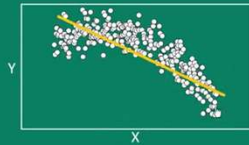


1. Linearity (Correct functional form)

Consider the following model:

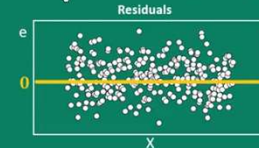
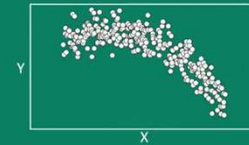
$$\text{Lung Function}_i = \beta_0 + \beta_1(\text{age})_i + \varepsilon_i$$



1. Linearity (Correct functional form)

Consider the following model:

$$\text{Lung Function}_i = \beta_0 + \beta_1(\text{age})_i + \beta_2(\text{age}^2)_i + \varepsilon_i$$



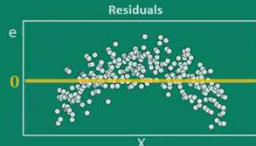
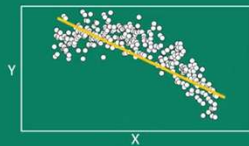
What's the issue?

- If functional form is incorrect, both the coefficients and standard errors in your output are unreliable

1. Linearity (Correct functional form)

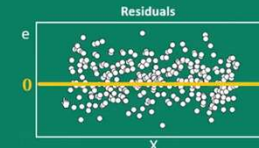
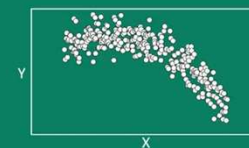
Consider the following model:

$$\text{Lung Function}_i = \beta_0 + \beta_1(\text{age})_i + \varepsilon_i$$



Consider the following model:

$$\text{Lung Function}_i = \beta_0 + \beta_1(\text{age})_i + \beta_2(\text{age}^2)_i + \varepsilon_i$$



What's the issue?

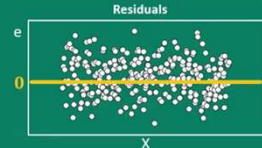
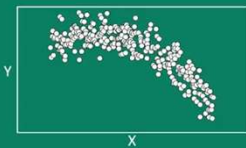
- If functional form is incorrect, both the coefficients and standard errors in your output are unreliable

Detection:

- Residual plots
- Likelihood ratio (LR) test

Consider the following model:

$$\text{Lung Function}_i = \beta_0 + \beta_1(\text{age})_i + \beta_2(\text{age}^2)_i + \varepsilon_i$$



What's the issue?

- If functional form is incorrect, both the coefficients and standard errors in your output are unreliable

Detection:

- Residual plots
- Likelihood ratio (LR) test

Remedies:

- Get the specification correct (trial and error)

Constant Variance
(no heteroskedasticity)

1. Linearity

2. Constant
Error Variance

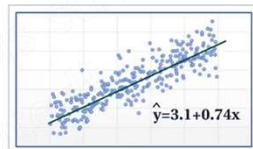
3. Independent
Error
Terms

4. Normal
errors

5. No multi-
collinearity

6. Exogeneity

$$Y = \beta_0 + \beta_1 X + \varepsilon$$



	Coef	Std Err	t Stat	P-value
Intercept	3.0605	1.4024	2.1823	0.02909
X	0.7393	0.2308	3.20339	0.00136

Regression Assumptions

Constant Variance
(no heteroskedasticity)

Consider the following model:

$$\text{Expenditure}_i = \beta_0 + \beta_1(\text{Income})_i + \varepsilon_i$$

Constant Variance (no heteroskedasticity)

Consider the following model:

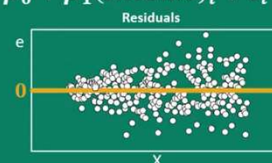
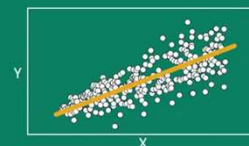
$$\text{Expenditure}_i = \beta_0 + \beta_1(\text{Income})_i + \varepsilon_i$$



Constant Variance (no heteroskedasticity)

Consider the following model:

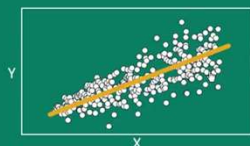
$$\text{Expenditure}_i = \beta_0 + \beta_1(\text{Income})_i + \varepsilon_i$$



Constant Variance (no heteroskedasticity)

Consider the following model:

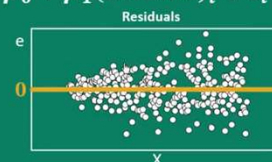
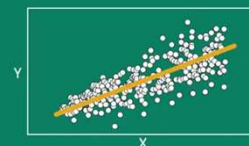
$$\text{Expenditure}_i = \beta_0 + \beta_1(\text{Income})_i + \varepsilon_i$$



Constant Variance (no heteroskedasticity)

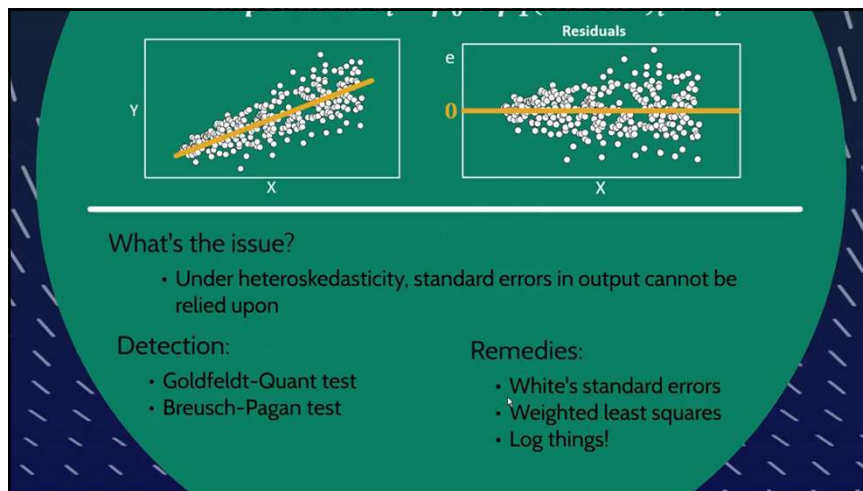
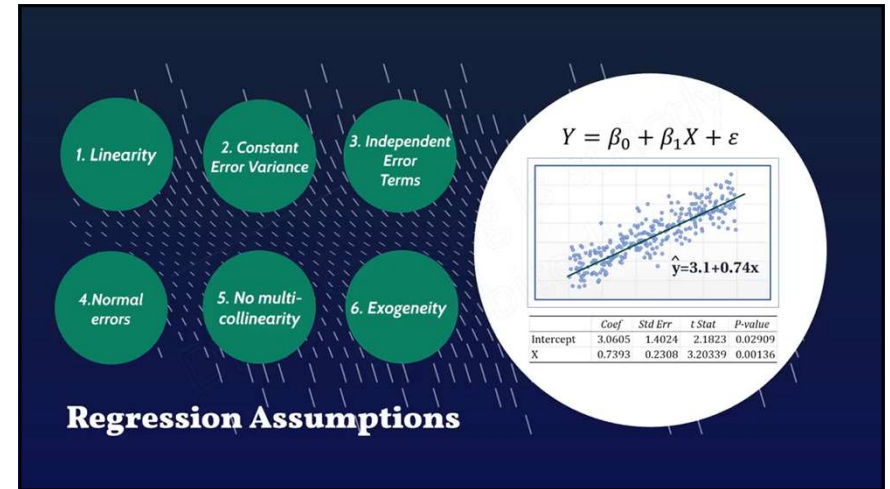
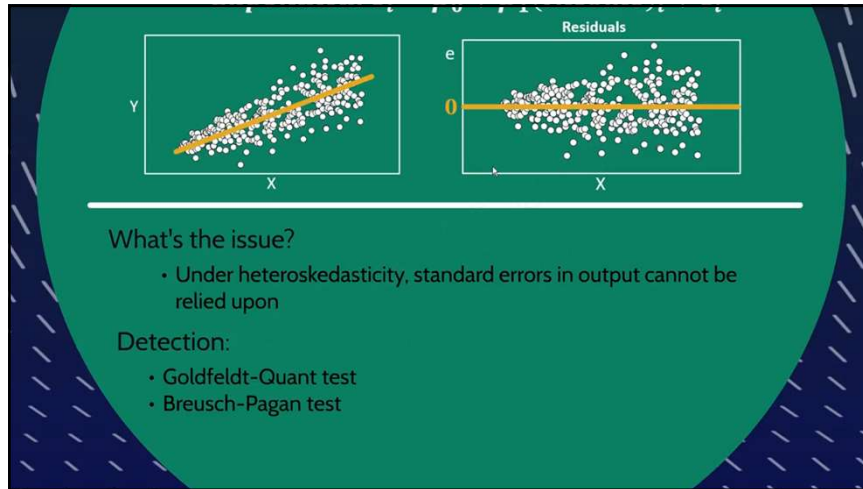
Consider the following model:

$$\text{Expenditure}_i = \beta_0 + \beta_1(\text{Income})_i + \varepsilon_i$$



What's the issue?

- Under heteroskedasticity, standard errors in output cannot be relied upon



Independent error terms (no autocorrelation)

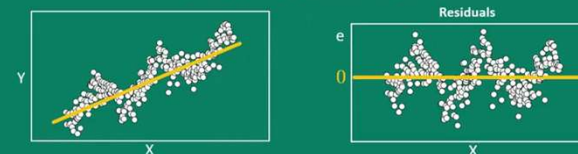
Consider the following model:

$$\text{Stock Index}_i = \beta_0 + \beta_1(\text{Time})_i + \varepsilon_i$$

Independent error terms (no autocorrelation)

Consider the following model:

$$\text{Stock Index}_i = \beta_0 + \beta_1(\text{Time})_i + \varepsilon_i$$



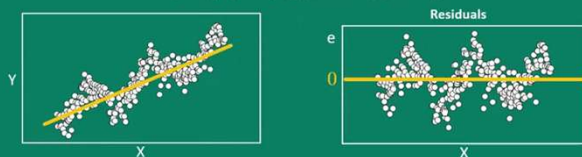
What's the issue?

- Under autocorrelation, standard errors in output cannot be relied upon

Independent error terms (no autocorrelation)

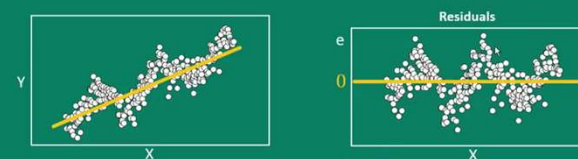
Consider the following model:

$$\text{Stock Index}_i = \beta_0 + \beta_1(\text{Time})_i + \varepsilon_i$$



Consider the following model:

$$\text{Stock Index}_i = \beta_0 + \beta_1(\text{Time})_i + \varepsilon_i$$

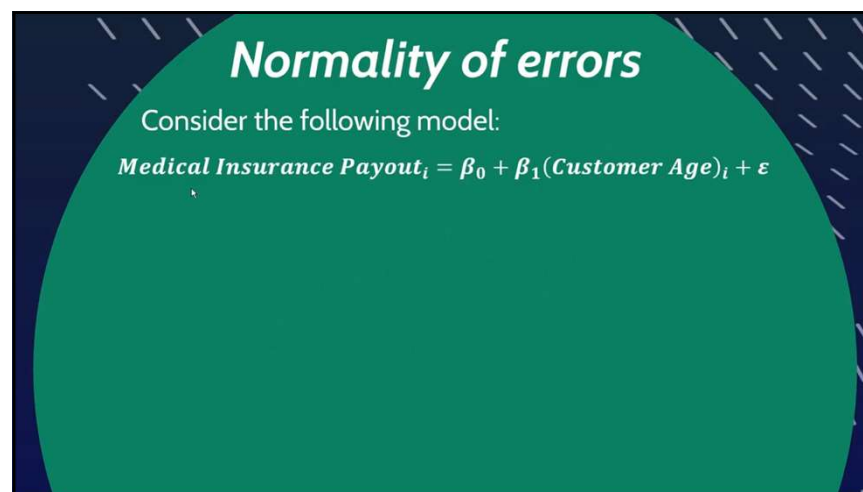
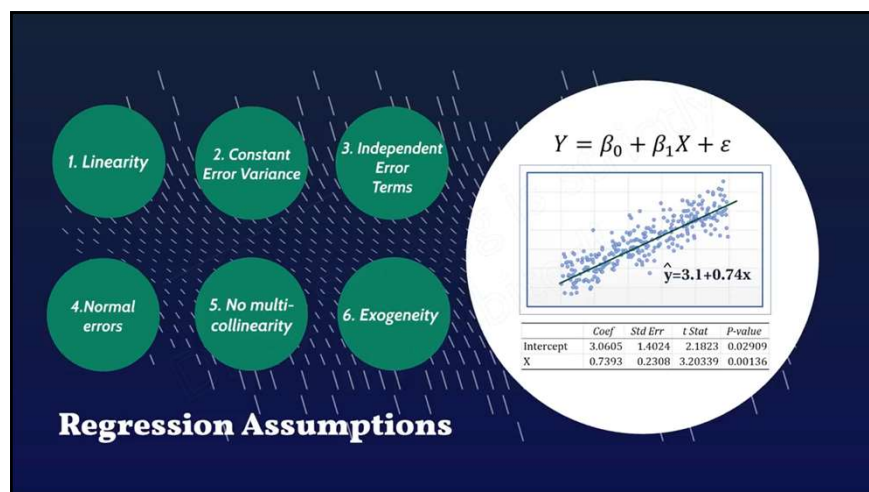
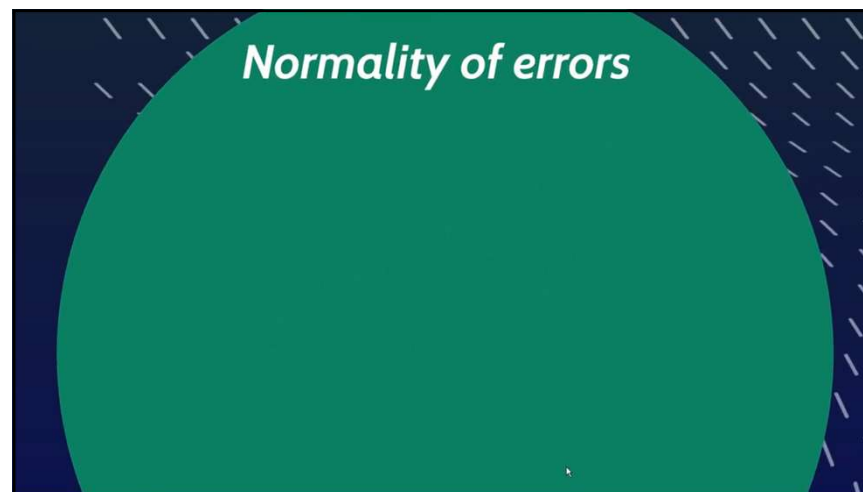
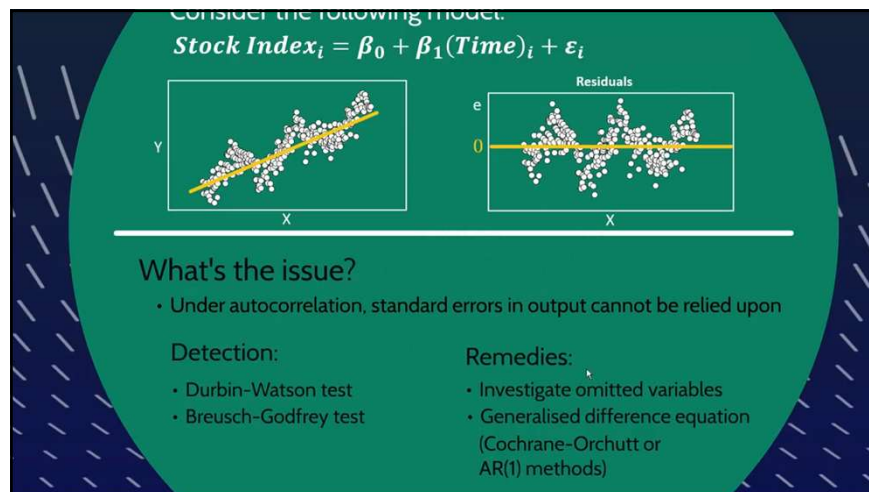


What's the issue?

- Under autocorrelation, standard errors in output cannot be relied upon

Detection:

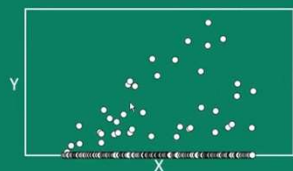
- Durbin-Watson test
- Breusch-Godfrey test



Normality of errors

Consider the following model:

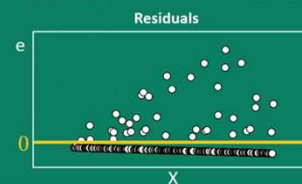
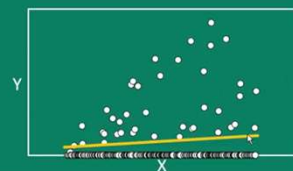
$$\text{Medical Insurance Payout}_i = \beta_0 + \beta_1(\text{Customer Age})_i + \varepsilon$$



Normality of errors

Consider the following model:

$$\text{Medical Insurance Payout}_i = \beta_0 + \beta_1(\text{Customer Age})_i + \varepsilon$$



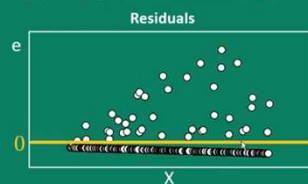
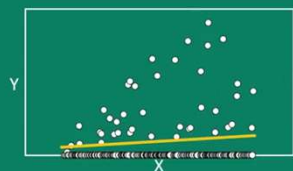
What's the issue?

- If normality is violated and n is small, standard errors in output are affected

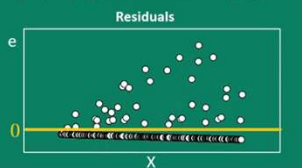
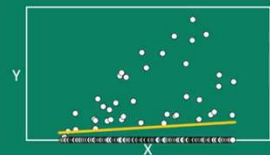
Normality of errors

Consider the following model:

$$\text{Medical Insurance Payout}_i = \beta_0 + \beta_1(\text{Customer Age})_i + \varepsilon$$



$$\text{Medical Insurance Payout}_i = \beta_0 + \beta_1(\text{Customer Age})_i + \varepsilon$$

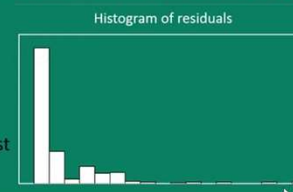


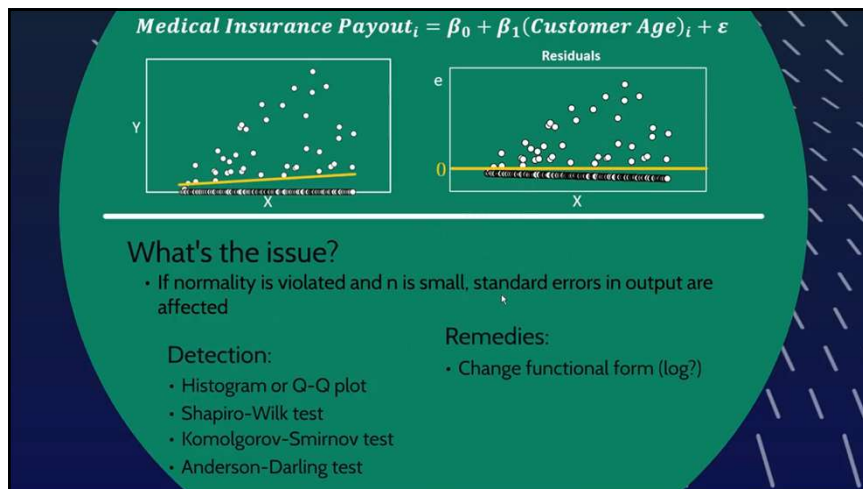
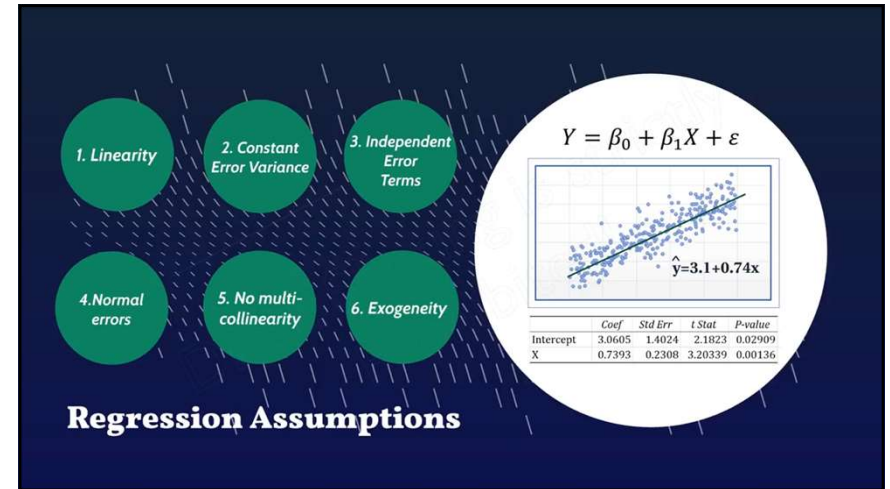
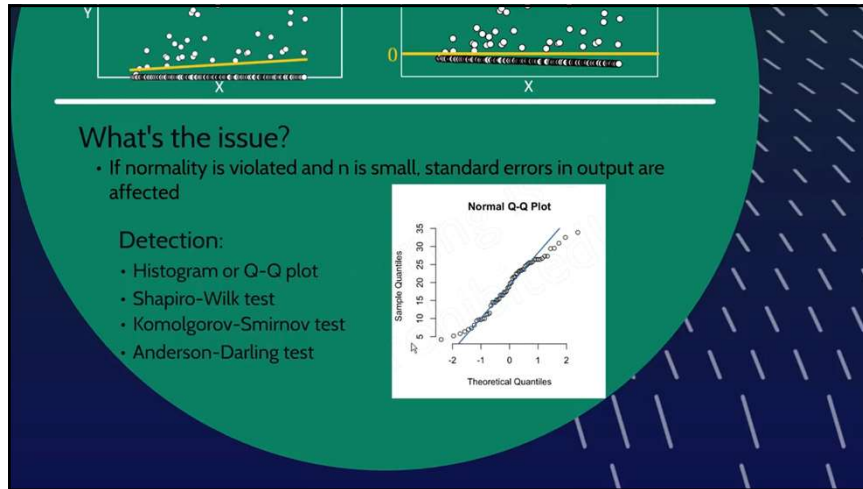
What's the issue?

- If normality is violated and n is small, standard errors in output are affected

Detection:

- Histogram or Q-Q plot
- Shapiro-Wilk test
- Kolmogorov-Smirnov test
- Anderson-Darling test





No multicollinearity

Consider the following model:

$$\begin{aligned} \text{Motor Accidents}_i &= \beta_0 + \beta_1(\text{Num cars})_i \\ &\quad + \beta_2(\text{Num residents})_i + \varepsilon \\ i &= \text{suburb } 1, 2, 3 \dots \end{aligned}$$

No multicollinearity

Consider the following model:

$$\begin{aligned} \text{Motor Accidents}_i &= \beta_0 + \beta_1(\text{Num cars})_i \\ &\quad + \beta_2(\text{Num residents})_i + \varepsilon \\ i &= \text{suburb } 1, 2, 3 \dots \end{aligned}$$

Multi-collinearity occurs where the X variables are themselves related

What's the issue?

- Coefficients and standard errors of affected variables are unreliable.

No multicollinearity

Consider the following model:

$$\begin{aligned} \text{Motor Accidents}_i &= \beta_0 + \beta_1(\text{Num cars})_i \\ &\quad + \beta_2(\text{Num residents})_i + \varepsilon \\ i &= \text{suburb } 1, 2, 3 \dots \end{aligned}$$

Multi-collinearity occurs where the X variables are themselves related

No multicollinearity

Consider the following model:

$$\begin{aligned} \text{Motor Accidents}_i &= \beta_0 + \beta_1(\text{Num cars})_i \\ &\quad + \beta_2(\text{Num residents})_i + \varepsilon \\ i &= \text{suburb } 1, 2, 3 \dots \end{aligned}$$

Multi-collinearity occurs where the X variables are themselves related

What's the issue?

- Coefficients and standard errors of affected variables are unreliable.

Detection:

- Look at correlation (ρ) between X variables
- Look at Variance Inflation Factors (VIF)

Multi-collinearity occurs where the X variables are themselves related

What's the issue?

- Coefficients and standard errors of affected variables are unreliable.

Detection:

- Look at correlation (ρ) between X variables
- Look at Variance Inflation Factors (VIF)

Remedies:

- Remove one of the variables

Exogeneity
(no omitted variable bias)

1. Linearity

2. Constant Error Variance

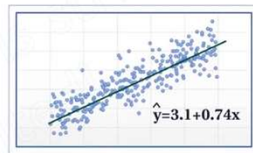
3. Independent Error Terms

4. Normal errors

5. No multi-collinearity

6. Exogeneity

$$Y = \beta_0 + \beta_1 X + \varepsilon$$



	Coef	Std Err	t Stat	P-value
Intercept	3.0605	1.4024	2.1823	0.02909
X	0.7393	0.2308	3.20339	0.00136

Regression Assumptions

Exogeneity
(no omitted variable bias)

Consider the following model:

$$\text{Salary}_i = \beta_0 + \beta_1(\text{Years of education})_i + \varepsilon_i$$

Exogeneity (no omitted variable bias)

Consider the following model:

$$\text{Salary}_i = \beta_0 + \beta_1(\text{Years of education})_i + \varepsilon_i$$

Socio-economic status affects **both** X and Y variables, thus would cause **omitted variable bias**.

Exogeneity (no omitted variable bias)

Consider the following model:

$$\text{Salary}_i = \beta_0 + \beta_1(\text{Years of education})_i + \varepsilon_i$$

Socio-economic status affects **both** X and Y variables, thus would cause **omitted variable bias**.

TECHNICALLY - Socio-economic status would affect ε_i in the model, thus, Education is no longer wholly exogenous as it can be explained in part by the error term.

What's the issue?

- Model can only be used for predictive purposes (can not infer causation)

Exogeneity (no omitted variable bias)

Consider the following model:

$$\text{Salary}_i = \beta_0 + \beta_1(\text{Years of education})_i + \varepsilon_i$$

Socio-economic status affects **both** X and Y variables, thus would cause **omitted variable bias**.

TECHNICALLY - Socio-economic status would affect ε_i in the model, thus, Education is no longer wholly exogenous as it can be explained in part by the error term.

Consider the following model:

$$\text{Salary}_i = \beta_0 + \beta_1(\text{Years of education})_i + \varepsilon_i$$

Socio-economic status affects **both** X and Y variables, thus would cause **omitted variable bias**.

TECHNICALLY - Socio-economic status would affect ε_i in the model, thus, Education is no longer wholly exogenous as it can be explained in part by the error term.

What's the issue?

- Model can only be used for predictive purposes (can not infer causation)

Detection:

- Intuition
- Checking correlations

Consider the following model:

$$\text{Salary}_i = \beta_0 + \beta_1(\text{Years of education})_i + \varepsilon_i$$

Socio-economic status affects both X and Y variables, thus would cause **omitted variable bias**.

TECHNICALLY - Socio-economic status would affect ε_i in the model, thus, Education is no longer wholly exogenous as it can be explained in part by the error term.

What's the issue?

- Model can only be used for predictive purposes (can not infer causation)

Detection:

- Intuition
- Checking correlations

Remedy:

- Using instrumental variables

THANK YOU!

1. Linearity

2. Constant Error Variance

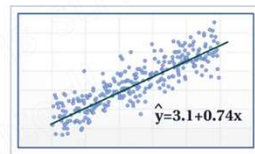
3. Independent Error Terms

4. Normal errors

5. No multi-collinearity

6. Exogeneity

$$Y = \beta_0 + \beta_1 X + \varepsilon$$



	Coef	Std Err	t Stat	P-value
Intercept	3.0605	1.4024	2.1823	0.02909
X	0.7393	0.2308	3.20339	0.00136

Regression Assumptions