

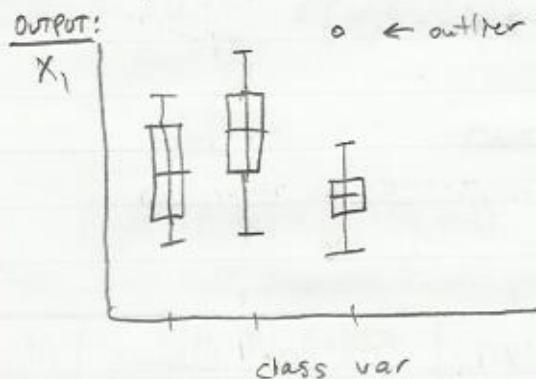
Chapter 3 Diagnostics

- Check for outlying predictor variables w/ boxplots

CODE: proc boxplot;

plot X_i * classVar / boxstyle = schematic; run;

OUTPUT:



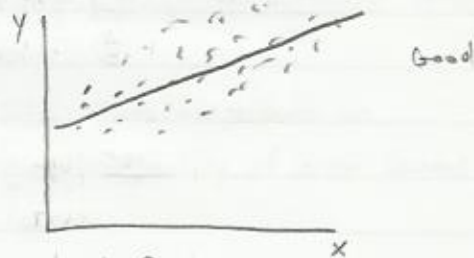
- Residuals: $e_i = Y - \hat{Y}_i$
 - $S^2(e_i) = \frac{\sum (e_i - \bar{e})^2}{n-2} = \frac{\sum e_i^2}{n-2} = \frac{SSE}{n-2} = MSE$
 - not independent Random Variables
- Studentized Residuals: $e_i^* = \frac{e_i - \bar{e}}{\sqrt{MSE}} = \frac{e_i}{\sqrt{MSE}}$
- We need to check for these problems
 - 1) The regression function isn't linear
 - 2) The error terms do not have constant variance
 - 3) The error terms are not independent
 - 4) The model fits all but one or a few outlier observations
 - 5) error terms not normally dist
 - 6) One or several important predictor variables have been omitted from the model

Problem One Non-linearity of regression function

Check using:

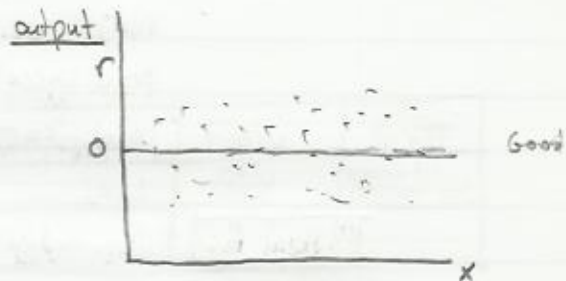
① Scatterplot: look for a linear pattern

CODE Proc sgscatter; output
plot y * x / loess;
run;



② Residual Plot: look for no pattern: balanced around 0

CODE Proc glm;
model y=x;
output out=glmData r=r;
run;
Proc sgscatter data=glmData;
plot r * x;
run;



Problem Two Nonconstant variance of errors

① Residual Plot: there should be no megaphone shape
- see code above

② Absolute value of residuals Plot: There should be no pattern or triangle shape

CODE: Proc Glim plots=all; output

model y=x;
output out=glmData r=r;
run;

data glmDataAbs;
set glmData;
absr = abs(r);
run;

proc sgscatter data=glmDataAbs;
plot absr * x;
run;

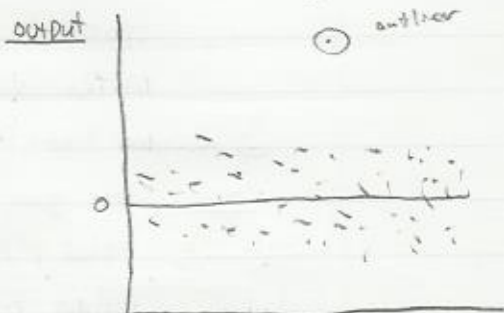


Problem Three Presence of outliers

check using

- ① Box plot - see earlier code
- ② Studentized Residuals vs. predictor: look for points far from cloud for outlying observations

```
CODE proc glm;
      model y=x;
      output out=glmdata rstudent=rs;
run;
proc sgscatter data=glmData;
  plot rs*x;
run;
```



Problem four

Non independence of error terms: plot residuals in order of time taken, there should be no patterns

```
CODE proc glm;
      model y=x;
      output out=glmData r=r;
run;
proc sgscatter;
  plot r*sequence-variable;
run;
```



Problem five

Non-normality of error terms

- ① Box plot - see earlier code: should be symmetric, no outliers
- ② QQ Plot - Given in glm diagnostics - should fit straight line
- ③ Histogram - Given in glm diagnostics - should look approx normal

If you don't like graphs there are tests! 😊

Brown Forsythe Test - Tests for statistically constant variance across groups

CODE proc glm;
 class class_var;
 model y = class_var;
 means class_var / hovtest = BF;
 run;

NOTE: Does not depend on normality of error terms

Output

Source	DF	Sum of Squares	Mean Square	F value	Pr > F
class_var				$\frac{MSC}{MSE}$	p value for constant var
error					

Breusch Pagan Test - Tests for constant variance

CODE; proc model;
 parms beta0 beta1;
 y = beta0 + beta1 * X;
 fit y / breusch = (1 x1);
 run;

Note: Assumes normal & independent errors

Equation	Test	Statistic	DF	Pr > ChiSq	Variables
	Breusch Pagan			test for constant variance	1, X

Normality tests

Shapiro-wilk's: dependent on sample size

Kolmogorov-Smirnov: dependent on vertical distance of $F_x(x)$ and $F_{normal}(x)$

Cramer-von Mises: uses empirical distribution functions > Quadratic EDF

Anderson Darling: Based on square Difference

CODE `proc univariate Normal;` output

`var x;`

`run;`

Test	Statistic	p value
Shapiro-Wilk	W	
Kolmogorov-Smirnov	D	
Cramer-von Mises	W_{sm}	
Anderson Darling	$A-g$	

F test for Lack of Fit: Test if model fits data well

- Assumes:
 - independent & normal errors
 - constant variance

CODE `proc reg;`
`model y=x/lackfit;`

`run;`

Output

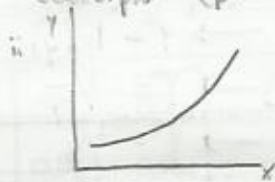
Source	DF	Sum of squares	Mean squares	F value	$P > F$
Model					All coeff = 0 p value
Error					
Lack of fit					Test regression function is linear
Pure error					
Corrected total					

- How to fix Nonlinearity of regression function
 - ① Add interaction terms
i.e. $Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \epsilon_i$
 - ② Transform data to make it linear, or more linear
- How to fix non constancy of error variance
 - ① Weighted least squares
 - ② Transform data to make it more constant
- How to fix Non independence of error terms
 - ① Use a model that allows this
- How to fix non normality of Error terms
 - ① Transform data to make errors more normal
- How to fix outlying observations
 - ① Disregard if they're not representative

Transformations

- To fix non linear relations; if error terms are approximately normal we want to transform x , not y

① Look at x^*y scatterplot (proc sgscatter)



② Transform x

i: $x^* = \sqrt{x}$
 $x^* = \ln(x)$

ii: $x^* = x^2$
 $x^* = e^x$

iii: $x^* = \frac{1}{x}$
 $x^* = e^{-x}$