# STAT 500

## Multiple Linear Regression Models
## More Examples

# Uncorrelated Predictors

Example: Yield of a chemical process (Myers)

$Y = $ Yield (%)
$X_1 = $ Temperature ($^o$F)
$X_2 = $ Time (hours)

Data:

| $Y$ | $X_1$ | $X_2$ |
|-----|-------|-------|
| 77 | 160 | 1 |
| 79 | 160 | 2 |
| 82 | 165 | 1 |
| 83 | 165 | 2 |
| 85 | 170 | 1 |
| 88 | 170 | 2 |
| 90 | 175 | 1 |
| 93 | 175 | 2 |

# Chemical Process Study

Full Factorial Design

$$r_{x_1, x_2} = \frac{\Sigma_{i=1}^{n}(x_{i1} - \bar{x}_1)(x_{i2} - \bar{x}_2)}{\sqrt{\Sigma_{i=1}^{n}(x_{i1} - \bar{x}_1)^2 \, \Sigma_{i=1}^{n}(x_{i2} - \bar{x}_2)^2}} = 0$$

# Estimated Models

Model 1: $\hat{Y}_i = -64.45 + 0.890x_{i1}$

$R^2 = 0.9435$

Model 2: $\hat{Y}_i = 81.25 + 2.25x_{i2}$

$R^2 = 0.0482$

Model 12: $\hat{Y}_i = -67.825 + 0.890x_{i1} + 2.250x_{i2}$

$R^2 = 0.9918$

# Chemical Process Study

| Source of variation | d.f. | SS | MS | F | p-val |
|---|---|---|---|---|---|
| reg on $x_1$ | 1 | 198.025 | 198.025 | 574.0 | .0001 |
| reg on $x_2$ after $x_1$ | 1 | 10.125 | 10.125 | 29.3 | .0029 |
| error | 5 | 1.725 | 0.345 | | |
| corrected total | 7 | 209.875 | | | |

| Source of variation | d.f. | SS | MS | F | p-val |
|---|---|---|---|---|---|
| reg on $x_2$ | 1 | 10.125 | 10.125 | 29.3 | .0029 |
| reg on $x_1$ after $x_2$ | 1 | 198.025 | 198.025 | 574.0 | .0001 |
| error | 5 | 1.725 | 0.345 | | |
| corrected total | 7 | 209.875 | | | |

# Complete Confounding

Example: Correlation between $X_1$ and $X_2$ is one.

| $Y$ | $X_1$ | $X_2$ |
|------|-------|-------|
| 1.95 | 1 | 5 |
| 6.25 | 2 | 10 |
| 9.85 | 3 | 15 |

# Estimated Models

Model 1: $\hat{Y}_i = -1.8833 + 3.95x_{i1}$

$R^2 = 0.9974$

Model 2: $\hat{Y}_i = -1.8833 + 0.79x_{i2}$

$R^2 = 0.9974$

Model 12:   Many choices for $b_1$ and $b_2$ in

$$\hat{Y}_i = b_0 + b_1 x_{i1} + b_2 x_{i2} = b_0 + b_1 x_{i1} + b_2(5x_{i1})$$

$$= b_0 + (b_1 + 5b_2)x_{i1}$$

$R^2 = 0.9974$

# Complete Confounding Example

| Source of variation | d.f. | SS | MS | F | p-val |
|---|---|---|---|---|---|
| reg on $x_1$ | 1 | 31.205 | 31.205 | 382.1 | .0325 |
| reg on $x_2$ after $x_1$ | 0 | 0.000 | 0.000 | NA | NA |
| error | 1 | 0.08167 | 0.08167 | | |
| corrected total | 2 | 31.28667 | | | |

| Source of variation | d.f. | SS | MS | F | p-val |
|---|---|---|---|---|---|
| reg on $x_2$ | 1 | 31.205 | 31.205 | 382.1 | .0325 |
| reg on $x_1$ after $x_2$ | 0 | 0.000 | 0.000 | NA | NA |
| error | 1 | 0.08167 | 0.08167 | | |
| corrected total | 2 | 31.28667 | | | |

# Partial Confounding

Example: Correlation between $X_1$ and $X_2$ is 0.95237

| $Y$ | $X_1$ | $X_2$ |
|-----|-------|-------|
| 1.8 | 1.0 | 5 |
| 1.7 | 1.1 | 6 |
| 5.4 | 1.8 | 11 |
| 6.1 | 2.0 | 10 |
| 7.0 | 2.1 | 9 |
| 9.6 | 3.0 | 15 |

# Estimated Models

Model 1: $\hat{Y}_i = -2.328 + 4.142 x_{i1}$

$R^2 = 0.978$

Model 2: $\hat{Y}_i = -2.114 + 0.791 x_{i2}$

$R^2 = 0.865$

Model 12: $\hat{Y}_i = -2.247 + 4.655 x_{i1} - 0.109 x_{i2}$

$R^2 = 0.980$

# Partial Confounding Example

| Source of variation | d.f. | SS | MS | F | p-val |
|---|---|---|---|---|---|
| reg on $x_1$ | 1 | 46.215 | 46.215 | 146.6 | 0.0012 |
| reg on $x_2$ after $x_1$ | 1 | 0.073 | 0.073 | 0.23 | 0.6639 |
| error | 3 | 0.946 | 0.315 | | |
| corrected total | 5 | 47.233 | | | |

| Source of variation | d.f. | SS | MS | F | p-val |
|---|---|---|---|---|---|
| reg on $x_2$ | 1 | 40.859 | 40.859 | 129.6 | 0.0015 |
| reg on $x_1$ after $x_2$ | 1 | 5.428 | 5.428 | 17.21 | 0.0254 |
| error | 3 | 0.946 | 0.315 | | |
| corrected total | 5 | 47.233 | | | |

# Multiple Regression
# Interpreting Regression Coefficients

$$Y_i = \beta_o + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik} + \epsilon_i$$

- $\beta_j$ is the $j$th regression coefficient or the $j$th *partial* regression coefficient

- $\beta_j$ is the change in the mean of $Y$ for a unit change in $X_j$ **with all other variables held constant**

- Sometimes this is not possible and the values of other explanatory variables change when $X_j$ changes: (e.g., polynomial terms $(X_j, X_j^2)$ or interaction terms $(X_i, X_j, X_i X_j)$ or other highly correlated predictors)

# Multiple Regression
# Interpreting Regression Coefficients

$$Y_i = \beta_o + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik} + \epsilon_i$$

- An alternative interpretation: $\beta_j$ is the linear effect of $X_j$ on $Y$ after adjusting for the linear effect of the other predictors on $Y$ and the linear effects of the other predictors on $X_j$

- Let $P_{-x_j}$ represent the projection matrix without variable $X_j$ (delete column $j+1$ from the model matrix $X$). Then, $\hat{\beta}_j$ is found from the regression of $(I - P_{-x_j})Y$ on $(I - P_{-x_j})X_j$

# Multiple Regression
# Interpreting Regression Coefficients

- Example: Brain size data (an observational study)
  *Question:* Do species with longer gestation times have bigger brains?

- Plots, biology $\Rightarrow$ linear in log variables

- Model 1: $log(brain)_i = \beta_0 + \beta_1 log(gest_i) + \epsilon_i$

  $\hat{\beta}_1 = 2.23 \Rightarrow$ Species differing by 1 unit log gestation time (e.g. log(gest) = 2 and log(gest)=1) differ in log(brain size) by 2.23 units, on average.

- Biology $\Rightarrow$ body size associated with both

# Multiple Regression
# Interpreting Regression Coefficients

- Model 2:

  $$log(brain_i) = \beta_0 + \beta_1 log(gest_i) + \beta_2 log(body_i) + \epsilon_i$$

  $$\hat{\beta}_1 = 0.668$$

  Two species with the same body size but differing by 1 unit log gestation time differ in log brain size by 0.668 units, on average.

- So, when is $\beta_j$ in multiple regression equal to $\beta_j$ from simple linear regression?

  Answer: When $X_j$ is uncorrelated with the rest of the explanatory variables.

# Multiple Regression
# Interpreting Regression Coefficients

- Consider the regression of one set of residuals
  $(I - P_{-x_j})Y$ on another set of residuals $(I - P_{-x_j})X_j$

  - Regress log(brain) on log(body):
    $$\text{residual} = e_i = (I - P_{-x_j})Y$$

  - Regress log(gest) on log(body):
    $$\text{residual} = g_i = (I - P_{-x_j})X_j$$

  - $\beta_2$ is regression coefficient for regression of $e_i$ on $g_i$:

  $$e_i = \beta_2 g_i + \eta_i$$