

Section 1

MANE 6313

Subsection 1

Week 12, Module A

Student Learning Outcome

- Select an appropriate experimental design with one or more factors,
- Select an appropriate model with one or more factors,
- Evaluate statistical analyses of experimental designs,
- Assess the model adequacy of any experimental design, and
- Interpret model results.

Module Learning Outcome

Describe linear regression.

Introduction to Linear Regression

- We are interested in a relationship between a single *dependent variable* or *response* y that depends on k *independent* or *regressor variables*.
- We assume that there is some mathematical function $y = \phi(x_1, x_2, \dots, x_k)$. In general, we don't know this function
- We'll use low order polynomial equations as an approximating function. This is called *empirical modeling*.
- What are methods that we can determine if there is a relationship between two (or more) variables?

Relationship between two or more variables

Example 12.8 Suppose that a scientist takes experimental data on the radius of a propellant grain Y as a function of powder temperature x_1 , extrusion rate

x_2 , and die temperature x_3 . Fit a linear regression model for predicting grain radius, and determine the effectiveness of each variable in the model. The data are as follows:

Grain radius	Powder temperature	Extrusion rate	Die temperature
82	150 (−1)	12 (−1)	220 (−1)
93	190 (1)	12 (−1)	220 (−1)
114	150 (−1)	24 (1)	220 (−1)
124	150 (−1)	12 (−1)	250 (1)
111	190 (1)	24 (1)	220 (−1)
129	190 (1)	12 (−1)	250 (1)
157	150 (−1)	24 (1)	250 (1)
164	190 (1)	24 (1)	250 (1)

Figure 1: Example 12.8

Linear regression models

- In general they look like

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{11} x_1^2 + \beta_{22} x_2^2 + \beta_{12} x_1 x_2 + \varepsilon$$

- This model is *linear* in the parameters β

- See graphical explanation from Ott.

measurement to another.

FIGURE 11.2

Theoretical distribution of y
in regression

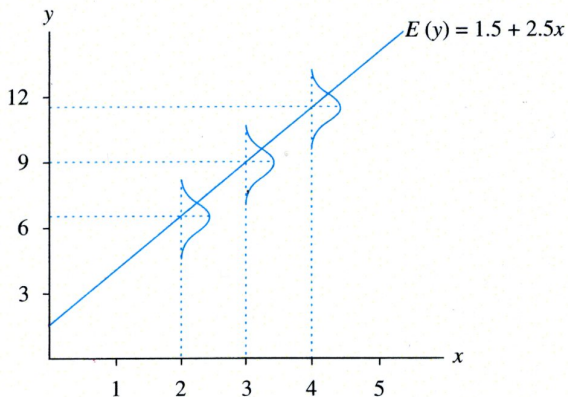


Figure 2: Figure 11.2

Estimation of Parameters

- Parameter estimates are derived using least squares. Goal is to minimize the squared error.
- Parameter estimation can be done algebraically or using linear algebra. Montgomery focuses on linear algebra formulation.
- In general, the matrix formulation is used. Model is defined to be

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon$$

- The least squares estimates are found by minimizing

$$L = \sum_i \varepsilon_i^2 = \varepsilon' \varepsilon = (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)$$

- The least squares estimates must satisfy

$$\frac{\partial L}{\partial \beta} \Big|_{\hat{\beta}} = -2\mathbf{X}'\mathbf{y} + 2\mathbf{X}'\mathbf{X}\hat{\beta} = 0$$

- Which leads to the solution

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

- We can define the predicted response to be

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}$$

- The residuals are defined to be

$$\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}}$$

- Thus the sum of squares errors can be shown to be

$$\begin{aligned} SS_E &= (\mathbf{y} - \mathbf{X}\hat{\beta})'(\mathbf{y} - \mathbf{X}\hat{\beta}) \\ &= \mathbf{y}'\mathbf{y} - \hat{\beta}'\mathbf{X}'\mathbf{y} \end{aligned}$$

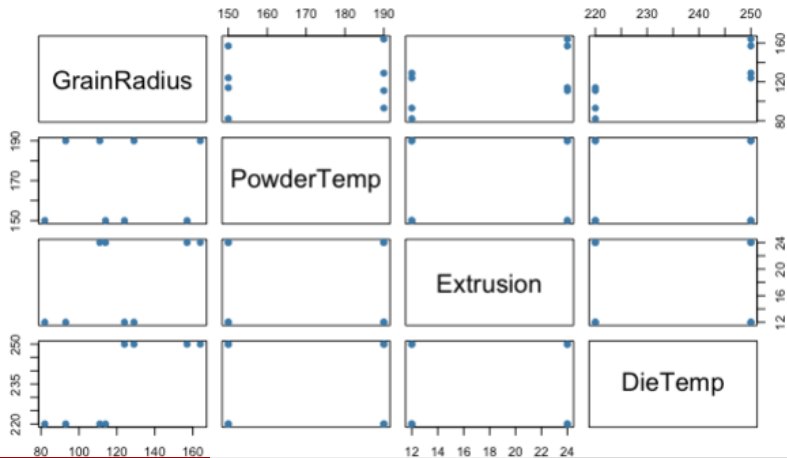
Coding Variables

- From the example, there were two ways to represent the same problem, coded and uncoded
- Why use coded variables^a
 - Computational ease and increased accuracy in estimating the model coefficients
 - Enhanced interpretability of the coefficient estimates in the model.
- Internally most statistical software codes for estimating parameters

^aKhuri, and Cornell (1987). Response Surfaces: Designs and Analyses. Dekker

Plot: Example 12.8

```
14 > ```{r}  
15 library(readxl)  
16 ex12_8.df <- read_excel("ex12_8.xlsx")  
17 plot(ex12_8.df,pch=20,cex=1.5,col='steelblue')  
18 >
```



Regression: Example 12.8

```

27 > ``{r}
28 ex12_8.model <- lm(GrainRadius~PowderTemp+Extrusion+DieTemp,data=ex12_8.df)
29 summary(ex12_8.model)
30 > ``

```

Call:

```
lm(formula = GrainRadius ~ PowderTemp + Extrusion + DieTemp,
    data = ex12_8.df)
```

Residuals:

1	2	3	4	5	6	7	8
-0.75	5.25	1.75	-2.25	-6.25	-2.25	1.25	3.25

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-284.50000	30.77729	-9.244	0.000761 ***
PowderTemp	0.12500	0.08501	1.470	0.215398
Extrusion	2.45833	0.28336	8.676	0.000972 ***
DieTemp	1.45000	0.11335	12.793	0.000215 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 4.809 on 4 degrees of freedom

Multiple R-squared: 0.9837, Adjusted R-squared: 0.9714

F-statistic: 80.36 on 3 and 4 DF, p-value: 0.0004967

Values from `lm()` function

Source: [https:](https://www.rdocumentation.org/packages/stats/versions/3.6.2/topics/lm)

[//www.rdocumentation.org/packages/stats/versions/3.6.2/topics/lm](https://www.rdocumentation.org/packages/stats/versions/3.6.2/topics/lm)

Value

`lm` returns an object of `class` `"lm"` or for multiple responses of class `c("mlm", "lm")`.

The functions `summary` and `anova` are used to obtain and print a summary and analysis of variance table of the results. The generic accessor functions `coefficients`, `effects`, `fitted.values` and `residuals` extract various useful features of the value returned by `lm`.

An object of class `"lm"` is a list containing at least the following components:

<code>coefficients</code>	a named vector of coefficients
<code>residuals</code>	the residuals, that is response minus fitted values.
<code>fitted.values</code>	the fitted mean values.
<code>rank</code>	the numeric rank of the fitted linear model.
<code>weights</code>	(only for weighted fits) the specified weights.
<code>df.residual</code>	the residual degrees of freedom.
<code>call</code>	the matched call.
<code>terms</code>	the <code>terms</code> object used.
<code>contrasts</code>	(only where relevant) the contrasts used.
<code>xlevels</code>	(only where relevant) a record of the levels of the factors used in fitting.
<code>offset</code>	the offset used (missing if none were used).
<code>y</code>	if requested, the response used.

