

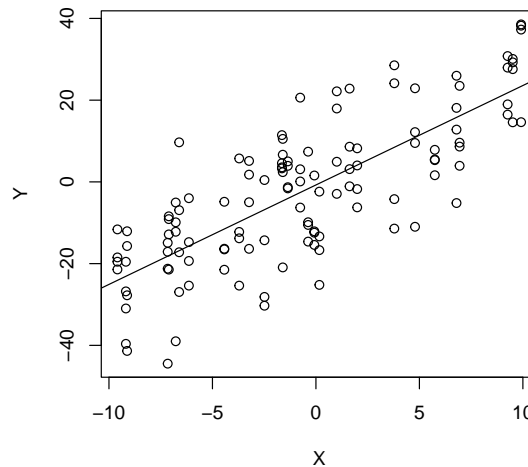
## Lesson 8. Conditions for a Simple Linear Regression Model – Part 1

### 1 Conditions for a simple linear regression model

- Recall that the **simple linear regression model** is

$$Y = \beta_0 + \beta_1 X + \varepsilon \quad \text{where } \varepsilon \sim \text{iid } N(0, \sigma_\varepsilon^2)$$

- The **errors**  $\varepsilon$  follow an identical normal distribution and are independent from one another



- When is a simple linear regression model reasonable?
  - Are we justified in using our model? How much can we trust predictions that come from the model?
- We check for the following conditions:

Condition	Explanation
<b>Linearity</b>	<ul style="list-style-type: none"> <li>• The overall relationship between the variables has a linear pattern</li> <li>• The average values of the response <math>Y</math> for each value of <math>X</math> fall on a common straight line</li> </ul>
<b>Independence</b>	<ul style="list-style-type: none"> <li>• The errors are independent from one another</li> <li>• The distance of one point from the line has no influence on the distance of another point</li> </ul>
<b>Normality</b>	<ul style="list-style-type: none"> <li>• The errors follow a normal distribution</li> </ul>
<b>Equal variance</b>	<ul style="list-style-type: none"> <li>• The variability in the errors is the same for all values of <math>X</math></li> <li>• In other words, the spread of the points around the line remains fairly constant</li> </ul>
<b>Randomness</b>	<ul style="list-style-type: none"> <li>• The data are obtained using a random process</li> </ul>

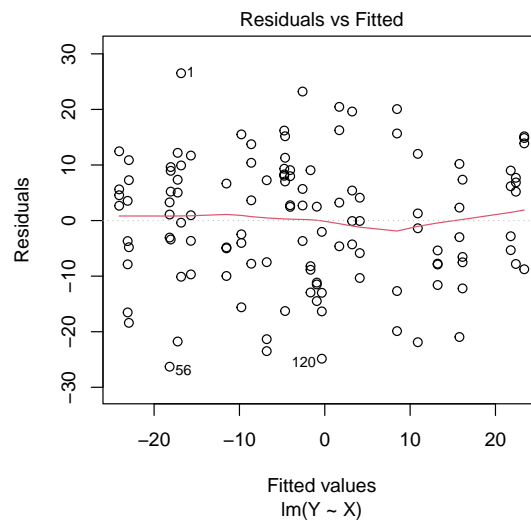
- Mnemonic: **LINER**
- Normality and randomness must be satisfied when we want to use the model for statistical inference (e.g., confidence intervals, hypothesis tests)

## 2 Assessing conditions for a simple linear regression model

- To easily assess the above conditions, we will use two diagnostic plots

### 2.1 Residuals vs. fitted values plot

- The **residuals vs. fitted values plot** reorients the axes so that the regression line is represented as a horizontal line through zero
  - $\left\{ \begin{array}{l} \text{Positive} \\ \text{Negative} \end{array} \right\}$  residuals are represented by points  $\left\{ \begin{array}{l} \text{above} \\ \text{below} \end{array} \right\}$  the regression line
  - This plot lets us focus on any clear patterns in the estimated errors (i.e., residuals)

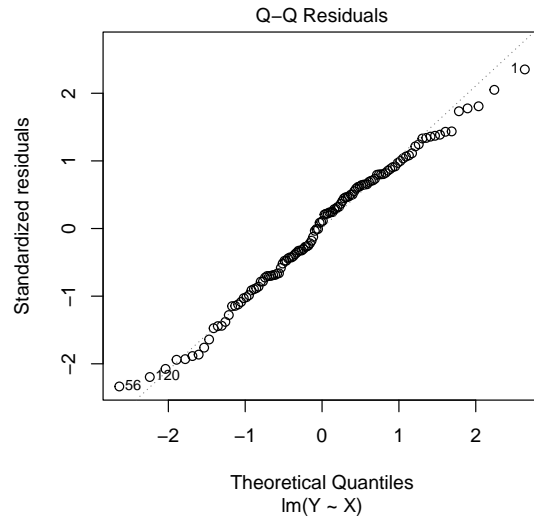


- The linearity condition is satisfied if

- The equal variance condition is satisfied if

## 2.2 Normal Q-Q plot of residuals

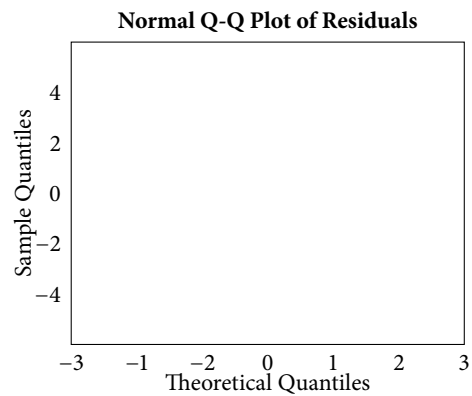
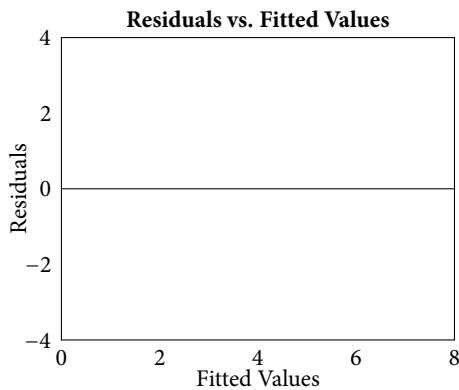
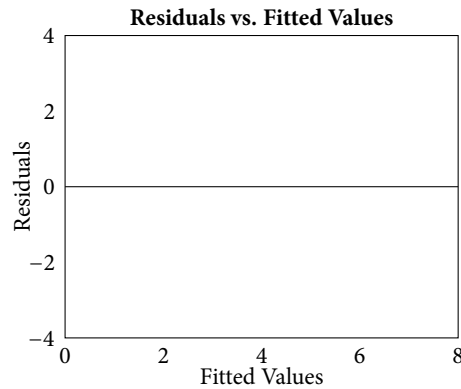
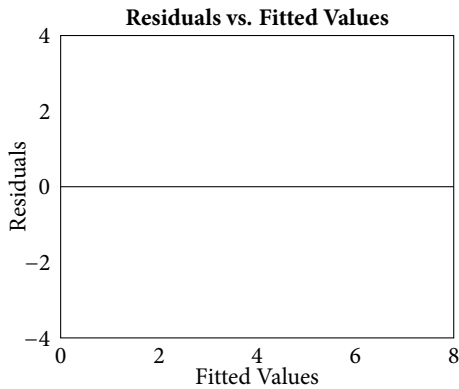
- An ideal Normal Q-Q plot of the residuals looks like



- The larger the sample size, the more lenient we can be about normality

**Example 1.** On the blank plots below, sketch the following:

- A residuals vs. fitted values plot where linearity is met, but equal variance is violated.
- A residuals vs. fitted values plot where equal variance is met, but linearity is violated.
- A residuals vs. fitted values plot where both equal variance and linearity are violated.
- A Normal Q-Q plot that shows dramatic violation of the normality condition.



### 2.3 Putting it all together...

Condition	Where to check	What we want
<b>Linearity</b>	Residuals vs. fitted values plot	Points randomly and evenly distributed above and below residual = 0 line, moving from left to right
<b>Independence</b>	Description of data collection	No indication that the errors influence each other
<b>Normality</b>	Normal Q-Q plot of residuals	Points in approximately straight line
<b>Equal variance</b>	Residuals vs. fitted values plot	Points span constant vertical width, moving from left to right
<b>Randomness</b>	Description of data collection	Data obtained using a random process, such as random sampling from a population or randomization in an experiment

- No model is perfect, and linear regression is fairly robust to slight violations.
- We will only be concerned with blatant violations.