

residuals are useful for detecting lack of fit and checking model assumptions

(Q: Why residuals can do the works?)

examine $\hat{\sigma}^2 \neq \sigma^2$

$\hat{\epsilon}$ for detecting over-fitting? (LNp 6-8)

① $Y = X\beta + \epsilon = \hat{Y} + \hat{\epsilon}$

② true: $Y = X_1\beta_1 + X_2\beta_2 + \epsilon = (X_1\beta_1 + H_f X_2\beta_2) + ((I - H_f)X_2\beta_2 + \epsilon) = \hat{Y}_{X_1} + \hat{\epsilon}_{X_1}^*$ (CF)

lack of fit

fitted: $Y = X_1\beta_1 + \epsilon^*$

Leverage

$X_i [\beta_1 + (x_i^T x_i)^{-1} x_i^T X_2 \beta_2]$ check LNp 5-9

leverage: $h_i \equiv H_{ii}$ (Note 1. $\text{var}(\hat{\epsilon}_i) = (1 - h_i)\sigma^2$. Note 2. h_i is known before observing Y)

(x_i, y_i)
 $\leftrightarrow \hat{\epsilon}_i$
 $\leftrightarrow h_i$

x_i whose h_i is large $\Rightarrow \text{var}(\hat{\epsilon}_i)$ small \Rightarrow fitted model has to force to fit close to y_i
 x_i whose h_i is small $\Rightarrow \text{var}(\hat{\epsilon}_i)$ large \Rightarrow in this x_i , model cannot fit so well

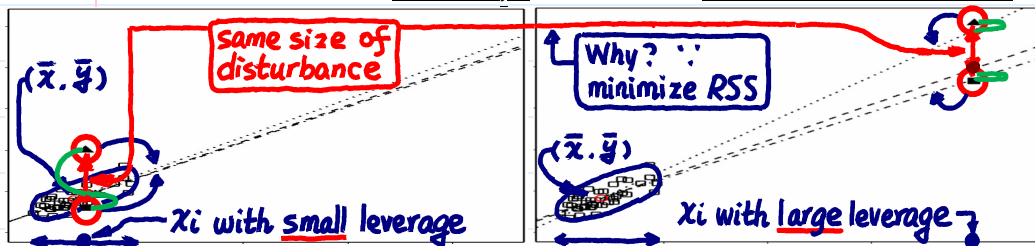
h_i roughly determines how close (x_i, y_i) to the regression surface (i.e., (x_i, \hat{y}_i))

observations with large h_i 's should be paid more attention. (Q: why?)

y_i not so close to \hat{y}_i

Recall. Over-fitting

Q: why x_i with large leverage has stronger influence on fit?



for linear model with an intercept, its fitted model must pass the point (\bar{x}, \bar{y})

$$E(y) = \beta_0 + \sum_{i=1}^{p-1} \beta_i x_i = \left(\beta_0 + \sum_{i=1}^{p-1} \beta_i \bar{x}_i \right) + \sum_{i=1}^{p-1} \beta_i (x_i - \bar{x}_i) = \beta'_0 + \sum_{i=1}^{p-1} \beta_i (x_i - \bar{x}_i)$$

orthogonality also, check LNp 3-7

centering each X

$\hat{\beta}'_0 = \bar{y} \Rightarrow E(y) = \bar{y} + \sum_{i=1}^{p-1} \hat{\beta}_i (x_i - \bar{x}_i)$

Q: why is it called leverage?

some facts about leverage: $X = [1 \ x_1 \ \dots \ x_{p-1}]$ $X^* = [x_1 - \bar{x}_1 \ \dots \ x_{p-1} - \bar{x}_{p-1}]$ p. 7-3

length of C.I. for prediction (LNp 5-4)

h_i corresponds to a Mahalanobis distance defined by X^* (X without intercept),

actually, $H = X(X^T X)^{-1} X^T \Rightarrow h_i = 1/n + [1/(n-1)] (x_i - \bar{x})^T \hat{\Sigma}_X^{-1} (x_i - \bar{x})$ $X^* = [1 \ X^*]$

where $\hat{\Sigma}_X$ is the estimated covariance of X^*

Recall. quadratic form in LNp 4-5 LNp 5-2

\Rightarrow "extreme" x_i has large leverage $\rightarrow (n-1)\hat{\Sigma}_X = X^{*T} X^*$

\Rightarrow a little change of y_i value on point with large leverage will change the fit a lot

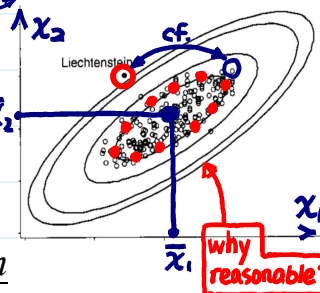
$\sum_{i=1}^n h_i = p$, and $1 \geq h_i \geq 1/n \ \forall i$

$\text{var}(\hat{\epsilon}_i) = (1 - h_i)\sigma^2 \geq 0$
 $R_i^2 \leq R_i (H^2 = H, R_i^2 + \sum_{j \neq i} H_{ij}^2 = R_i)$

trace(H) = sum of eigenvalues of H

(Q: what h_i 's are too large? Note: an average leverage is p/n)

\Rightarrow large leverage $\gg p/n \Rightarrow$ "rule of thumb": if $h_i > 2p/n$, look closer



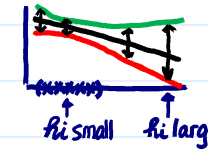
$\text{var}(\hat{Y}) = \text{var}(HY) = H\sigma^2 \Rightarrow \text{var}(\hat{y}_i) = h_i\sigma^2$ Recall C.I. of prediction (LNp 5-4-5)

$\text{Var}(\hat{\epsilon}) = \sigma^2 I$

(internally) studentized residuals r_i 's: usually different (When can they be identical?)

because $\text{var}(\hat{\epsilon}_i) = (1 - h_i)\sigma^2$, let $r_i = \hat{\epsilon}_i / [(1 - h_i)^{1/2} \hat{\sigma}]$, then $\text{var}(r_i) \approx 1$ (if model assumptions are correct)

$\text{cor}(r_i, r_j) \approx 0$



- non-constant variance removed
- dependence is very small in practice
- sum of r_i 's is not zero ($\sum \hat{\epsilon}_i = 0$)
- r_i is slightly correlated with \hat{y}_i .
- studentized residuals are preferred in residual plots ($\hat{\epsilon} \perp \hat{Y}$ and $\text{cov}(\hat{\epsilon}, \hat{Y}) = 0$)
- if there is some underlying heteroscedasticity (i.e., violation of $\text{var}(\epsilon) = \sigma^2 I$) in the errors, studentization cannot correct it

← (check LNp.7-1) unusual observations → **Outlier** → could carry misleading
useful information p. 7-4

• an outlier is a point that does not fit the current model (Q: possible cause?)

⇒ usually, large residual (Q: why?) → ①, ② in LNp.7-2.

• Q: is there a problem if (raw or studentized) residuals are used to detect outliers?

⇒ outliers may affect the fit (see plot)

⇒ cross-validated, leave-one-out

• **idea**: exclude i^{th} observation and re-compute the estimates to get $\hat{\beta}_{(i)}$ and $\hat{\sigma}_{(i)}$,

where (i) denotes that the i^{th} case has been excluded. Then, $\hat{\beta}_{(i)}$ & $\hat{\sigma}_{(i)}$ not influenced by the i^{th} observation

consider $y_i - \hat{y}_{(i)}$, where $\hat{y}_{(i)} = x_i^T \hat{\beta}_{(i)}$. (Q: why is it better in detecting outliers?)

$x_i^T \beta + \epsilon_i \rightarrow x_i^T \hat{\beta}_{(i)}$
 $var(y_i - \hat{y}_{(i)}) = \sigma^2(1 + x_i^T(X_{(i)}^T X_{(i)})^{-1}x_i)$ (Hint. prediction of future observation)

(x_i, y_i) can be regarded as a future obs. for the fitted model

• jackknife (or externally studentized, or crossvalidated) residuals developed without using the i^{th} observation

bias correction
 $t_i = (y_i - \hat{y}_{(i)}) / [(1 + x_i^T(X_{(i)}^T X_{(i)})^{-1}x_i)^{1/2} \hat{\sigma}_{(i)}]$

which are distributed as $t_{(n-1)-p}$ under null, if model is correct and $\epsilon \sim N(0, \sigma^2 I)$

∴ i^{th} obs excluded

• a simpler way to calculate t_i (avoid doing n regression)

$H_0^{(i)}$: i^{th} obs. not an outlier
 $H_A^{(i)}$: i^{th} obs. is an outlier
 $i = 1, 2, \dots, n$

raw residual $y_i - \hat{y}_i$
 $t_i = \hat{\epsilon}_i / [(1 - h_i)^{1/2} \hat{\sigma}_{(i)}] = r_i \cdot ((n-p-1)/(n-p-r_i^2))^{1/2}$

• test for outliers

externally → $\hat{\sigma}$
 internally → $\hat{\sigma}$
 studentized residual → 2-sided test

➤ given a specific case i , conclude an outlier if $|t_i| > t_{n-1-p}^{(\alpha/2)}$

$(X_{(i)}^T X_{(i)})^{-1} = (X^T X)^{-1} + \frac{(X^T X)^{-1} x_i x_i^T (X^T X)^{-1}}{1 - r_i}$

➤ in practice, a few (or all) t_i 's will be tested ⇒ problem of multiple testing → LNp.4-16

need to adjust the significance level of the test accordingly $\alpha \leq P(\cup RR_i | H_0) = \alpha^*$

Bonferroni correction
 H_0 : no outlier in the n observations against H_1 : at least one outlier

$1 - \alpha^* = 1 - \text{Prob}(\text{Type I error} | H_0) = \text{Prob}(\text{all tests accept} | H_0)$



$= 1 - \text{Prob}(\text{at least one rejected} | H_0) \geq 1 - \sum_i \text{Prob}(\text{test } i \text{ rejects} | H_0) = 1 - n\alpha$

⇒ conclude an outlier if $|t_i| > t_{n-p-1}^{(\alpha/2n)}$ if $n\alpha = 0.05$ (↔ $\alpha = 0.05/n$), then $\alpha^* \leq 0.05$
 ⇒ it's conservative, tends not to label points as outliers (especially when n large)

• some problems about outliers:

➤ two or more outliers next to each other can hide each other (see an example in Lab)

➤ change model
 an outlier in one model may not be an outlier in another when the variables have been changed or transformed ⇒ reinvestigate outliers when model changed

➤ the error distribution may not be Normal ⇒ larger residuals may be expected.
 large n. → e.g. heavy-tailed error like Cauchy.

➤ for large datasets, individual outliers are usually much less of a problem from the viewpoint of fit. In this case, it's still worth identifying outliers if these types of points are worth knowing about in the particular application. For large datasets, we need only worry about clusters of outliers. these few units/subjects may carry different information or phenomena

• **Q:** What should be done if some observations are identified as outliers?

- check for a data entry error first ← *data cleaning*, **隨** ⇒ outlier | **規** ⇒ outlier ←
- Examine the physical context (sometimes, outliers may have physical significance)
- exclude the point from the analysis
 - try to re-include later if model changed
 - if exclude permanently, report
- dangerous to exclude them in an automatic manner ← e.g. NASA, 1985.

Keep model
change data

Keep data
change model

For heavy-tailed error distribution

❖ Reading: Faraway (2005, 1st ed.), 4.2.2 ❖ Further reading: D&S, 8.1

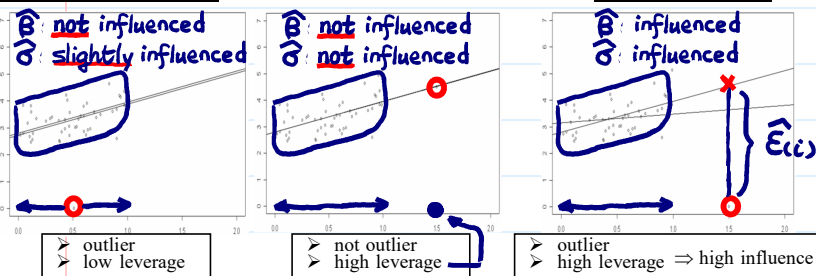
∴ [leverage residual] → **Influential observation**

large or medium leverages large or medium residuals

individual pattern

$\hat{\beta}, \hat{y}, \hat{\sigma}$

- Each observations have different influence/contribution to the fitted model. Our fitted model should not change too much (i.e., robust) just because of adding/dropping a specific observation.
- influential point: one whose removal from data would cause large change in the fit.
- an influential point may or may not be an outlier and may or may not have large leverage but it will tend to have at least one of those two properties.



- Q₁:** what points are more influential?
Q₂: what information should be included in the detection of influential points?

• **measures of influence** (**Q:** how to numerically characterize “large change in fit”?) p. 7-7

- Cross-validation** ➢ change in coefficients: $\hat{\beta} - \hat{\beta}_{(i)}$ (**Q:** how large is large?)
 ➢ change in fit: $\hat{Y} - \hat{Y}_{(i)}$ (**Q:** how large is large?)

Note.
 ① $\hat{\beta} - \hat{\beta}_{(i)}, \hat{Y} - \hat{Y}_{(i)}$ have units and scales. We may use physical standard to examine whether the change is physically significant.
 ② $(\hat{\beta} - \hat{\beta}_{(i)}) / \hat{\beta}, (\hat{Y} - \hat{Y}_{(i)}) / \hat{Y}$
 ③ compare $\hat{\beta}_{(i)} / \text{se}(\hat{\beta})$ to $\hat{\beta} / \text{se}(\hat{\beta})$

➢ Cook's statistics/distances (scale and unit free):

Note. $\text{Var}(\hat{\beta}) = (X^T X)^{-1} \sigma^2$ $D_i = (\hat{\beta} - \hat{\beta}_{(i)})^T (X^T X) (\hat{\beta} - \hat{\beta}_{(i)}) / (p \hat{\sigma}^2)$

standardization. Why? ∴ $\hat{\beta}$ have different units
 $= (\hat{Y} - \hat{Y}_{(i)})^T (\hat{Y} - \hat{Y}_{(i)}) / (p \hat{\sigma}^2)$
 $= (1/p) r_i^2 (h_i / (1 - h_i))$ ← **leverage part**

cf. Mahalanobis distance

⇒ it's a combination of residual and leverage. (**Q:** what are the effects of

residual and leverage on Cook's statistic?) **like in sampling model (LNp.5-7)** null: F-dist.

$D_i > 1$ has been suggested

⇒ **Q:** how large is large? If assume X is multivariate Normal, can do a test on D_i . However, normality may not be a reasonable assumption in practice.

➢ Others: DFFITS, Atkinson's modified Cook's statistics

- suggestion for identifying influential points: use relatively large D_i as an indication of a possible problem, then examine their $\hat{\beta} - \hat{\beta}_{(i)}$ and/or $\hat{Y} - \hat{Y}_{(i)}$.

$\begin{matrix} \uparrow XB \\ Y_X \sim N(\mu_X, \sigma_X^2) \end{matrix}$

❖ Reading: Faraway (2005, 1st ed.), 4.2.3 ❖ Further reading: D&S, 8.3, 8.4

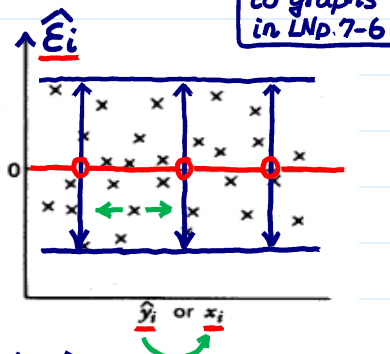
Residual plots (better to (i) remove outliers & influential obs (ii) use studentized residuals)

$(\sigma_X \leftrightarrow \mu_X)$ $\mu(\sigma_X \leftrightarrow X)$

- residual plots: plot residuals (or absolute values of residuals) against (i) \hat{y} , (ii) x_k (for predictors in model and not in model), (iii) combination (or transformation) of x_k 's, (iv) time order (if available), (v) any other quantities relevant to residuals

(**Q:** why draw residual plots?) → check ①② in LNp.7-2

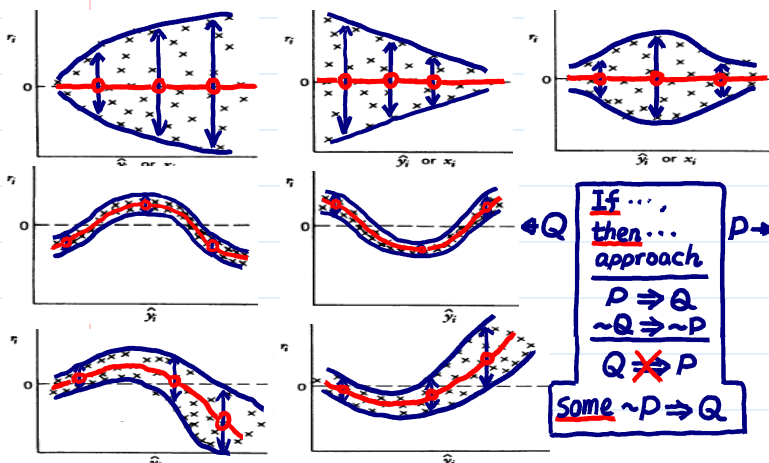
- in residual plots,
 - find overall patterns from the shape of all points (cf., residuals used in checking outliers or influential points \Rightarrow identifying individually unusual point)
 - check assumptions: (i) non-constant variance; (ii) incorrectly specified mean structure (i.e., $E(Y)=X\beta$ too simple (lack of fit))
 - rather subjective
- a satisfactory residual plot (null plot)
 - constant variance \checkmark
 - no curvature in the mean of residuals \checkmark



Note: one satisfactory residual plot cannot guarantee the

cf. residual plots for other variables will be satisfactory

- some possible patterns in unsatisfactory residual plots: (alternative)

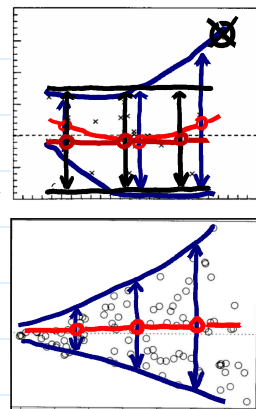


evidence of non-constant variance

curvature in the mean of residuals \Rightarrow evidence of incorrectly specified mean structure

evidence of non-constant variance and incorrectly specified mean structure

- unfortunately, in real data set, it's rare the pattern is so clear (Q: what will you conclude from the residual plot on the right?)
- in models with many terms or models with complex non-linear mean structure, cannot necessarily associate shapes in a residual plot with a particular problem with the assumptions, e.g.,



(Lnp. 7-2) true model: $E(Y) = |x_1| / [2 + (1.5 + x_2)^2]$ with constant variance
 fitted model: $E(Y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$ approximate

nonlinear mean structure

- possible remedies for unsatisfactory residual plots

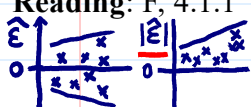
unsatisfactory residual plot	plot residuals against ...		
	\hat{y}	x_k	time order
non-constant variance	1. weighted least square 2. transform y	1. weighted least square 2. transform y	weighted least square
curvature in mean structure	1. add extra term 2. transform y	1. add extra term of x_k 2. transform y	add term of time in model

❖ Reading: F, 4.1.1

❖ Further reading: D&S, 2.5

change current model

11/24



① (Lnp. 2) \rightarrow Non-constant variance \leftarrow overall pattern

- if not sure, plot absolute values of residuals against \hat{y} , x_k 's, time order $\hat{\beta}_{WLS}, \hat{\beta}_{GLS}$
- when non-constant variance exists, $\hat{\beta}_{OLS}$ will be more variable than the best estimates ($\hat{\beta}_{OLS}$ unbiased but not BLUE) and $\hat{\sigma}$ wrong (\Rightarrow test and C.I. inaccurate)