

# “MULTIPLE REGRESSION AND ISSUES IN REGRESSION ANALYSIS”

MSR = Mean Regression Sum of Squares  
 MSE = Mean Squared Error  
 RSS = Regression Sum of Squares  
 SSE = Sum of Squared Errors/Residuals  
 $\alpha$  = Level of Significance  
 ML = Machine Learning

$F_c$  = Critical F taken from F Distribute Table  
 $H_0$  = Null Hypothesis  
 $H_{\alpha}$  = Alternative Hypothesis  
 X = Independent Variable  
 Y = Dependent Variable  
 F = F Statistic (calculated)

**1. INTRODUCTION**

- Multiple linear regression models are more sophisticated.
- They incorporate more than one independent variable.

**2. MULTIPLE LINEAR REGRESSIONS**

- Allows determining effects of more than one independent variable on a particular dependent variable
- $Y_i = b_0 + b_1X_{1i} + b_2X_{2i} + \dots + b_nX_n + E_i$
- Tells the impact on Y by changing  $X_1$  by 1 unit keeping other independent variables same.
- Individual slope coefficients (e.g.  $b_1$ ) in multiple regressions known as partial regression/slope coefficients.

**2.1 Assumption of the Multiple Linear Regression Model**

**2.2 Predicting the Dependent Variable in a Multiple Regression Model**

**2.3 Testing Whether All Population Regression Coefficients Equals Zero**

**2.4 Adjusted  $R^2$**

- Relationship b/w Y and  $X_1, X_2, X_3, \dots, X_n$  is linear.
- Independent variables are not random and no exact linear relationship exists b/w 2 or more independent variables.
- Expected value of error terms is 0.
- Variance of error term is same for all observations.
- Error term is uncorrelated across observations.
- Error term is normally distributed.

- Obtain estimates of regression parameters.
  - *estimates* =  $b_0^{\wedge}, b_1^{\wedge}, b_2^{\wedge}, \dots, b_n^{\wedge}$
  - *regression parameters* =  $b_0, b_1, b_2, \dots, b_k$
- Determine assumed values of  $\hat{X}_{1i}, \hat{X}_{2i}, \dots, \hat{X}_{ki}$
- Compute predicted value of  $\hat{Y}$  using  $\hat{Y}_i = \hat{b}_0 + \hat{b}_1\hat{X}_{1i} + \hat{b}_2\hat{X}_{2i} + \dots + \hat{b}_k\hat{X}_{ki}$
- To predict dependent variable:
  - Be confident that assumptions of the regression are met.
  - Predictions regarding X must be within reliable range of data used to estimate the model.

- $H_0 \Rightarrow$  All slope coefficients are simultaneously = 0, none of the X variable helps explain Y.
- To test  $H_0$  F-test is used.
- T-test cannot be used.

$$F = \frac{MSR}{MSE} = \frac{RSS/k}{SSE/(n-(k+1))}$$

Where

$$RSS = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$$

$$SSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

Decision rule  $\Rightarrow$  reject  $H_0$  if  $F > F_c$  (for given  $\alpha$ ).

- It is a one-tailed test.
- $d_f$  numerator = k
- $d_f$  denominator =  $n-(k+1)$ .
- For k and n the test statistic representing  $H_0$ , all slope coefficients are equal to 0, is  $F_{k, n-(k+1)}$
- In F-distribution table  $f_0 F_{k, n-(k+1)}$  where K represents column and  $n-(k+1)$  represents row.
- Significance of F in ANOVA table represents 'p value'.
- $\uparrow$  F-statistic  $\downarrow$  chances of Type I error.

- $R^2 \uparrow$  with addition of independent variables (X) in regression
- *Adjusted  $R^2$*  ( $\bar{R}^2$ ) =  $1 - \left(\frac{n-1}{n-k-1}\right)(1 - R^2)$ .
- When  $k \geq 1 \Rightarrow R^2 > \bar{R}^2$
- $\bar{R}^2$  can be -ve but  $R^2$  is always +ve.
- If  $\bar{R}^2$  is used for comparing regression models.
  - Sample size must be the same
  - Dependent variable is defined in the same way.
- $\uparrow \bar{R}^2$  Does not necessarily indicate regression is well specified.

3. USING DUMMY VARIABLES IN REGRESSION

- Dummy variable  $\Rightarrow$  takes 1 if particular condition is true & 0 when it is false.
- Diligence is required in choosing no. of dummy variables.
- Usually  $n-1$  dummy variables are used where  $n =$  no. of categories.

4. VIOLATIONS OF REGRESSION ASSUMPTIONS

4.1 Heteroskedasticity

- Variance of errors differs across observations  $\Rightarrow$  heteroskedastic
- Variance of errors is similar across observations  $\Rightarrow$  homoskedastic
- Usually no systematic relationship exists b/w X & regression residuals.
- If systematic relationship is present  $\Rightarrow$  heteroskedasticity can exist.

4.2 Serial Correlation

- Regression errors correlated across observations.
- Usually arises in time-series regression.

4.3 Multicollinearity

- Occurs when two or more independent variables (X) are highly correlated with each other.
- Regression can be estimated but result becomes problematic.
- Serious practical concern due to commonly found approximate linear relation among financial variables.

4.4 Summarizing the Issues

On page 3

On page 4

4.1.1 The Consequence of Heteroskedasticity

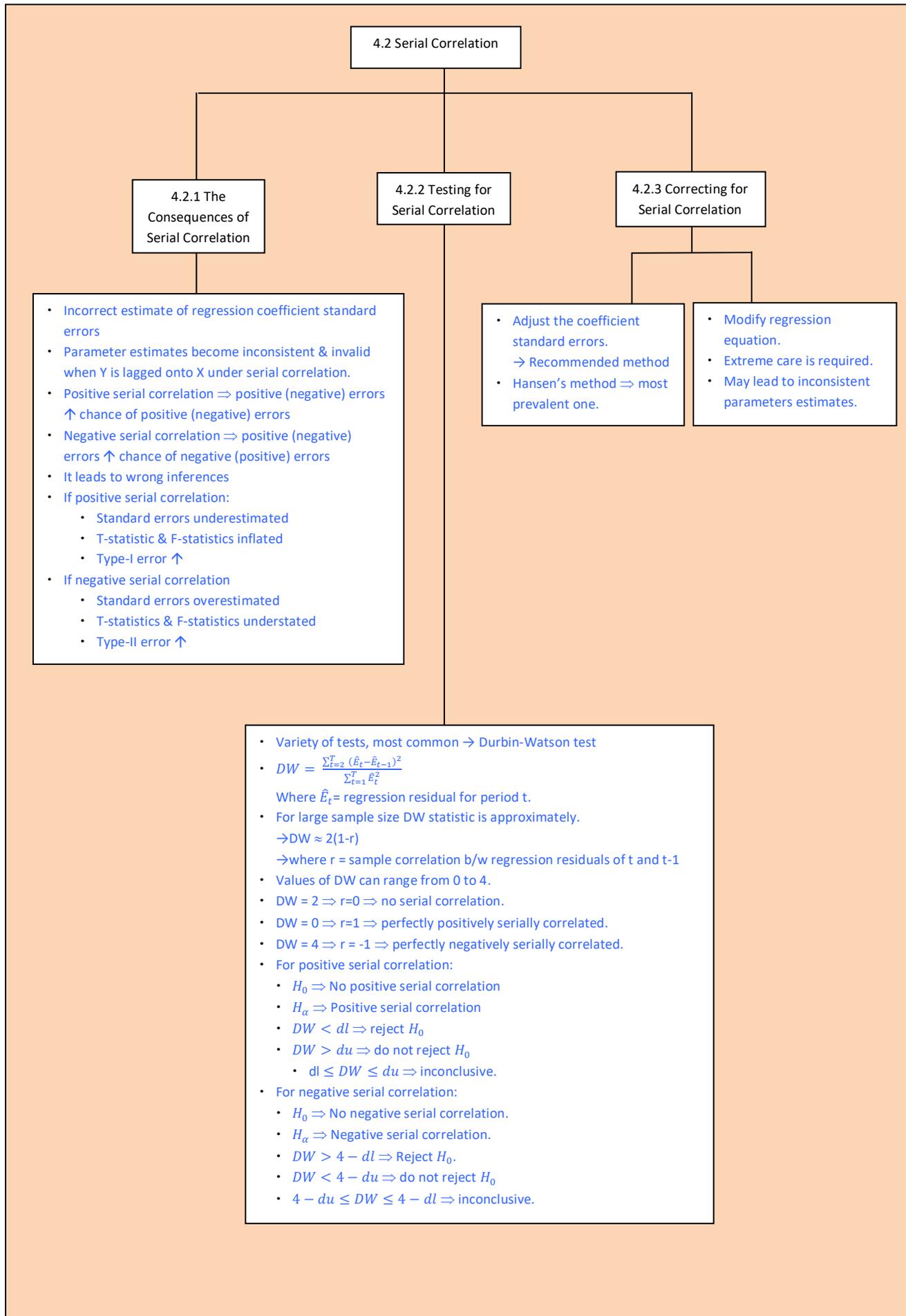
- It can lead to mistake in inference. Does not affect consistency.
- F-test becomes unreliable.
- Due to biased estimators of standard errors, t-test also becomes unreliable.
- Heteroskedasticity effects may include:
  - underestimation of estimated standard errors
  - inflated t-statistic
- Ignoring heteroskedasticity leads to significant relationship that does not exist actually.
- It becomes more serious while developing investment strategy using regression analysis.
- **Unconditional heteroskedasticity**  $\Rightarrow$  when heteroskedasticity of error variance is not correlated with independent variables in the multiple regression.
  - Create major problems for statistical inference.
- **Conditional heteroskedasticity**  $\Rightarrow$  when heteroskedasticity of error variance is correlated with the independent variables.
  - It causes most problems.
  - Can be tested & corrected easily through many statistically software packages.

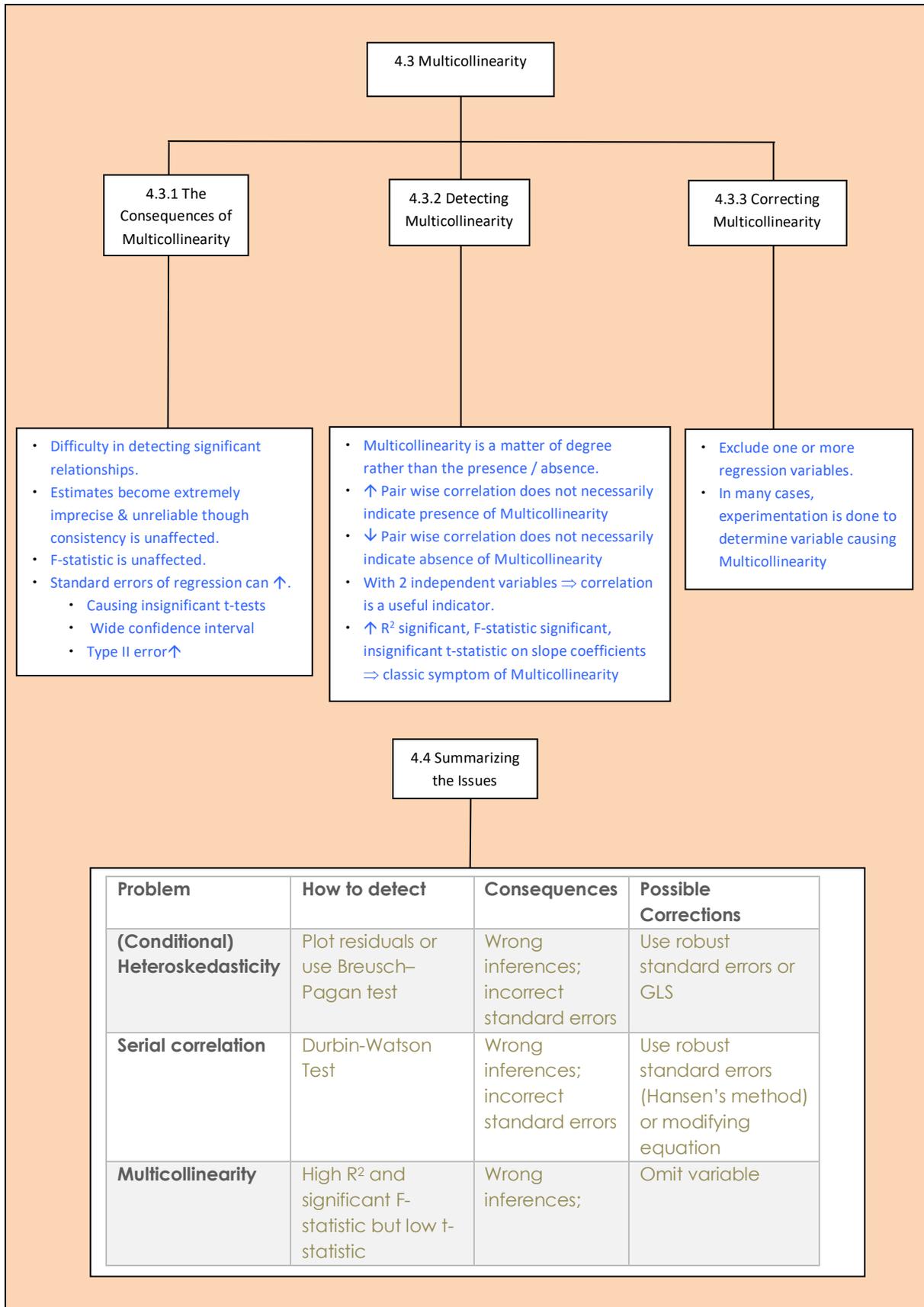
4.1.2 Testing for Heteroskedasticity

- **Breush-Pagan test** is widely used.
- Regression squared residuals of regression on independent variables.
  - Independent variables explain much of the variation of errors  $\Rightarrow$  conditional heteroskedasticity exists.
- $H_0 =$  no conditional heteroskedasticity exists.
- $H_a =$  conditional heteroskedasticity exist
- Under Breush-pagan test statistic =  $nR^2$
- $R^2$ : from regression of squared residuals on X
- Critical value  $\Rightarrow$  calculated  $\chi^2$  distribution.
- $d_f =$  no. of independent variables
- Reject  $H_0$  if test-static > critical value.

4.1.3 Correcting for Heteroskedasticity

- |  |  |
|--|--|
| <p><b>Robust Standard Errors</b></p> <ul style="list-style-type: none"> <li>• Corrects standard error of estimated coefficients.</li> <li>• Also known as heteroskedasticity consistent standards errors or white-corrected standards errors.</li> </ul> | <p><b>Generalized Least Squares</b></p> <ul style="list-style-type: none"> <li>• Modify original equation.</li> <li>• Requires economic expertise to implement correctly on financial data.</li> </ul> |
|--|--|





- Difficulty in detecting significant relationships.
- Estimates become extremely imprecise & unreliable though consistency is unaffected.
- F-statistic is unaffected.
- Standard errors of regression can ↑
  - Causing insignificant t-tests
  - Wide confidence interval
  - Type II error ↑

- Multicollinearity is a matter of degree rather than the presence / absence.
- ↑ Pair wise correlation does not necessarily indicate presence of Multicollinearity
- ↓ Pair wise correlation does not necessarily indicate absence of Multicollinearity
- With 2 independent variables ⇒ correlation is a useful indicator.
- ↑ R<sup>2</sup> significant, F-statistic significant, insignificant t-statistic on slope coefficients ⇒ classic symptom of Multicollinearity

- Exclude one or more regression variables.
- In many cases, experimentation is done to determine variable causing Multicollinearity

4.4 Summarizing the Issues

Problem	How to detect	Consequences	Possible Corrections
<b>(Conditional) Heteroskedasticity</b>	Plot residuals or use Breusch-Pagan test	Wrong inferences; incorrect standard errors	Use robust standard errors or GLS
<b>Serial correlation</b>	Durbin-Watson Test	Wrong inferences; incorrect standard errors	Use robust standard errors (Hansen's method) or modifying equation
<b>Multicollinearity</b>	High R <sup>2</sup> and significant F-statistic but low t-statistic	Wrong inferences;	Omit variable

