

Ch 6

All of these
can be done
in proc reg
and glm
we just
change
model statement

First-Order Model w/ two predictors: $Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i$
 $E(Y) = \beta_0 + \beta_1 X_1 + \beta_2 X_2$ ↑ response plane

First-Order Model w/ more than 2 predictors: $Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_{p-1} X_{i,p-1} + \epsilon_i$

$$E(Y) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_{p-1} X_{p-1}$$

Assumptions: $\beta_0, \beta_1, \dots, \beta_{p-1} \rightarrow$ parameters

$X_1, X_2, \dots, X_{p-1} \rightarrow$ known constants

$\epsilon_i \sim N(0, \sigma^2)$

Qualitative Variables

Dummy variables: $X = \begin{cases} 0 & \text{if group 1} \\ 1 & \text{if group 2} \end{cases}$

Polynomial Regression $Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i1}^2 + \beta_4 X_{i2}^2 + \beta_5 X_{i1} X_{i2} + \epsilon_i$

Transformed Regression $\ln(Y_i) = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i$

Generalized Linear Regression $Y_i = c_{i0} \beta_0 + c_{i1} \beta_1 + \dots + c_{i,p-1} \beta_{p-1} + \epsilon_i$

In matrix terms:
$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & X_{11} & X_{12} & \dots & X_{1,p-1} \\ 1 & X_{21} & X_{22} & \dots & X_{2,p-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & \dots & X_{n,p-1} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

least squares criterion: $Q = \sum (Y_i - \beta_0 - \beta_1 X_{i1} - \dots - \beta_{p-1} X_{i,p-1})^2$

coefficients $b = (X'X)^{-1} (X'Y)$

predictions $\hat{Y} = Xb = HY$ where $H = X(X'X)^{-1}X'$

residuals $e = Y - \hat{Y} = Y - Xb = (I - H)Y$

cov $\sigma^2(b) = \sigma^2(I - H)$

s^2 $s^2(e) = \text{MSE}(I - H)$

$$\cdot \text{SSTO} = \underbrace{Y'Y}_{\uparrow n-1 \text{ dof}} - \frac{1}{n} Y'JY = Y' \left[I - \left(\frac{1}{n}\right) J \right] Y \quad J = \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix}$$

$$\cdot \text{SSE} = \underbrace{e'e}_{\uparrow n-p \text{ dof}} = (Y - Xb)'(Y - Xb) = Y'Y - b'Y'Y = Y'(I - H)Y$$

$$\cdot \text{SSR} = \underbrace{b'X'Y'}_{\uparrow p-1 \text{ dof}} - \left(\frac{1}{n}\right) Y'JY = Y' \left[H - \frac{1}{n} J \right] Y$$

$$\cdot \text{MSR} = \frac{\text{SSR}}{p-1}$$

NOTE: $E(\text{MSR}) = \sigma^2$ when coefficients = 0

$$\cdot \text{MSE} = \frac{\text{SSE}}{n-p}$$

otherwise $E(\text{MSR}) > \sigma^2$

$$\cdot F^* = \frac{\text{MSR}}{\text{MSE}}$$

Tests $H_0: \beta_0 = \beta_1 = \dots = \beta_{p-1} = 0$

$$\cdot R^2 = \frac{\text{SSR}}{\text{SSTO}} = 1 - \frac{\text{SSE}}{\text{SSTO}} \quad \text{Note. Adding predictors, even just noise, can only increase } R$$

$$\cdot R_a^2 = \frac{\frac{\text{SSR}}{n-p}}{\frac{\text{SSTO}}{n-1}} = 1 - \left(\frac{n-1}{n-p}\right) \frac{\text{SSE}}{\text{SSTO}} \quad \text{Note this } R_a^2 \text{ makes up for the issue w/ } R^2 \text{ by penalizing additional predictors}$$

$$\cdot \sigma^2(b) = \sigma^2(X'X)^{-1} = \begin{bmatrix} \sigma^2(b_0) & \sigma(b_0, b_1) & \dots & \sigma(b_0, b_{p-1}) \\ \sigma(b_1, b_0) & \sigma^2(b_1) & \dots & \sigma(b_1, b_{p-1}) \\ \vdots & \vdots & \ddots & \vdots \\ \sigma(b_{p-1}, b_0) & \sigma(b_{p-1}, b_1) & \dots & \sigma^2(b_{p-1}) \end{bmatrix}$$

* covb option *

$$\cdot S^2(b) = \text{MSE}(X'X)^{-1} = \begin{bmatrix} s^2(b_0) & s(b_0, b_1) & \dots & s(b_0, b_{p-1}) \\ s(b_1, b_0) & s^2(b_1) & \dots & s(b_1, b_{p-1}) \\ \vdots & \vdots & \ddots & \vdots \\ s(b_{p-1}, b_0) & s(b_{p-1}, b_1) & \dots & s^2(b_{p-1}) \end{bmatrix}$$

$$\cdot \frac{b_k - \beta_k}{\text{se}(b_k)} \sim t_{(n-p)}$$

* clb option *

$$\circ \circ \text{ } 1 - \alpha \text{ CI for } b_k: b_k \pm t(1 - \frac{\alpha}{2}; n-p) \text{se}(b_k)$$

• Test $b_k = 0$

$$t^* = \frac{b_k}{\text{se}(b_k)}$$

Reject when $|t^*| > t(1 - \frac{\alpha}{2}; n-p)$

(2)

• Simultaneous joint confidence intervals

• g parameters to estimate jointly

• $b_k \pm t(1-\alpha/2g; n-p) s \hat{\epsilon} b_k$

Bonferroni CI

↑ Done by choosing $\alpha = \frac{\alpha}{2g}$

SAS CODE

proc reg;

model $y = x_1 x_2 x_3 x_4$ / clb covb alpha=.05; * $\alpha = \frac{\alpha}{2g}$ for bonferroni

run;

OUTPUT

| Source | DF | Sum of Sq | Mean Sq | F val | P val |
|-----------------|-------------------|-----------|---------|-------------------|---------------|
| Model | # of Coeff | SSR | MSR | $\frac{MSR}{MSE}$ | Test Stat = 0 |
| Error | $n - \#$ of Coeff | SSE | MSE | | |
| Corrected Total | n | SSTO | | | |

| Root MSE | \sqrt{MSE} | R-Square | R^2 |
|----------------|--|-------------|---------|
| Dependent Mean | \bar{y} | Adj-Rsquare | R_a^2 |
| Coeff Var | $(\frac{\sqrt{MSE}}{\text{Dep Mean}}) * 100$ | | |

1 for cont. vars
cl for class vars

| Var | DF | Parameter estimate | standard error | t value | P val | $(1-\alpha) * 100$ CI |
|-------|----|--------------------|----------------|---------------|--------------|-----------------------|
| int | | b_0 | $se(b_0)$ | $b_0/se(b_0)$ | Test $b_0=0$ | |
| x_1 | | b_1 | $se(b_1)$ | $b_1/se(b_1)$ | ' | |
| x_2 | | b_2 | $se(b_2)$ | $b_2/se(b_2)$ | ' | |
| x_3 | | b_3 | $se(b_3)$ | $b_3/se(b_3)$ | ' | |
| x_4 | | b_4 | $se(b_4)$ | $b_4/se(b_4)$ | ' | |

cov of estimates

| Var | b_0 | b_1 | b_2 | b_3 | b_4 |
|-------|---------------|---------------|---------------|---------------|---------------|
| b_0 | $s^2(b_0)$ | $s(b_0, b_1)$ | $s(b_0, b_2)$ | $s(b_0, b_3)$ | $s(b_0, b_4)$ |
| b_1 | $s(b_1, b_0)$ | $s^2(b_1)$ | $s(b_1, b_2)$ | $s(b_1, b_3)$ | $s(b_1, b_4)$ |
| b_2 | $s(b_2, b_0)$ | $s(b_2, b_1)$ | $s^2(b_2)$ | $s(b_2, b_3)$ | $s(b_2, b_4)$ |
| b_3 | $s(b_3, b_0)$ | $s(b_3, b_1)$ | $s(b_3, b_2)$ | $s^2(b_3)$ | $s(b_3, b_4)$ |
| b_4 | $s(b_4, b_0)$ | $s(b_4, b_1)$ | $s(b_4, b_2)$ | $s(b_4, b_3)$ | $s^2(b_4)$ |

How to get cov matrix?

```
SAS CODE Proc glm;
model y = x1 x2 x3 x4 / solution clparm alpha = .05;
run;
```

Output

| Source | DF | Sum of squares | Mean Square | F-value | Pr > F |
|-----------------|--------------------|----------------|-------------|-------------------|-------------------------|
| Model | # of predictors | SSR | MSR | $\frac{MSR}{MSE}$ | Tests all $\beta_i = 0$ |
| Error | $n - \text{pred.}$ | SSE | MSE | | |
| Corrected Total | n | SSTO | | | |

| R square | Coeff Var | Root MSE | Y Mean |
|----------|---|--------------|----------------------------------|
| R^2 | $\frac{\text{Root MSE}}{Y \text{ mean}} \times 100$ | \sqrt{MSE} | $\frac{1}{n} \sum Y_i = \bar{Y}$ |

Type I: sequential sum of squares: Sum of squares obtained by adding the next variable to the model given previous are a part of the model

| Source | DF | Type I SS | Mean Square | F value | Pr > F |
|--------|-----------------------|----------------------------|----------------------------|---------|---|
| x_1 | 1 for continuous | SSR(x_1) | MSR(x_1) | | Tests $\beta_i = 0$ given previous β_j in the model |
| x_2 | c-1 for factor levels | SSR($x_2 x_1$) | MSR($x_2 x_1$) | | |
| x_3 | | SSR($x_3 x_2, x_1$) | MSR($x_3 x_2, x_1$) | | |
| x_4 | | SSR($x_4 x_3, x_2, x_1$) | MSR($x_4 x_3, x_2, x_1$) | | |

Type III: Marginal sum of squares: Sum of squares obtained by adding this variable given all others are part of the model

| Source | DF | Type III SS | Mean Square | F value | Pr > F |
|--------|----------------------------|----------------------------|----------------------------|---------|--|
| x_1 | 1 for continuous variables | SSR($x_1 x_2, x_3, x_4$) | MSR($x_1 x_2, x_3, x_4$) | | Tests $\beta_i = 0$ given all other β_j in the model |
| x_2 | c-1 for factor levels | SSR($x_2 x_1, x_3, x_4$) | MSR($x_2 x_1, x_3, x_4$) | | |
| x_3 | | SSR($x_3 x_1, x_2, x_4$) | MSR($x_3 x_1, x_2, x_4$) | | |
| x_4 | | SSR($x_4 x_1, x_2, x_3$) | MSR($x_4 x_1, x_2, x_3$) | | |

| Parameter | Estimate | Standard error | t value | $P_r = t $ | $1 - \alpha$ % | Conf int |
|-----------|----------|----------------|-----------------------|---------------------|----------------|----------|
| int | b_0 | $se(b_0)$ | $\frac{b_0}{se(b_0)}$ | Tests $\beta_i = 0$ | | |
| x_1 | b_1 | $se(b_1)$ | $\frac{b_1}{se(b_1)}$ | | | |
| x_2 | b_2 | $se(b_2)$ | $\frac{b_2}{se(b_2)}$ | | | |
| x_3 | b_3 | $se(b_3)$ | $\frac{b_3}{se(b_3)}$ | | | |
| x_4 | b_4 | $se(b_4)$ | $\frac{b_4}{se(b_4)}$ | | | |

① Estimation & CI of mean of \hat{y} given observation vector \underline{x}

Interval Estimation of Mean Response

- Given $\underline{X}_n = \begin{bmatrix} 1 \\ x_{n1} \\ x_{n2} \\ \vdots \\ x_{np-1} \end{bmatrix}$, $E(Y_n) = \underline{X}_n' \underline{\beta}$ and $\hat{Y}_n = \underline{X}_n' \underline{b}$
 - $E(\hat{Y}_n) = E(Y_n)$
 - $\sigma^2(\hat{Y}_n) = \sigma^2 \underline{X}_n' (\underline{X}' \underline{X})^{-1} \underline{X}_n$
 - $s^2(\hat{Y}_n) = \text{MSE}(\underline{X}_n' (\underline{X}' \underline{X})^{-1} \underline{X}_n) = \underline{X}_n' s^2(\underline{b}) \underline{X}_n$
- $1 - \alpha$ Confidence Interval for $E\{\hat{Y}_n\}$

$$\hat{Y}_n \pm t(1 - \alpha/2, n-p) s(\hat{Y}_n)$$

$p = \# \text{ of predictors} + 1$

② Estimation & CI of prediction \hat{Y}_n given new observation

Interval Estimation of Prediction

- Given $\underline{X}_n = \begin{bmatrix} 1 \\ x_{n1} \\ x_{n2} \\ \vdots \\ x_{np-1} \end{bmatrix}$, $s^2(\hat{Y}_n) = \text{MSE}(\underline{X}_n' (\underline{X}' \underline{X})^{-1} \underline{X}_n) = \underline{X}_n' s^2(\underline{b}) \underline{X}_n$
 - $s^2(\text{pred}) = \text{MSE} + s^2(\hat{Y}_n) = \text{MSE}(1 + \underline{X}_n' (\underline{X}' \underline{X})^{-1} \underline{X}_n)$
- $1 - \alpha$ Confidence Interval for Prediction of new observation

$$\hat{Y}_n \pm t(1 - \alpha/2; n-p) s\{\text{pred}\}$$

③ Estimation & CI of mean of m predictions \hat{Y}_n given new observation

Interval Estimation of Mean of m Predictions

- Given $\underline{X}_n = \begin{bmatrix} 1 \\ x_{n1} \\ x_{n2} \\ \vdots \\ x_{np-1} \end{bmatrix}$, $s^2(\hat{Y}_n) = \text{MSE}(\underline{X}_n' (\underline{X}' \underline{X})^{-1} \underline{X}_n) = \underline{X}_n' s^2(\underline{b}) \underline{X}_n$
 - $s^2(\text{pred mean}) = \frac{\text{MSE}}{m} + s^2(\hat{Y}_n) = \text{MSE}(\frac{1}{m} + \underline{X}_n' (\underline{X}' \underline{X})^{-1} \underline{X}_n)$
- $1 - \alpha$ Confidence Interval for Mean of m Predictions

$$\hat{Y}_n \pm t(1 - \alpha/2; n-p) s\{\text{pred mean}\}$$

SAS CODE Reg

```
proc reg;
  model y = x1 x2 x3 x4 / clm cli;
run;
```

SAS CODE GLM

```
proc reg;
  model y = x1 x2 x3 x4 / clm cli;
run;
```

Diagnostics & Remedial Measures

Scatter Plot Matrix

| | | |
|---|----------------|----------------|
| Y | | |
| | X ₁ | |
| | | X ₂ |

SAS CODE

```
proc sgscatter;
  matrix y x1 x2;
run;
```

| | | | |
|----------------|---------------------------------------|---|---|
| | Y | X ₁ | X ₂ |
| Y | r_{yy} | r_{yx_1} p-value for $\beta_1=0$ | r_{yx_2} p-value for $\beta_2=0$ |
| X ₁ | r_{x_1y} p-value for $\beta_1=0$ | $r_{x_1x_1}$ | $r_{x_1x_2}$ p-value for $\beta_{x_2}=0$ |
| X ₂ | r_{x_2y} p-value for $\beta_2=0$ | $r_{x_2x_1}$ p-value for $\beta_{x_1}=0$ | $r_{x_2x_2}$ |

SAS CODE

```
proc corr pearson;
  var y x1 x2;
run;
```

Note: We use these to check for multicollinearity among predictors

Residual Plots

Code, same as before page 11

- Residuals vs. predictors - same as before, check marginal plots to check if curvature effect is needed
- Residuals vs. new predictors - See correlation with response to gauge worthiness of model
- Residuals vs. Fitted Values
Abs(Residuals) vs. Fitted Values → Check for megaphone shape → non constant variance
- If this is an issue, plot abs(residuals) vs predictor and find those causing the issue marginally

(2)

- Same Normality tests on pg 14 apply
- Same constant variance tests apply pg 13
- Same F Test for Lack of fit applies pg 14
- Same remedies on pages 15-16 apply