

# Lesson 7: SLR: Checking model assumptions

Nicky Wakim

2025-01-27

# Learning Objectives

1. Describe the model assumptions made in linear regression using ordinary least squares
2. Determine if the relationship between our sampled  $X$  and  $Y$  is linear
3. Use QQ plots to determine if our fitted model holds the normality assumption
4. Use residual plots to determine if our fitted model holds the equality of variance assumption

# Let's remind ourselves of one model we have been working with

- We have been looking at the association between life expectancy and female literacy rate
- We used OLS to find the coefficient estimates of our best-fit line

Population model:

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

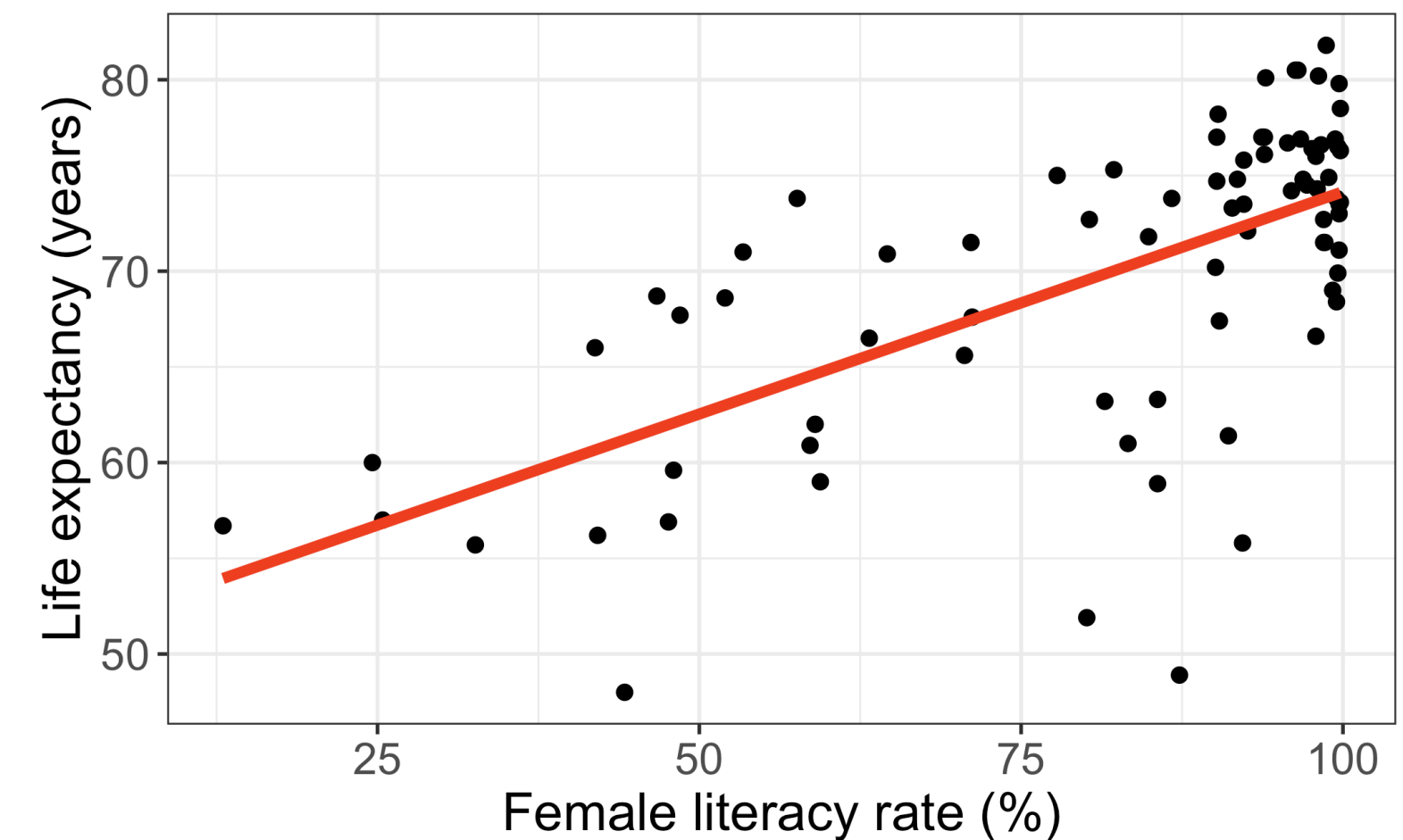
Estimated model:

term	estimate	std.error	statistic	p.value
(Intercept)	50.93	2.66	19.14	0.00
FemaleLiteracyRate	0.23	0.03	7.38	0.00

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 \cdot X$$

$$\widehat{\text{LE}} = 50.9 + 0.232 \cdot \text{FLR}$$

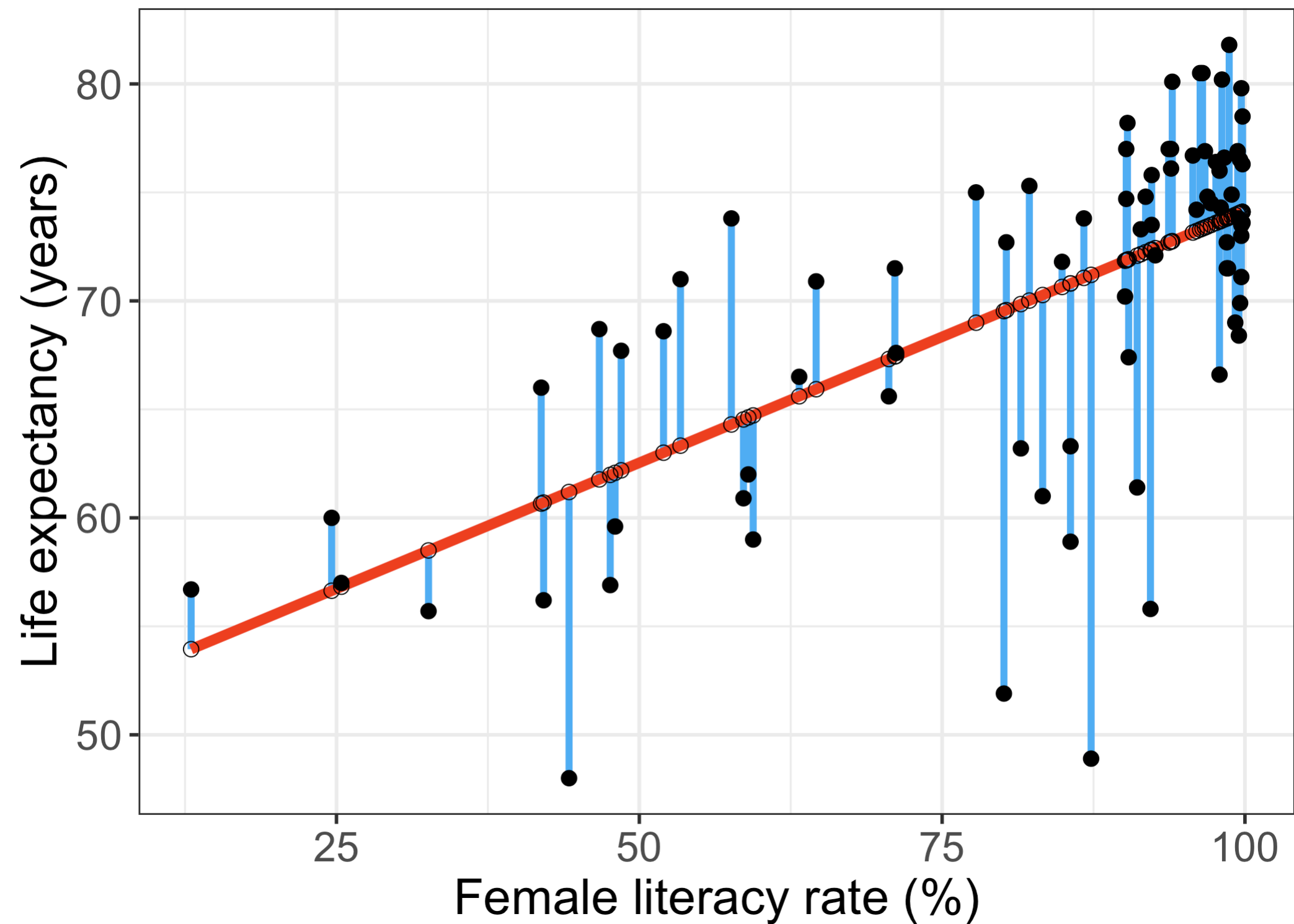
Relationship between life expectancy and the female literacy rate in 2011



# Our residuals will help us a lot in our diagnostics and assumptions!

- The **residuals**  $\hat{\epsilon}_i$  are the vertical distances between
  - the observed data  $(X_i, Y_i)$
  - the fitted values (regression line)  
 $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$

$$\hat{\epsilon}_i = Y_i - \hat{Y}_i, \text{ for } i = 1, 2, \dots, n$$



# Learning Objectives

1. Describe the model assumptions made in linear regression using ordinary least squares
2. Determine if the relationship between our sampled  $X$  and  $Y$  is linear
3. Use QQ plots to determine if our fitted model holds the normality assumption
4. Use residual plots to determine if our fitted model holds the equality of variance assumption

# Least-squares model assumptions: LINE

These are the model assumptions made in ordinary least squares:

[L] Linearity of relationship between variables

[I] Independence of the  $Y$  values

[N] Normality of the  $Y$ 's given  $X$  (or residuals)

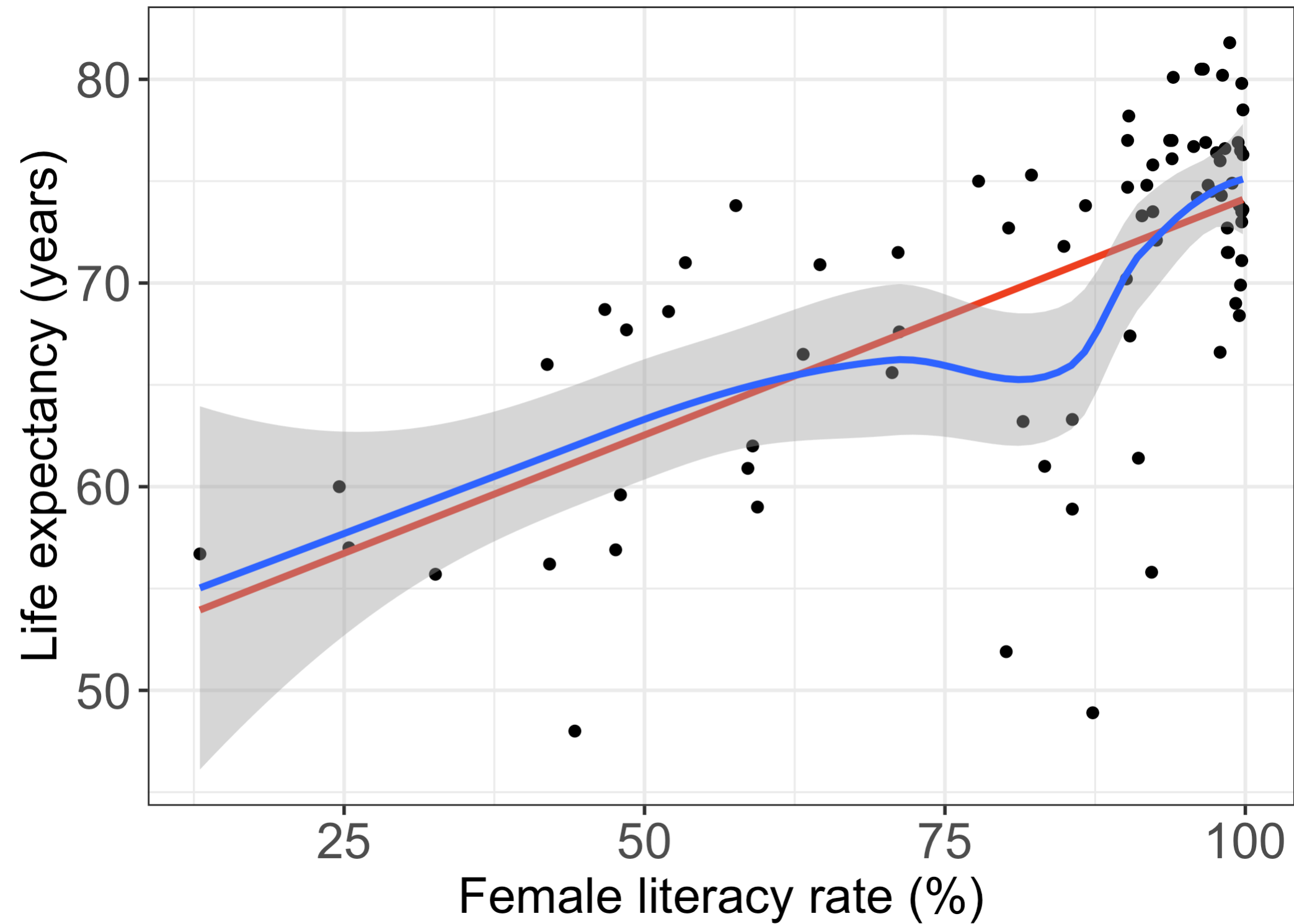
[E] Equality of variance of the residuals (homoscedasticity)

**Note:** These assumptions are baked into the *population model*. We look at the *population parameters* when we discuss these assumptions, but we use the *estimated model* with our data to check *if the assumptions are held up*.

# L: Linearity

- The relationship between the variables is linear (a straight line):
  - The mean value of  $Y$  given  $X$ ,  $\mu_{y|x}$  or  $E[Y|X]$ , is a straight-line function of  $X$

$$\mu_{y|x} = \beta_0 + \beta_1 \cdot X$$



# I: Independence of observations

- The  $Y$ -values are statistically independent of one another
- Examples of when they are *not* independent, include
  - repeated measures (such as baseline, 3 months, 6 months)
  - data from clusters, such as different hospitals or families
- This condition is checked by reviewing the study *design* and not by inspecting the data
- How to analyze data using regression models when the  $Y$ -values are not independent is covered in BSTA 519 (Longitudinal data)

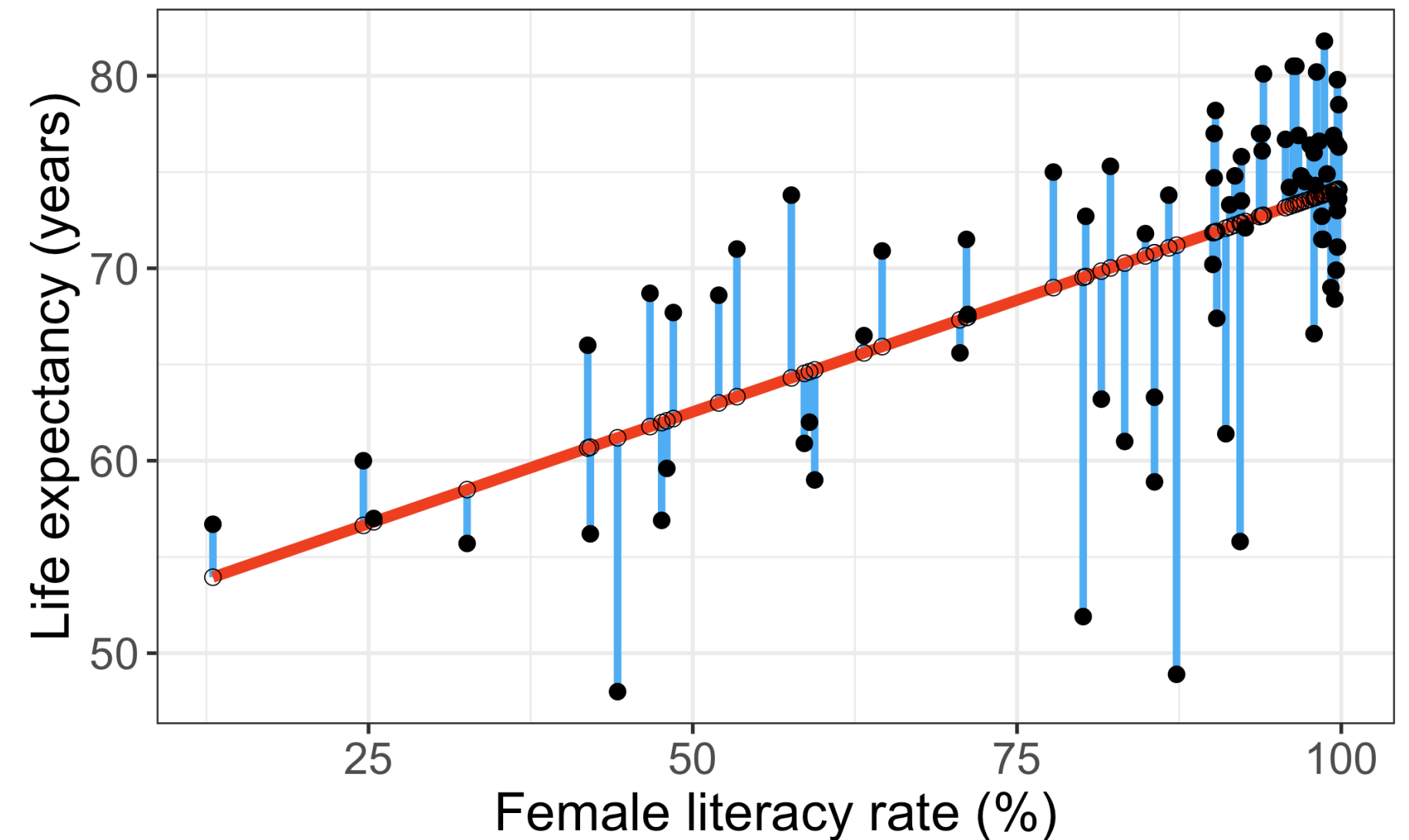


# Poll Everywhere Question 1

# N: Normality

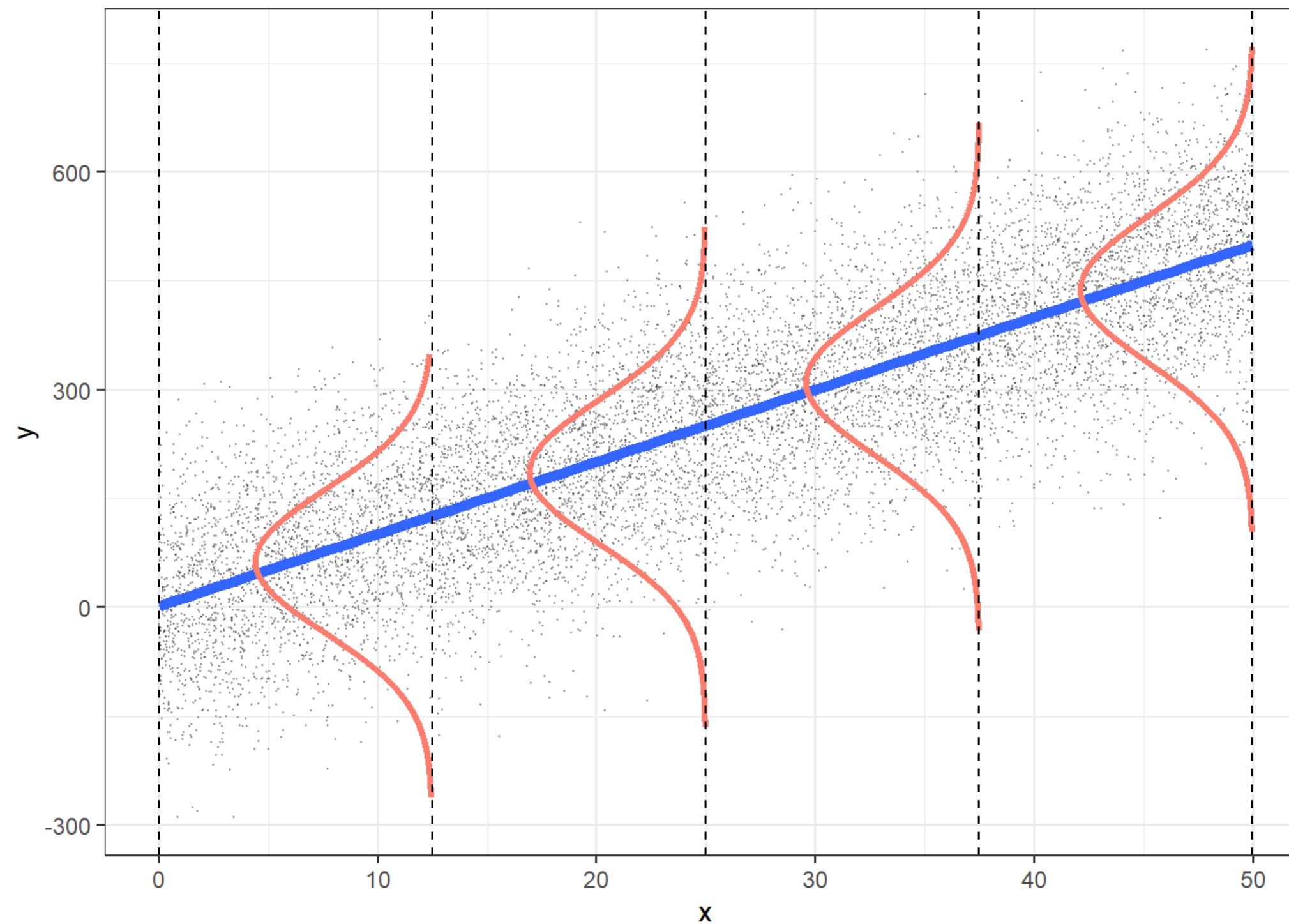
- For any fixed value of  $X$ ,  $Y$  has normal distribution.
  - Note: This is not about  $Y$  alone, but  $Y|X$
- Equivalently, the measurement (random) errors  $\epsilon_i$  's normally distributed
  - This is more often what we check

Relationship between life expectancy and the female literacy rate in 2011



# E: Equality of variance of the residuals

- The variance of  $Y$  given  $X$  ( $\sigma^2_{Y|X}$ ), is the same for any  $X$ 
  - We use just  $\sigma^2$  to denote the common variance
- This is also called **homoscedasticity**



# Summary of LINE model assumptions

- $Y$  values are independent (check study design!)

The distribution of  $Y$  given  $X$  is

- normal
- with mean  $\mu_{y|x} = \beta_0 + \beta_1 \cdot X$
- and common variance  $\sigma^2$

This means that the residuals are

- normal
- with mean = 0
- and common variance  $\sigma^2$

In mathematical form:

$$Y|X \stackrel{\text{i.i.d.}}{\sim} N(\beta_0 + \beta_1 X, \sigma^2)$$

$$\epsilon \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma^2)$$

# How do we determine if our model follows the LINE assumptions?

## [L] Linearity of relationship between variables

Check if there is a linear relationship between the mean response ( $Y$ ) and the explanatory variable ( $X$ )

## [I] Independence of the $Y$ values

Check that the observations are independent

## [N] Normality of the $Y$ 's given $X$ (residuals)

Check that the responses (at each level  $X$ ) are normally distributed

- Usually measured through the residuals

## [E] Equality of variance of the residuals (homoscedasticity)

Check that the variance (or standard deviation) of the responses is equal for all levels of  $X$

- Usually measured through the residuals

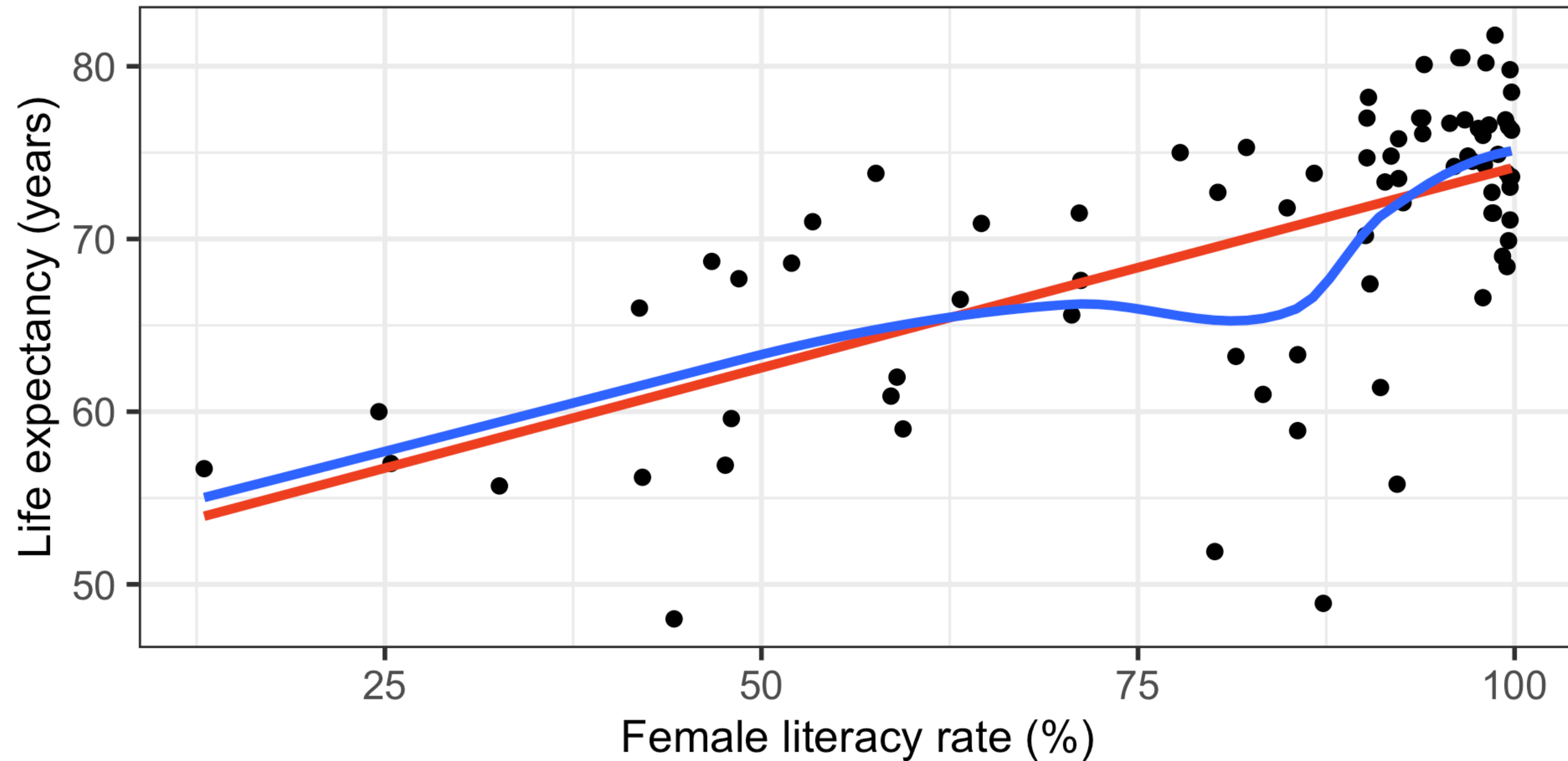
# Learning Objectives

1. Describe the model assumptions made in linear regression using ordinary least squares
2. Determine if the relationship between our sampled  $X$  and  $Y$  is linear
3. Use QQ plots to determine if our fitted model holds the normality assumption
4. Use residual plots to determine if our fitted model holds the equality of variance assumption

# L: Linearity of relationship between variables

Is the association between the variables linear?

- Diagnostic tool: Scatterplot of  $X$  vs.  $Y$



# Poll Everywhere Question 2



# I: Independence of the residuals ( $Y$ values)

- Are the data points independent of each other?
- **Diagnostic tool:** reviewing the study *design* and not by inspecting the data

# Learning Objectives

1. Describe the model assumptions made in linear regression using ordinary least squares
2. Determine if the relationship between our sampled X and Y is linear
3. Use QQ plots to determine if our fitted model holds the normality assumption
4. Use residual plots to determine if our fitted model holds the equality of variance assumption

# N: Normality of the residuals

- We need to check if the errors/residuals ( $\epsilon_i$ 's) are normally distributed
- **Diagnostic tools:**
  - Distribution plots of residuals
  - QQ plots of residuals
- Extra resource on how QQ plots are made

# N: Extract model's residuals in R

- First extract the residuals' values from the model output using the `augment()` function from the `broom` package.
- Get a tibble with the original data, as well as the residuals and some other important values.

```
1 model1 <- lm(LifeExpectancyYrs ~ FemaleLiteracyRate,  
2             data = gapm)  
3 aug1 <- augment(model1)  
4  
5 glimpse(aug1)
```

Rows: 80

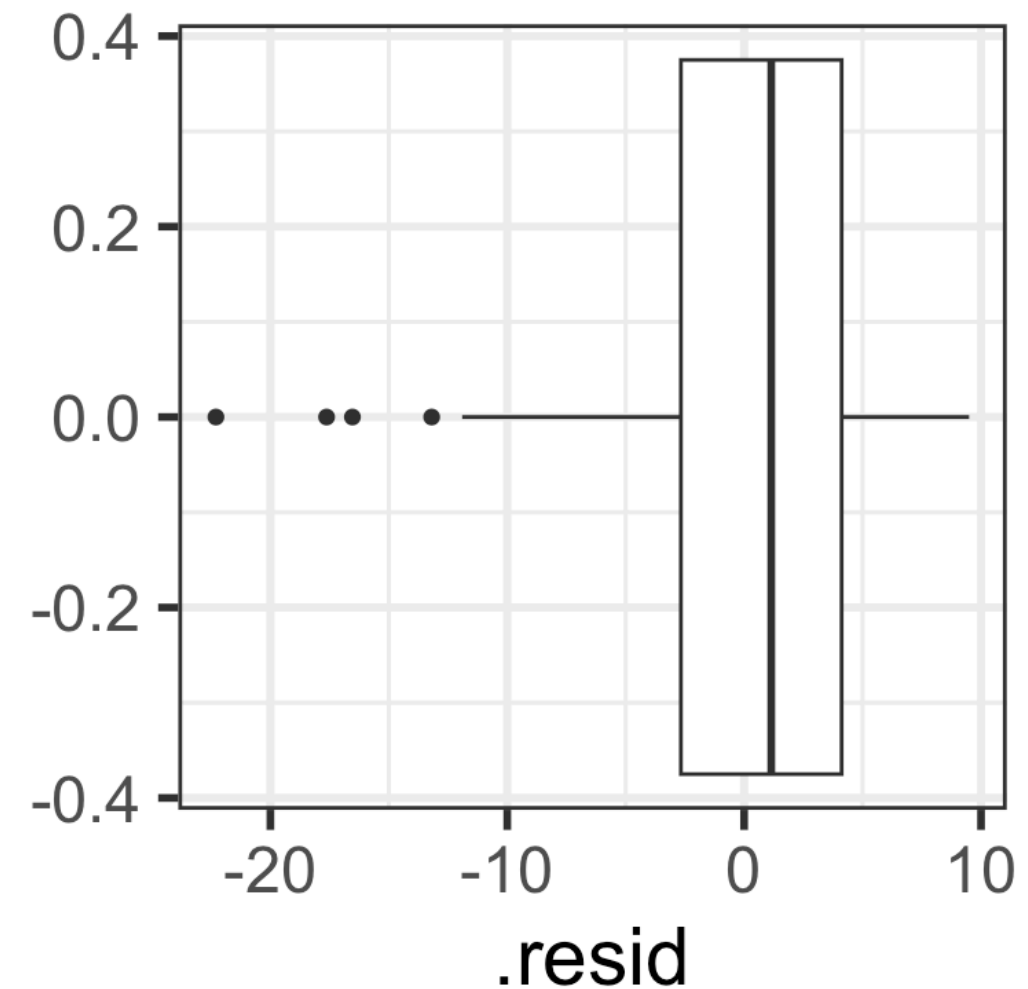
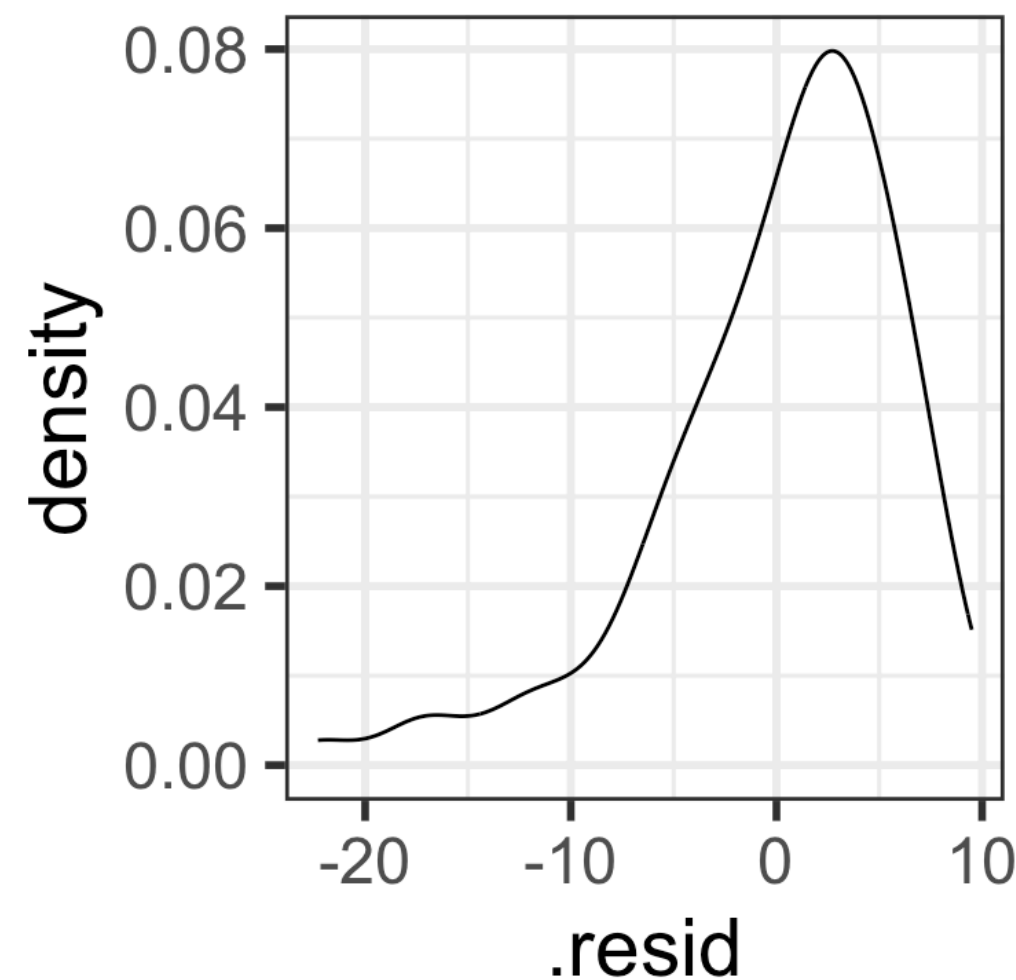
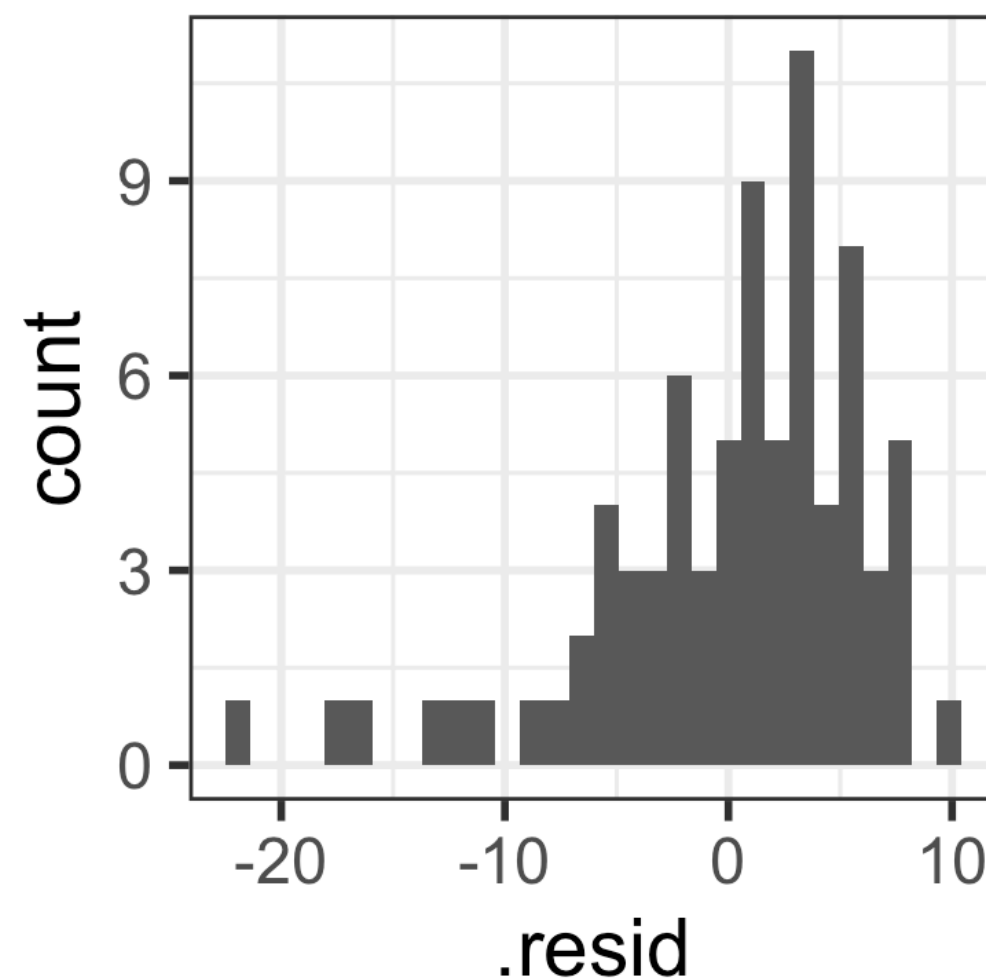
Columns: 8

```
$ LifeExpectancyYrs <dbl> 56.7, 76.7, 60.9, 76.9, 76.0, 73.8, 71.0, 76.9, 58...  
$ FemaleLiteracyRate <dbl> 13.0, 95.7, 58.6, 99.4, 97.9, 99.5, 53.4, 96.7, 85...  
$ .fitted           <dbl> 53.94643, 73.14897, 64.53453, 74.00809, 73.65980, 7...  
$ .resid            <dbl> 2.7535654, 3.5510294, -3.6345319, 2.8919074, 2.3402...  
$ .hat              <dbl> 0.13628996, 0.01768176, 0.02645854, 0.02077123, 0.0...  
$ .sigma            <dbl> 6.172684, 6.168414, 6.167643, 6.172935, 6.176043, 6...  
$ .cooks            <dbl> 1.835891e-02, 3.062372e-03, 4.887448e-03, 2.400993e...  
$ .std.resid        <dbl> 0.48238134, 0.58332052, -0.59972251, 0.47579667, 0...
```

# N: Check normality with distribution plots of residuals (1/2)

Note that below I save each figure as an object, and then combine them together in one row of output using `grid.arrange()` from the `gridExtra` package

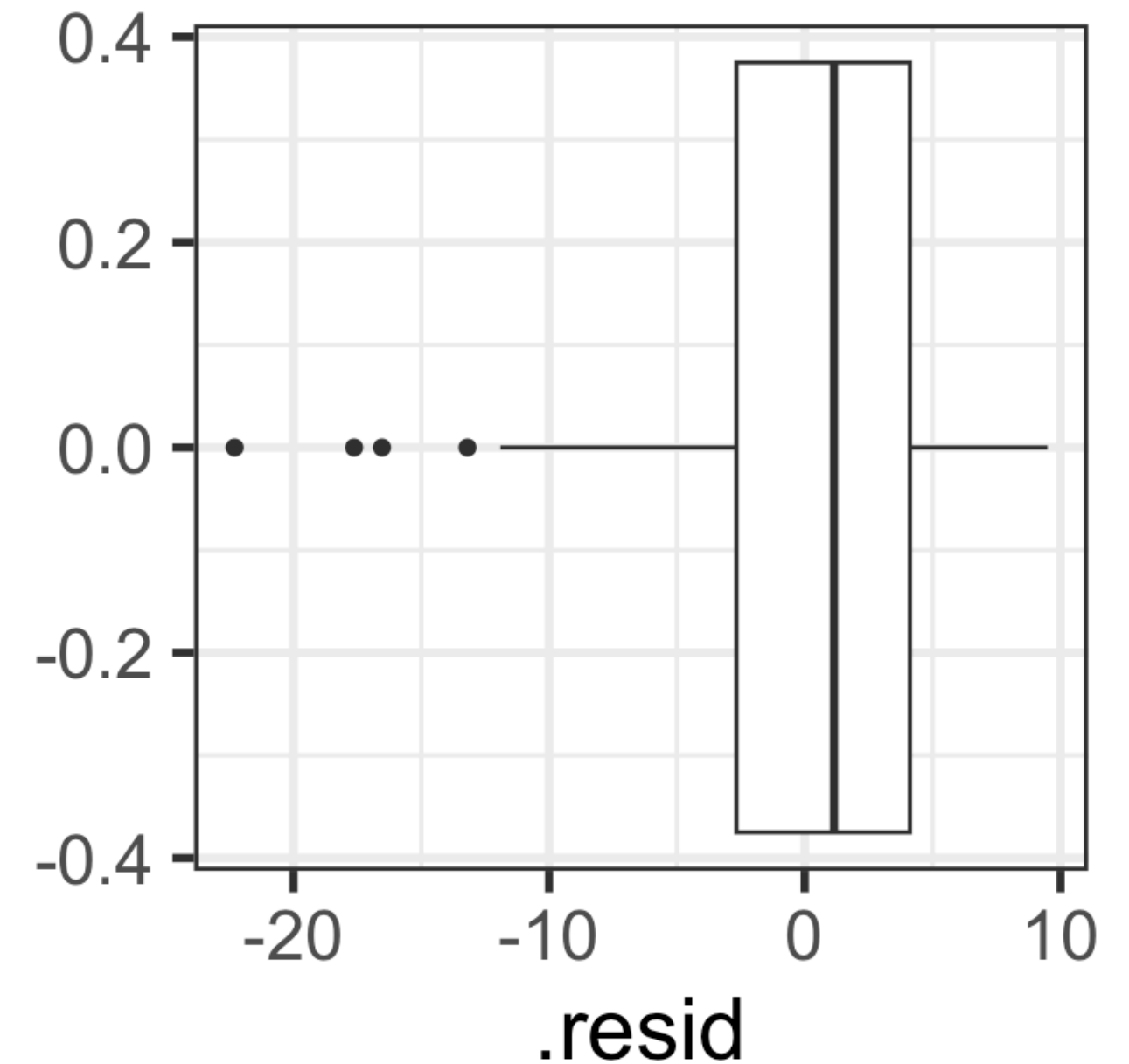
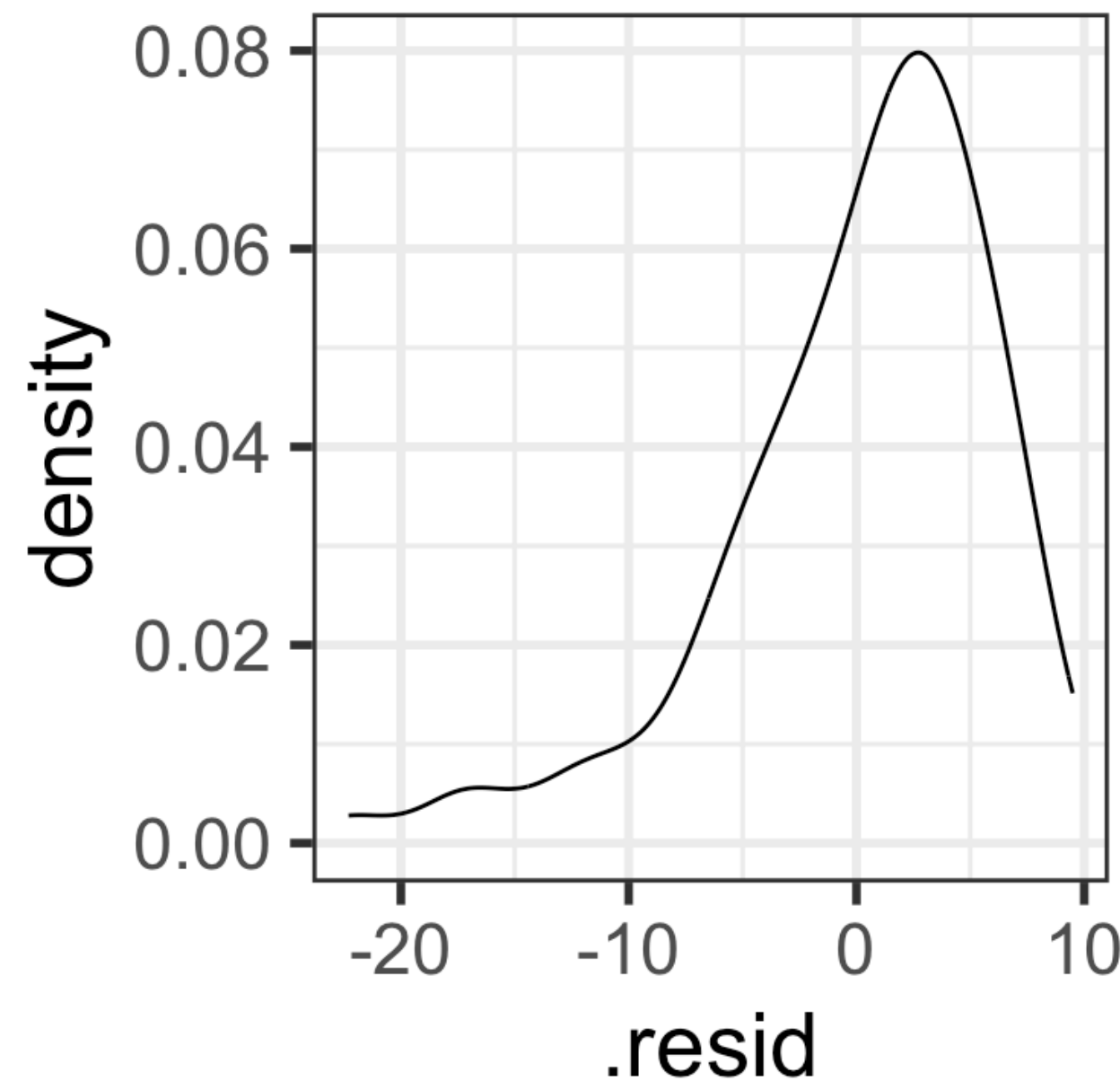
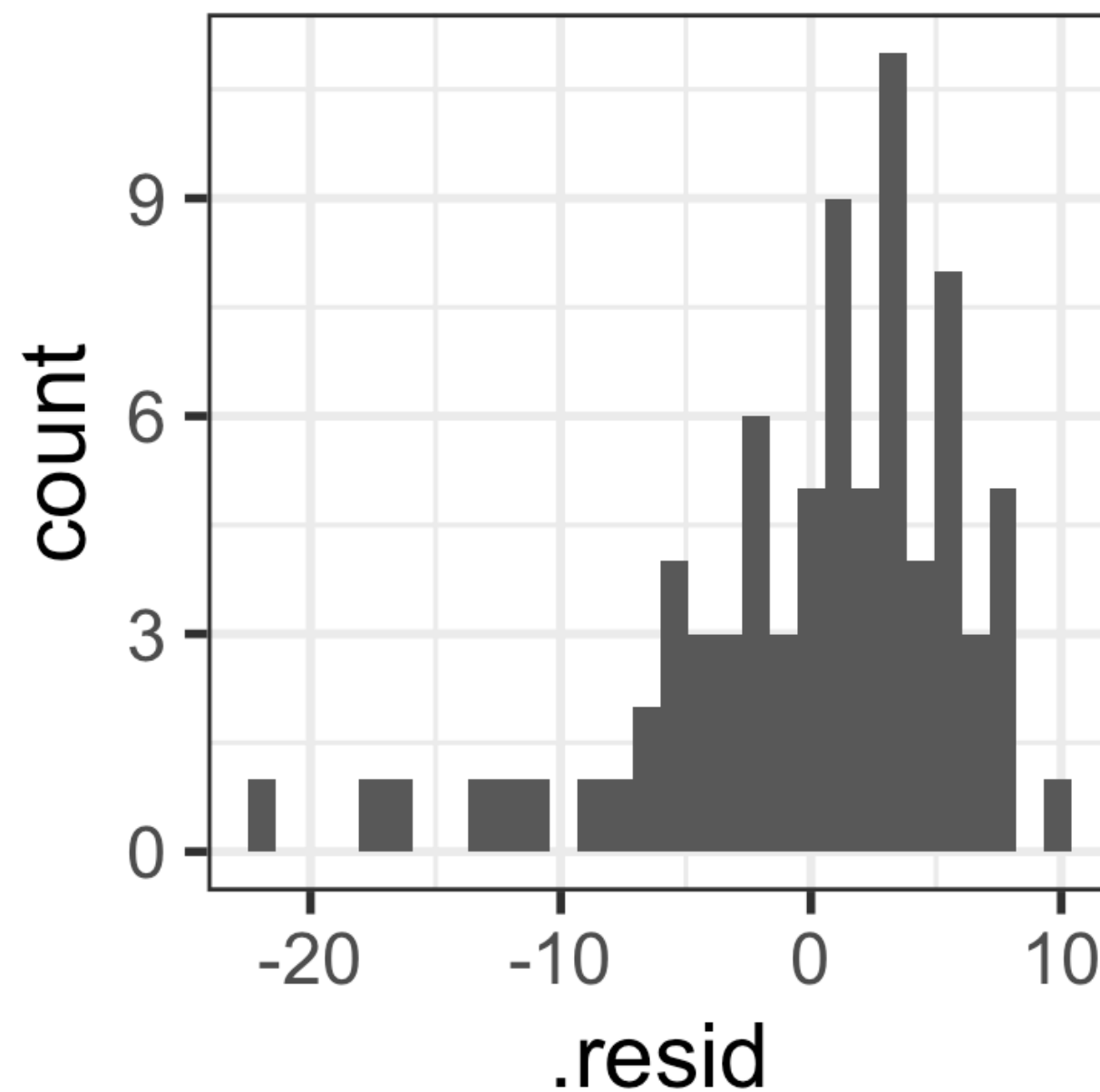
```
1 hist1 <- ggplot(aug1, aes(x = .resid)) + geom_histogram()  
2  
3 density1 <- ggplot(aug1, aes(x = .resid)) + geom_density()  
4  
5 box1 <- ggplot(aug1, aes(x = .resid)) + geom_boxplot()  
6  
7 grid.arrange(hist1, density1, box1, nrow = 1)
```



# N: Check normality with distribution plots of residuals (2/2)

- So do these plots of the residuals look normal?

```
1 grid.arrange(hist1, density1, box1, nrow = 1)
```



- My assessment: Looks like our residuals could be normal if we did not have those values around -20

# N: Normal QQ plots (QQ = quantile-quantile)

- It can be tricky to eyeball with a histogram or density plot whether the residuals are normal or not
- QQ plots are often used to help with this

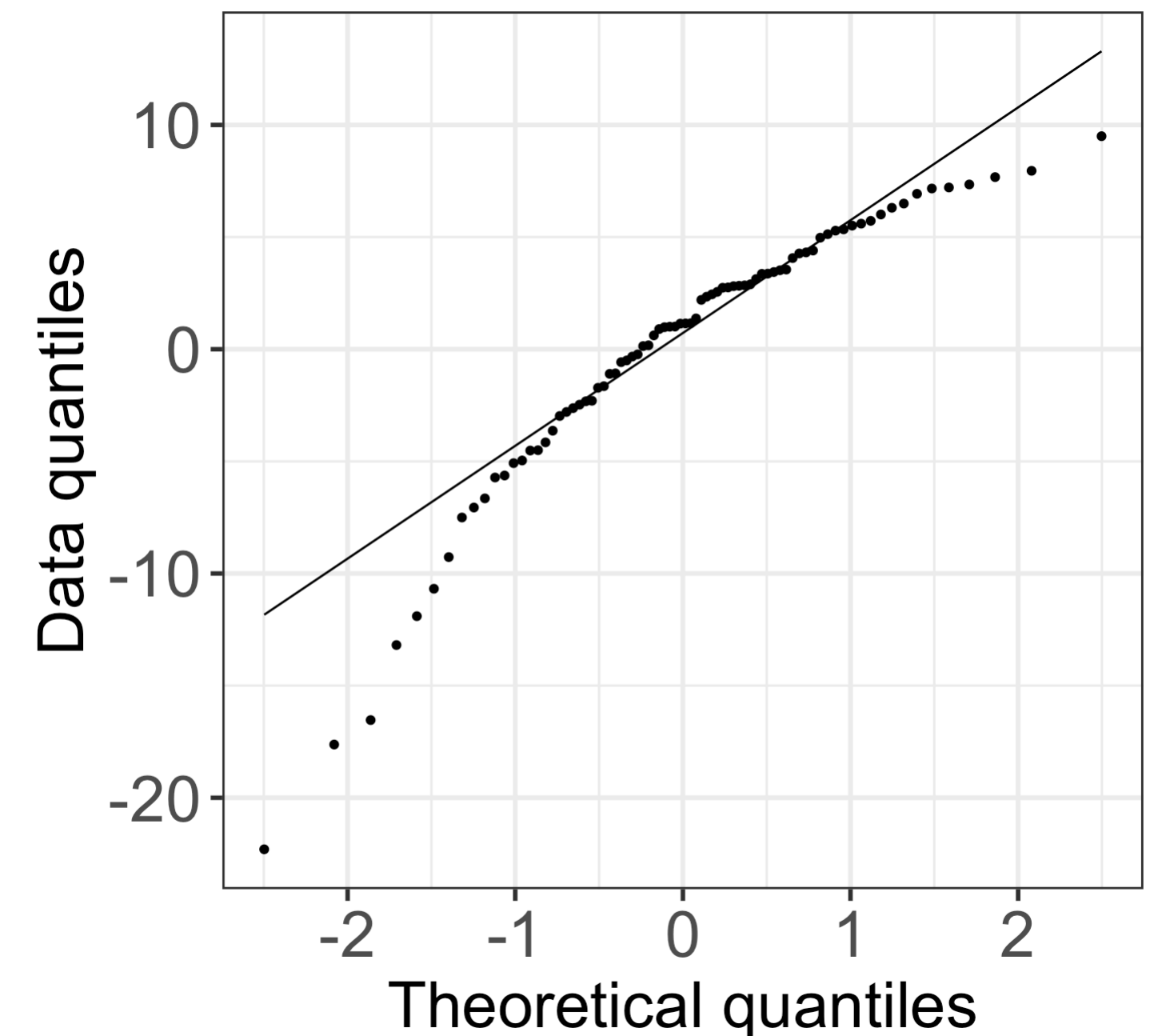
- *Vertical axis: data quantiles*

- data points are sorted in order and
  - assigned quantiles based on how many data points there are

- *Horizontal axis: theoretical quantiles*

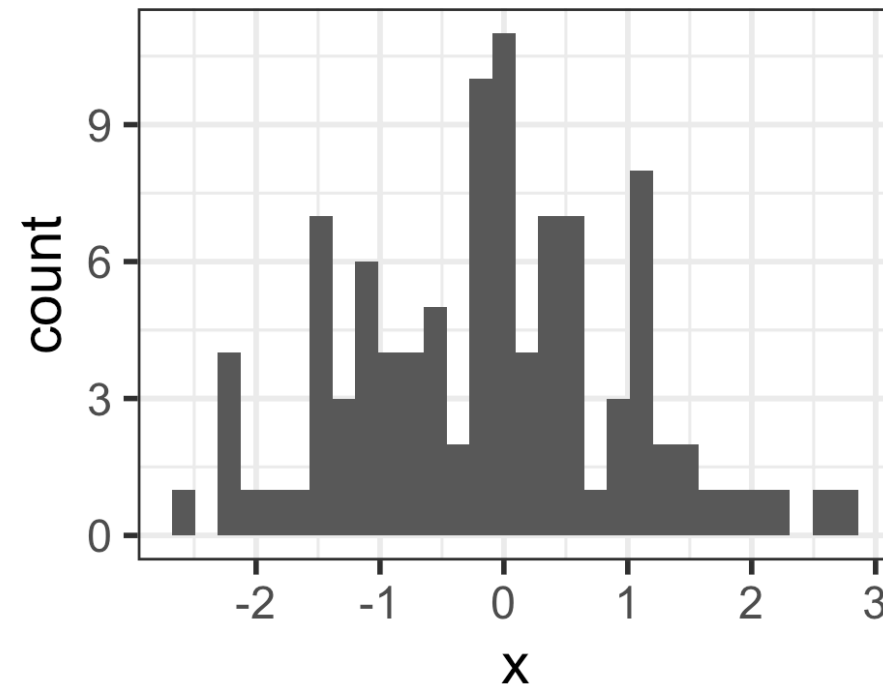
- mean and standard deviation (SD) calculated from the data points
  - theoretical quantiles are calculated for each point, assuming the data are modeled by a normal distribution with the mean and SD of the data

- **Data are approximately normal if points fall on a line.**

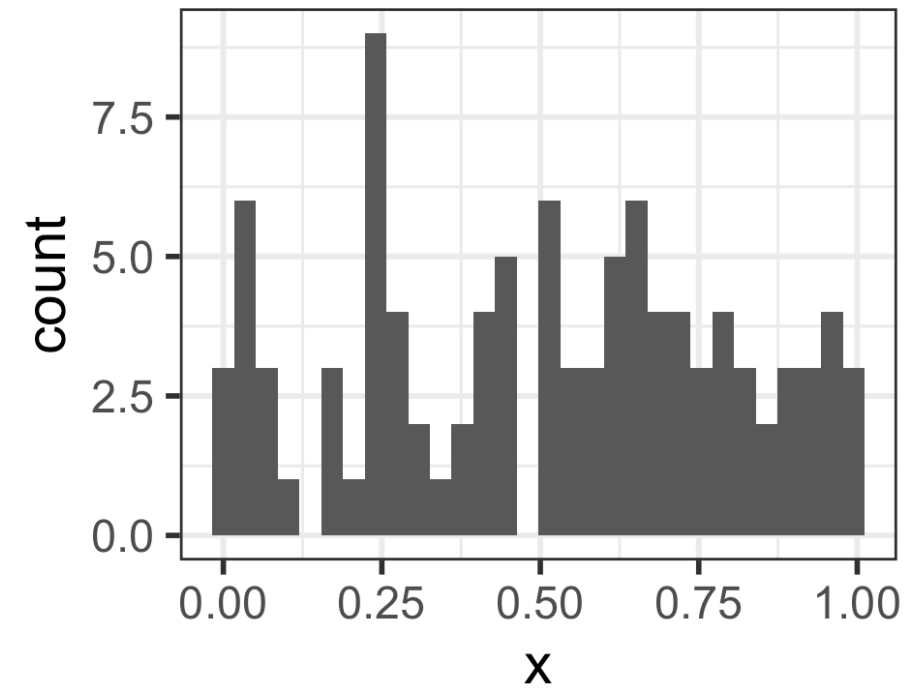


# N: Examples of Normal QQ plots (from $n = 100$ observations)

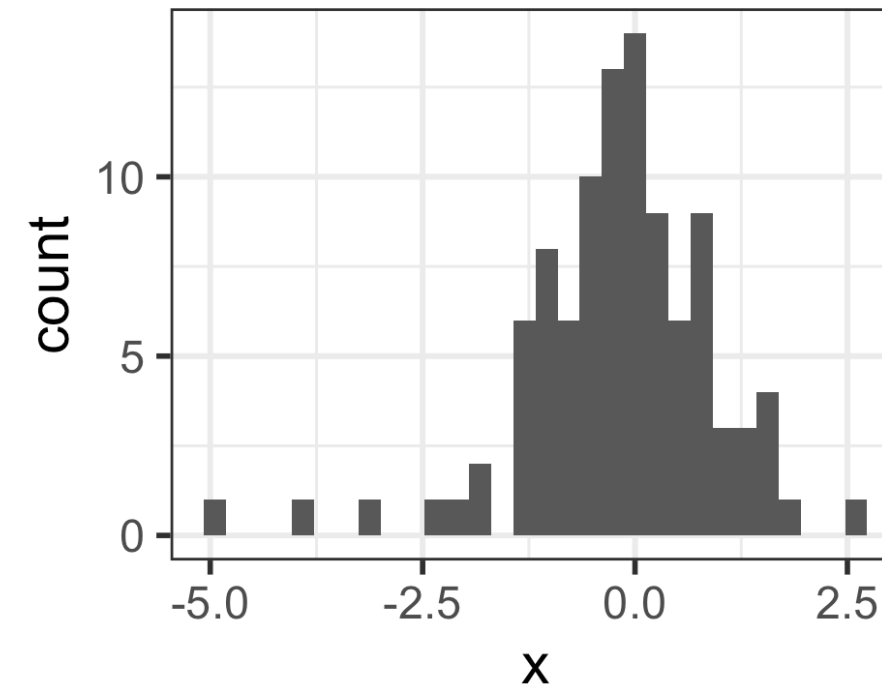
Normal



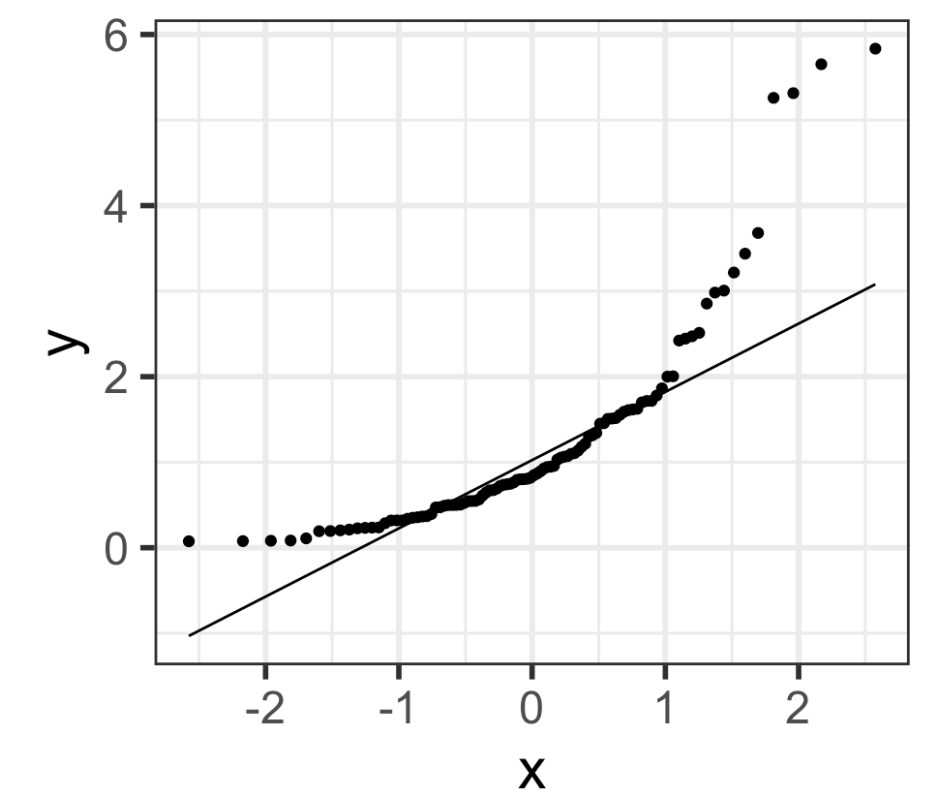
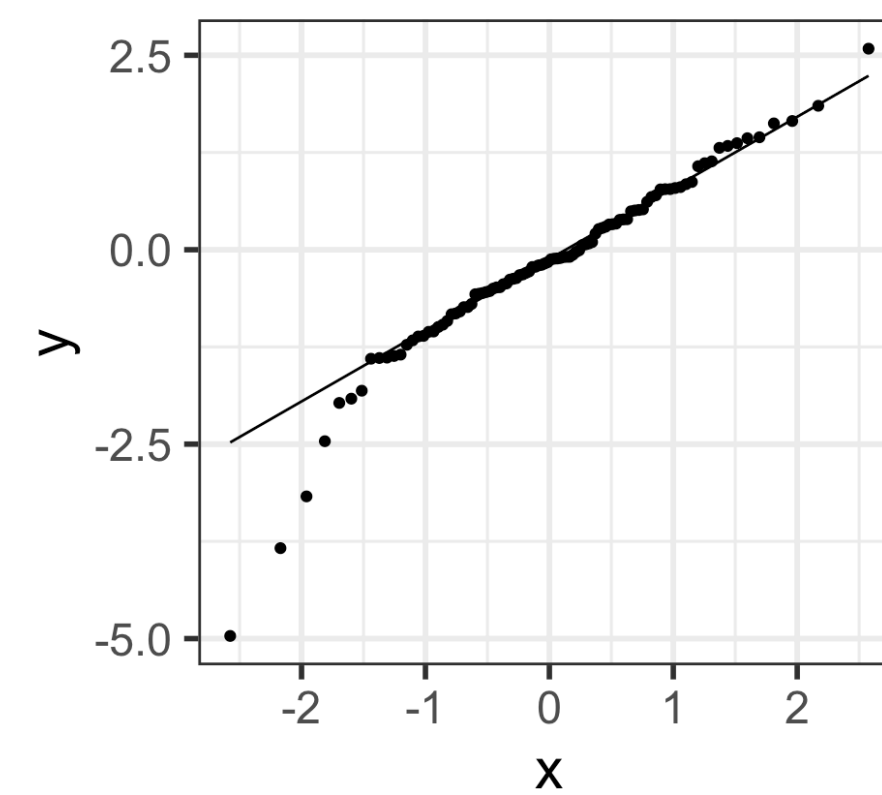
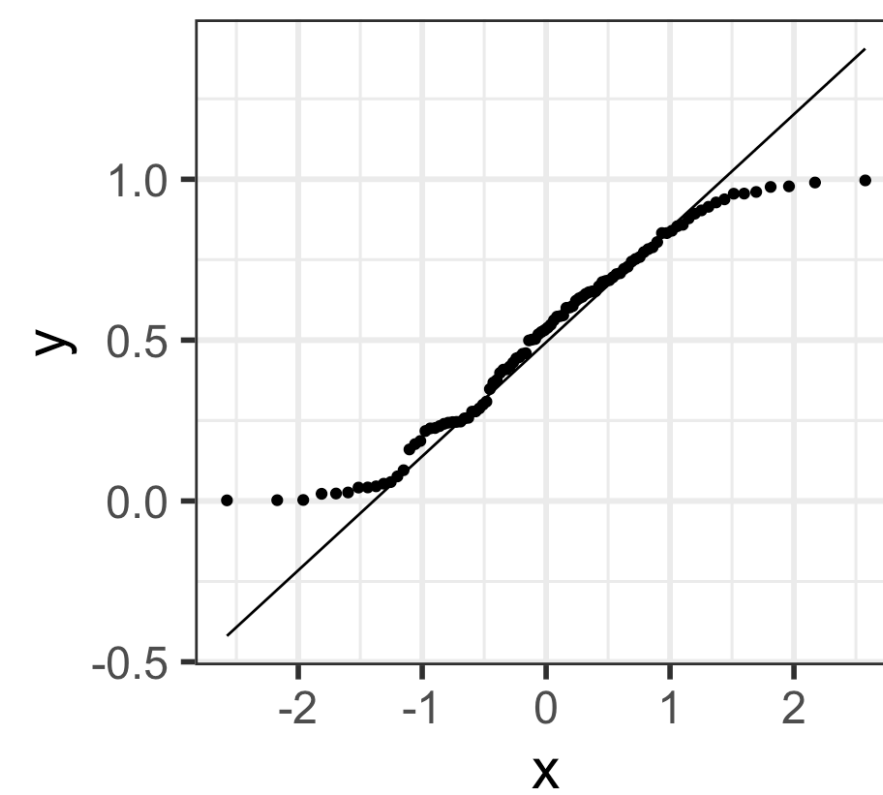
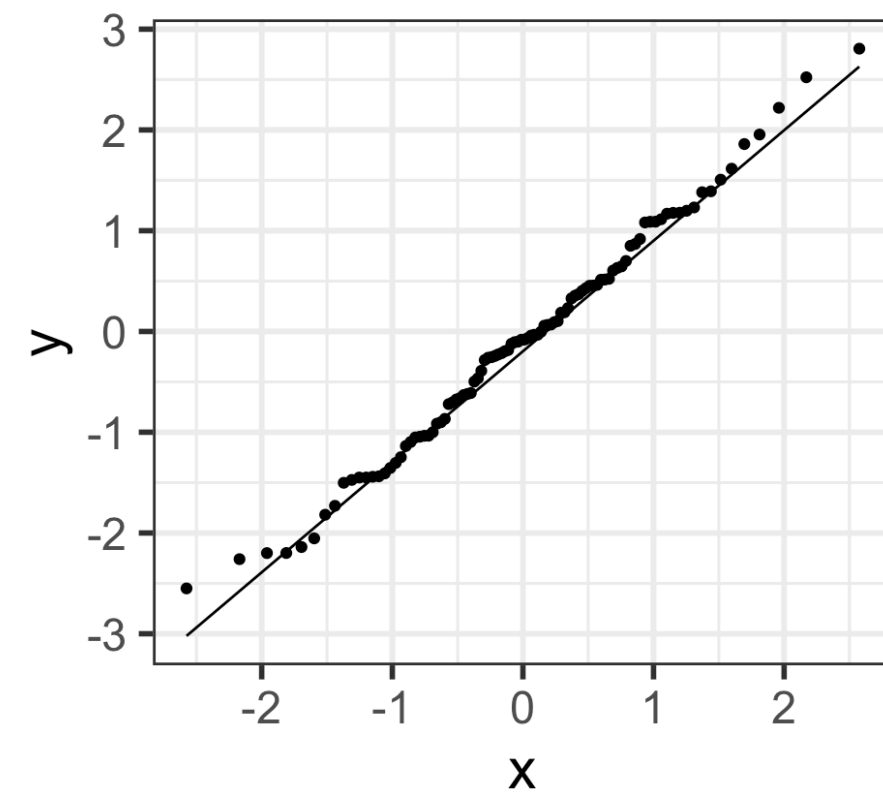
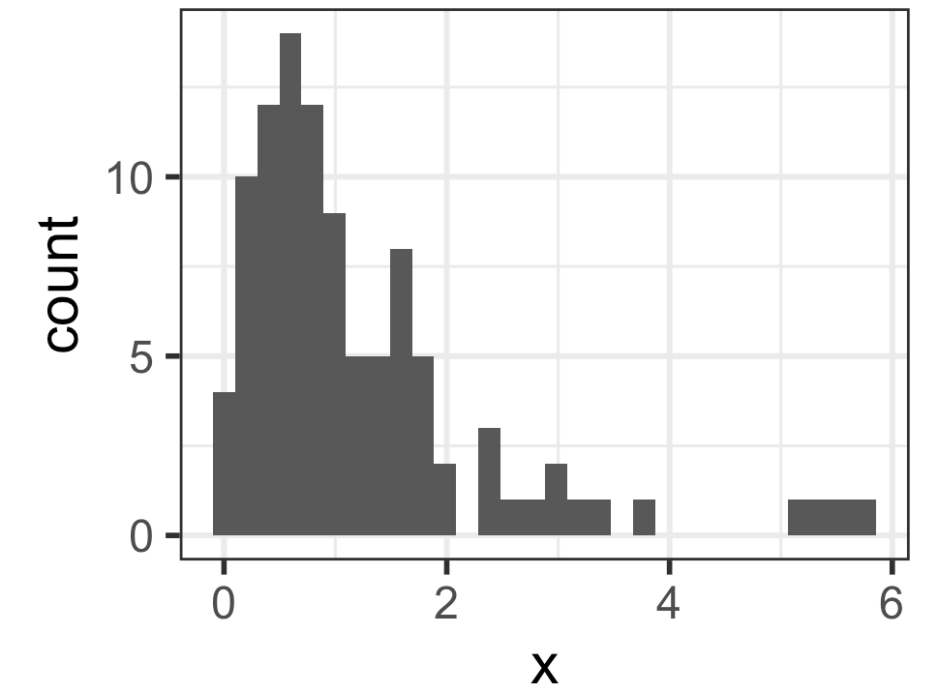
Uniform



T



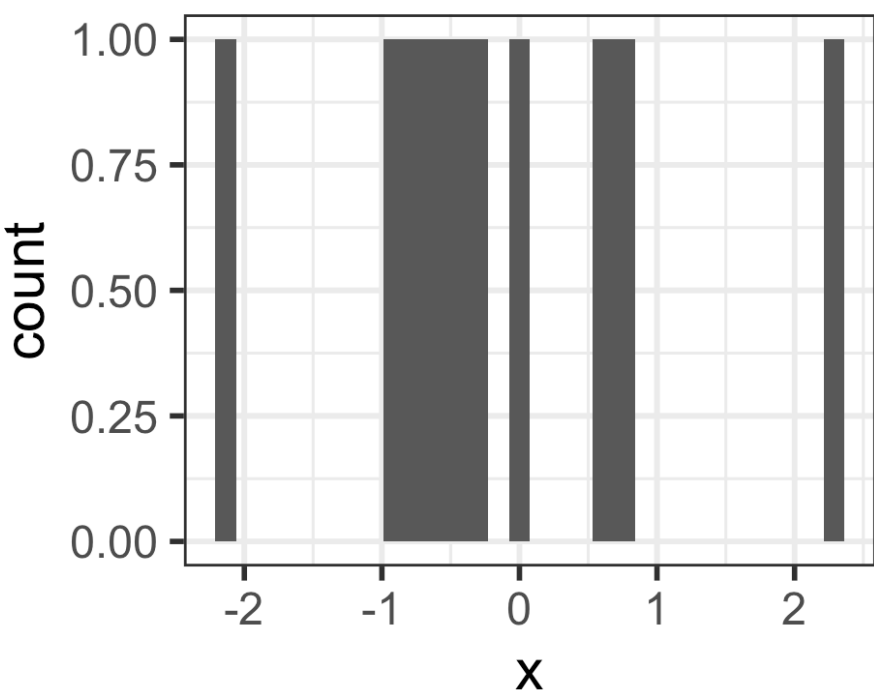
Skewed



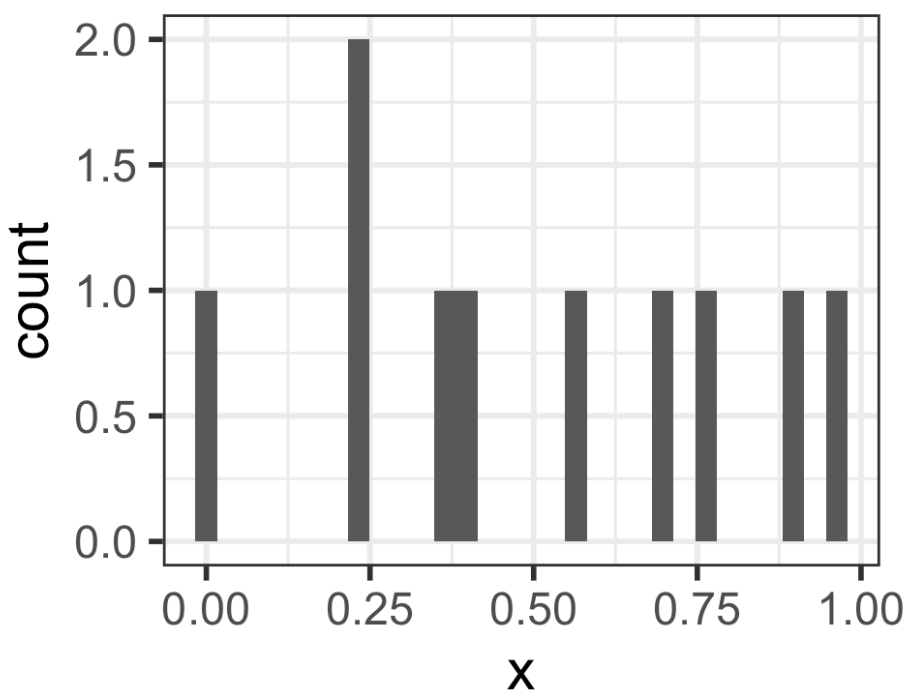


# N: Examples of Normal QQ plots (from $n = 10$ observations)

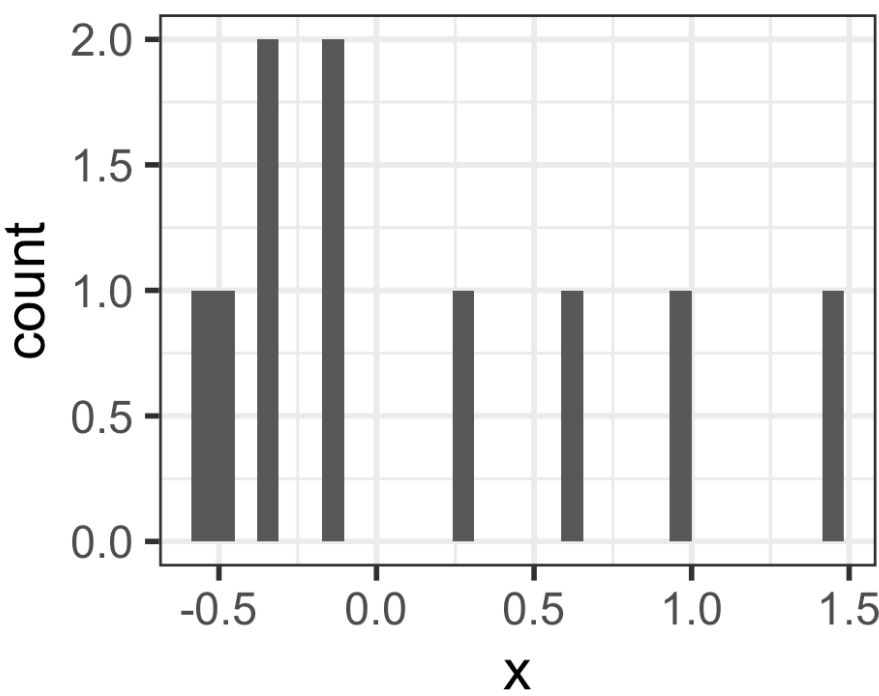
Normal



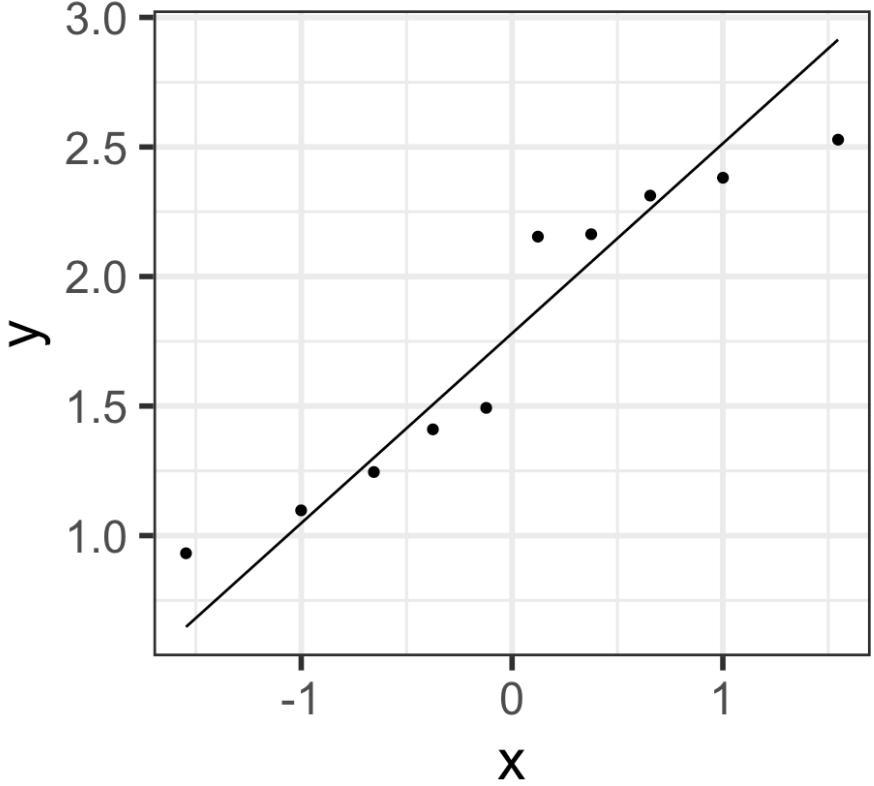
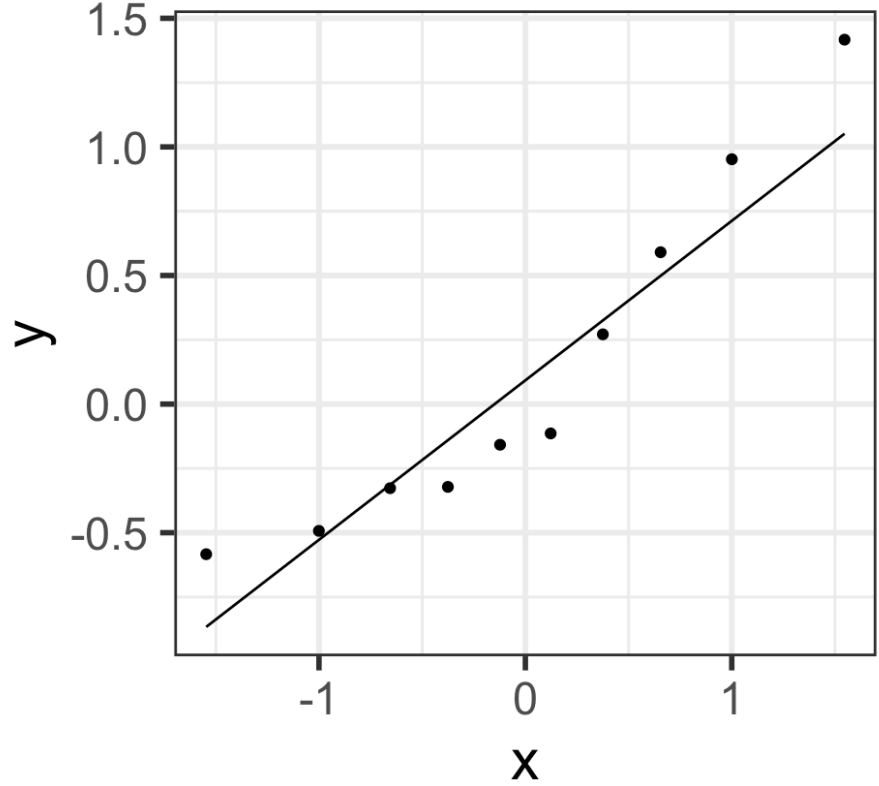
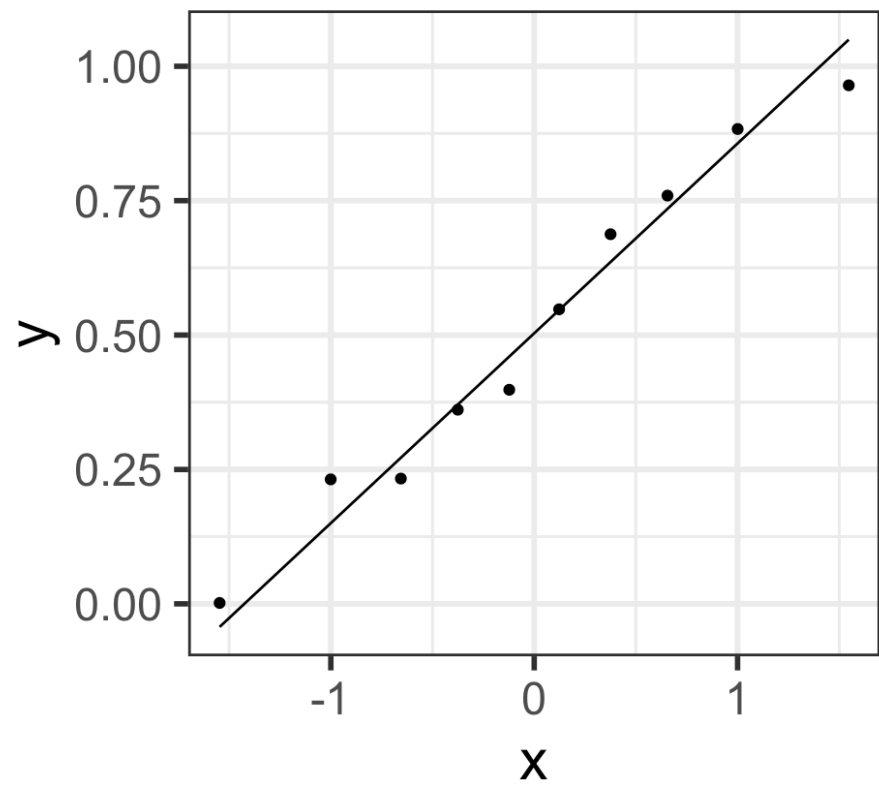
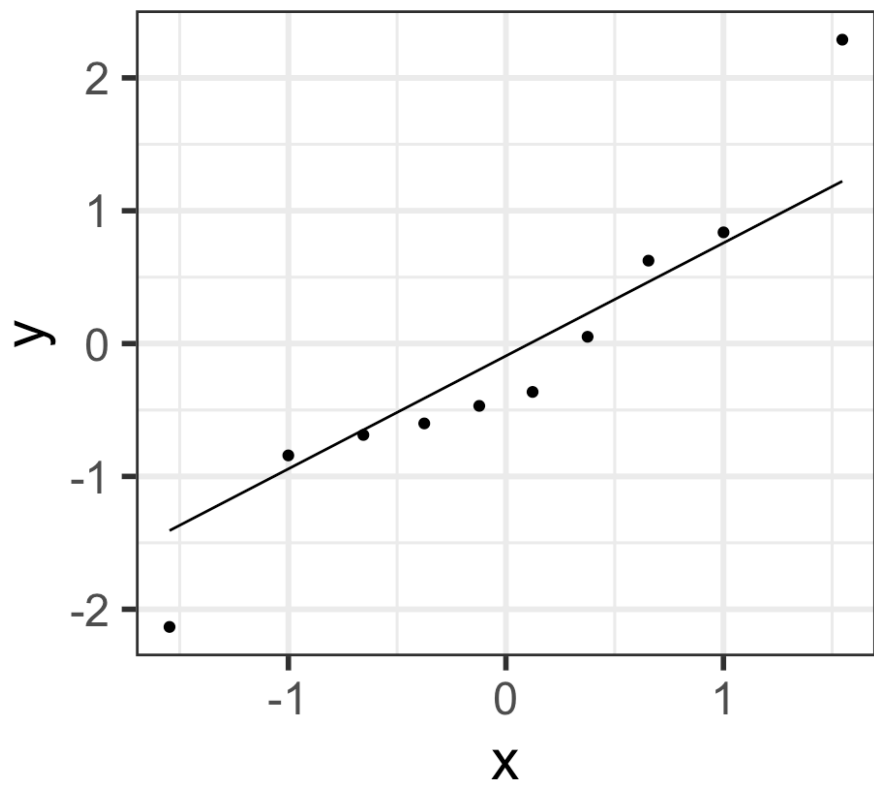
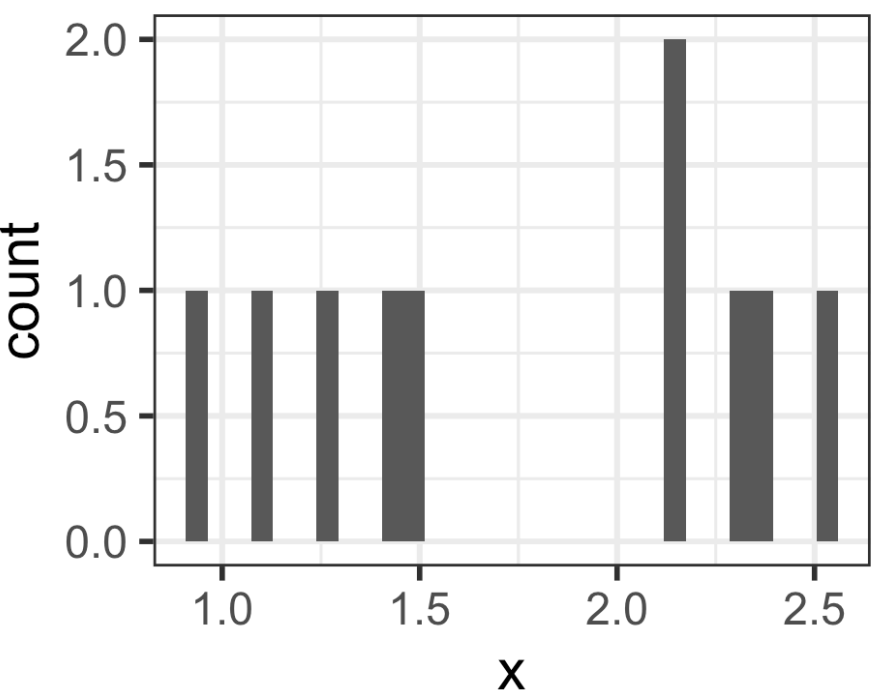
Uniform



T

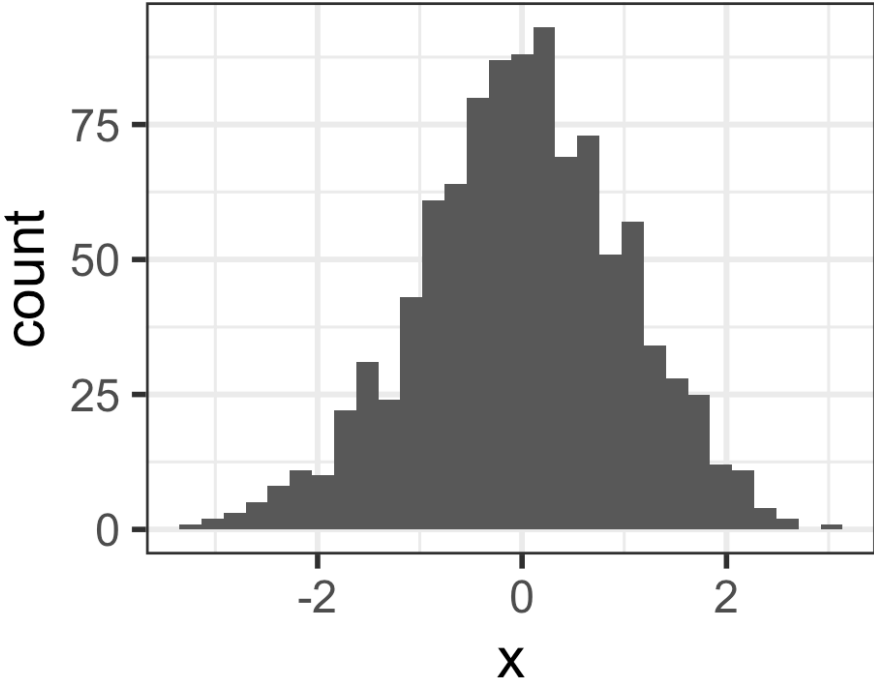


Skewed

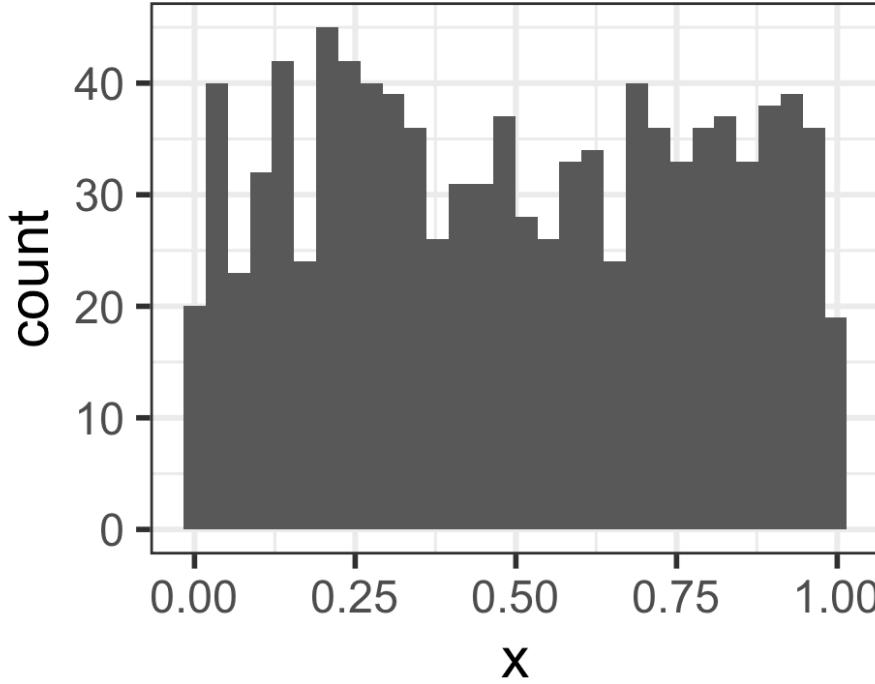


# N: Examples of Normal QQ plots (from $n = 1000$ observations)

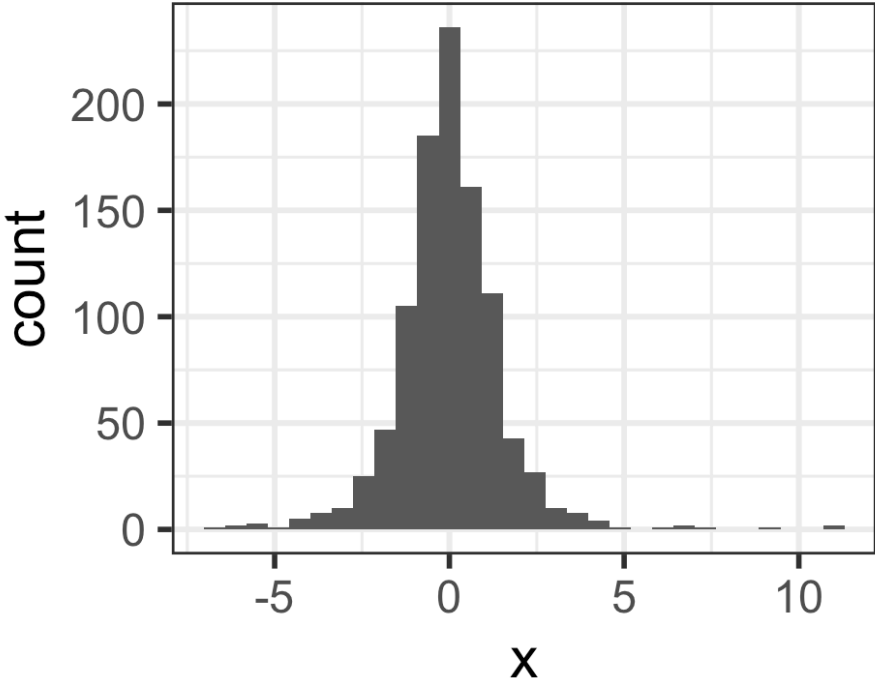
Normal



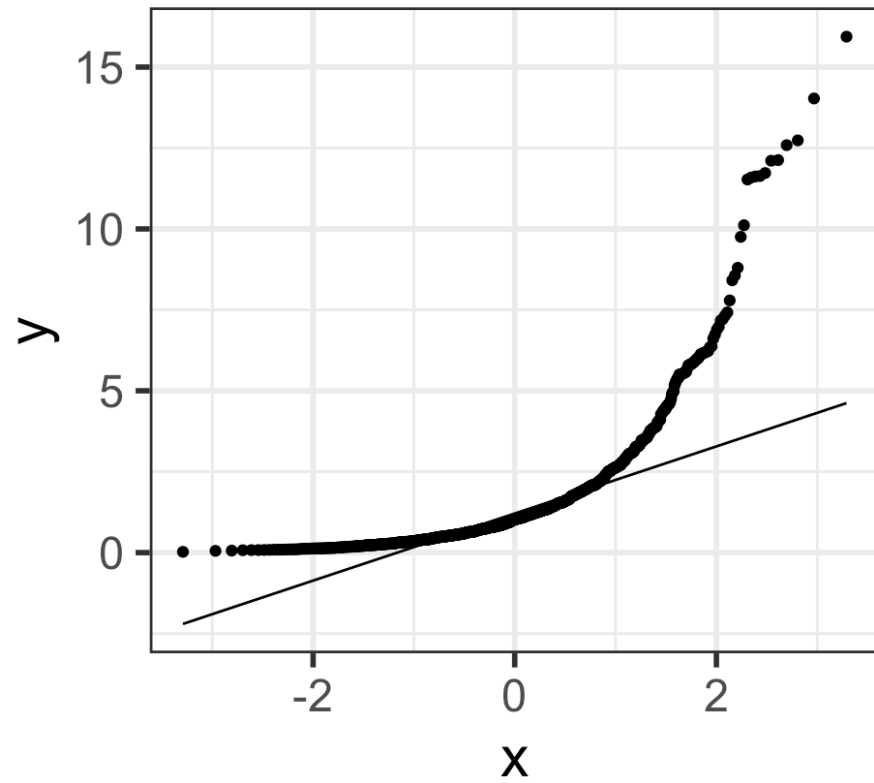
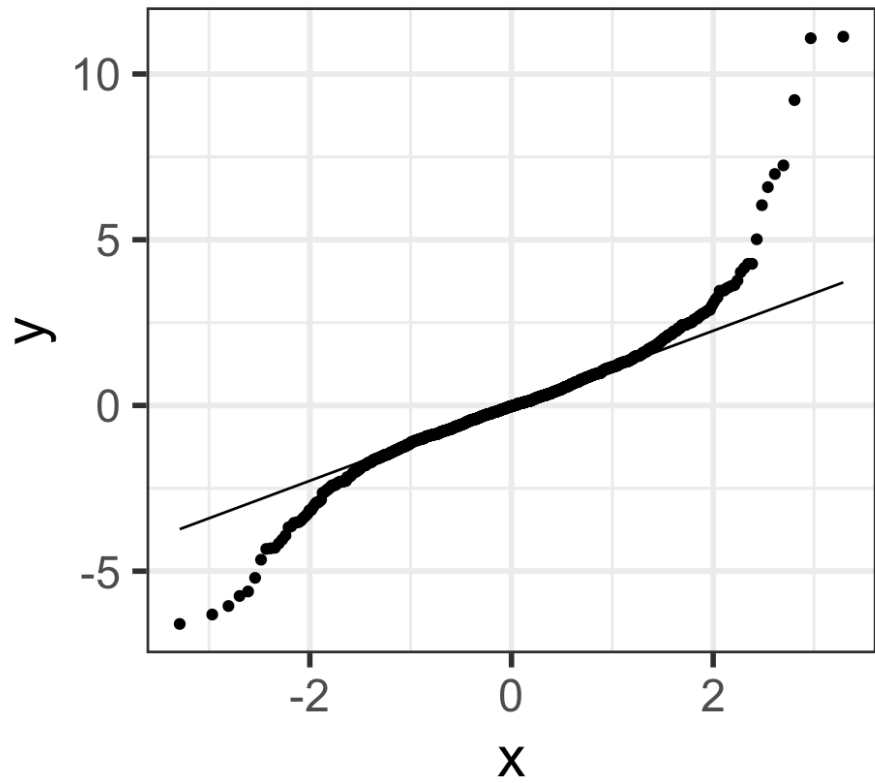
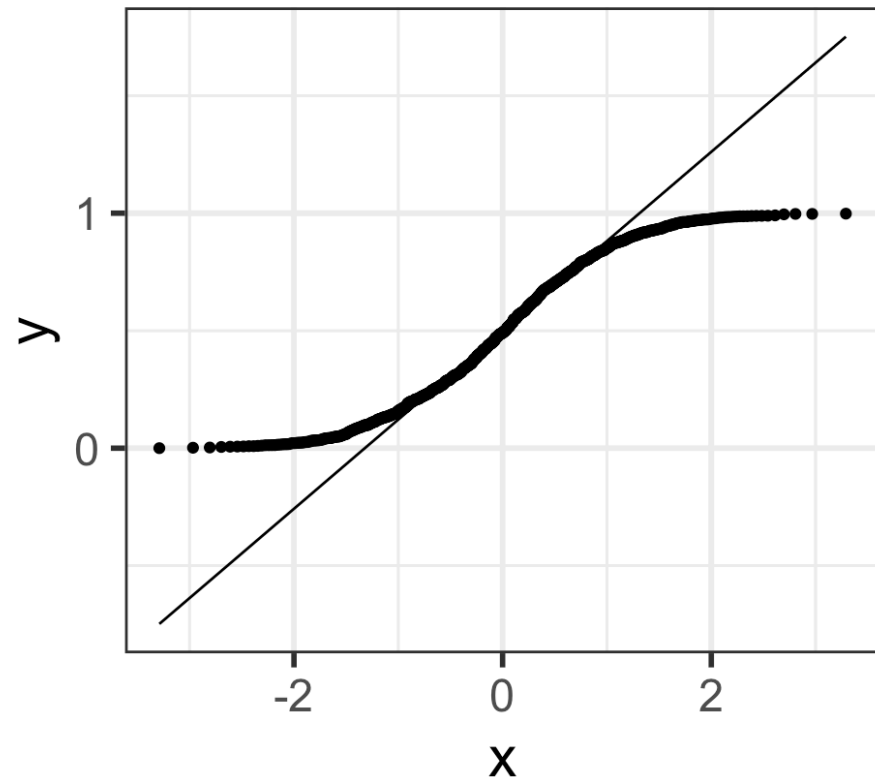
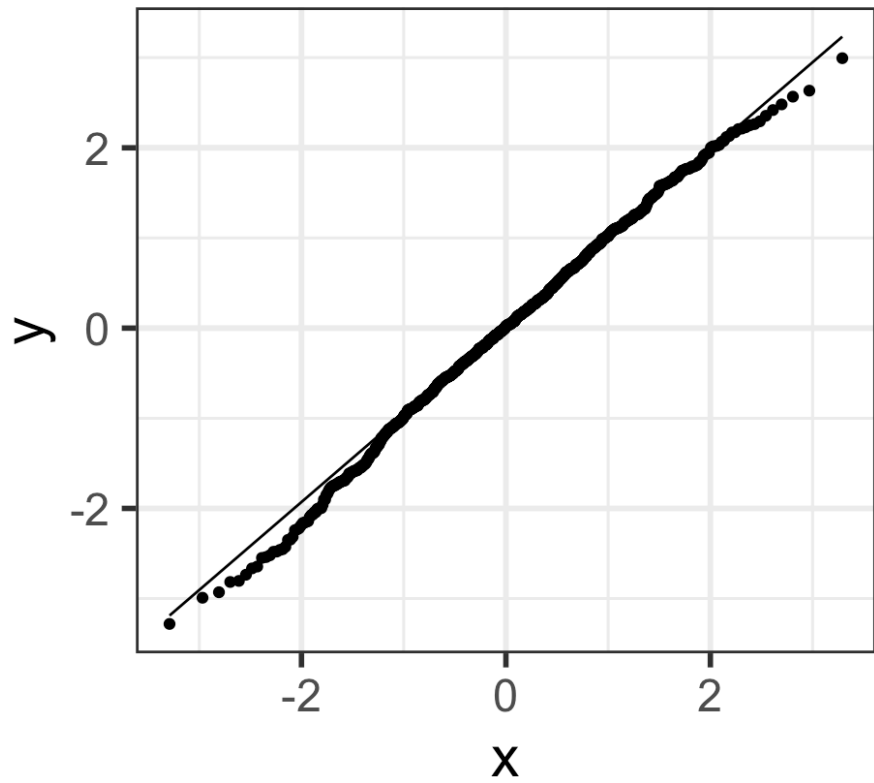
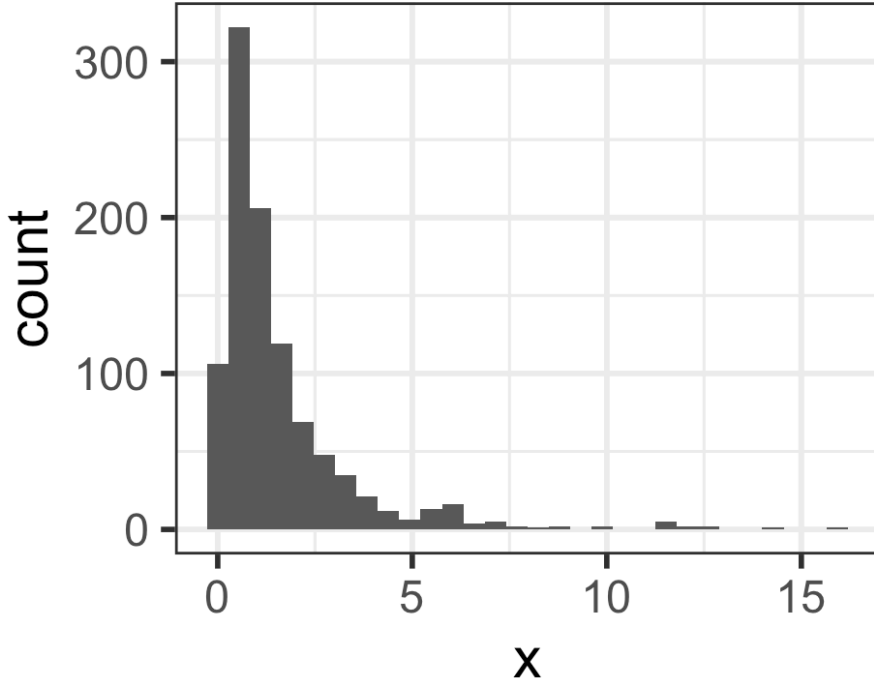
Uniform



T



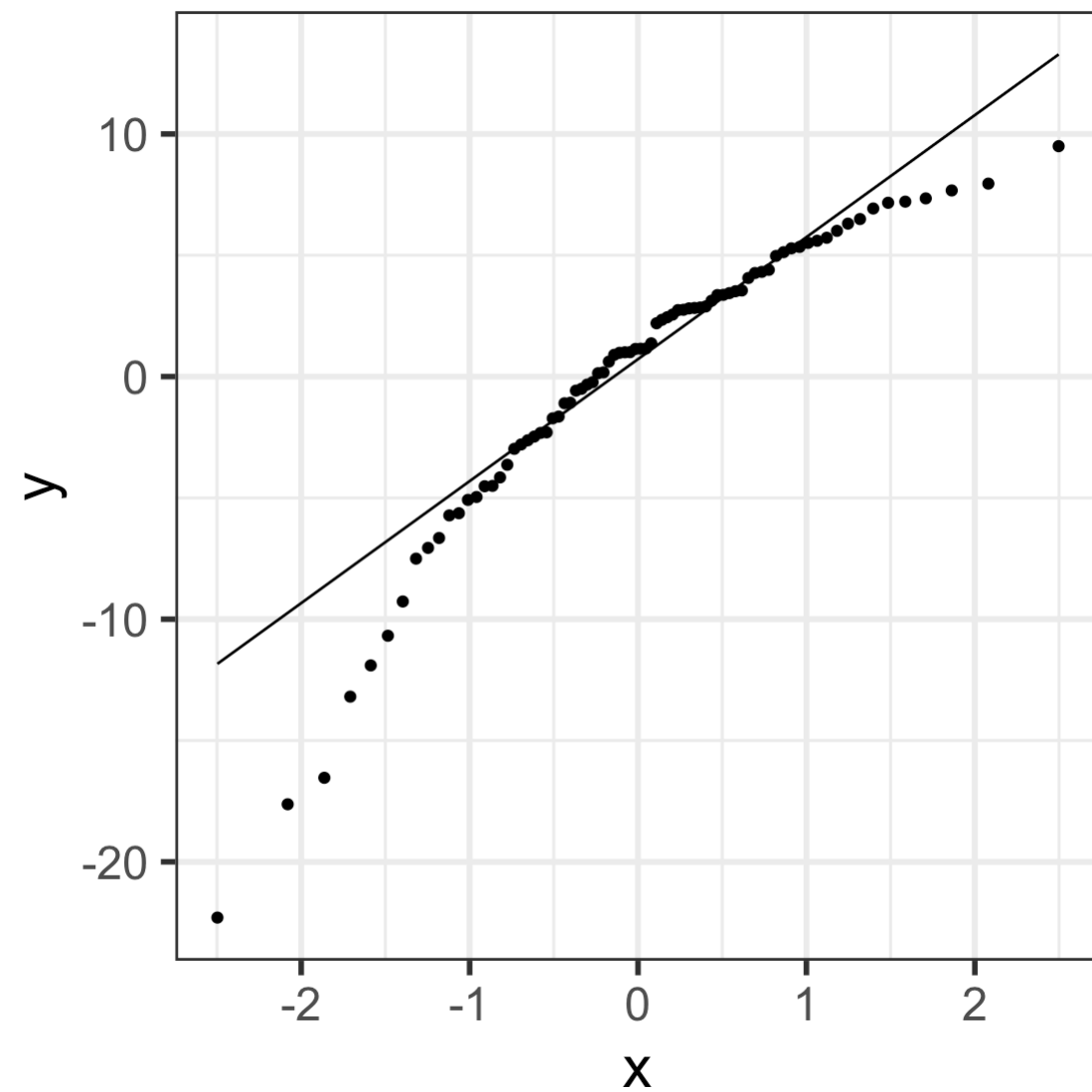
Skewed



# N: We can compare the QQ plots: model vs. theoretical

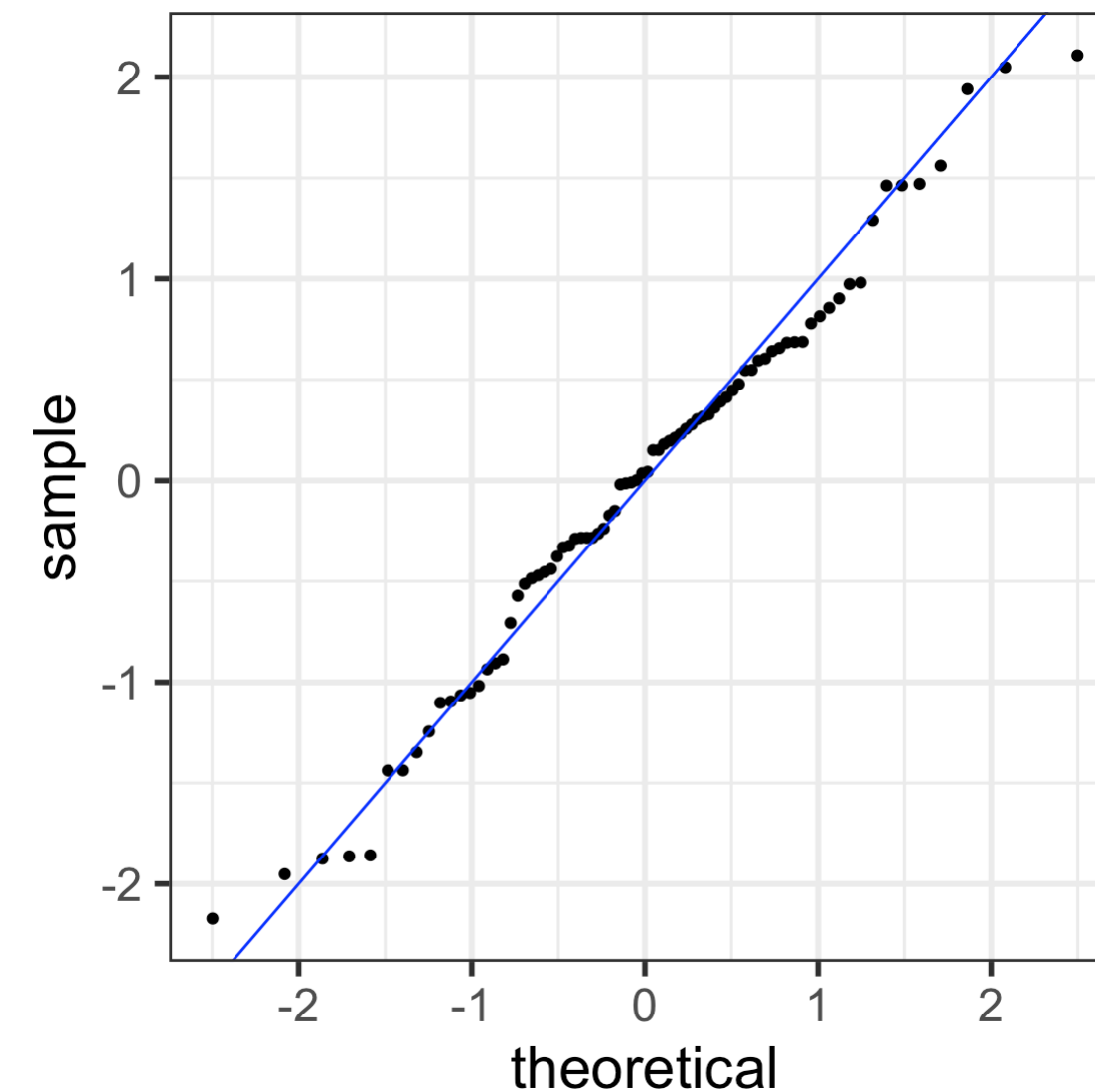
- QQ plot from Life Expectancy vs. Female Literacy Rate Regression

```
1 ggplot(aug1,  
2       aes(sample = .resid)) +  
3   stat_qq() +  
4   stat_qq_line()
```



- Simulated QQ plot of Normal Residuals with  $n = 80$

```
1 ggplot() +  
2   stat_qq(aes(  
3     sample = rnorm(80))) +  
4   geom_abline(  
5     intercept = 0, slope = 1,  
6     color = "blue")
```



# N: Shapiro-Wilk Test of Normality

- Goodness-of-fit test for the normal distribution: Is there evidence that our residuals are from a normal distribution?
- Hypothesis test:

$H_0$  : data are from a normally distributed population

$H_1$  : data are NOT from a normally distributed population

```
1 shapiro.test(aug1$.resid)
```

Shapiro-Wilk normality test

data: aug1\$.resid

W = 0.90575, p-value = 2.148e-05

## Conclusion

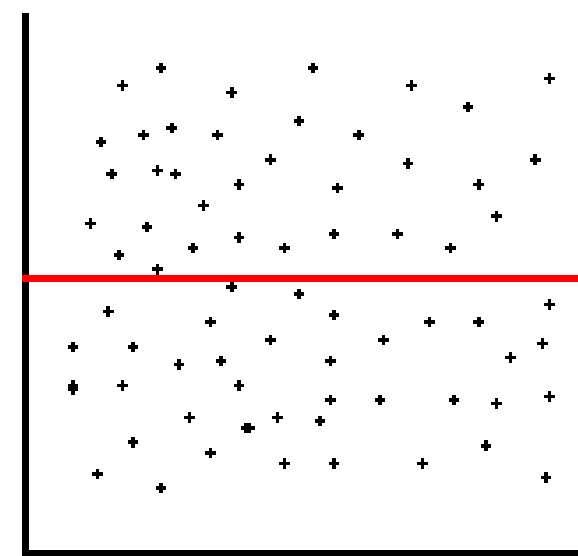
Reject the null. Data are not from a normal distribution.

# Learning Objectives

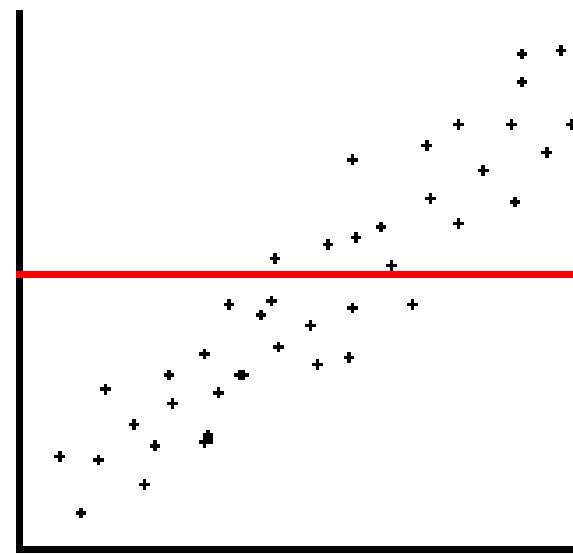
1. Describe the model assumptions made in linear regression using ordinary least squares
2. Determine if the relationship between our sampled  $X$  and  $Y$  is linear
3. Use QQ plots to determine if our fitted model holds the normality assumption
4. Use residual plots to determine if our fitted model holds the equality of variance assumption

# E: Equality of variance of the residuals

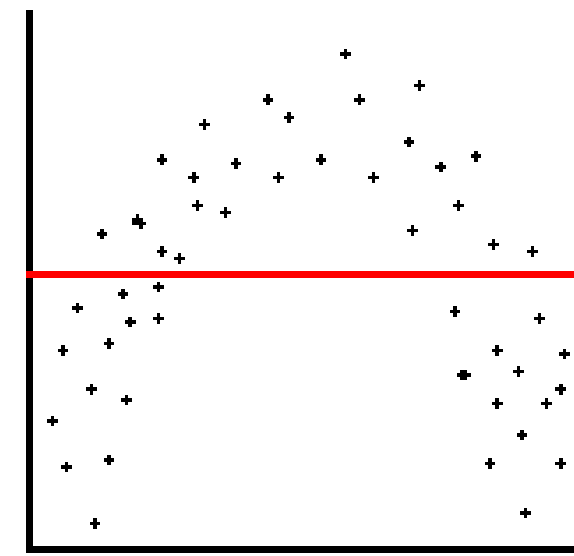
- Homoscedasticity: How do we determine if the variance across X values is constant?
- Diagnostic tool: **residual plot**



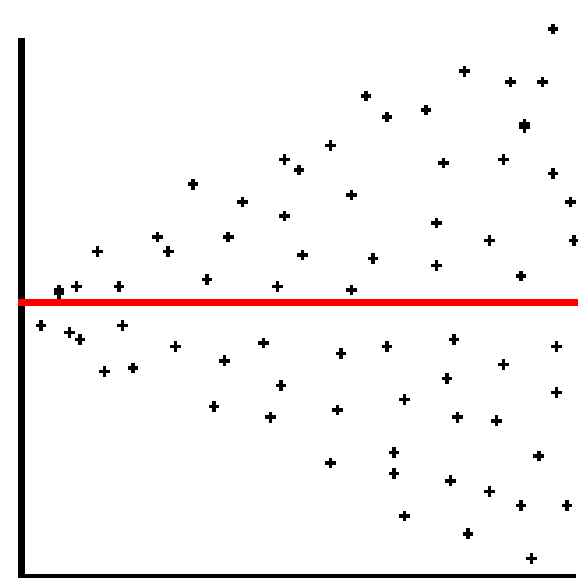
(a) Unbiased and Homoscedastic



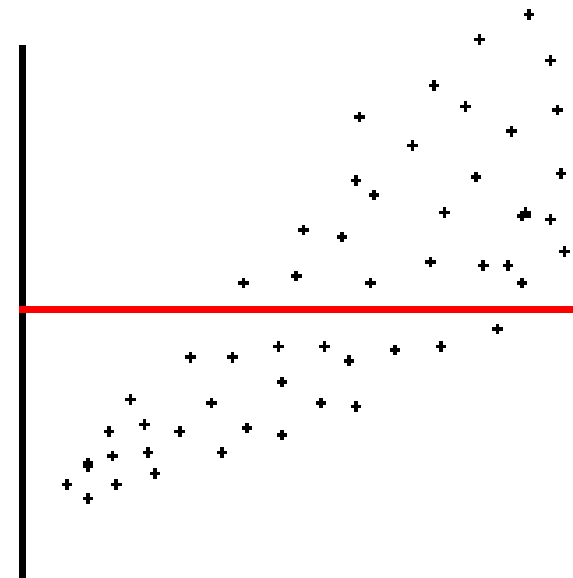
(b) Biased and Homoscedastic



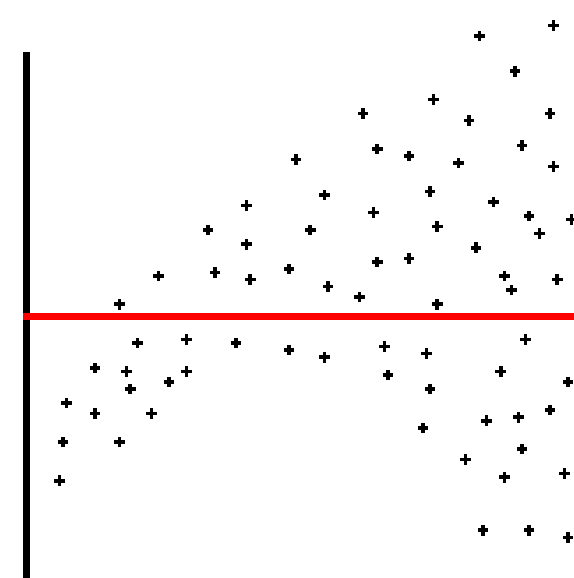
(c) Biased and Homoscedastic



(d) Unbiased and Heteroscedastic



(e) Biased and Heteroscedastic



(f) Biased and Heteroscedastic

## E: Creating a residual plot

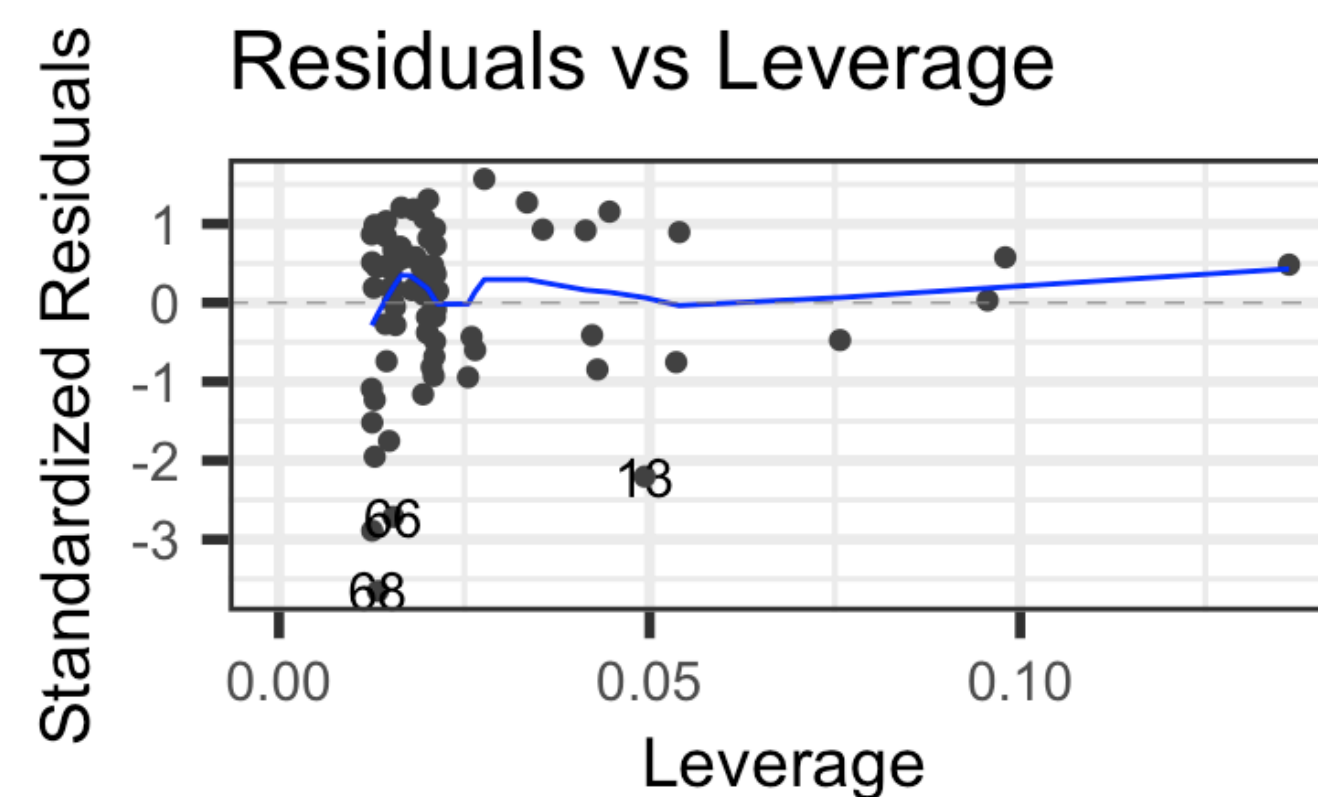
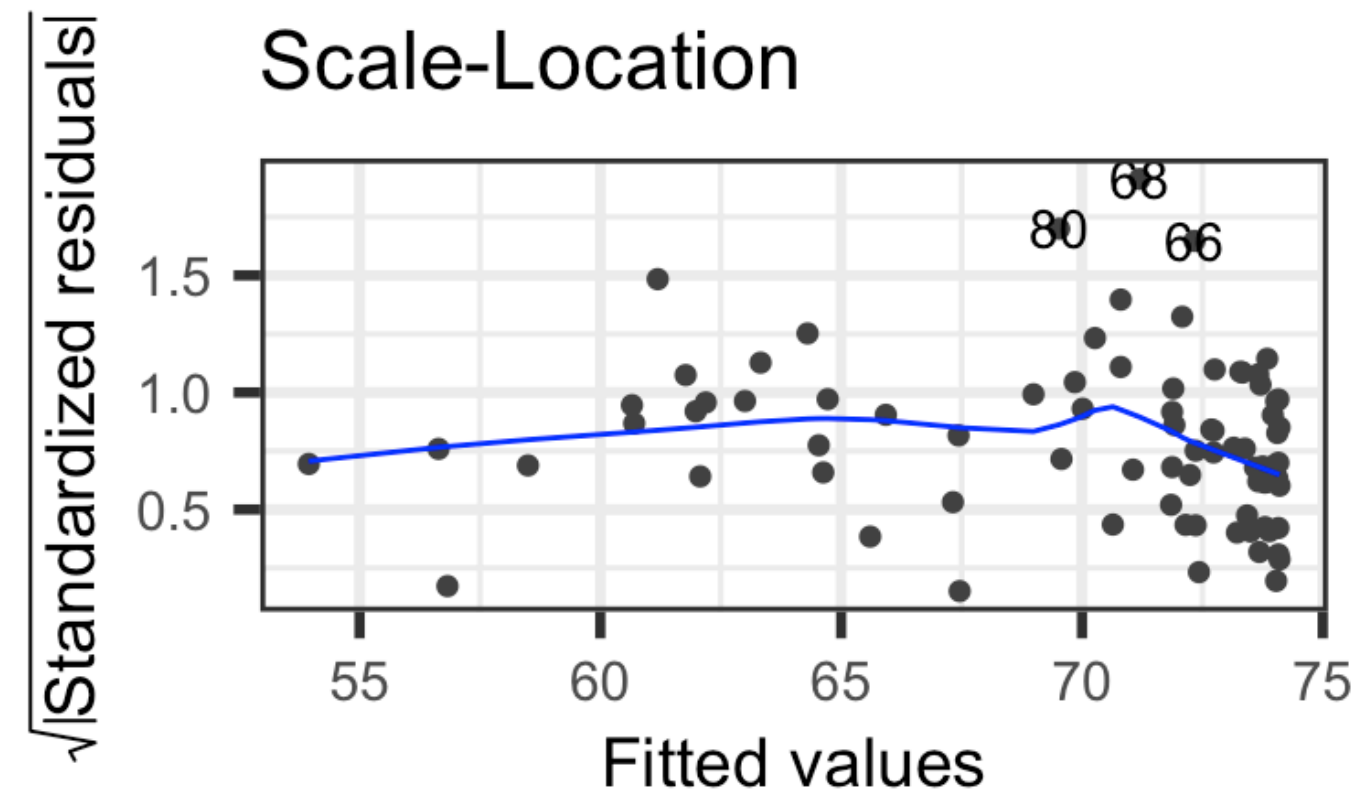
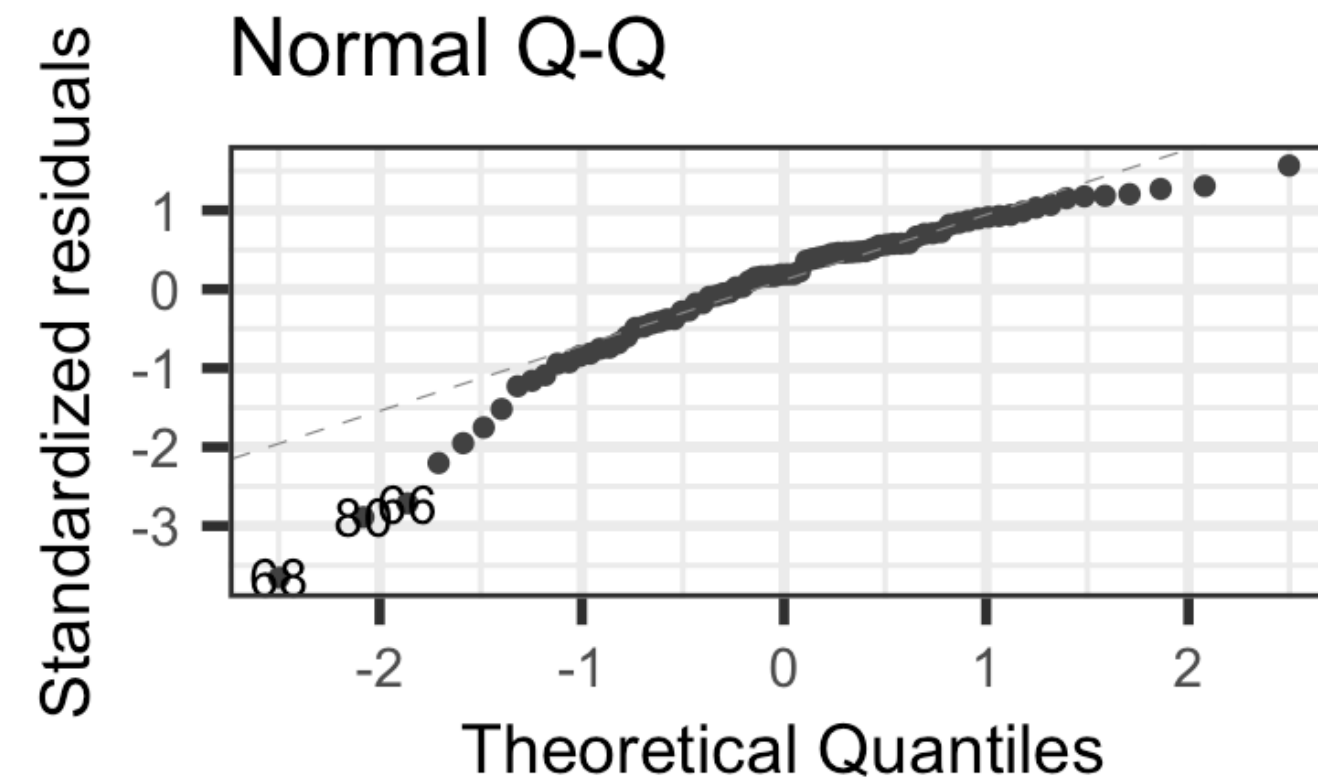
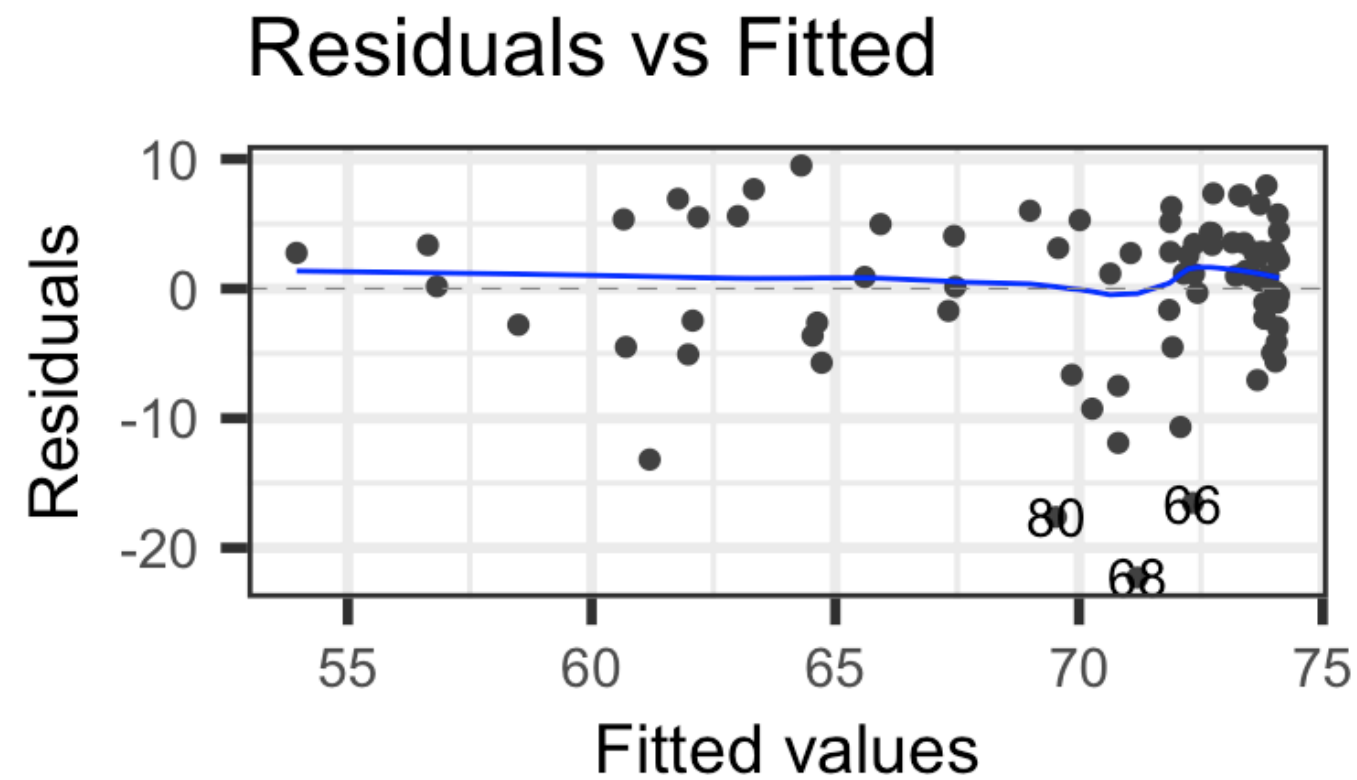
- $x$  = explanatory variable from regression model
  - (or the fitted values for a multiple regression)
- $y$  = residuals from regression model

```
1 ggplot(aug1,  
2       aes(x = FemaleLiteracyRate,  
3           y = .resid)) +  
4   geom_point(size = 2) +  
5   geom_abline(intercept = 0, slope = 0,  
6             size = 2, color = "#FF8021") +  
7   labs(title = "Residual plot") +  
8   theme(axis.title = element_text(size = 30),  
9         axis.text = element_text(size = 30))
```



# autoplot() can be a helpful tool

```
1 library(ggfortify)
2 autoplot(model1) + theme(text=element_text(size=14))
```





# Summary of the assumptions and their diagnostic tool

Assumption	What needs to hold?	Diagnostic tool
Linearity	<ul style="list-style-type: none"><li>Relationship between <math>X</math> and <math>Y</math> is linear</li></ul>	<ul style="list-style-type: none"><li>Scatterplot of <math>Y</math> vs. <math>X</math></li></ul>
Independence	<ul style="list-style-type: none"><li>Observations are independent from each other</li></ul>	<ul style="list-style-type: none"><li>Study design</li></ul>
Normality	<ul style="list-style-type: none"><li>Residuals (and thus <math>Y X</math>) are normally distributed</li></ul>	<ul style="list-style-type: none"><li>QQ plot of residuals</li><li>Distribution of residuals</li></ul>
Equality of variance	<ul style="list-style-type: none"><li>Variance of residuals (and thus <math>Y X</math>) is same across <math>X</math> values (homoscedasticity)</li></ul>	<ul style="list-style-type: none"><li>Residual plot</li></ul>

# We didn't really go over our options when these assumptions do not hold

- We will consider this more once we get into multiple linear regression
- For now, with SLR, when assumptions do not hold, I conclude we need to add more variables in the model
- Another note: I did not make these plots very presentable
  - Axes were left with whatever names were given to them
  - These plots are usually just for us!
  - Not really something that you include in a formal report

