4.5.1.2.]
 2. [How can I use my process model for calibration?](#) [4.5.2.]
 1. [Single-Use Calibration Intervals](#) [4.5.2.1.]
 3. [How can I optimize my process using the process model?](#) [4.5.3.]
6. [Case Studies in Process Modeling](#) [4.6.]
 1. [Load Cell Calibration](#) [4.6.1.]
 1. [Background & Data](#) [4.6.1.1.]
 2. [Selection of Initial Model](#) [4.6.1.2.]
 3. [Model Fitting - Initial Model](#) [4.6.1.3.]
 4. [Graphical Residual Analysis - Initial Model](#) [4.6.1.4.]
 5. [Interpretation of Numerical Output - Initial Model](#) [4.6.1.5.]
 6. [Model Refinement](#) [4.6.1.6.]
 7. [Model Fitting - Model #2](#) [4.6.1.7.]
 8. [Graphical Residual Analysis - Model #2](#) [4.6.1.8.]
 9. [Interpretation of Numerical Output - Model #2](#) [4.6.1.9.]
 10. [Use of the Model for Calibration](#) [4.6.1.10.]
 11. [Work This Example Yourself](#) [4.6.1.11.]
 2. [Alaska Pipeline](#) [4.6.2.]
 1. [Background and Data](#) [4.6.2.1.]
 2. [Check for Batch Effect](#) [4.6.2.2.]
 3. [Initial Linear Fit](#) [4.6.2.3.]
 4. [Transformations to Improve Fit and Equalize Variances](#) [4.6.2.4.]
 5. [Weighting to Improve Fit](#) [4.6.2.5.]
 6. [Compare the Fits](#) [4.6.2.6.]
 7. [Work This Example Yourself](#) [4.6.2.7.]
 3. [Ultrasonic Reference Block Study](#) [4.6.3.]
 1. [Background and Data](#) [4.6.3.1.]
 2. [Initial Non-Linear Fit](#) [4.6.3.2.]
 3. [Transformations to Improve Fit](#) [4.6.3.3.]
 4. [Weighting to Improve Fit](#) [4.6.3.4.]
 5. [Compare the Fits](#) [4.6.3.5.]
 6. [Work This Example Yourself](#) [4.6.3.6.]
 4. [Thermal Expansion of Copper Case Study](#) [4.6.4.]
 1. [Background and Data](#) [4.6.4.1.]

2. [Rational Function Models](#) [4.6.4.2.]
 3. [Initial Plot of Data](#) [4.6.4.3.]
 4. [Quadratic/Quadratic Rational Function Model](#) [4.6.4.4.]
 5. [Cubic/Cubic Rational Function Model](#) [4.6.4.5.]
 6. [Work This Example Yourself](#) [4.6.4.6.]
-
7. [References For Chapter 4: Process Modeling](#) [4.7.]
-
8. [Some Useful Functions for Process Modeling](#) [4.8.]
 1. [Univariate Functions](#) [4.8.1.]
 1. [Polynomial Functions](#) [4.8.1.1.]
 1. [Straight Line](#) [4.8.1.1.1.]
 2. [Quadratic Polynomial](#) [4.8.1.1.2.]
 3. [Cubic Polynomial](#) [4.8.1.1.3.]
 2. [Rational Functions](#) [4.8.1.2.]
 1. [Constant / Linear Rational Function](#) [4.8.1.2.1.]
 2. [Linear / Linear Rational Function](#) [4.8.1.2.2.]
 3. [Linear / Quadratic Rational Function](#) [4.8.1.2.3.]
 4. [Quadratic / Linear Rational Function](#) [4.8.1.2.4.]
 5. [Quadratic / Quadratic Rational Function](#) [4.8.1.2.5.]
 6. [Cubic / Linear Rational Function](#) [4.8.1.2.6.]
 7. [Cubic / Quadratic Rational Function](#) [4.8.1.2.7.]
 8. [Linear / Cubic Rational Function](#) [4.8.1.2.8.]
 9. [Quadratic / Cubic Rational Function](#) [4.8.1.2.9.]
 10. [Cubic / Cubic Rational Function](#) [4.8.1.2.10.]
 11. [Determining m and n for Rational Function Models](#) [4.8.1.2.11.]

[4. Process Modeling](#)

4.1. Introduction to Process Modeling

*Overview
of Section
4.1*

The goal for this section is to give the big picture of function-based process modeling. This includes a discussion of what process modeling is, the goals of process modeling, and a comparison of the different statistical methods used for model building. Detailed information on how to collect data, construct appropriate models, interpret output, and use process models is covered in the following sections. The final section of the chapter contains case studies that illustrate the general information presented in the first five sections using data from a variety of scientific and engineering applications.

*Contents
of Section
4.1*

1. [What is process modeling?](#)
2. [What terminology do statisticians use to describe process models?](#)
3. [What are process models used for?](#)
 1. [Estimation](#)
 2. [Prediction](#)
 3. [Calibration](#)
 4. [Optimization](#)
4. [What are some of the statistical methods for model building?](#)
 1. [Linear Least Squares Regression](#)
 2. [Nonlinear Least Squares Regression](#)
 3. [Weighted Least Squares Regression](#)
 4. [LOESS \(aka LOWESS\)](#)

[4. Process Modeling](#)[4.1. Introduction to Process Modeling](#)

4.1.1. What is process modeling?

Basic Definition

Process modeling is the concise description of the total variation in one quantity, y , by partitioning it into

1. a deterministic component given by a mathematical function of one or more other quantities, x_1, x_2, \dots , plus
2. a random component that follows a particular probability distribution.

Example

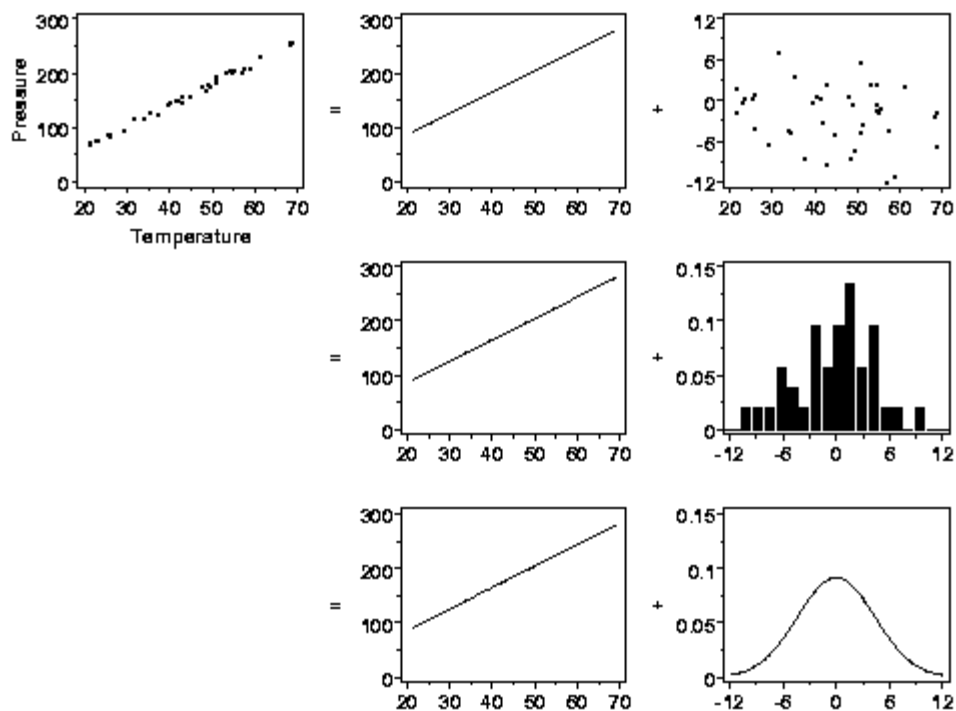
For example, the total variation of the measured pressure of a fixed amount of a gas in a tank can be described by partitioning the variability into its deterministic part, which is a function of the temperature of the gas, plus some left-over random error. Charles' Law states that the pressure of a gas is proportional to its temperature under the conditions described here, and in this case most of the variation will be deterministic. However, due to measurement error in the pressure gauge, the relationship will not be purely deterministic. The random errors cannot be characterized individually, but will follow some probability distribution that will describe the relative frequencies of occurrence of different-sized errors.

Graphical Interpretation

Using the example above, the definition of process modeling can be graphically depicted like this:

*Click Figure
for Full-Sized
Copy*

4.1.1. What is process modeling?

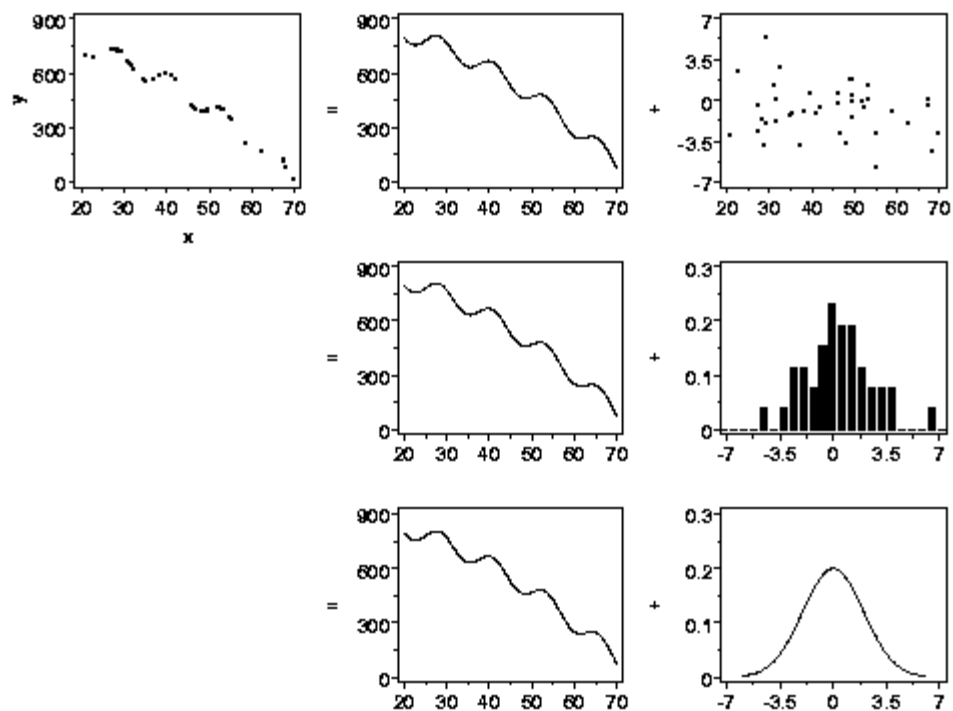


The top left plot in the figure shows pressure data that vary deterministically with temperature except for a small amount of random error. The relationship between pressure and temperature is a straight line, but not a perfect straight line. The top row plots on the right-hand side of the equals sign show a partitioning of the data into a perfect straight line and the remaining "unexplained" random variation in the data (note the different vertical scales of these plots). The plots in the middle row of the figure show the deterministic structure in the data again and a [histogram](#) of the random variation. The histogram shows the relative frequencies of observing different-sized random errors. The bottom row of the figure shows how the relative frequencies of the random errors can be summarized by a (normal) probability distribution.

An Example from a More Complex Process

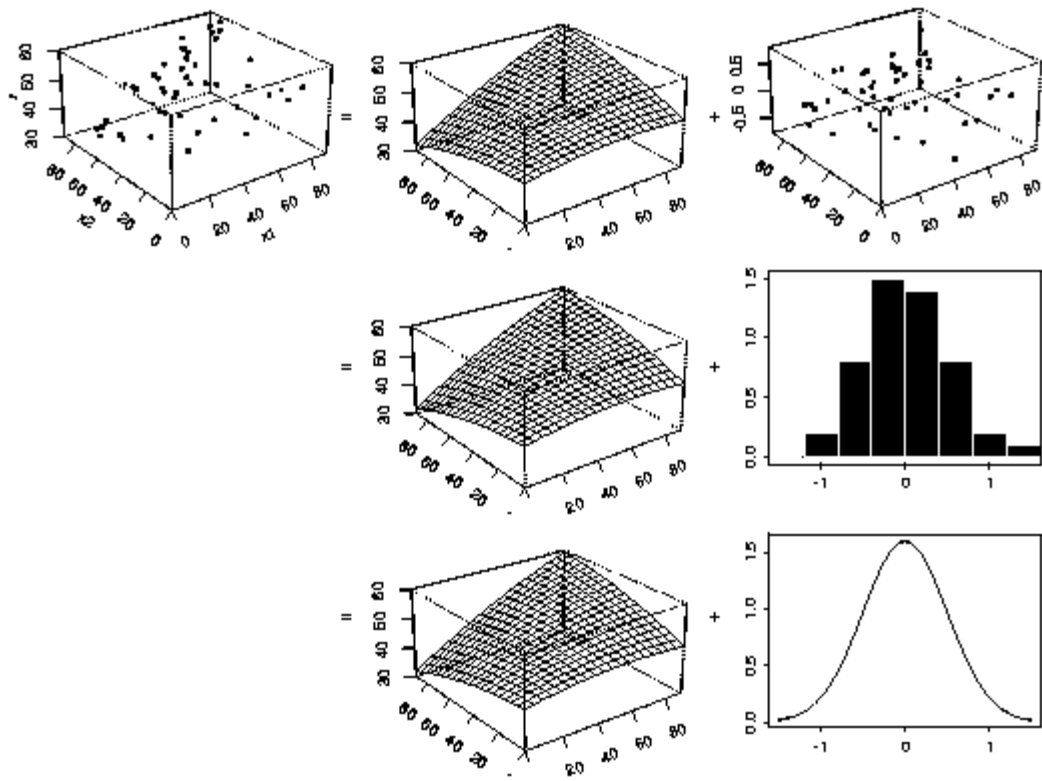
Of course, the straight-line example is one of the simplest functions used for process modeling. Another example is shown below. The concept is identical to the straight-line example, but the structure in the data is more complex. The variation in y is partitioned into a deterministic part, which is a function of another variable, x , plus some left-over random variation. (Again note the difference in the vertical axis scales of the two plots in the top right of the figure.) A probability distribution describes the leftover random variation.

4.1.1. What is process modeling?



An Example with Multiple Explanatory Variables

The examples of process modeling shown above have only one explanatory variable but the concept easily extends to cases with more than one explanatory variable. The three-dimensional perspective plots below show an example with two explanatory variables. Examples with three or more explanatory variables are exactly analogous, but are difficult to show graphically.





[4. Process Modeling](#)

[4.1. Introduction to Process Modeling](#)

4.1.2. What terminology do statisticians use to describe process models?

Model Components

There are three main parts to every process model. These are

1. the response variable, usually denoted by y ,
2. the mathematical function, usually denoted as $f(\vec{x}; \vec{\beta})$, and
3. the random errors, usually denoted by ϵ .

Form of Model

The general form of the model is

$$y = f(\vec{x}; \vec{\beta}) + \epsilon.$$

All process models discussed in this chapter have this general form. As alluded to [earlier](#), the random errors that are included in the model make the relationship between the response variable and the predictor variables a "statistical" one, rather than a perfect deterministic one. This is because the functional relationship between the response and predictors holds only on average, not for each data point.

Some of the details about the different parts of the model are discussed below, along with alternate terminology for the different components of the model.

Response Variable

The response variable, y , is a quantity that varies in a way that we hope to be able to summarize and exploit via the modeling process. Generally it is known that the variation of the response variable is systematically related to the values of one or more other variables before the modeling process is begun, although testing the existence and nature of this dependence is part of the modeling process itself.

Mathematical Function

The mathematical function consists of two parts. These parts are the predictor variables, x_1, x_2, \dots and the parameters, β_0, β_1, \dots . The predictor variables are observed along with the response variable. They are the quantities described on the previous page as inputs to the mathematical function, $f(\vec{x}; \vec{\beta})$. The collection of all of the predictor variables is denoted by \vec{x}

for short.

$$\vec{x} \equiv (x_1, x_2, \dots)$$

The parameters are the quantities that will be estimated during the modeling process. Their true values are unknown and unknowable, except in simulation experiments. As for the predictor variables, the collection of all of the parameters is denoted by $\vec{\beta}$ for short.

$$\vec{\beta} \equiv (\beta_0, \beta_1, \dots)$$

The parameters and predictor variables are combined in different forms to give the function used to describe the deterministic variation in the response variable. For a straight line with an unknown intercept and slope, for example, there are two parameters and one predictor variable

$$f(x; \vec{\beta}) = \beta_0 + \beta_1 x.$$

For a straight line with a known slope of one, but an unknown intercept, there would only be one parameter

$$f(x; \vec{\beta}) = \beta_0 + x.$$

For a quadratic surface with two predictor variables, there are six parameters for the full model.

$$f(\vec{x}; \vec{\beta}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2 + \beta_{11} x_1^2 + \beta_{22} x_2^2$$

Random Error

Like the parameters in the mathematical function, the random errors are unknown. They are simply the difference between the data and the mathematical function. They are assumed to follow a particular probability distribution, however, which is used to describe their aggregate behavior. The probability distribution that describes the errors has a mean of zero and an unknown standard deviation, denoted by σ , that is another parameter in the model, like the β 's.

Alternate Terminology

Unfortunately, there are no completely standardized names for the parts of the model discussed [above](#). Other publications or software may use different terminology. For example, another common name for the response variable is "dependent variable". The response variable is also simply called "the response" for short. Other names for the predictor variables include "explanatory variables", "independent variables", "predictors" and "regressors". The mathematical function used to describe the deterministic variation in the response variable is sometimes called the "regression function", the "regression equation", the "smoothing function", or the "smooth".

*Scope of
"Model"*

In its correct usage, the term "model" refers to the equation [above](#) and also includes the underlying assumptions made about the probability distribution used to describe the variation of the random errors. Often, however, people will also use the term "model" when referring specifically to the mathematical function describing the deterministic variation in the data. Since the function is part of the model, the more limited usage is not wrong, but it is important to remember that the term "model" might refer to more than just the mathematical function.



[4. Process Modeling](#)

[4.1. Introduction to Process Modeling](#)

4.1.3. What are process models used for?

Three Main Purposes Process models are used for four main purposes:

1. estimation,
2. prediction,
3. calibration, and
4. optimization.

The rest of this page lists brief explanations of the different uses of process models. More detailed explanations of the uses for process models are given in the subsections of this section listed at the bottom of this page.

Estimation The goal of estimation is to determine the value of the [regression function](#) (i.e., the average value of the response variable), for a particular combination of the values of the predictor variables. Regression function values can be estimated for any combination of predictor variable values, including values for which no data have been measured or observed. Function values estimated for points within the observed space of predictor variable values are sometimes called interpolations. Estimation of regression function values for points outside the observed space of predictor variable values, called extrapolations, are sometimes necessary, but require caution.

Prediction The goal of prediction is to determine either

1. the value of a new observation of the response variable, or
2. the values of a specified proportion of all future observations of the response variable

for a particular combination of the values of the predictor variables. Predictions can be made for any combination of predictor variable values, including values for which no data have been measured or observed. As in the case of estimation, predictions made outside the observed space of predictor variable values are sometimes necessary, but require caution.

Calibration The goal of calibration is to quantitatively relate measurements made using one measurement system to those of another measurement system. This is done so that

measurements can be compared in common units or to tie results from a relative measurement method to absolute units.

Optimization Optimization is performed to determine the values of process inputs that should be used to obtain the desired process output. Typical optimization goals might be to maximize the yield of a process, to minimize the processing time required to fabricate a product, or to hit a target product specification with minimum variation in order to maintain specified tolerances.

*Further
Details*

1. [Estimation](#)
2. [Prediction](#)
3. [Calibration](#)
4. [Optimization](#)



[4. Process Modeling](#)

[4.1. Introduction to Process Modeling](#)

4.1.4. What are some of the different statistical methods for model building?

*Selecting an
Appropriate
Stat
Method:
General
Case*

For many types of data analysis problems there are no more than a couple of general approaches to be considered on the route to the problem's solution. For example, there is often a dichotomy between highly-efficient methods appropriate for data with noise from a normal distribution and more general methods for data with other types of noise. Within the different approaches for a specific problem type, there are usually at most a few competing statistical tools that can be used to obtain an appropriate solution. The bottom line for most types of data analysis problems is that selection of the best statistical method to solve the problem is largely determined by the goal of the analysis and the nature of the data.

*Selecting an
Appropriate
Stat
Method:
Modeling*

Model building, however, is different from most other areas of statistics with regard to method selection. There are more general approaches and more competing techniques available for model building than for most other types of problems. There is often more than one statistical tool that can be effectively applied to a given modeling application. The large menu of methods applicable to modeling problems means that there is both more opportunity for effective and efficient solutions and more potential to spend time doing different analyses, comparing different solutions and mastering the use of different tools. The remainder of this section will introduce and briefly discuss some of the most popular and well-established statistical techniques that are useful for different model building situations.

*Process
Modeling
Methods*

1. [Linear Least Squares Regression](#)
2. [Nonlinear Least Squares Regression](#)
3. [Weighted Least Squares Regression](#)
4. [LOESS \(aka LOWESS\)](#)



[4. Process Modeling](#)

4.2. Underlying Assumptions for Process Modeling

- Implicit Assumptions Underlie Most Actions* Most, if not all, thoughtful actions that people take are based on ideas, or assumptions, about how those actions will affect the goals they want to achieve. The actual assumptions used to decide on a particular course of action are rarely laid out explicitly, however. Instead, they are only implied by the nature of the action itself. Implicit assumptions are inherent to process modeling actions, just as they are to most other types of action. It is important to understand what the implicit assumptions are for any process modeling method because the validity of these assumptions affect whether or not the goals of the analysis will be met.
- Checking Assumptions Provides Feedback on Actions* If the implicit assumptions that underlie a particular action are not true, then that action is not likely to meet expectations either. Sometimes it is abundantly clear when a goal has been met, but unfortunately that is not always the case. In particular, it is usually not possible to obtain immediate feedback on the attainment of goals in most process modeling applications. The goals of process modeling, such as answering a scientific or engineering question, depend on the correctness of a process model, which can often only be directly and absolutely determined over time. In lieu of immediate, direct feedback, however, indirect information on the effectiveness of a process modeling analysis can be obtained by checking the validity of the underlying assumptions. Confirming that the underlying assumptions are valid helps ensure that the methods of analysis were appropriate and that the results will be consistent with the goals.
- Overview of Section 4.2* This section discusses the specific underlying assumptions associated with most model-fitting methods. In discussing the underlying assumptions, some background is also provided on the consequences of stopping the modeling process short of completion and leaving the results of an analysis at odds with the underlying assumptions. Specific data analysis methods that can be used to check whether or not the assumptions hold in a particular case are discussed in [Section 4.4.4](#).
- Contents of Section 4.2*
1. [What are the typical underlying assumptions in process modeling?](#)

1. [The process is a *statistical* process.](#)
2. [The means of the random errors are zero.](#)
3. [The random errors have a constant standard deviation.](#)
4. [The random errors follow a normal distribution.](#)
5. [The data are randomly sampled from the process.](#)
6. [The explanatory variables are observed without error.](#)

[4. Process Modeling](#)[4.2. Underlying Assumptions for Process Modeling](#)

4.2.1. What are the typical underlying assumptions in process modeling?

*Overview of
Section
4.2.1*

This section lists the typical assumptions underlying most process modeling methods. On each of the following pages, one of the six major assumptions is described individually; the reasons for its importance are also briefly discussed; and any methods that are not subject to that particular assumption are noted. As discussed on the [previous page](#), these are implicit assumptions based on properties inherent to the process modeling methods themselves. Successful use of these methods in any particular application hinges on the validity of the underlying assumptions, whether their existence is acknowledged or not. [Section 4.4.4](#) discusses methods for checking the validity of these assumptions.

*Typical
Assumptions
for Process
Modeling*

1. [The process is a *statistical* process.](#)
2. [The means of the random errors are zero.](#)
3. [The random errors have a constant standard deviation.](#)
4. [The random errors follow a normal distribution.](#)
5. [The data are randomly sampled from the process.](#)
6. [The explanatory variables are observed without error.](#)

[4. Process Modeling](#)

4.3. Data Collection for Process Modeling

*Collecting
Good Data*

This section lays out some general principles for collecting data for construction of process models. Using well-planned data collection procedures is often the difference between successful and unsuccessful experiments. In addition, well-designed experiments are often less expensive than those that are less well thought-out, regardless of overall success or failure.

Specifically, this section will answer the question:

What can the analyst do even prior to collecting the data (that is, at the experimental design stage) that would allow the analyst to do an optimal job of modeling the process?

*Contents:
Section 3*

This section deals with the following five questions:

1. [What is design of experiments \(DOE\)?](#)
2. [Why is experimental design important for process modeling?](#)
3. [What are some general design principles for process modeling?](#)
4. [I've heard some people refer to "optimal" designs, shouldn't I use those?](#)
5. [How can I tell if a particular experimental design is good for my application?](#)



[4. Process Modeling](#)

[4.3. Data Collection for Process Modeling](#)

4.3.1. What is design of experiments (DOE)?

Systematic Approach to Data Collection

Design of experiments (DOE) is a systematic, rigorous approach to engineering problem-solving that applies principles and techniques at the data collection stage so as to ensure the generation of valid, defensible, and supportable engineering conclusions. In addition, all of this is carried out under the constraint of a minimal expenditure of engineering runs, time, and money.

DOE Problem Areas

There are four general engineering problem areas in which DOE may be applied:

1. Comparative
2. Screening/Characterizing
3. Modeling
4. Optimizing

Comparative

In the first case, the engineer is interested in assessing whether a change in a single factor has in fact resulted in a change/improvement to the process as a whole.

Screening Characterization

In the second case, the engineer is interested in "understanding" the process as a whole in the sense that he/she wishes (after design and analysis) to have in hand a ranked list of important through unimportant factors (most important to least important) that affect the process.

Modeling

In the third case, the engineer is interested in functionally modeling the process with the output being a good-fitting (= high predictive power) mathematical function, and to have good (= maximal accuracy) estimates of the coefficients in that function.

Optimizing

In the fourth case, the engineer is interested in determining optimal settings of the process factors; that is, to determine for each factor the level of the factor that optimizes the process response.

In this section, we focus on case 3: modeling.

[4. Process Modeling](#)

[4.3. Data Collection for Process Modeling](#)

4.3.2. Why is experimental design important for process modeling?

Output from Process Model is Fitted Mathematical Function

The output from process modeling is a fitted mathematical function with estimated coefficients. For example, in modeling resistivity, y , as a function of dopant density, x , an analyst may suggest the function

$$y = \beta_0 + \beta_1 x + \beta_{11} x^2 + \epsilon$$

in which the coefficients to be estimated are β_0 , β_1 , and β_{11} . Even for a given functional form, there is an infinite number of potential coefficient values that potentially may be used. Each of these coefficient values will in turn yield predicted values.

What are Good Coefficient Values?

Poor values of the coefficients are those for which the resulting predicted values are considerably different from the observed raw data y . Good values of the coefficients are those for which the resulting predicted values are close to the observed raw data y . The best values of the coefficients are those for which the resulting predicted values are close to the observed raw data y , and the statistical uncertainty connected with each coefficient is small.

There are two considerations that are useful for the generation of "best" coefficients:

1. Least squares criterion
2. Design of experiment principles

Least Squares Criterion

For a given data set (e.g., 10 (x, y) pairs), the most common procedure for obtaining the coefficients for

$$y = f(x; \vec{\beta}) + \epsilon$$

is the [least squares estimation criterion](#). This criterion yields coefficients with predicted values that are closest to the raw data y in the sense that the sum of the squared differences between the raw data and the predicted values is as small as possible.

The overwhelming majority of regression programs today use the least squares criterion for estimating the model coefficients. Least squares estimates are popular because

1. the estimators are statistically optimal (BLUEs: Best Linear Unbiased Estimators);
2. the estimation algorithm is mathematically tractable, in closed form, and therefore easily programmable.

How then can this be improved? For a given set of x values it cannot be; but frequently the choice of the x values is under our control. If we can select the x values, the coefficients will have less variability than if the x are not controlled.

Design of Experiment Principles

As to what values should be used for the x 's, we look to established experimental design principles for guidance.

Principle 1: Minimize Coefficient Estimation Variation

The first principle of experimental design is to control the values within the x vector such that after the y data are collected, the subsequent model coefficients are as good, in the sense of having the smallest variation, as possible.

The key underlying point with respect to design of experiments and process modeling is that even though (for simple (x,y) fitting, for example) the least squares criterion may yield optimal (minimal variation) estimators for a given distribution of x values, some distributions of data in the x vector may yield better (smaller variation) coefficient estimates than other x vectors. If the analyst can specify the values in the x vector, then he or she may be able to drastically change and reduce the noisiness of the subsequent least squares coefficient estimates.

Five Designs

To see the effect of experimental design on process modeling, consider the following simplest case of fitting a line:

$$y = \beta_0 + \beta_1 x + \epsilon$$

Suppose the analyst can afford 10 observations (that is, 10 (x,y) pairs) for the purpose of determining optimal (that is, minimal variation) estimators of β_0 and β_1 . What 10 x values should be used for the purpose of collecting the corresponding 10 y values? Colloquially, where should the 10 x values be sprinkled along the horizontal axis so as to minimize the variation of the least squares estimated coefficients for β_0 and β_1 ? Should the 10 x values be:

1. ten equi-spaced values across the range of interest?
2. five replicated equi-spaced values across the range of

interest?

3. five values at the minimum of the x range and five values at the maximum of the x range?
4. one value at the minimum, eight values at the mid-range, and one value at the maximum?
5. four values at the minimum, two values at mid-range, and four values at the maximum?

or (in terms of "quality" of the resulting estimates for β_0 and β_1) perhaps it doesn't make any difference?

For each of the above five experimental designs, there will of course be y data collected, followed by the generation of least squares estimates for β_0 and β_1 , and so each design will in turn yield a fitted line.

Are the Fitted Lines Better for Some Designs?

But are the fitted lines, i.e., the fitted process models, better for some designs than for others? Are the coefficient estimator variances smaller for some designs than for others? For given estimates, are the resulting predicted values better (that is, closer to the observed y values) than for other designs? The answer to all of the above is YES. It DOES make a difference.

The most popular answer to the above question about which design to use for linear modeling is design #1 with ten equi-spaced points. It can be shown, however, that the variance of the estimated slope parameter depends on the design according to the relationship

$$\text{Var}(\hat{\beta}_1) \propto \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

Therefore to obtain minimum variance estimators, one maximizes the denominator on the right. To maximize the denominator, it is (for an arbitrarily fixed \bar{x}), best to position the x 's as far away from \bar{x} as possible. This is done by positioning half of the x 's at the lower extreme and the other half at the upper extreme. This is design #3 above, and this "dumbbell" design (half low and half high) is in fact the best possible design for fitting a line. Upon reflection, this is intuitively arrived at by the adage that "2 points define a line", and so it makes the most sense to determine those 2 points as far apart as possible (at the extremes) and as well as possible (having half the data at each extreme). Hence the design of experiment solution to model processing when the model is a line is the "dumbbell" design--half the X's at each extreme.

What is the Worst

What is the worst design in the above case? Of the five designs, the worst design is the one that has maximum

Design? variation. In the mathematical expression above, it is the one that minimizes the denominator, and so this is design #4 above, for which almost all of the data are located at the mid-range. Clearly the estimated line in this case is going to chase the solitary point at each end and so the resulting linear fit is intuitively inferior.

Designs 1, 2, and 5 What about the other 3 designs? Designs 1, 2, and 5 are useful only for the case when we think the model may be linear, but we are not sure, and so we allow additional points that permit fitting a line if appropriate, but build into the design the "capacity" to fit beyond a line (e.g., quadratic, cubic, etc.) if necessary. In this regard, the ordering of the designs would be

- design 5 (if our worst-case model is quadratic),
- design 2 (if our worst-case model is quartic)
- design 1 (if our worst-case model is quintic and beyond)



[4. Process Modeling](#)

[4.3. Data Collection for Process Modeling](#)

4.3.3. What are some general design principles for process modeling?

Experimental Design Principles Applied to Process Modeling

There are six principles of experimental design as applied to process modeling:

1. Capacity for Primary Model
2. Capacity for Alternative Model
3. Minimum Variance of Coefficient Estimators
4. Sample where the Variation Is
5. Replication
6. Randomization

We discuss each in detail below.

Capacity for Primary Model

For your best-guess model, make sure that the design has the capacity for estimating the coefficients of that model. For a simple example of this, if you are fitting a quadratic model, then make sure you have at least three distinct horizontal axis points.

Capacity for Alternative Model

If your best-guess model happens to be inadequate, make sure that the design has the capacity to estimate the coefficients of your best-guess back-up alternative model (which means implicitly that you should have already identified such a model). For a simple example, if you suspect (but are not positive) that a linear model is appropriate, then it is best to employ a globally robust design (say, four points at each extreme and three points in the middle, for a ten-point design) as opposed to the locally optimal design (such as five points at each extreme). The locally optimal design will provide a best fit to the line, but have no capacity to fit a quadratic. The globally robust design will provide a good (though not optimal) fit to the line and additionally provide a good (though not optimal) fit to the quadratic.

Minimum Variance of Coefficient Estimators

For a given model, make sure the design has the property of minimizing the variation of the least squares estimated coefficients. This is a general principle that is always in effect but which in practice is hard to implement for many models beyond the simpler 1-factor

$$y = f(x; \vec{\beta}) + \epsilon$$

models. For more complicated 1-factor models, and for most multi-factor

$$y = f(\vec{x}; \vec{\beta}) + \epsilon$$

models, the expressions for the variance of the least squares estimators, although available, are complicated and assume more than the analyst typically knows. The net result is that this principle, though important, is harder to apply beyond the simple cases.

Sample Where the Variation Is (Non Constant Variance Case)

Regardless of the simplicity or complexity of the model, there are situations in which certain regions of the curve are noisier than others. A simple case is when there is a linear relationship between x and y but the recording device is proportional rather than absolute and so larger values of y are intrinsically noisier than smaller values of y . In such cases, sampling where the variation is means to have more replicated points in those regions that are noisier. The practical answer to how many such replicated points there should be is

$$n_i = \frac{1}{\sigma_i^2}$$

with σ_i denoting the theoretical standard deviation for that given region of the curve. Usually σ_i is estimated by a-priori guesses for what the local standard deviations are.

Sample Where the Variation Is (Steep Curve Case)

A common occurrence for non-linear models is for some regions of the curve to be steeper than others. For example, in fitting an exponential model (small x corresponding to large y , and large x corresponding to small y) it is often the case that the y data in the steep region are intrinsically noisier than the y data in the relatively flat regions. The reason for this is that commonly the x values themselves have a bit of noise and this x -noise gets translated into larger y -noise in the steep sections than in the shallow sections. In such cases, when we know the shape of the response curve well enough to identify steep-versus-shallow regions, it is often a good idea to sample more heavily in the steep regions than in the shallow regions. A practical rule-of-thumb for where to position the x values in such situations is to

1. sketch out your best guess for what the resulting curve will be;
2. partition the vertical (that is the y) axis into n equi-spaced points (with n denoting the total number of

- data points that you can afford);
3. draw horizontal lines from each vertical axis point to where it hits the sketched-in curve.
4. drop a vertical projection line from the curve intersection point to the horizontal axis.

These will be the recommended x values to use in the design.

The above rough procedure for an exponentially decreasing curve would thus yield a logarithmic preponderance of points in the steep region of the curve and relatively few points in the flatter part of the curve.

Replication

If affordable, replication should be part of every design. Replication allows us to compute a model-independent estimate of the process standard deviation. Such an estimate may then be used as a criterion in an objective [lack-of-fit test](#) to assess whether a given model is adequate. Such an objective lack-of-fit F-test can be employed only if the design has built-in replication. Some replication is essential; replication at every point is ideal.

Randomization

Just because the x 's have some natural ordering does not mean that the data should be collected in the same order as the x 's. Some aspect of randomization should enter into every experiment, and experiments for process modeling are no exception. Thus if you are sampling ten points on a curve, the ten y values should not be collected by sequentially stepping through the x values from the smallest to the largest. If you do so, and if some extraneous drifting or wear occurs in the machine, the operator, the environment, the measuring device, etc., then that drift will unwittingly contaminate the y values and in turn contaminate the final fit. To minimize the effect of such potential drift, it is best to randomize (use random number tables) the sequence of the x values. This will not make the drift go away, but it will spread the drift effect evenly over the entire curve, realistically inflating the variation of the fitted values, and providing some mechanism after the fact (at the residual analysis model validation stage) for uncovering or discovering such a drift. If you do not randomize the run sequence, you give up your ability to detect such a drift if it occurs.



[4. Process Modeling](#)

[4.3. Data Collection for Process Modeling](#)

4.3.4. I've heard some people refer to "optimal" designs, shouldn't I use those?

Classical Designs Heavily Used in Industry

The most heavily used designs in industry are the "classical designs" (full factorial designs, fractional factorial designs, Latin square designs, Box-Behnken designs, etc.). They are so heavily used because they are optimal in their own right and have served superbly well in providing efficient insight into the underlying structure of industrial processes.

Reasons Classical Designs May Not Work

Cases do arise, however, for which the tabulated classical designs do not cover a particular practical situation. That is, user constraints preclude the use of tabulated classical designs because such classical designs do not accommodate user constraints. Such constraints include:

1. Limited maximum number of runs:

User constraints in budget and time may dictate a maximum allowable number of runs that is too small or too "irregular" (e.g., "13") to be accommodated by classical designs--even fractional factorial designs.

2. Impossible factor combinations:

The user may have some factor combinations that are impossible to run. Such combinations may at times be specified (to maintain balance and orthogonality) as part of a recommended classical design. If the user simply omits this impossible run from the design, the net effect may be a reduction in the quality and optimality of the classical design.

3. Too many levels:

The number of factors and/or the number of levels of some factors intended for use may not be included in tabulations of classical designs.

4. Complicated underlying model:

The user may be assuming an underlying model that is too complicated (or too non-linear), so that classical designs would be inappropriate.

*What to Do If
Classical
Designs Do Not
Exist?*

If user constraints are such that classical designs do not exist to accommodate such constraints, then what is the user to do?

The previous section's list of design criteria (capability for the primary model, capability for the alternate model, minimum variation of estimated coefficients, etc.) is a good passive target to aim for in terms of desirable design properties, but provides little help in terms of an active formal construction methodology for generating a design.

*Common
Optimality
Criteria*

To satisfy this need, an "optimal design" methodology has been developed to generate a design when user constraints preclude the use of tabulated classical designs. Optimal designs may be optimal in many different ways, and what may be an optimal design according to one criterion may be suboptimal for other criteria. Competing criteria have led to a literal alphabet-soup collection of optimal design methodologies. The four most popular ingredients in that "soup" are:

- D-optimal designs: minimize the generalized variance of the parameter estimators.
- A-optimal designs: minimize the average variance of the parameter estimators.
- G-optimal designs: minimize the maximum variance of the predicted values.
- V-optimal designs: minimize the average variance of the predicted values.

Need 1: a Model

The motivation for optimal designs is the practical constraints that the user has. The advantage of optimal designs is that they do provide a reasonable design-generating methodology when no other mechanism exists. The disadvantage of optimal designs is that they require a model from the user. The user may not have this model.

All optimal designs are model-dependent, and so the quality of the final engineering conclusions that result from the ensuing design, data, and analysis is dependent on the correctness of the analyst's assumed model. For example, if the responses from a particular process are actually being drawn from a cubic model and the analyst assumes a linear model and uses the corresponding optimal design to generate data and perform the data

analysis, then the final engineering conclusions will be flawed and invalid. Hence one price for obtaining an in-hand generated design is the designation of a model. All optimal designs need a model; without a model, the optimal design-generation methodology cannot be used, and general design principles must be reverted to.

*Need 2: a
Candidate Set of
Points*

The other price for using optimal design methodology is a user-specified set of candidate points. Optimal designs will not generate the best design points from some continuous region--that is too much to ask of the mathematics. Optimal designs will generate the best subset of n points from a larger superset of candidate points. The user must specify this candidate set of points. Most commonly, the superset of candidate points is the full factorial design over a fine-enough grid of the factor space with which the analyst is comfortable. If the grid is too fine, and the resulting superset overly large, then the optimal design methodology may prove computationally challenging.

*Optimal
Designs are
Computationally
Intensive*

The optimal design-generation methodology is computationally intensive. Some of the designs (e.g., D-optimal) are better than other designs (such as A-optimal and G-optimal) in regard to efficiency of the underlying search algorithm. Like most mathematical optimization techniques, there is no iron-clad guarantee that the result from the optimal design methodology is in fact the true optimum. However, the results are usually satisfactory from a practical point of view, and are far superior than any ad hoc designs.

For further details about optimal designs, the analyst is referred to [Montgomery \(2001\)](#).



4. [Process Modeling](#)

4.3. [Data Collection for Process Modeling](#)

4.3.5. How can I tell if a particular experimental design is good for my application?

*Assess
Relative to
the Six
Design
Principles*

If you have a design, generated by whatever method, in hand, how can you assess its after-the-fact goodness? Such checks can potentially parallel the list of the [six general design principles](#). The design can be assessed relative to each of these six principles. For example, does it have capacity for the primary model, does it have capacity for an alternative model, etc.

Some of these checks are quantitative and complicated; other checks are simpler and graphical. The graphical checks are the most easily done and yet are among the most informative. We include two such graphical checks and one quantitative check.

*Graphically
Check for
Univariate
Balance*

If you have a design that claims to be globally good in k factors, then generally that design should be locally good in each of the individual k factors. Checking high-dimensional global goodness is difficult, but checking low-dimensional local goodness is easy. Generate k counts plots, with the levels of factors x_i plotted on the horizontal axis of each plot and the number of design points for each level in factor x_i on the vertical axis. For most good designs, these counts should be about the same (= balance) for all levels of a factor. Exceptions exist, but such balance is a low-level characteristic of most good designs.

*Graphically
Check for
Bivariate
Balance*

If you have a design that is purported to be globally good in k factors, then generally that design should be locally good in all pairs of the individual k factors. Graphically check for such 2-way balance by generating plots for all pairs of factors, where the horizontal axis of a given plot is x_i and the vertical axis is x_j . The response variable y does NOT come into play in these plots. We are only interested in characteristics of the design, and so only the x variables are involved. The 2-way plots of most good designs have a certain symmetric and balanced look about them--all combination points should be covered and each combination point should have about the same number of points.

Check for For optimal designs, metrics exist (D-efficiency, A-

*Minimal
Variation*

efficiency, etc.) that can be computed and that reflect the quality of the design. Further, relative ratios of standard deviations of the coefficient estimators and relative ratios of predicted values can be computed and compared for such designs. Such calculations are commonly performed in computer packages which specialize in the generation of optimal designs.



[HOME](#)

[TOOLS & AIDS](#)

[SEARCH](#)

[BACK](#) [NEXT](#)



4. [Process Modeling](#)

4.4. Data Analysis for Process Modeling

Building a Good Model This section contains detailed discussions of the necessary steps for developing a good process model after data have been collected. A general model-building framework, applicable to multiple statistical methods, is described with method-specific points included when necessary.

Contents: *Section 4*

1. [What are the basic steps for developing an effective process model?](#)
2. [How do I select a function to describe my process?](#)
 1. [Incorporating Scientific Knowledge into Function Selection](#)
 2. [Using the Data to Select an Appropriate Function](#)
 3. [Using Methods that Do Not Require Function Specification](#)
3. [How are estimates of the unknown parameters obtained?](#)
 1. [Least Squares](#)
 2. [Weighted Least Squares](#)
4. [How can I tell if a model fits my data?](#)
 1. [How can I assess the sufficiency of the functional part of the model?](#)
 2. [How can I detect non-constant variation across the data?](#)
 3. [How can I tell if there was drift in the measurement process?](#)
 4. [How can I assess whether the random errors are independent from one to the next?](#)
 5. [How can I test whether or not the random errors are normally distributed?](#)
 6. [How can I test whether any significant terms are missing or misspecified in the functional part of the model?](#)
 7. [How can I test whether all of the terms in the functional part of the model are necessary?](#)
5. [If my current model does not fit the data well, how can I improve it?](#)
 1. [Updating the Function Based on Residual Plots](#)
 2. [Accounting for Non-Constant Variation Across the Data](#)
 3. [Accounting for Errors with a Non-Normal Distribution](#)



[HOME](#)

[TOOLS & AIDS](#)

[SEARCH](#)

[BACK](#) [NEXT](#)



[4. Process Modeling](#)

[4.4. Data Analysis for Process Modeling](#)

4.4.1. What are the basic steps for developing an effective process model?

Basic Steps Provide Universal Framework The basic steps used for model-building are the same across all modeling methods. The details vary somewhat from method to method, but an understanding of the common steps, combined with the typical [underlying assumptions](#) needed for the analysis, provides a framework in which the results from almost any method can be interpreted and understood.

Basic Steps of Model Building The basic steps of the model-building process are:

1. model selection
2. model fitting, and
3. model validation.

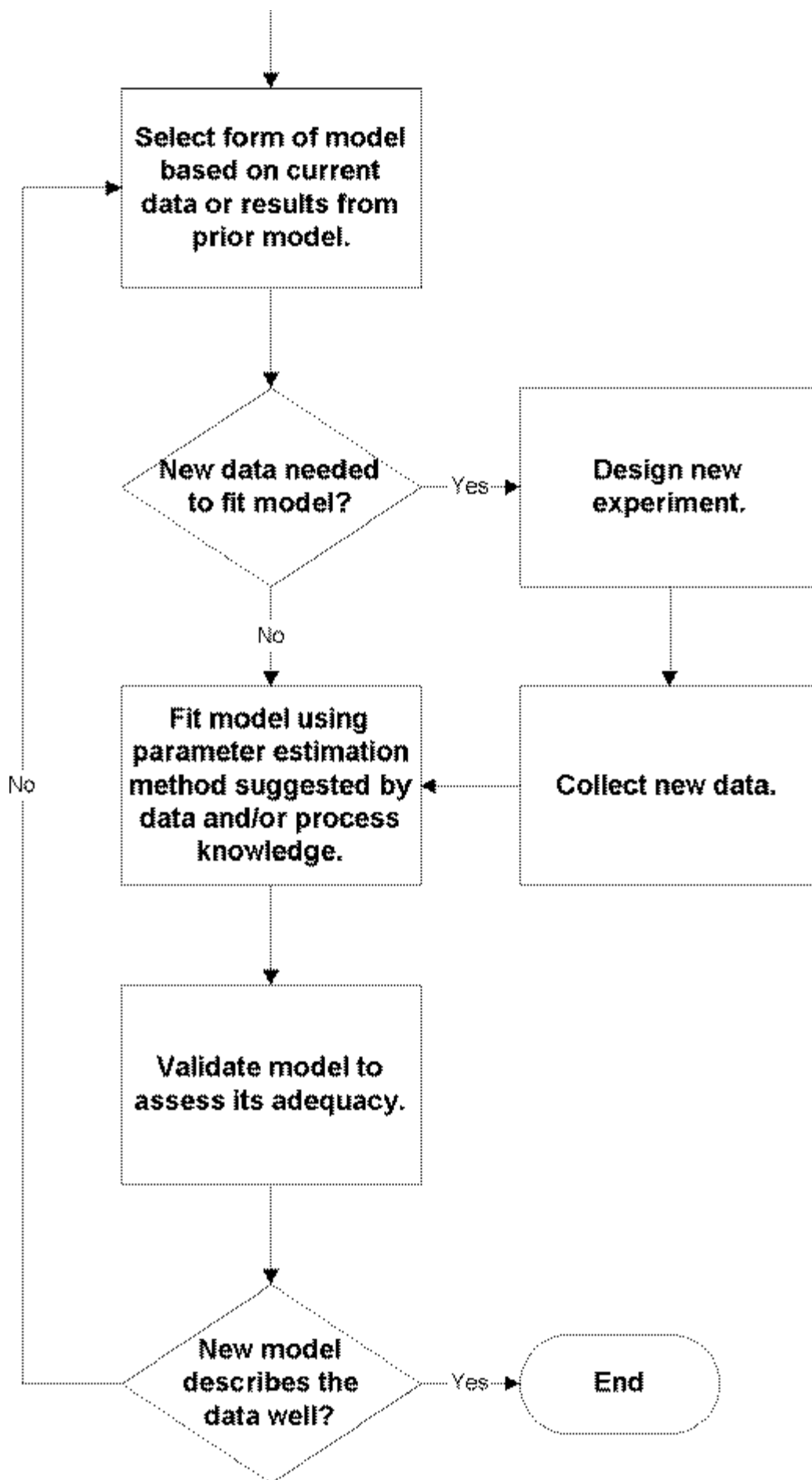
These three basic steps are used iteratively until an appropriate model for the data has been developed. In the model selection step, plots of the data, process knowledge and assumptions about the process are used to determine the form of the model to be fit to the data. Then, using the selected model and possibly information about the data, an appropriate model-fitting method is used to estimate the unknown parameters in the model. When the parameter estimates have been made, the model is then carefully assessed to see if the underlying assumptions of the analysis appear plausible. If the assumptions seem valid, the model can be used to answer the scientific or engineering questions that prompted the modeling effort. If the model validation identifies problems with the current model, however, then the modeling process is repeated using information from the model validation step to select and/or fit an improved model.

A Variation on the Basic Steps The three basic steps of process modeling described in the paragraph above assume that the data have already been collected and that the same data set can be used to fit all of the candidate models. Although this is often the case in model-building situations, one variation on the basic model-building sequence comes up when additional data are needed to fit a newly hypothesized model based on a model fit to the initial data. In this case two additional steps, [experimental design](#) and data collection, can be added to the basic sequence between model selection and model-fitting. The flow chart below shows the basic model-fitting sequence with the integration of the related data collection steps into the model-building process.



4.4.1. What are the basic steps for developing an effective process model?

Model Building Sequence



Examples illustrating the model-building sequence in real applications can be found in the case studies in [Section 4.6](#). The specific tools and techniques used in the basic model-building steps are described in the remainder of this section.

*Design of
Initial
Experiment*

Of course, considering the model selection and fitting before collecting the initial data is also a good idea. Without data in hand, a hypothesis about what the data will look like is needed in order to guess what the initial model should be. Hypothesizing the outcome of an experiment is not always possible, of course, but efforts made in the earliest stages of a project often maximize the efficiency of the whole model-building process and result in the best possible models for the process. More details about experimental design can be found in [Section 4.3](#) and in [Chapter 5: Process Improvement](#).



[4. Process Modeling](#)

[4.4. Data Analysis for Process Modeling](#)

4.4.2. How do I select a function to describe my process?

Synthesis of Process Information Necessary Selecting a model of the right form to fit a set of data usually requires the use of empirical evidence in the data, knowledge of the process and some trial-and-error experimentation. As mentioned on the previous page, model building is always an iterative process. Much of the need to iterate stems from the difficulty in initially selecting a function that describes the data well. Details about the data are often not easily visible in the data as originally observed. The fine structure in the data can usually only be elicited by use of model-building tools such as residual plots and repeated refinement of the model form. As a result, it is important not to overlook any of the sources of information that indicate what the form of the model should be.

Answer Not Provided by Statistics Alone Sometimes the different sources of information that need to be integrated to find an effective model will be contradictory. An open mind and a willingness to think about what the data are saying is important. Maintaining balance and looking for alternate sources for unusual effects found in the data are also important. For example, in the [load cell calibration case study](#) the statistical analysis pointed out that the model initially thought to be appropriate did not account for all of the structure in the data. A refined model was developed, but the appearance of an unexpected result brings up the question of whether the original understanding of the problem was inaccurate, or whether the need for an alternate model was due to experimental artifacts. In the load cell problem it was easy to accept that the refined model was closer to the truth, but in a more complicated case additional experiments might have been needed to resolve the issue.

Knowing Function Types Helps Another helpful ingredient in model selection is a wide knowledge of the shapes that different mathematical functions can assume. Knowing something about the models that have been found to work well in the past for different application types also helps. A menu of different functions on the next page, Section 4.4.2.1. (links provided below), provides one way to learn about the function shapes and flexibility. Section 4.4.2.2. discusses how general function features and qualitative scientific information can be combined to help with model selection. Finally, Section 4.4.2.3. points to

methods that don't require specification of a particular function to be fit to the data, and how models of those types can be refined.

1. [Incorporating Scientific Knowledge into Function Selection](#)
2. [Using the Data to Select an Appropriate Function](#)
3. [Using Methods that Do Not Require Function Specification](#)



[HOME](#)

[TOOLS & AIDS](#)

[SEARCH](#)

[BACK](#) [NEXT](#)



[4. Process Modeling](#)

[4.4. Data Analysis for Process Modeling](#)

4.4.3. How are estimates of the unknown parameters obtained?

Parameter Estimation in General

After selecting the basic form of the functional part of the model, the next step in the model-building process is estimation of the unknown parameters in the function. In general, this is accomplished by solving an optimization problem in which the objective function (the function being minimized or maximized) relates the response variable and the functional part of the model containing the unknown parameters in a way that will produce parameter estimates that will be close to the true, unknown parameter values. The unknown parameters are, loosely speaking, treated as variables to be solved for in the optimization, and the data serve as known coefficients of the objective function in this stage of the modeling process.

In theory, there are as many different ways of estimating parameters as there are objective functions to be minimized or maximized. However, a few principles have dominated because they result in parameter estimators that have good statistical properties. The two major methods of parameter estimation for process models are maximum likelihood and least squares. Both of these methods provide parameter estimators that have many good properties. Both maximum likelihood and least squares are sensitive to the presence of outliers, however. There are also many newer methods of parameter estimation, called robust methods, that try to balance the efficiency and desirable properties of least squares and maximum likelihood with a lower sensitivity to outliers.

Overview of Section 4.3

Although robust techniques are valuable, they are not as well developed as the more traditional methods and often require specialized software that is not readily available. Maximum likelihood also requires specialized algorithms in general, although there are important special cases that do not have such a requirement. For example, for data with normally distributed random errors, the least squares and maximum likelihood parameter estimators are identical. As a result of these software and developmental issues, and the coincidence of maximum likelihood and least squares in many applications, this section currently focuses on parameter estimation only by least squares methods. The remainder of this section offers some intuition into how least squares works

and illustrates the effectiveness of this method.

*Contents
of Section
4.3*

1. [Least Squares](#)
2. [Weighted Least Squares](#)



[HOME](#)

[TOOLS & AIDS](#)

[SEARCH](#)

[BACK](#) [NEXT](#)

[4. Process Modeling](#)
[4.4. Data Analysis for Process Modeling](#)

4.4.4. How can I tell if a model fits my data?

R² Is Not Enough!

Model validation is possibly the most important step in the model building sequence. It is also one of the most overlooked. Often the validation of a model seems to consist of nothing more than quoting the R^2 statistic from the fit (which measures the fraction of the total variability in the response that is accounted for by the model). Unfortunately, a high R^2 value does not guarantee that the model fits the data well. Use of a model that does not fit the data well cannot provide good answers to the underlying engineering or scientific questions under investigation.

Main Tool: Graphical Residual Analysis

There are many statistical tools for model validation, but the primary tool for most process modeling applications is graphical residual analysis. Different types of plots of the residuals ([see definition below](#)) from a fitted model provide information on the adequacy of different aspects of the model. Numerical methods for model validation, such as the R^2 statistic, are also useful, but usually to a lesser degree than graphical methods. Graphical methods have an advantage over numerical methods for model validation because they readily illustrate a broad range of complex aspects of the relationship between the model and the data. Numerical methods for model validation tend to be narrowly focused on a particular aspect of the relationship between the model and the data and often try to compress that information into a single descriptive number or test result.

Numerical Methods' Forte

Numerical methods do play an important role as confirmatory methods for graphical techniques, however. For example, the [lack-of-fit test](#) for assessing the correctness of the functional part of the model can aid in interpreting a borderline residual plot. There are also a few modeling situations in which graphical methods cannot easily be used. In these cases, numerical methods provide a fallback position for model validation. One common situation when numerical validation methods take precedence over graphical methods is when the number of parameters being estimated is relatively close to the size of the data set. In this situation residual plots are often difficult to interpret due to constraints on the residuals imposed by the estimation of the unknown parameters. One area in which this typically happens is in optimization applications using designed experiments. Logistic regression

with binary data is another area in which graphical residual analysis can be difficult.

Residuals The residuals from a fitted model are the differences between the responses observed at each combination values of the explanatory variables and the corresponding prediction of the response computed using the regression function.

Mathematically, the definition of the residual for the i^{th} observation in the data set is written

$$e_i = y_i - f(\vec{x}_i; \vec{\beta}),$$

with y_i denoting the i^{th} response in the data set and \vec{x}_i represents the list of explanatory variables, each set at the corresponding values found in the i^{th} observation in the data set.

Example The data listed below are from the [Pressure/Temperature example](#) introduced in [Section 4.1.1](#). The first column shows the order in which the observations were made, the second column indicates the day on which each observation was made, and the third column gives the ambient temperature recorded when each measurement was made. The fourth column lists the temperature of the gas itself (the explanatory variable) and the fifth column contains the observed pressure of the gas (the response variable). Finally, the sixth column gives the corresponding values from the fitted straight-line regression function.

$$\hat{P} = 7.749695 + 3.930123T$$

and the last column lists the residuals, the difference between columns five and six.

*Data,
Fitted
Values &
Residuals*

Run Fitted Order	Day	Ambient Temperature	Temperature	Pressure	Residual
1	1	23.820	54.749	225.066	
222.920		2.146			
2	1	24.120	23.323	100.331	
99.411		0.920			
3	1	23.434	58.775	230.863	
238.744		-7.881			
4	1	23.993	25.854	106.160	
109.359		-3.199			
5	1	23.375	68.297	277.502	
276.165		1.336			
6	1	23.233	37.481	148.314	
155.056		-6.741			
7	1	24.162	49.542	197.562	
202.456		-4.895			
8	1	23.667	34.101	138.537	
141.770		-3.232			
9	1	24.056	33.901	137.969	
140.983		-3.014			
10	1	22.786	29.242	117.410	
122.674		-5.263			
11	2	23.785	39.506	164.442	
163.013		1.429			
12	2	22.987	43.004	181.044	

4.4.4. How can I tell if a model fits my data?

176.759		4.285		
13	2	23.799	53.226	222.179
216.933		5.246		
14	2	23.661	54.467	227.010
221.813		5.198		
15	2	23.852	57.549	232.496
233.925		-1.429		
16	2	23.379	61.204	253.557
248.288		5.269		
17	2	24.146	31.489	139.894
131.506		8.388		
18	2	24.187	68.476	273.931
276.871		-2.940		
19	2	24.159	51.144	207.969
208.753		-0.784		
20	2	23.803	68.774	280.205
278.040		2.165		
21	3	24.381	55.350	227.060
225.282		1.779		
22	3	24.027	44.692	180.605
183.396		-2.791		
23	3	24.342	50.995	206.229
208.167		-1.938		
24	3	23.670	21.602	91.464
92.649		-1.186		
25	3	24.246	54.673	223.869
222.622		1.247		
26	3	25.082	41.449	172.910
170.651		2.259		
27	3	24.575	35.451	152.073
147.075		4.998		
28	3	23.803	42.989	169.427
176.703		-7.276		
29	3	24.660	48.599	192.561
198.748		-6.188		
30	3	24.097	21.448	94.448
92.042		2.406		
31	4	22.816	56.982	222.794
231.697		-8.902		
32	4	24.167	47.901	199.003
196.008		2.996		
33	4	22.712	40.285	168.668
166.077		2.592		
34	4	23.611	25.609	109.387
108.397		0.990		
35	4	23.354	22.971	98.445
98.029		0.416		
36	4	23.669	25.838	110.987
109.295		1.692		
37	4	23.965	49.127	202.662
200.826		1.835		
38	4	22.917	54.936	224.773
223.653		1.120		
39	4	23.546	50.917	216.058
207.859		8.199		
40	4	24.450	41.976	171.469
172.720		-1.251		

Why Use Residuals?

If the model fit to the data were correct, the residuals would approximate the random errors that make the relationship between the explanatory variables and the response variable a [statistical relationship](#). Therefore, if the residuals appear to behave randomly, it suggests that the model fits the data well. On the other hand, if non-random structure is evident in the residuals, it is a clear sign that the model fits the data poorly. The subsections listed below detail the types of plots to use to test different aspects of a model and give guidance on the correct interpretations of different results that could be observed for each type of plot.

Model Validation

1. [How can I assess the sufficiency of the functional part of the model?](#)

Specifics

2. [How can I detect non-constant variation across the data?](#)
3. [How can I tell if there was drift in the process?](#)
4. [How can I assess whether the random errors are independent from one to the next?](#)
5. [How can I test whether or not the random errors are distributed normally?](#)
6. [How can I test whether any significant terms are missing or misspecified in the functional part of the model?](#)
7. [How can I test whether all of the terms in the functional part of the model are necessary?](#)

[4. Process Modeling](#)

[4.4. Data Analysis for Process Modeling](#)

4.4.5. If my current model does not fit the data well, how can I improve it?

What Next?

Validating a model using residual plots, formal hypothesis tests and descriptive statistics would be quite frustrating if discovery of a problem meant restarting the modeling process back at square one. Fortunately, however, there are also techniques and tools to remedy many of the problems uncovered using residual analysis. In some cases the model validation methods themselves suggest appropriate changes to a model at the same time problems are uncovered. This is especially true of the graphical tools for model validation, though tests on the parameters in the regression function also offer insight into model refinement. Treatments for the various model deficiencies that were diagnosed in [Section 4.4.4.](#) are demonstrated and discussed in the subsections listed below.

Methods for Model Improvement

1. [Updating the Function Based on Residual Plots](#)
2. [Accounting for Non-Constant Variation Across the Data](#)
3. [Accounting for Errors with a Non-Normal Distribution](#)

[4. Process Modeling](#)

4.5. Use and Interpretation of Process Models

*Overview
of Section
4.5*

This section covers the interpretation and use of the models developed from the collection and analysis of data using the procedures discussed in [Section 4.3](#) and [Section 4.4](#). Three of the main uses of such models, estimation, prediction and calibration, are discussed in detail. Optimization, another important use of this type of model, is primarily discussed in [Chapter 5: Process Improvement](#).

*Contents
of Section
4.5*

1. [What types of predictions can I make using the model?](#)
 1. [How do I estimate the average response for a particular set of predictor variable values?](#)
 2. [How can I predict the value and estimate the uncertainty of a single response?](#)
2. [How can I use my process model for calibration?](#)
 1. [Single-Use Calibration Intervals](#)
3. [How can I optimize my process using the process model?](#)



[4. Process Modeling](#)

[4.5. Use and Interpretation of Process Models](#)

4.5.1. What types of predictions can I make using the model?

Detailed Information on Prediction

This section details some of the different types of predictions that can be made using the various process models whose development is discussed in [Section 4.1](#) through [Section 4.4](#). Computational formulas or algorithms are given for each different type of estimation or prediction, along with simulation examples showing its probabilistic interpretation. An introduction to the different types of estimation and prediction can be found in [Section 4.1.3.1](#). A brief description of estimation and prediction versus the other uses of process models is given in [Section 4.1.3](#).

Different Types of Predictions

1. [How do I estimate the average response for a particular set of predictor variable values?](#)
2. [How can I predict the value and estimate the uncertainty of a single response?](#)

[4. Process Modeling](#)

[4.5. Use and Interpretation of Process Models](#)

4.5.2. How can I use my process model for calibration?

*Detailed
Calibration
Information*

This section details some of the different types of calibrations that can be made using the various process models whose development was discussed in previous sections. Computational formulas or algorithms are given for each different type of calibration, along with simulation examples showing its probabilistic interpretation. An introduction to calibration can be found in [Section 4.1.3.2](#). A brief comparison of calibration versus the other uses of process models is given in [Section 4.1.3](#). Additional information on calibration is available in [Section 3](#) of [Chapter 2: Measurement Process Characterization](#).

*Calibration
Procedures*

1. [Single-Use Calibration Intervals](#)

[4. Process Modeling](#)

[4.5. Use and Interpretation of Process Models](#)

4.5.3. How can I optimize my process using the process model?

Detailed Information on Process Optimization

Process optimization using models fit to data collected using [response surface designs](#) is primarily covered in [Section 5.5.3](#) of [Chapter 5: Process Improvement](#). In that section detailed information is given on how to determine the correct process inputs to hit a target output value or to maximize or minimize process output. Some background on the use of process models for optimization can be found in [Section 4.1.3.3](#) of this chapter, however, and information on the basic analysis of data from optimization experiments is covered along with that of other types of models in [Section 4.1](#) through [Section 4.4](#) of this chapter.

Contents of Chapter 5 Section 5.5.3.

1. [Optimizing a Process](#)
 1. [Single response case](#)
 1. [Path of steepest ascent](#)
 2. [Confidence region for search path](#)
 3. [Choosing the step length](#)
 4. [Optimization when there is adequate quadratic fit](#)
 5. [Effect of sampling error on optimal solution](#)
 6. [Optimization subject to experimental region constraints](#)
 2. [Multiple response case](#)
 1. [Path of steepest ascent](#)
 2. [Desirability function approach](#)
 3. [Mathematical programming approach](#)



[4. Process Modeling](#)

4.6. Case Studies in Process Modeling

*Detailed,
Realistic
Examples*

The general points of the first five sections are illustrated in this section using data from physical science and engineering applications. Each example is presented step-by-step in the text and is often cross-linked with the relevant sections of the chapter describing the analysis in general. Each analysis can also be repeated using a worksheet linked to the appropriate Dataplot macros. The worksheet is also linked to the step-by-step analysis presented in the text for easy reference.

*Contents:
Section 6*

1. [Load Cell Calibration](#)
 1. [Background & Data](#)
 2. [Selection of Initial Model](#)
 3. [Model Fitting - Initial Model](#)
 4. [Graphical Residual Analysis - Initial Model](#)
 5. [Interpretation of Numerical Output - Initial Model](#)
 6. [Model Refinement](#)
 7. [Model Fitting - Model #2](#)
 8. [Graphical Residual Analysis - Model #2](#)
 9. [Interpretation of Numerical Output - Model #2](#)
 10. [Use of the Model for Calibration](#)
 11. [Work this Example Yourself](#)
2. [Alaska Pipeline Ultrasonic Calibration](#)
 1. [Background and Data](#)
 2. [Check for Batch Effect](#)
 3. [Initial Linear Fit](#)
 4. [Transformations to Improve Fit and Equalize Variances](#)
 5. [Weighting to Improve Fit](#)
 6. [Compare the Fits](#)
 7. [Work This Example Yourself](#)
3. [Ultrasonic Reference Block Study](#)
 1. [Background and Data](#)
 2. [Initial Non-Linear Fit](#)
 3. [Transformations to Improve Fit](#)
 4. [Weighting to Improve Fit](#)
 5. [Compare the Fits](#)
 6. [Work This Example Yourself](#)
4. [Thermal Expansion of Copper Case Study](#)
 1. [Background and Data](#)
 2. [Exact Rational Models](#)
 3. [Initial Plot of Data](#)
 4. [Fit Quadratic/Quadratic Model](#)
 5. [Fit Cubic/Cubic Model](#)

6. [Work This Example Yourself](#)



[HOME](#)

[TOOLS & AIDS](#)

[SEARCH](#)

[BACK](#) [NEXT](#)



[4. Process Modeling](#)

[4.6. Case Studies in Process Modeling](#)

4.6.1. Load Cell Calibration

Quadratic Calibration This example illustrates the construction of a linear regression model for load cell data that relates a known load applied to a load cell to the deflection of the cell. The model is then used to calibrate future cell readings associated with loads of unknown magnitude.

1. [Background & Data](#)
2. [Selection of Initial Model](#)
3. [Model Fitting - Initial Model](#)
4. [Graphical Residual Analysis - Initial Model](#)
5. [Interpretation of Numerical Output - Initial Model](#)
6. [Model Refinement](#)
7. [Model Fitting - Model #2](#)
8. [Graphical Residual Analysis - Model #2](#)
9. [Interpretation of Numerical Output - Model #2](#)
10. [Use of the Model for Calibration](#)
11. [Work This Example Yourself](#)



[4. Process Modeling](#)

[4.6. Case Studies in Process Modeling](#)

4.6.2. Alaska Pipeline

Non-Homogeneous Variances

This example illustrates the construction of a linear regression model for Alaska pipeline ultrasonic calibration data. This case study demonstrates the use of transformations and weighted fits to deal with the violation of the assumption of [constant standard deviations](#) for the random errors. This assumption is also called homogeneous variances for the errors.

1. [Background and Data](#)
2. [Check for a Batch Effect](#)
3. [Fit Initial Model](#)
4. [Transformations to Improve Fit and Equalize Variances](#)
5. [Weighting to Improve Fit](#)
6. [Compare the Fits](#)
7. [Work This Example Yourself](#)

[4. Process Modeling](#)

[4.6. Case Studies in Process Modeling](#)

4.6.3. Ultrasonic Reference Block Study

*Non-Linear
Fit with Non-
Homogeneous
Variances*

This example illustrates the construction of a non-linear regression model for ultrasonic calibration data. This case study demonstrates fitting a non-linear model and the use of transformations and weighted fits to deal with the violation of the assumption of [constant standard deviations](#) for the errors. This assumption is also called homogeneous variances for the errors.

1. [Background and Data](#)
2. [Fit Initial Model](#)
3. [Transformations to Improve Fit](#)
4. [Weighting to Improve Fit](#)
5. [Compare the Fits](#)
6. [Work This Example Yourself](#)

[4. Process Modeling](#)

[4.6. Case Studies in Process Modeling](#)

4.6.4. Thermal Expansion of Copper Case Study

Rational Function Models This case study illustrates the use of a class of nonlinear models called rational function models. The data set used is the thermal expansion of copper related to temperature.

This data set was provided by the NIST scientist Thomas Hahn.

Contents

1. [Background and Data](#)
2. [Rational Function Models](#)
3. [Initial Plot of Data](#)
4. [Fit Quadratic/Quadratic Model](#)
5. [Fit Cubic/Cubic Model](#)
6. [Work This Example Yourself](#)



[4. Process Modeling](#)

4.7. References For Chapter 4: Process Modeling

Handbook of Mathematical Functions with Formulas, Graphs and Mathematical Tables (1964) Abramowitz M. and Stegun I. (eds.), U.S. Government Printing Office, Washington, DC, 1046 p.

Berkson J. (1950) "Are There Two Regressions?," *Journal of the American Statistical Association*, Vol. 45, pp. 164-180.

Carroll, R.J. and Ruppert D. (1988) *Transformation and Weighting in Regression*, Chapman and Hall, New York.

Cleveland, W.S. (1979) "Robust Locally Weighted Regression and Smoothing Scatterplots," *Journal of the American Statistical Association*, Vol. 74, pp. 829-836.

Cleveland, W.S. and Devlin, S.J. (1988) "Locally Weighted Regression: An Approach to Regression Analysis by Local Fitting," *Journal of the American Statistical Association*, Vol. 83, pp. 596-610.

Fuller, W.A. (1987) *Measurement Error Models*, John Wiley and Sons, New York.

Graybill, F.A. (1976) *Theory and Application of the Linear Model*, Duxbury Press, North Scituate, Massachusetts.

Graybill, F.A. and Iyer, H.K. (1994) *Regression Analysis: Concepts and Applications*, Duxbury Press, Belmont, California.

Harter, H.L. (1983) "Least Squares," *Encyclopedia of Statistical Sciences*, Kotz, S. and Johnson, N.L., eds., John Wiley & Sons, New York, pp. 593-598.

Montgomery, D.C. (2001) *Design and Analysis of Experiments*, 5th ed., Wiley, New York.

Neter, J., Wasserman, W., and Kutner, M. (1983) *Applied Linear Regression Models*, Richard D. Irwin Inc., Homewood, IL.

Ryan, T.P. (1997) *Modern Regression Methods*, Wiley, New York

Seber, G.A.F and Wild, C.F. (1989) *Nonlinear Regression*, John Wiley and Sons, New York.

Stigler, S.M. (1978) "Mathematical Statistics in the Early States," *The Annals of Statistics*, Vol. 6, pp. 239-265.

Stigler, S.M. (1986) *The History of Statistics: The Measurement of Uncertainty Before 1900*, The Belknap Press of Harvard University Press, Cambridge, Massachusetts.

NIST
SEMATECH

[HOME](#)

[TOOLS & AIDS](#)

[SEARCH](#)

[BACK](#) [NEXT](#)



[4. Process Modeling](#)

4.8. Some Useful Functions for Process Modeling

*Overview
of Section
4.8*

This section lists some functions commonly-used for process modeling. Constructing an exhaustive list of useful functions is impossible, of course, but the functions given here will often provide good starting points when an empirical model must be developed to describe a particular process.

Each function listed here is classified into a family of related functions, if possible. Its statistical type, linear or nonlinear in the parameters, is also given. Special features of each function, such as asymptotes, are also listed along with the function's domain (the set of allowable input values) and range (the set of possible output values). Plots of some of the different shapes that each function can assume are also included.

*Contents
of Section
4.8*

1. [Univariate Functions](#)
 1. [Polynomials](#)
 2. [Rational Functions](#)