## References

- G&S Chapter 4: Do the data fit the assumptions?
- KKMN Chapter 12: Regression diagnostics
- NKNW Chapter 9

## 3 Sets of issues...

1   At the level of the <u>individual observation</u>...

   - outliers in X    ---> "influence"

   - outliers in Y    ---> residuals, rough & refined

   - overall effect    ---> Cooks' "distance"

                  --->   Fit

                  --->    's

2   At the level of the <u>model</u>

   - model mis-specified

     - variables omitted

     - have important variables, but mis-specify form of equation

  Analysis of Residuals key in -1- and -2-

-3- At the level of the <u>variables</u>

   - joint distribution of variables less than optimal
     (multi)collinearity  (Chapter 5)

## 1. AT LEVEL OF THE INDIVIDUAL OBSERVATION...

### Potential to Influence fitted model
*(see G&S Chapter 4 pages 130-134)*

   • Leverage ( h )

     • function of x's (cf. [ X - $\bar{x}$ ] term in SE for projection
     • $0 < h \le 1$;

       **average** value: $\dfrac{\#\ terms\ \text{in equation}}{n}$ [# includes $_0$ ]

       **"watch out for"** value: $> 2 \times$ **average** value

### Residuals

*(see G&S Chapter 4 pages 119-130  and 134-136)*

- *unscaled* ("raw"):       $e = y - \hat{y}$

- *standardized:*       $e\ /\ RMSE$

- *refined*: e's at "X-edges" are less variable: $\text{var}(e) = {}^2(1 - h)$

       *studentized:*

       $e\ /\ \{\ RMSE\ \sqrt{1 - h}\ \}$

- *refined further*    to make outliers "stand out",
               re-estimate RMSE by deleting observation

       *studentized deleted:*

       $e\ /\ \{\ RMSE_{(-i)}\ \sqrt{1 - h}\ \}$

## 1. AT LEVEL OF THE INDIVIDUAL OBSERVATION...

**Residuals** (see comment re studentized residuals at top of p 135 of G&S)

| G&S text name | symbol | definition | SAS INSIGHT name | prefix | SAS Proc REG Keywords in Output and Plot subcommands | NKNW text |
|---|---|---|---|---|---|---|
| unscaled / "raw" | $e$ | $y - \hat{y}$ | residual | **R_Y** | residual | $e$ |
| standardized | $e_s$ | $\dfrac{y - \hat{y}}{RMSE}$ | | | | $e^*$ |
| studentized/ internally standardized | $r$ | $\dfrac{y - \hat{y}}{RMSE \sqrt{1 - h}}$ | standardized | **RS_Y** | student | $r$ |
| studentized deleted / jackknife/ externally standardized | $r_{(-i)}$ | $\dfrac{y - \hat{y}}{RMSE_{(-i)} \sqrt{1 - h}}$ | studentized | **RT_Y** | rstudent | $t$ |

# 1. AT LEVEL OF THE INDIVIDUAL OBSERVATION...

**Actual Influence on fitted model**

### Cook's Distance (D)

• (Scaled) Distance of <u>regression coefficients</u> [b's] obtained <u>without</u> the observation from those obtained with all n observations.

• Also expressible as..
(Scaled) Distance of <u>fitted y's</u> obtained without the observation from those obtained with all n observations.

• D's "tend to follow"  $F_{(\text{\# terms}, n - \text{\# terms})}$  distribution, so

**D > 1**          **D > 4**

see also Ch 232 of KKMN

• D is a function of studentized residual (r) <u>and</u> leverage (h)

$$D = \frac{r^2}{\text{\# of terms}} \times \frac{h}{1 - h}$$

so larger if larger r <u>and</u> larger h

### Change in Fitted (predicted) Values (DFFitsS)

• the amount by which the predicted Y value for the observation in question changes when the observation itself is excluded from the analysis.

• a standardized measure

$$\text{DFFitsS} > 2 \sqrt{\frac{\text{\# terms}}{n}}$$

### Dfbeta's

• (Standardized) measure of amount by which each term in the model changes when the observation itself is excluded from the analysis.

**DFbeta > 2**

## 2. AT THE LEVEL OF THE MODEL

### Model evaluation...
*(see G&S Chapter 4 pages 145-151 and 170)*

The right variables, in the appropriate form?

- Residual Plots

    Worry more about bimodal distributions of residuals
    than lack of Gaussian-ness (e.g long tails) per se. A
    bimodal distribution might hint at an omitted binary
    covariate

### In <u>Simple</u> Linear Regression

- can directly identify need for quadratic term in X, other
    transformation etc.

### In <u>Multiple</u> Linear regression...

- plot residuals vs predicted and against each X

- plot Partial Y residuals vs. partial X residuals
    Called **Partial Leverage** Plots in INSIGHT
    (cf Annotated guide**)**

## 3. AT THE LEVEL OF THE VARIABLES...

### Multi-collinearity *cf Chapter 5 and annotated guide to Output...*

#### how to diagnose it

Pairwise correlations of X's

Variance Inflation Factors (better),together with
collinearity diagnostics

If estimates of b's "flip" (change sign)
(remember the hammock or trampoline!! )

or have very large Standard Errors

#### when does Multicollinearity matter?

if seek reasonably uncorrelated estimates of 2 or more 's
but unfavourable distribution of corresponding X's

e.g. if in a study of hearing (lung function) loss, wish to
"separate" the    for age from the    for years worked
in noisy (dusty) jobs

#### what to do about it

- drop one X (or drop the project!)
- get outside estimates for some of the    's
- increase sample size, and study more favourable X's

- if adding powers or products of other variables,
    *center* X variables first!

#### when does Multicollinearity matter less?

if merely interested in prediction even then, may want to
reduce the "dimensionality" of the X's using a technique
such as principal components