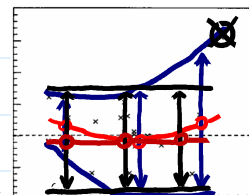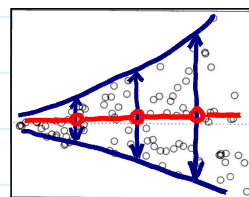p. 7-9

➤ unfortunately, in <u>real data set</u>, it's <u>rare</u> the <u>pattern is so clear</u>
  (**Q**: <u>what will you conclude</u> from the <u>residual plot on the right</u>?)

➤ in models with <u>many terms</u> or models with <u>complex non-linear</u>
  <u>mean structure</u>, <u>cannot necessarily associate shapes</u> in a residual
  <u>plot</u> with a <u>particular problem</u> with the assumptions, e.g.,

(LNp.7-2) ① ②  true model:  $E(Y)=|x_1|/[2+(1.5+x_2)^2]$ with <u>constant variance</u>
  fitted model:  $E(Y)=\underline{\beta_0}+\underline{\beta_1}\underline{x_1}+\underline{\beta_2}\underline{x_2}$ → approximate   [nonlinear mean structure]

• <u>possible remedies</u> for <u>unsatisfactory</u> residual plots

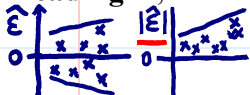| unsatisfactory residual plot | plot residuals <u>against</u> ... | | |
|---|---|---|---|
| | $\hat{y}$ | $x_k$ | time order |
| non-constant variance | 1. weighted least square<br>2. transform $y$ | 1. weighted least square<br>2. transform $y$ | weighted least square |
| curvature in mean structure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

*different philoso-phy*

❖ **Reading**: F, 4.1.1      ❖ **Further reading**: D&S, 2.5     → change current model.
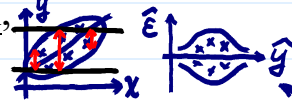
11/24

① (LNp.2) → **Non-constant variance** ← overall pattern

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

---

p. 7-10

• It's <u>better</u> try to <u>understand</u> the <u>cause</u> of <u>non-constant variance</u> before
  taking any <u>remedies</u>, e.g., (1) <u>larger response</u> have <u>more "room"</u> to <u>vary</u>,
  (2) <u>response constrained</u> to lie between a <u>maximum</u> and a <u>minimum</u>,
  (3) <u>response from Poisson distribution</u> or <u>binomial distribution</u>, …
  → count data. ①          ②
  $\Rightarrow$ <u>discovering reasons</u> to <u>support</u> the <u>remedies you are going to take</u>

$\sigma_x = \sqrt{\mu_x}$

• remedies for non-constant variance  ① $Y_x \sim P(\mu_x)$, $E(Y_x)=Var(Y_x) \Rightarrow \mu_x = \sigma_x^2 \rightarrow \hat{\varepsilon}$
  ➤ <u>weighted least squares</u>  ② $Y_x \sim Bin(n, P_x)$, $E(Y_x)=nP_x$, $Var(Y_x)=nP_x(1-P_x) \rightarrow \hat{\varepsilon}$
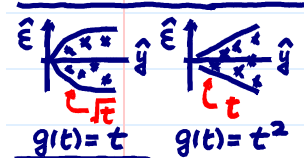    ▪ need <u>weighting information</u> (may from plotting <u>residual vs. $x_k$</u>) or
    ▪ model the <u>form of $\Sigma$</u> and using <u>IRWLS</u> ← check LNp 6-5 ↳ $(\sigma_x \leftrightarrow \chi_k)$ cf. $\boxed{\sigma_x^2 \propto g(\mu_x)}$ g-variance function
  ➤ transform $Y$ (may use <u>information</u> from plotting <u>residual vs. $\hat{y}$</u>) $(\sigma_x \leftrightarrow \mu_x)$
    ▪ idea: find a <u>transformation $h$</u> such that <u>var$(h(y_x))$ is a constant</u>, (**Q**: how? Hint:
      δ-method)

[variance stabilizing transfor-mation]
$Y=X\beta+\varepsilon$
$h(Y)=X\beta'+\varepsilon'$

$$h(y_x) \underset{=}{\approx} h(E(y))+(y-E(y))h'(E(y)) \cancel{\otimes} \Rightarrow var(h(y)) \approx var(y)[h'(E(y))]^2=c$$
  └ Taylor expansion
hope var$(h(y))$ to be a constant $c \Rightarrow h'(E(y_x)) \propto 1/(var(y_x))^{1/2}$
$$h(E(y)) \equiv \int 1/(var(y))^{1/2}\, d(E(y)) = \int \frac{1}{\sqrt{g(\mu_x)}} d\mu_x \propto \frac{1}{\sqrt{g(\mu_x)}}$$

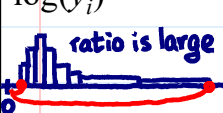$g(t)=t$   $g(t)=t^2$
$v(\sigma_x \leftrightarrow \mu_x)$   □ Example 1: $var(y_x) \propto [E(y_x)]^2 \Rightarrow$ suggest $h(y)=\log(y)$  $h(t)=\int \frac{1}{\sqrt{t^2}} dt = ln(t)$
  $g(t)=t \quad \sigma_x^2 \propto \mu_x$
$x(\sigma_x^2 \leftrightarrow \mu_x)$   □ Example 2: $var(y_x) \propto E(y_x) \Rightarrow$ suggest $h(y)=y^{1/2}$  $h(t)=\int \frac{1}{\sqrt{t}} dt \propto t^{1/2}$
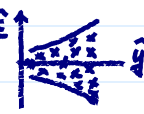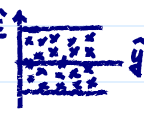
□ Note: in <u>residual plot</u>, tend to see $[var(y_x)]^{1/2}$ rather than var$(y_x)$ (example?)

**p. 7-11**

- practical **problems**:
  - □ if $y_i \leq 0$, for <u>some</u> $i$, <u>square root</u> or <u>log</u> transformations <u>fail</u> $\Rightarrow$ can do transformation on $y_i + d$, where $d$ is some small amount s.t. $y_i + d > 0$ for all $i$
  - □ transformation may make interpretation difficult — interpretation of $\hat{\beta}'$ (unit=?) cf. interpretation of $\hat{\beta}$
- <u>example of transformations</u> ← Weisberg (2005). Applied Linear Regression, Sec.8.3.1

| | | |
|---|---|---|
| sqrt($y_i$) | $var(y_i) \propto E(y_i)^{\underline{1}}$ | useful for count data from Poisson distribution |
| log($y_i$) | $var(y_i) \propto [E(y_i)]^{\underline{2}}$ | very common, good candidate if the range of $Y$ is very broad, say $y_x$'s value ranges from 10 to $10^5, 10^6$ |
| $Y|x$ ratio is large $\to y_x \to \log(y_x)$ | $\sigma_x/u_x$ = a constant (coefficient of variation) | |
| $1/y_i$ $\hat{\varepsilon}$ $t^2$ $\hat{y}$ | $var(y_i) \propto [E(y_i)]^{\underline{4}}$ $g(t)=t^4$ $h(t)=\int \frac{1}{t^4}dt = \int \frac{1}{t^2}dt \propto t^{-1}$ | appropriate when responses are "bunched" near zero, but, in markedly decreasing numbers, large responses do occur. most data / very few / heavy-tailed dist. |
| $\sin^{-1}(sqrt(y_i))$ | $var(y_i) \propto E(y_i)(1-E(y_i))$ $g(t)=t(1-t)$ | for binomial proportions $h(t)=\int \frac{1}{\sqrt{t(1-t)}}dt$ |

- ➢ <u>do nothing</u> $\Rightarrow$ because (i) $\hat{\beta}_{OLS}$ is still <u>unbiased</u>, although <u>not BLUE</u>; (ii) tests and C.I. inaccurate, but bootstrap may be used to get <u>more accurate results</u>
- ➢ use <u>generalized linear model</u> (e.g., Poisson/binomial $y \Rightarrow var(y_x)$: function of $E(y_x)$)
- • <u>formal test</u> for <u>non-constant variance</u> (or $\hat{\varepsilon}^2$)
  - ➢ <u>regressing absolute residuals</u> on $\hat{y}$ or $x_k$'s, slope>0, slope=0

---

**p. 7-12**

- ➢ <u>data with replication</u> $\Rightarrow$ can <u>estimate variances</u> of distinct $x_i$'s and <u>test</u> their <u>homogeneity</u> (see D&S, 2.2)
  - $H_0: \sigma_{x_1}^2 = \sigma_{x_2}^2 = \cdots = \sigma_{x_k}^2$ ← Bartlett's test (Recall $\hat{\sigma}_{pe}^2$ in LNp.6-9)
- ➢ <u>data without replication</u> $\Rightarrow$ assign variance a <u>model</u>, test whether parameters in the model equal zero (see Weinberg (2005), 8.3.2) — e.g. $\sigma_z^2 = \sigma^2 \cdot \exp(\lambda z)$ parameters
- ➢ <u>formal test</u> may be <u>good</u> at detecting a <u>particular kind</u> of non-constant variance $\lambda \neq 0$ (depending on the <u>alternative hypothesis</u>), but <u>always do the residual plots</u>

❖ **Reading**: Faraway (2005, 1st ed.), 4.1.1　❖ **Further reading**: D&S, 2.2, 13.6

test $H_0$ vs. $H_1$. $H_0 \cup H_1$ = all possible models

② (LNp.2) → **Curvature in the mean of residuals** ← overall pattern

- • <u>related</u> to the concept of <u>lack-of-fit</u> (tests for <u>lack-of-fit</u> can be <u>used</u> if possible), i.e., the <u>current model</u>, $E(Y)=X\beta$, may need to be <u>modified</u> for *achieving <u>better fitting</u>* — Weisberg (2005). Sec.8.2　← too simple
- ⊙ A <u>simple test</u> for <u>curvature</u>: test whether a <u>plot of residuals</u> versus a quantity $U$ (e.g., $\hat{y}$ or $x_k$'s) is a <u>null plot</u> or has <u>curvature</u> (can be any reasonable known function of $U$) $ax^2+bx+c$
  - $\Rightarrow$ <u>refit</u> the original mean structure with an <u>additional term</u> $U^2$ added
  - $\Rightarrow$ <u>significant</u> $t$-test for $U^2$ suggests <u>curvature</u> (be aware of <u>collinearity</u> between $U^2$ and <u>other terms</u> in original mean structure)
- • **Q**: how to identify *why* the <u>non-linearity</u> happened? high cor$(x,x^2)$ $\Rightarrow$ can use <u>orthogonal polynomial</u> to remove collinearity (future lecture)
  - ➢ plot <u>residuals</u> against $\hat{y}$ $\Rightarrow$ can tell you <u>whether</u> some <u>problems</u> exist, but <u>cannot tell</u> you <u>why</u>

➤ plot <u>residuals</u> against $x_k$'s or $y$ against $x_k$'s ⇒ <u>may tell</u> you <u>why</u> <span style="color:red">possible variables/effects to be added into model matrix $X$</span> this problem happened, but in multivariate regression there may <span style="color:red">e.g. $cor(x_k, x_i) \approx 1$</span> exist <u>correlation</u> between <u>predictors</u>, then it's <u>difficult</u> to <u>find why</u> <span style="color:red">$x_i$ already appears in the model matrix $X$. Then, $\hat{\varepsilon} \perp x_i$ and $cor(\hat{\varepsilon}, x_k) \approx 0$</span>

➤ *added variable* (*partial regression*) plots

**Fitted model:** $Y = X_1\beta_1 + \varepsilon$

- <u>recall:</u> <span style="color:red">$\hat{\varepsilon}_Y(x_k) = \beta_2 \hat{\varepsilon}_{x_k} + \varepsilon$</span>

$$Y = X_1\beta_1 + X_2\beta_2 + \varepsilon = (X_1\beta_1 + H_1X_2\beta_2) + ((I-H_1)X_2\beta_2 + \varepsilon)$$

1. regress $Y$ on <u>all predictors except</u> $x_k$ ⇒ get <u>residuals</u> $\hat{\varepsilon}_Y(x_k)$
2. regress $x_k$ on <u>all predictors except</u> $x_k$ ⇒ get <u>residuals</u> $\hat{\varepsilon}_{x_k}$

<span style="color:red">the estimated coefficient of $x_k$ when $x_k$ is included in model matrix</span>

<span style="color:red">remove collinearity btw $x_k$ & the other predictors</span>

▫ $\hat{\varepsilon}_Y(x_k)$: <u>part of $Y$ not explained</u> by <u>all predictors except</u> $x_k$
▫ $\hat{\varepsilon}_{x_k}$ : <u>part of $x_k$ not explained</u> by <u>all predictors except</u> $x_k$

3. plot $\hat{\varepsilon}_Y(x_k)$ versus $\hat{\varepsilon}_{x_k}$ <span style="color:red">Fit simple linear regression: $\hat{\varepsilon}_Y(x_k) = \beta_0 + \beta_k\hat{\varepsilon}_{x_k} + \varepsilon$</span>

- the <u>slope</u> of a <u>fitted line</u> to the added variable plot is $\hat{\beta}_k$ and <u>intercept=0</u> (the line passes $(0, 0)$) <span style="color:blue">check LNp.7.2</span>

<span style="color:blue">∵ $\sum \hat{\varepsilon}_Y(x_k) = 0$ $\sum \hat{\varepsilon}_{x_k} = 0$ if $X_1$ contains the intercept</span>

- a <u>strong relationship</u> between the <u>plotted quantities</u> corresponds to a <u>strong adjusted relationship</u> between $y$ and $x_k$

**Fitted model:** $Y = X_1\beta_1 + X_2\beta_2 + \varepsilon$

- can be used to check if <u>new predictors</u> should be <u>included</u>

<span style="color:red">slope $\hat{\beta}_k$     slope $\hat{\beta}_k$</span>

➤ <u>partial residual plots</u> <span style="color:red">fitted linear line (Note: $\hat{\varepsilon} \perp x_k$): slope $\hat{\beta}_k$, pass $(\bar{x}_k, \hat{\beta}_k\bar{x}_k)$</span>



- plot $\hat{\varepsilon} + \hat{\beta}_k x_k$ versus $x_k$ ⇒ <u>same interpretation</u> as <u>added variable plots</u>
- $y - \sum_{j \neq k} \hat{\beta}_j x_j = \hat{y} + \hat{\varepsilon} - \sum_{j \neq k}\hat{\beta}_j x_j = \hat{\varepsilon} + \hat{\beta}_k x_k$ <span style="color:red">can use <u>Lowess</u> (future lecture) to find how $x_k$ affect $y$</span>

---

• <u>remedies</u> for <u>curvature</u> ⇒ <u>adjust</u> the <u>mean structure</u>, E($Y$)=$X\beta$, for <u>better fitting</u>

➤ <u>many many modeling techniques</u> in addition to <u>linear regression</u> can be <span style="color:blue">link function</span> adopted (<u>GLM</u>, <u>additive model</u>, <u>nonparametric regression</u>, <u>ACE</u>, <u>AVAS</u>, <u>regression trees</u>, <u>regression spline</u>, <u>MARS</u>)

<span style="color:red">additivity & variance stabiliza-tion</span>

<span style="color:red">multivariate adaptive regression spline</span> <span style="color:red">alternating coditional expectation</span> <span style="color:red">Hastie & Tibshirani (1990). Generalized additive models</span>

➤ <u>add more</u> (<u>polynomial</u> or <u>cross product</u>) <u>terms</u>
<span style="color:red">interaction</span>

- <u>may identify required terms</u> from <u>residual plot</u>, <u>added variable plot</u>, or <u>partial residual plot</u> (<u>polynomial model</u> will be <u>introduced</u> in <u>further lecture</u>)

<span style="color:red">Linear model use <u>linear structure</u> $X\beta$ to approximate the functional relationship between $y_x \longleftrightarrow x$ 規律</span>

➤ <u>transformation</u> of <u>response</u> or <u>predictors</u>. idea behind the approach:

(i) a <u>statistical model</u> is a <u>local approximation</u> of the <u>underlying system</u>

(ii) when the <u>mean structure</u> of the <u>underlying system</u> is <u>non-linear</u> and <u>complex</u>, a <u>linear approximation</u> over a <u>relatively wide range</u> of $X$ may be <u>inadequate</u> (e.g., <span style="color:red">$E[\log(Y)]$</span>

<span style="color:red">a good fitted model would require a lot of effects (e.g., check LNp 7-9)</span>

$$E(Y) = \beta_0 x_1^{\beta_1} x_2^{\beta_2}) \quad \Rightarrow \quad \log(E(Y)) = \log(\beta_0) + \beta_1\log(x_1) + \beta_2\log(x_2)$$

<span style="color:blue">cf. variance stablizing transfor-mation (LNp.10)</span>
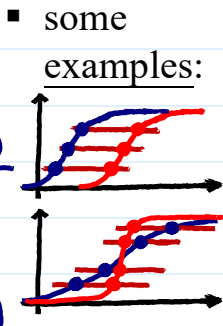
(iii) we <u>sometimes</u> can find <u>suitable transformations</u> of data that will permit a <u>non-linear model</u> to be <u>better approximated</u> (after transformation) by a <u>linear one</u> (e.g., <span style="color:red">use <u>fewer terms</u></span>

<span style="color:red">a model with only a few terms</span>

$$E(\log(Y)) \approx \log(\beta_0) + \beta_1\log(x_1) + \beta_2\log(x_2))$$

*(left margin handwritten)*
$\dfrac{z_1 - \mu_1}{\sigma_1}$
$\sim \dfrac{z_2 - \mu_2}{\sigma_2}$
$\sim N(0,1)$
If $(z_1, z_2)$
such that
$\Phi\left(\dfrac{z_1 - \mu_1}{\sigma_1}\right)$
$= \Phi\left(\dfrac{z_2 - \mu_2}{\sigma_2}\right)$
$\Rightarrow (z_1 - \mu_1)/\sigma_1$
$= (z_2 - \mu_2)/\sigma_2$

■ some examples:

| transformation | | non-linear model |
|---|---|---|
| log(y) | log(x) | $E(y) = \alpha \prod x_j^{\beta_j}$ |
| log(y) | x | $E(y) = \alpha \exp(\Sigma \beta_j x_j)$ |
| y | log(x) | $E(y) = \alpha + \Sigma \beta_j \log(x_j)$ |
| 1/y | 1/x | $E(y) = 1/[\alpha + \Sigma (\beta_j/x_j)]$ |
| 1/y | x | $E(y) = 1/(\alpha + \Sigma \beta_j x_j)$ |
| y | 1/x | $E(y) = \alpha + \Sigma \beta_j (1/x_j)$ |

■ There exists <u>numerical method</u> for finding a <u>suitable transformation</u> to <u>improve the fit</u> and/or <u>remedy non-constant variance</u> (e.g., <u>Box-Cox transformation</u>, future lectures)
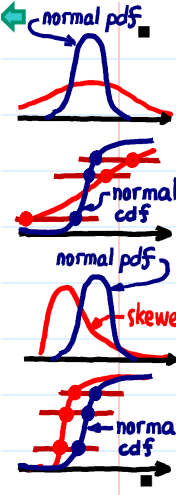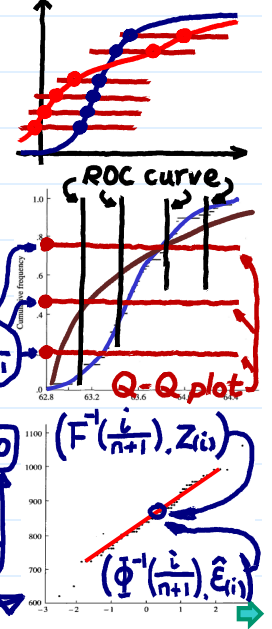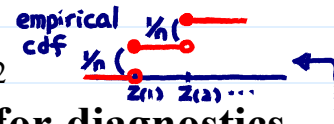
❖ **Reading**: F, 4.3, 7.2.4　　❖ **Further reading**: D&S, 8.2

*(handwritten)* empirical cdf $\frac{1}{n}$( $\frac{1}{n}$( $z_{(1)} \; z_{(2)} \cdots$

## **Various plots and tests for diagnostics**

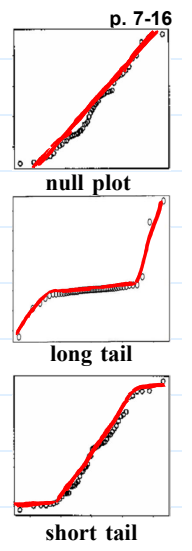*(handwritten)* quantile vs. quantile

• <u>Q-Q plot</u>

➤ **Q**: we often see the statement "$z_1, z_2, \ldots, z_m$ are i.i.d. from a <u>cdf</u> $F$", how to examine if $F$ is an <u>appropriate</u> distribution assumption for $z_j$'s? (<u>Hint</u>: examine the <u>similarity</u> btw <u>cdf</u> and <u>empirical cdf</u>)

*(handwritten)* $z_{(1)}, \cdots, z_{(n)}$ ← each probability $\frac{1}{n}$

➤ <u>normal (probability) plot</u>: assessing <u>normality assumption</u> of $\underline{\varepsilon}$ *(handwritten $\frac{i}{n+1}$)* (Note: <u>tests</u> and <u>C.I.</u> depend on <u>normality</u> assumption)

*(handwritten)* ROC curve, Q-Q plot

1. <u>sort the data</u> $\hat{\varepsilon}_{(1)} \le \hat{\varepsilon}_{(2)} \le \cdots \le \hat{\varepsilon}_{(n)}$　$\hat{\varepsilon} \sim N(0, (I-H)\sigma^2)$ if $\varepsilon \sim N(0, \sigma^2 I)$

*(handwritten)* slope = 1 intercept = 0 ; $(F^{-1}(\frac{i}{n+1}), z_{(i)})$

2. plot $\hat{\varepsilon}_{(i)}$ against $\Phi^{-1}(i/(n+1))$, where $\Phi$ is the <u>cdf of N(0, 1)</u>

■ If the residuals are <u>normally distributed</u>, an approximately *(cf)* <u>straight-line relationship</u> will be observed (<u>null plot</u>)

*(handwritten)* $(\Phi^{-1}(\frac{i}{n+1}), \hat{\varepsilon}_{(i)})$

*(left margin handwritten)* normal pdf; normal cdf; normal pdf — skewed; normal cdf

■ <u>non-normality</u>: <u>long-tail</u>, <u>short-tail</u>, <u>asymmetric</u> *(handwritten)* tend to generate data with large residuals ⇒ can remove outliers then plot again

□ <u>worst case</u> is <u>long-tail</u>; <u>mild non-normality</u> can safely be <u>ignored</u>; the <u>larger the sample size</u>, the <u>less troublesome</u> the <u>non-normality</u> *(handwritten)* Why? CLT for $\hat{\beta}$

□ for <u>long-tail</u>, (i) use <u>test</u> based on <u>other distributions</u>, or <u>bootstrap</u>, or <u>permutation tests</u> (ii) for <u>estimation</u>, use <u>robust methods</u> (e.g., <u>least absolute deviation</u> instead of <u>least square</u>)

□ <u>asymmetric</u>, <u>transform</u> $Y$ (e.g., <u>Box-Cox method</u>)

□ <u>short-tail</u> can be reasonably <u>ignored</u>

■ <u>formal tests</u> exists (such as <u>Kolmogorov-Smirnov test</u>), but <u>not</u> as <u>flexible</u> as the <u>Q-Q plot</u>

**null plot**

**long tail**

**asymmetric**　　**short tail**

➤ <u>normal plot</u> can be applied to identify <u>extreme values</u> (e.g., in <u>residuals</u>, *(handwritten)* outliers <u>leverages</u>, <u>Cook's statistics</u>, ...): in the case, <u>not interested in a straight line</u> *(handwritten: overall pattern)* relationship, but rather looking for points *that depart from the straight line* *(handwritten: individual pattern)*

*(left margin handwritten)* influential observations

• <u>half-normal plot</u> *(handwritten)* ← absolute values

*(handwritten)* cdf of N(0,1)

1. <u>sort the absolute data</u> $|\hat{\varepsilon}|_{(1)} \le |\hat{\varepsilon}|_{(2)} \le \cdots \le |\hat{\varepsilon}|_{(n)}$

2. plot $|\hat{\varepsilon}|_{(i)}$ against $\Phi^{-1}((n+i)/(2n+1))$

*(handwritten)* $\hat{\varepsilon}_i$, $h_i$, $c_i$ ; $|\hat{\varepsilon}|_{(i)}$
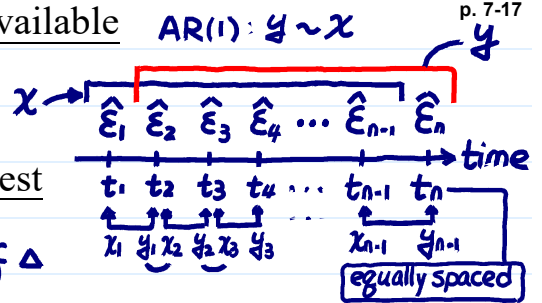
➤ usually used to identify "<u>extreme</u>" values

➤ can be used to examine <u>residuals</u>, <u>leverages</u>, <u>Cook's statistics</u>, <u>treatment effects</u> (especially for <u>experimental data without replicates</u>)

*(handwritten)* cdf of $|X|$, where $X \sim N(0,1)$ ; $|\hat{\varepsilon}_i|$, $|h_i|$, $|c_i|$, $|\hat{\beta}_i|$ ; $|\hat{\beta}_i|$'s

- diagnostic of correlated errors when a <u>time order</u> is <u>available</u>          AR(1): $y \sim x$          p. 7-17
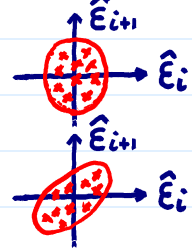
    ➢ plot $\hat{\varepsilon}$ against <u>time</u>   *lag 2, lag 3,* $\cdots$ ← *lag 1*
    
    ➢ plot $\hat{\varepsilon}_{i+1}$ against $\hat{\varepsilon}_i$, when $i$ related to <u>time</u>
    
    ➢ use <u>formal tests</u> like the <u>Durbin-Watson</u> or <u>runs test</u>

$x \to$ $\hat{\varepsilon}_1$ $\hat{\varepsilon}_2$ $\hat{\varepsilon}_3$ $\hat{\varepsilon}_4$ $\cdots$ $\hat{\varepsilon}_{n-1}$ $\hat{\varepsilon}_n$

$t_1$ $t_2$ $t_3$ $t_4$ $\cdots$ $t_{n-1}$ $t_n$ → *time*

$x_1$ $y_1$ $x_2$ $y_2$ $x_3$ $y_3$   $x_{n-1}$ $y_{n-1}$

*equally spaced*

$\Sigma(\hat{\varepsilon}_i - \hat{\varepsilon}_{i-1})^2$
$= \Sigma \hat{\varepsilon}_i^2$
$+ \Sigma \hat{\varepsilon}_{i-1}^2$
$- 2\Sigma \hat{\varepsilon}_i \hat{\varepsilon}_{i-1}$

$$DW = \frac{\sum_{i=2}^{n} (\hat{\epsilon}_i - \hat{\epsilon}_{i-1})^2}{\sum_{i=1}^{n} \hat{\epsilon}_i^2}$$

*Information of* Δ

a *Correlation matrix of* $\varepsilon_i$'s $(y_i$'s$)$

*block-diagonal correlation matrix (LNp.6-1)*

cf. $\varepsilon_1$ $\varepsilon_2$ $\varepsilon_3$ $\varepsilon_4$ $\cdots$ Δ

| | | | | |
|---|---|---|---|---|
| $\varepsilon_1$ | 1 | Δ | × | □ |
| $\varepsilon_2$ | Δ | 1 | Δ | × |
| $\varepsilon_3$ | × | Δ | 1 | Δ |
| $\varepsilon_4$ | □ | × | Δ | 1 |

- $0 \leq DW \leq 4$

- positively correlated $\Rightarrow DW \to 0$

- negatively correlated $\Rightarrow DW \to 4$

- under <u>null</u> (i.e., <u>correlation=0</u>) $\Rightarrow DW \approx 2$

- <u>null distribution</u> depends on $\underline{X}$ ← $\because \hat{\varepsilon} = (I-H)\varepsilon$ and $\varepsilon \sim N(0, \sigma^2 I)$ *under null.*

    ➢ use <u>GLS</u> when you have <u>correlated errors</u>

$1 > |Δ| > |×| > |□| > \cdots$

*cor* ↓ *as time duration* ↑

❖ **Reading**: F, 4.1.2, 4.1.3          ❖ **Further reading**: D&S, 2.4, 2.7, chapter 7

$\hat{\varepsilon}_i$   $\hat{Δ} = cor(\hat{\varepsilon}_{i-1}, \hat{\varepsilon}_i) \approx 1$   $\approx$ "$y = x$"   $(0,0)$   $0 \approx$

$\hat{\varepsilon}_i$   $\hat{Δ} = cor(\hat{\varepsilon}_{i-1}, \hat{\varepsilon}_i) \approx -1$   $\approx$ "$y = -x$"   $(0,0)$   $0 \approx$

"$y \approx x$"   $\hat{\varepsilon}_t$   $0$ → $t$

"$y \approx -x$"   $\hat{\varepsilon}_t$   $0$ → $t$