

# Multiple Regression

**Dan Frey**

**Associate Professor of Mechanical Engineering and Engineering Systems**



# Plan for Today

- Multiple Regression
  - Estimation of the parameters
  - Hypothesis testing
  - Regression diagnostics
  - Testing lack of fit
- Case study
- Next steps

# The Model Equation

For a single variable

$$Y = \alpha + \beta x + \varepsilon$$

For multiple variables

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

*$\alpha$  is renamed  $\beta_0$*

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix} \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

*$p = k + 1$*

*These 1's allow  $\beta_0$  to enter the equation without being mult by  $x$ 's*

# The Model Equation $y = X\beta + \varepsilon$

Each row of  $X$   
is paired with  
an observation

Each column of  $X$   
is paired with a  
coefficient

$$E(\varepsilon_i) = 0$$

$$Var(\varepsilon_i) = \sigma^2$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix} \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

There are  $n$   
observations of  
the response

$$\begin{bmatrix} \beta_0 & \beta_1 & \cdots & \beta_k \end{bmatrix}$$

There are  $k$  coefficients

Each observation  
is affected by an  
independent  
homoscedastic  
normal variates

# Accounting for Indices

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

$n \times 1$

$n \times p$   $p \times 1$

$n \times 1$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix} \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

*Note: A red box highlights the element  $x_{1k}$  in the matrix  $\mathbf{X}$ , with a red arrow pointing to it from the label  $p = k + 1$  above.*

# Concept Question

Which of these is a valid  $\mathbf{X}$  matrix?

$$\mathbf{X} = \begin{bmatrix} 1 & 5.0m & 0.3\text{sec} \\ 1 & 7.1m & 0.2\text{sec} \\ 1 & 3.2m & 0.7\text{sec} \\ 1 & 5.4m & 0.4\text{sec} \end{bmatrix}$$

A

$$\mathbf{X} = \begin{bmatrix} 1 & 5.0m & 0.3m \\ 1 & 7.1V & 0.2V \\ 1 & 3.2\text{sec} & 0.7\text{sec} \\ 1 & 5.4A & 0.4A \end{bmatrix}$$

B

$$\mathbf{X} = \begin{bmatrix} 1 & 5.0m & 0.1\text{sec} \\ 1 & 7.1m & 0.3\text{sec} \end{bmatrix}$$

C

1) A only

2) B only

3) C only

4) A and B

5) B and C

6) A and C

7) A, B, & C

8) None

9) I don't know

# Adding h.o.t. to the Model Equation

Each row of  $\mathbf{X}$  is paired with an observation

You can add interactions

You can add curvature

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{11}x_{12} & x_{11}^2 \\ 1 & x_{21} & x_{22} & x_{21}x_{22} & x_{21}^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n1}x_{n2} & x_{n1}^2 \end{bmatrix}$$

There are  $n$  observations of the response

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_0 & \beta_1 & \beta_2 & \beta_{12} & \beta_{11} \end{bmatrix}$$

$\boldsymbol{\beta}$

# Estimation of the Parameters $\beta$

Assume the model equation

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

We wish to minimize the sum squared error

$$L = \boldsymbol{\varepsilon}^T \boldsymbol{\varepsilon} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

To minimize, we take the derivative and set it equal to zero

$$\left. \frac{\partial L}{\partial \boldsymbol{\beta}} \right|_{\hat{\boldsymbol{\beta}}} = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X}\hat{\boldsymbol{\beta}}$$

The solution is

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

And we define the fitted model

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

# MathCad Demo Montgomery Example 10-1

Montgomery, D. C., 2001, *Design and Analysis of Experiments*, John Wiley & Sons.

Done in MathCad:

## DCM Example 10-1

ORIGIN := 1

	Viscosity of a Polymer	Reaction temperature	Catalyst feed rate
$y :=$	$\begin{pmatrix} 2256 \\ 2340 \\ 2426 \\ 2293 \\ 2330 \\ 2368 \\ 2250 \\ 2409 \\ 2364 \\ 2379 \\ 2440 \\ 2364 \\ 2404 \\ 2317 \\ 2309 \\ 2328 \end{pmatrix}$	$X :=$	$\begin{pmatrix} 1 & 80 & 8 \\ 1 & 93 & 9 \\ 1 & 100 & 10 \\ 1 & 82 & 12 \\ 1 & 90 & 11 \\ 1 & 99 & 8 \\ 1 & 81 & 8 \\ 1 & 96 & 10 \\ 1 & 94 & 12 \\ 1 & 93 & 11 \\ 1 & 97 & 13 \\ 1 & 95 & 11 \\ 1 & 100 & 8 \\ 1 & 85 & 12 \\ 1 & 86 & 9 \\ 1 & 87 & 12 \end{pmatrix}$

$$X := \text{augment}\left[X, \begin{matrix} \longrightarrow \\ \left(X^{(2)}\right)^2 \end{matrix}\right] \quad \text{disabled} \quad +$$
  

$$p := \text{cols}(X) \quad p = 3$$

$$n := \text{rows}(X) \quad n = 16$$

$$k := p - 1$$

$$\beta_{\text{hat}} := (X^T \cdot X)^{-1} \cdot X^T \cdot y \quad \beta_{\text{hat}} = \begin{pmatrix} 1.566 \times 10^3 \\ 7.621 \\ 8.585 \end{pmatrix}$$

$$y_{\text{hat}} := X \cdot \beta_{\text{hat}} \quad \underline{\underline{e}} := y - y_{\text{hat}}$$

# Breakdown of Sum Squares

“Grand Total  
Sum of Squares”

$$GTSS = \sum_{i=1}^n y_i^2$$

$SS$  due to mean  
 $= n\bar{y}^2$

$$SS_T = \sum_{i=1}^n (y_i - \bar{y})^2$$

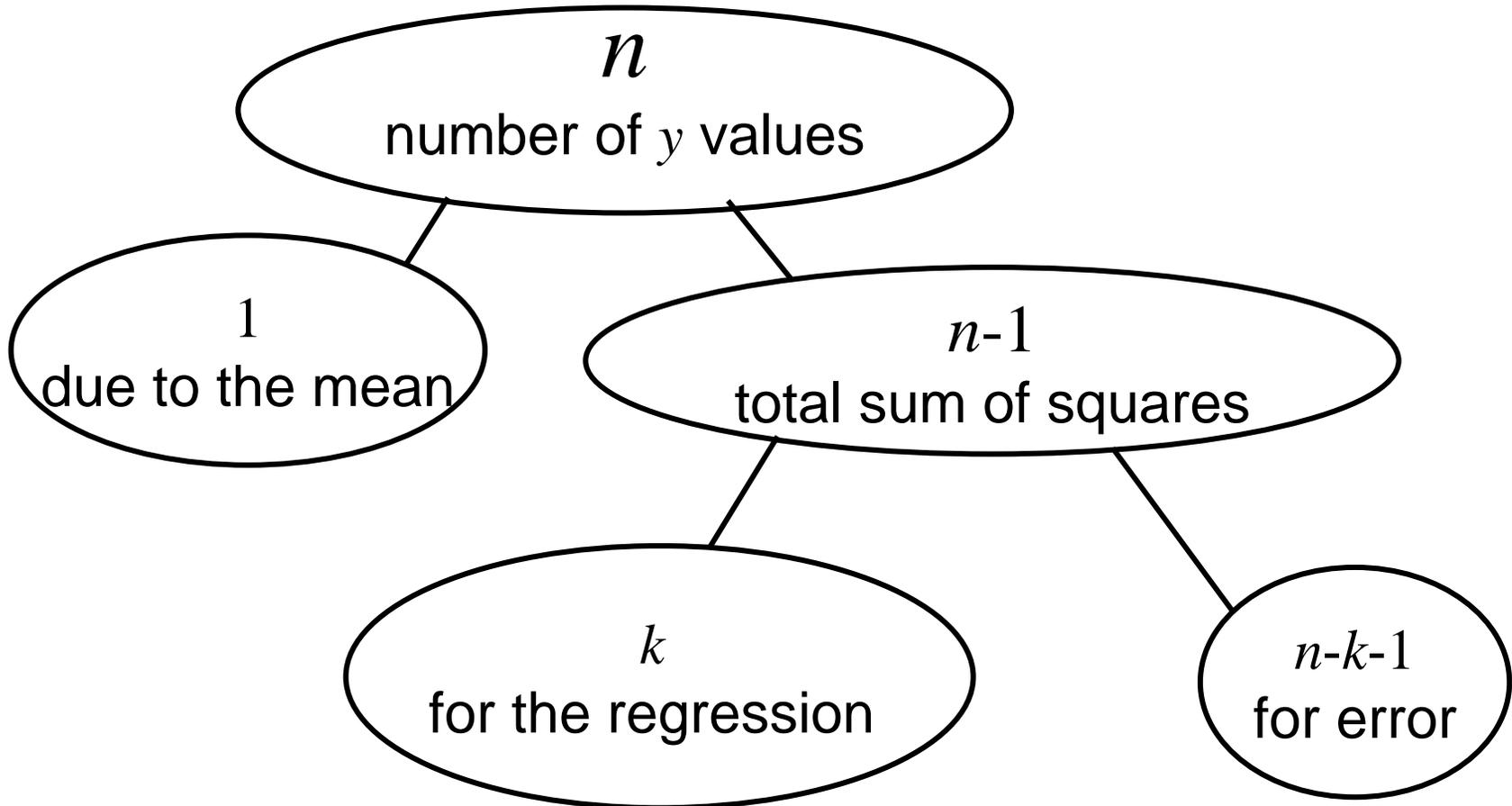
$$SS_R = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

$$SS_E = \sum_{i=1}^n \mathbf{e}_i^2$$

$SS_{PE}$

$SS_{LOF}$

# Breakdown of DOF



# Estimation of the Error Variance $\sigma^2$

Remember the the model equation  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$   $\boldsymbol{\varepsilon} \sim N(0, \sigma)$

If assumptions of the model equation hold, then

$$E(SS_E / (n - k - 1)) = \sigma^2$$

So an **unbiased** estimate of  $\sigma^2$  is

$$\hat{\sigma}^2 = SS_E / (n - k - 1)$$

a.k.a. “coefficient of multiple determination”

## $R^2$ and Adjusted $R^2$

What fraction of the total sum of squares ( $SS_T$ ) is accounted for jointly by all the parameters in the fitted model?

$$R^2 \equiv \frac{SS_R}{SS_T} = 1 - \frac{SS_E}{SS_T} \quad R^2 \text{ can only rise as parameters are added}$$

$$R_{adj}^2 \equiv 1 - \frac{SS_E / (n - p)}{SS_T / (n - 1)} = 1 - \left( \frac{n - 1}{n - p} \right) (1 - R^2)$$

$R_{adj}^2$  can rise or drop as parameters are added

# Back to MathCad Demo Montgomery Example 10-1

$$F_{\text{ww}} := \frac{SS_R}{SS_E} \cdot \frac{n - k - 1}{k}$$

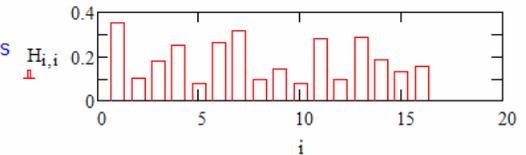
$$F = 82.505 \quad qF(0.05, k, n - k - 1) = 0.051$$

"reject Ho" if  $F > qF(0.05, k, n - k - 1)$  = "reject Ho"  
"accept Ho" otherwise

$$H_{\text{ww}} := X \cdot (X^T \cdot X)^{-1} \cdot X^T$$

$$2 \cdot \frac{p}{n} = 0.375 \quad i := 1..n$$

Influence of the observations



Covariance of the residuals

	1	2	3	4	5
1	174.074	-25.35	17.472	-32.12	-11.55
2	-25.35	240.182	-25.095	3.015	-10.307
3	17.472	-25.095	220.326	17.077	-11.703
4	-32.12	3.015	17.077	200.413	-28.569
5	-11.55	-10.307	-11.703	-28.569	247.028
6	-13.233	-39.278	-47.08	34.612	-3.194
7	-89.303	-26.083	14.074	-28.608	-11.11
8	0.567	-22.163	-33.688	3.028	-13.462
9	18.594	-5.781	-23.693	-25.044	-23.292
10	1.128	-12.506	-21.896	-18.033	-19.257
11	44.511	-0.523	-32.286	-25.032	-26.447
12	9.581	-13.972	-28.691	-11.008	-18.377
13	-9.007	-40.011	-50.478	38.124	-2.754
14	-19.441	0.816	6.884	-56.654	...

$$\frac{SS_E}{n - k - 1} \cdot (\text{identity}(n) - H) =$$

Montgomery, D. C., 2001, *Design and Analysis of Experiments*, John Wiley & Sons.

# Why Hypothesis Testing is Important in Multiple Regression

- Say there are 10 regressor variables
- Then there are 11 coefficients in a linear model
- To make a fully 2<sup>nd</sup> order model requires
  - 10 curvature terms in each variable
  - 10 choose 2 = 45 interactions
- You'd need 68 samples just to get the matrix  $\mathbf{X}^T\mathbf{X}$  to be invertible
- You need a way to discard insignificant terms

# Test for Significance of Regression

The hypotheses are

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0$$

$$H_1 : \beta_j \neq 0 \text{ for at least one } j$$

The test statistic is

$$F_0 = \frac{SS_R/k}{SS_E/(n-k-1)}$$

Reject  $H_0$  if  $F_0 > F_{\alpha, k, n-k-1}$

# Test for Significance Individual Coefficients

The hypotheses are  $H_0 : \beta_j = 0$

$$H_1 : \beta_j \neq 0$$

The test statistic is

$$t_0 = \frac{\hat{\beta}_j}{\sqrt{\hat{\sigma}^2 C_{jj}}}$$

$$C = (\mathbf{X}^T \mathbf{X})^{-1}$$

← Standard error

$$\sqrt{\hat{\sigma}^2 C_{jj}}$$

Reject  $H_0$  if  $|t_0| > t_{\alpha/2, n-k-1}$

# Test for Significance of Groups of Coefficients

Partition the coefficients into two groups  $\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$  to be removed  
to remain

Reduced model  $y = \mathbf{X}_2 \beta_2 + \varepsilon$

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix} \quad \rightarrow \quad \mathbf{X}_2$$

$$H_0 : \beta_1 = \mathbf{0}$$

$$H_1 : \beta_1 \neq \mathbf{0}$$

Basically, you form  $\mathbf{X}_2$  by removing the columns associated with the coefficients you are testing for significance

# Test for Significance Groups of Coefficients

Reduced model  $\mathbf{y} = \mathbf{X}_2\boldsymbol{\beta}_2 + \boldsymbol{\varepsilon}$

The regression sum of squares for the reduced model is

$$SS_R(\boldsymbol{\beta}_2) = \mathbf{y}^T \mathbf{H}_2 \mathbf{y} - n\bar{y}^2$$

Define the sum squares of the removed set given the other coefficients are in the model

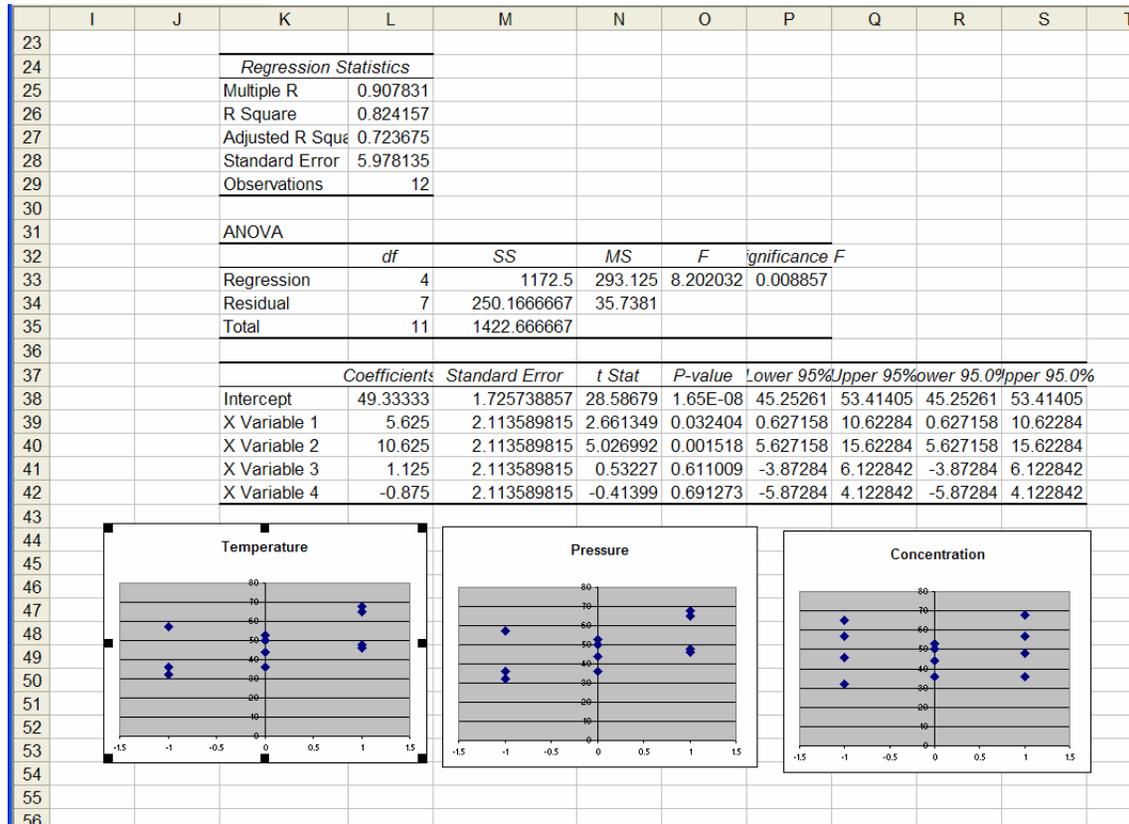
$$SS_R(\boldsymbol{\beta}_1 | \boldsymbol{\beta}_2) \equiv SS_R(\boldsymbol{\beta}) - SS_R(\boldsymbol{\beta}_2)$$

The partial  
*F* test

$$F_0 = \frac{SS_R(\boldsymbol{\beta}_1 | \boldsymbol{\beta}_2) / r}{SS_E / (n - p)}$$

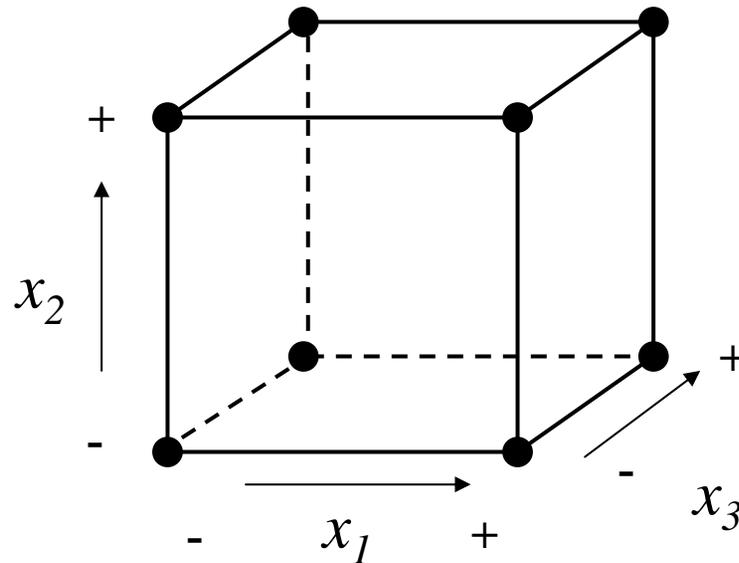
Reject  $H_0$  if  $F_0 > F_{\alpha, r, n-p}$

# Excel Demo -- Montgomery Ex10-2



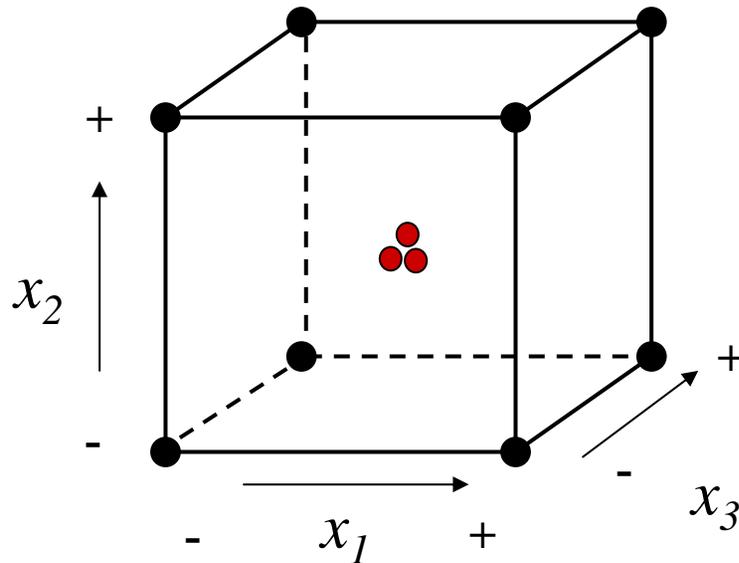
# Factorial Experiments

## Cuboidal Representation



Exhaustive search of the space of discrete 2-level factors is the **full factorial  $2^3$  experimental design**

# Adding Center Points



Center points allow an experimenter to check for curvature and, if replicated, allow for an estimate of **pure experimental error**

# Plan for Today

- Mud cards
- Multiple Regression
  - Estimation of the parameters
  - Hypothesis testing
-  Regression diagnostics
  - Testing lack of fit
- Case study
- Next steps

# The “Hat” Matrix

Since 
$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

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

therefore 
$$\hat{\mathbf{y}} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

So we define 
$$\mathbf{H} \equiv \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

Which maps from  
observations  $\mathbf{y}$  to  
predictions  $\hat{\mathbf{y}}$

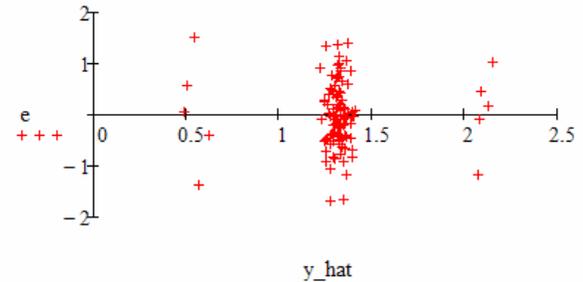
$$\hat{\mathbf{y}} = \mathbf{H} \mathbf{y}$$

# Influence Diagnostics

- The relative disposition of points in  $x$  space determines their effect on the coefficients
- The hat matrix  $\mathbf{H}$  gives us an ability to check for leverage points
- $h_{ij}$  is the amount of leverage exerted by point  $\mathbf{y}_j$  on  $\hat{\mathbf{y}}_i$
- Usually the diagonal elements  $\sim p/n$  and it is good to check whether the diagonal elements within 2X of that

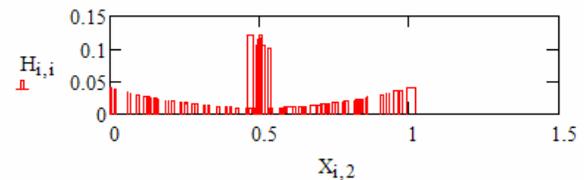
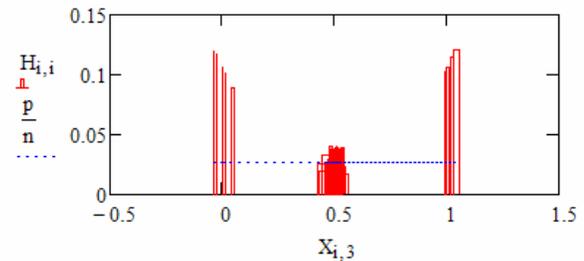
# MathCad Demo on Distribution of Samples and Its Effect on Regression

## Plot the residuals



$$H := X \cdot (X^T \cdot X)^{-1} \cdot X^T \quad i := 1..110$$

## Influence of the observations



# Standardized Residuals

The residuals are defined as  $\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}}$

So an unbiased estimate of  $\sigma^2$  is  $\hat{\sigma}^2 = SS_E / (n - p)$

The **standardized residuals** are defined as  $\mathbf{d} = \frac{\mathbf{e}}{\hat{\sigma}}$

If these elements were  $z$ -scores then with probability 99.7%

$$-3 < d_i < 3$$

# Studentized Residuals

The residuals are defined as  $\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}}$

therefore  $\mathbf{e} = \mathbf{y} - \mathbf{H}\mathbf{y} = (\mathbf{I} - \mathbf{H})\mathbf{y}$

So the **covariance matrix** of the residuals is  $\text{Cov}(\mathbf{e}) = \sigma^2 \text{Cov}(\mathbf{I} - \mathbf{H})$

The **studentized residuals** are defined as  $r_i = \frac{e_i}{\sqrt{\hat{\sigma}^2 (1 - h_{ii})}}$

If these elements were  $z$ -scores then with probability 99.7%

$$\text{????} \quad -3 < r_i < 3$$

# Testing for Lack of Fit

(Assuming a Central Composite Design)

- Compute the standard deviation of the center points and assume that represents the  $MS_{PE}$

$$MS_{PE} = \frac{\sum (y_i - \bar{y})}{\text{center points} \over n_C - 1}$$

$$MS_{LOF} = \frac{SS_{LOF}}{p}$$

$$SS_{PE} = (n - 1)MS_{PE}$$

$$SS_{PE} + SS_{LOF} = SS_E$$

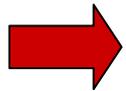
$$F_0 = \frac{MS_{LOF}}{MS_{PE}}$$

# Concept Test

- You perform a linear regression of 100 data points ( $n=100$ ). There are two independent variables  $x_1$  and  $x_2$ . The regression  $R^2$  is 0.72. Both  $\beta_1$  and  $\beta_2$  pass a  $t$  test for significance. You decide to add the interaction  $x_1x_2$  to the model. Select all the things that cannot happen:
  - 1) Absolute value of  $\beta_1$  decreases
  - 2)  $\beta_1$  changes sign
  - 3)  $R^2$  decreases
  - 4)  $\beta_1$  fails the  $t$  test for significance

# Plan for Today

- Mud cards
- Multiple Regression
  - Estimation of the parameters
  - Hypothesis testing
  - Regression diagnostics
  - Testing lack of fit



Case study

- Next steps

# Scenario

- The FAA and EPA are interested in reducing CO<sub>2</sub> emissions
- Some parameters of airline operations are thought to effect CO<sub>2</sub> (e.g., Speed, Altitude, Temperature, Weight)
- Imagine flights have been made with special equipment that allowed CO<sub>2</sub> emission to be measured (data provided)
- You will report to the FAA and EPA on your analysis of the data and make some recommendations

# Phase One

- Open a Matlab window
- Load the data (load FAAcase3.mat)
- Explore the data

# Phase Two

- Do the regression
- Examine the betas and their intervals
- Plot the residuals

```
y=[CO2./ground_speed];  
ones(1:3538)=1;  
X=[ones' TAS alt temp weight];  
[b,bint,r,rint,stats] = regress(y,X,0.05);  
yhat=X*b;  
plot(yhat,r,'+')
```

```
dims=size(X);
i=2:dims(1)-1;
climb(1)=1;
climb(dims(1))=0;
des(1)=0;
des(dims(1))=1;
climb(i)=(alt(i)>(alt(i-1)+100))|(alt(i+1)>(alt(i)+100));
des(i)=(alt(i)<(alt(i-1)-100))|(alt(i+1)<(alt(i)-100));
for i=dims(1):-1:1
if climb(i)|des(i)
            y(i,:)=[]; X(i,:)=[]; yhat(i,:)=[]; r(i,:)=[];
end
end
hold off
plot(yhat,r,'or')
```

This code will remove the points at which the aircraft is climbing or descending

# Try The Regression Again on Cruise Only Portions

- What were the effects on the residuals?
- What were the effects on the betas?

```
hold off  
[b,bint,r,rint,stats] = regress(y,X,0.05);  
yhat=X*b;  
plot(yhat,r,'+')
```

# See What Happens if We Remove Variables

- Remove weight & temp
- Do the regression (CO2 vs TAS & alt)
- Examine the betas and their intervals

```
[b,bint,r,rint,stats] = regress(y,X(:,1:3),0.05);
```

# Phase Three

- Try different data (flight34.mat)
- Do the regression
- Examine the betas and their intervals
- Plot the residuals

```
y=[fuel_burn];  
ones(1:34)=1;  
X=[ones' TAS alt temp];  
[b,bint,r,rint,stats] = regress(y,X,0.05);  
yhat=X*b;  
plot(yhat,r,'+')
```

# Adding Interactions

```
X(:,5)=X(:,2).*X(:,3);
```

This line will add a  
interaction

What's the effect  
on the  
regression?

# Case Wrap-Up

- What were the recommendations?
- What other analysis might be done?
- What were the key lessons?

# Next Steps

- Wednesday 25 April
  - Design of Experiments
  - Please read "Statistics as a Catalyst to Learning"
- Friday 27 April
  - Recitation to support the term project
- Monday 30 April
  - Design of Experiments
- Wednesday 2 May
  - Design of Computer Experiments
- Friday 4 May?? Exam review??
- Monday 7 May – Frey at NSF
- Wednesday 9 May – Exam #2