

STA 302 / 100!

Note Title

11/10/2011

My new mantra: Answers must be
unambiguous, accurate
deep — for any marks
— for a good mark

Back to Breukelen example - unequal variance

Untransformed data, no weights

$$\text{fitted equation: } \hat{y} = 1886.2 - 54.0 x \quad (\text{Volts})$$

(time)

Untransformed data, with weights $\propto \frac{1}{\text{var at each } x_i}$

$$\text{fitted equation: } \hat{y} = 119.7 - 3.13 x \quad \text{Volts.}$$

time

Equation is completely — because large values that occurred

at values of x at which there is large variance are down weighted
 Looked at standardized residuals for weighted LS.
 ↳ use these always for weighted LS rather than raw residuals.

Student. vs. p. $\frac{r_i}{\sqrt{v_i}}$ $\frac{1}{\sqrt{n}}$ $\frac{1}{\sqrt{y}}$ - increasing, possibly curvature, approx. constant variance.

Normal quantile plot - distribution of standardized residuals is right-skewed.

Weighted LS fixed non-constant variance problem but other problems

Need a transformation to fix other problems

If assumptions of normality and independence hold, the same general procedure for t tests and CIs (t.i.F) apply as in unweighted least squares

$$\text{df } t_{obs} = \frac{b_{1W}}{s_w}$$

where

$$s_w = \sqrt{\frac{\sum w_i (x_i - \bar{x})^2}{n-1}}$$

where $\sum w_i (x_i - \bar{x})^2 = \sum w_i (x_i - \bar{x})^2$

$$\text{and } S_M^2 = \frac{\sum w_i (y_i - b_{0M} - b_{1M} x_i)^2}{n-2}$$

Under $H_0: \beta_1 = 0$, t_{obs} is an observation from t_{n-2} distribution

MULTIPLE REGRESSION

- more than one explanatory variable

Why do this? - multiple x 's arise naturally

- want to control for some x 's to consider effect on y of other x 's
over and above control variables

- want to fit a polynomial

- want to compare regression line for
2 or more groups

Multiple Regression

p = number of explanatory variables in the model

$$\text{Model } \underline{Y}_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_p X_{ip} + e_i$$

$i = 1, \dots, n$

Need to estimate $p+2$ parameters ($p+1$ β 's + σ^2)

Need $n > p+2$ observations (minimum)

In matrix notation,

$$Y = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}, \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}, \quad e = \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix}$$

$Y = X\beta + e$

"Design matrix"

$$X = \begin{bmatrix} 1 & X_{11} & X_{12} & \dots & X_{1p} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & X_{21} & X_{22} & \dots & X_{2p} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & X_{n1} & X_{n2} & \dots & X_{np} \end{bmatrix} \quad n \times (p+1)$$

Gauss-Markov Assumptions:

$$E(\tilde{e} | X) = 0$$

$$\text{Var}(\tilde{e} | X) = \sigma^2 I$$

For inference we need \tilde{e} have a multivariate normal distribution:

Least squares estimator:
$$\tilde{b} = \frac{X'X}{X'X} X'Y$$

$$\tilde{b} = \begin{pmatrix} b_0 \\ b_1 \\ \vdots \\ b_p \end{pmatrix}$$

Residuals
$$\tilde{e} = Y - \hat{Y} = Y - X\tilde{b} = Y - X \frac{X'X}{X'X} X'Y$$

$$= (I - H)Y$$

Estimate of σ^2 : $S^2 = MSE = \frac{\sum \hat{e}_i^2}{df \text{ error}} = \frac{\hat{e}'\hat{e}}{df \text{ error}}$

$df \text{ error} = n - p - 1$

$RSS = \hat{e}'\hat{e} = Y'(I-H)Y$ $\text{rank}(I-H)$

$= \text{rank}(I) - \text{rank}(H)$
 $= n - (p+1)$

assuming columns
of X are linearly
independent

Need $E(RSS) = (n - p - 1)\sigma^2$ to show S^2 is unbiased

estimator for σ^2

Proof: same as SLR proof

$$\text{except } \text{trace}(I - H) = \text{trace}(I) - \text{trace}(H)$$

$$H = X(X'X)^{-1}X'$$

$$\text{trace}(H) = \text{trace}(X'X(X'X)^{-1})$$

$$= n - \text{trace}(\underbrace{X'X(X'X)^{-1}}_{I_{p+1}})$$

//

$$= n - (p+1)$$

Rainfall Example

y_i : corn yield (bushels / acre) in 6 US States, 1890 - 1927

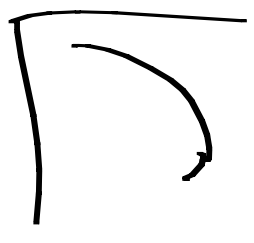
Explanatory variable: rainfall (inches)

Straight line not adequate

Non-monotonic relationship, curvature

Fit quadratic model:

$$y = \beta_0 + \beta_1 \text{rain} + \beta_2 \text{rain}^2 + e$$



- multiple regression: 2 predictors: rain, rain²
= still linear regression (linear in β 's)

$$p = 2$$

X matrix has 3 columns:

$$n = 38$$

$$\text{df for error} = 38 - 3 = 35$$

$$\text{df for SST} = n - 1 = 37$$

$$\text{df for SSR}_{\text{reg}} = 2 = p$$

Fitted equation:

$$\widehat{\text{yield}} = -5.015 + b_0 + b_1 \text{rain} - 0.229 \text{rain}^2$$

In general, in multiple regression interpret coefficient of a predictor variable (β_j or its est. b_j) as change in mean value of Y associated with a one unit change in the predictor variable (X_j) with all other variables held constant.

For rainfall example, this is impossible.

Can say for this example

If rain increases from 8" to 9", estimated change in mean of yield is 6.004 = .229 (81-64)
 ~ 2.1 bn/acre.

If rain increases from 14" to 15", estimated change in mean of yield is

$$6.504 - .229(15^2 - 14^2) = -.637 \text{ bn/acre}$$

Is coefficient of quadratic term = 0?

Test: $H_0: \beta_2 = 0$ vs $H_a: \beta_2 \neq 0$

Test statistic: $t_{obs} = \frac{b_2}{\text{s.e.}(\hat{\beta}_2)}$

$$\text{Var}(\hat{\beta}) = \sigma^2 (X'X)^{-1}$$

2nd diagonal entry of $\text{Var}(\hat{\beta})$ is $\text{Var}(\hat{\beta}_2)$
Subst s^2 (=MSE) for s^2 and take square
root to get S.E. ($\hat{\beta}_2$)

$$t_{obs} = \frac{-1.229}{\sqrt{1.0886}} = -2.16$$

→ from SAS

Under H_0 , this is an observation from
a t distribution with 35 df
 $= n - p - 1$

$$p = .0140$$

Conclude that we have some evidence that coefficient of rain² is not 0 so we should keep it in the model; even with linear term in model.

Note: Test of $H_0: \beta_j = 0$ gives an indication of whether or not j th predictor variable statistically significantly contributes to the estimation/prediction of Y over and above the other predictor variables.
That is, the assumes the other variables

are in the model

R^2 in Multiple Regression

$$R^2 = \frac{SSR}{SST}$$

For rainfall example, quadratic regression $R^2 = .209 \approx .21$

Called the "coefficient of multiple determination"
It's not the square of a correlation anymore.

It's guaranteed to increase if you add more predictor variables.
Why?