

17-803 Empirical Methods

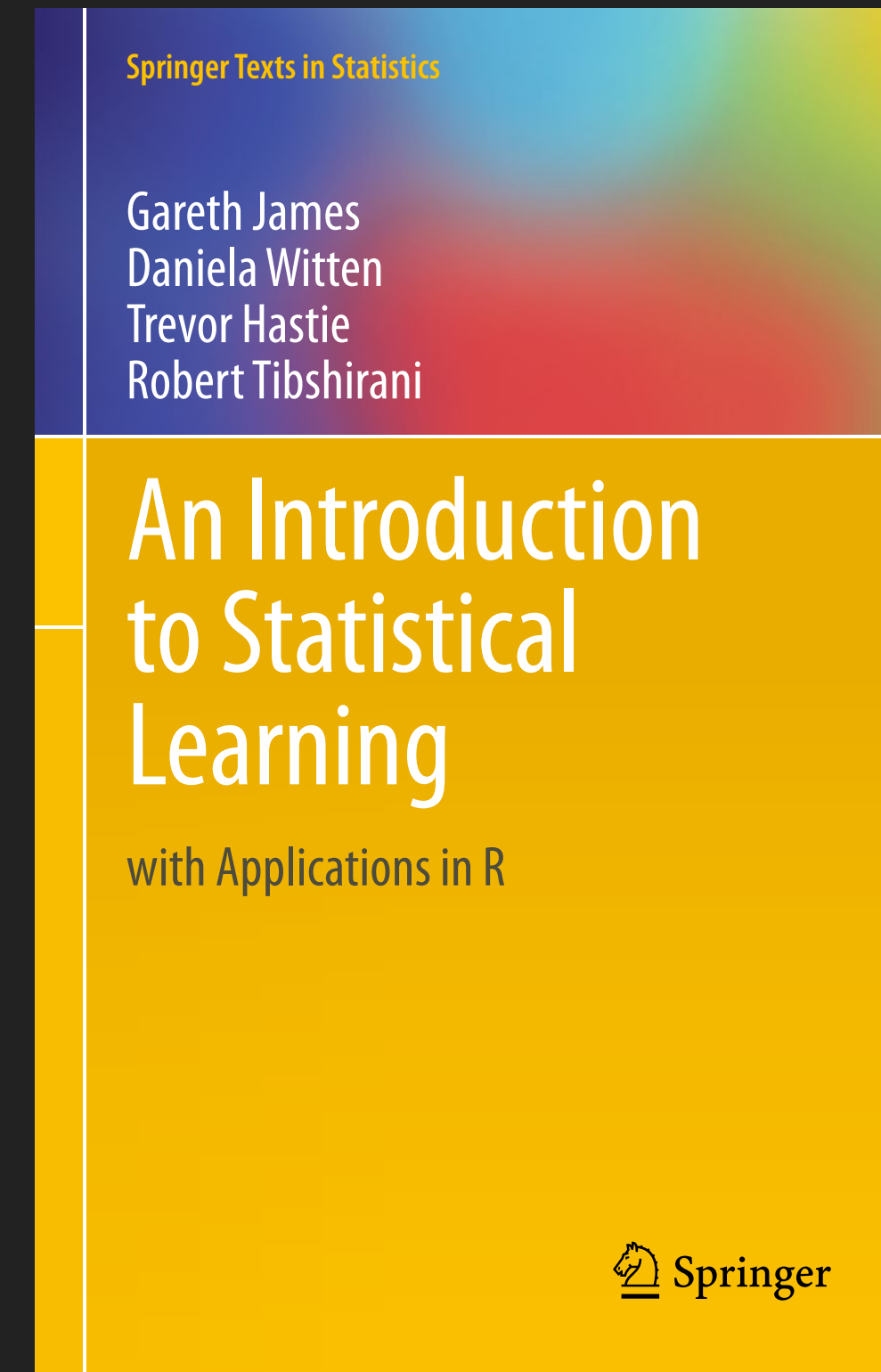
Bogdan Vasilescu, S3D

Regression Modeling (Part 1)

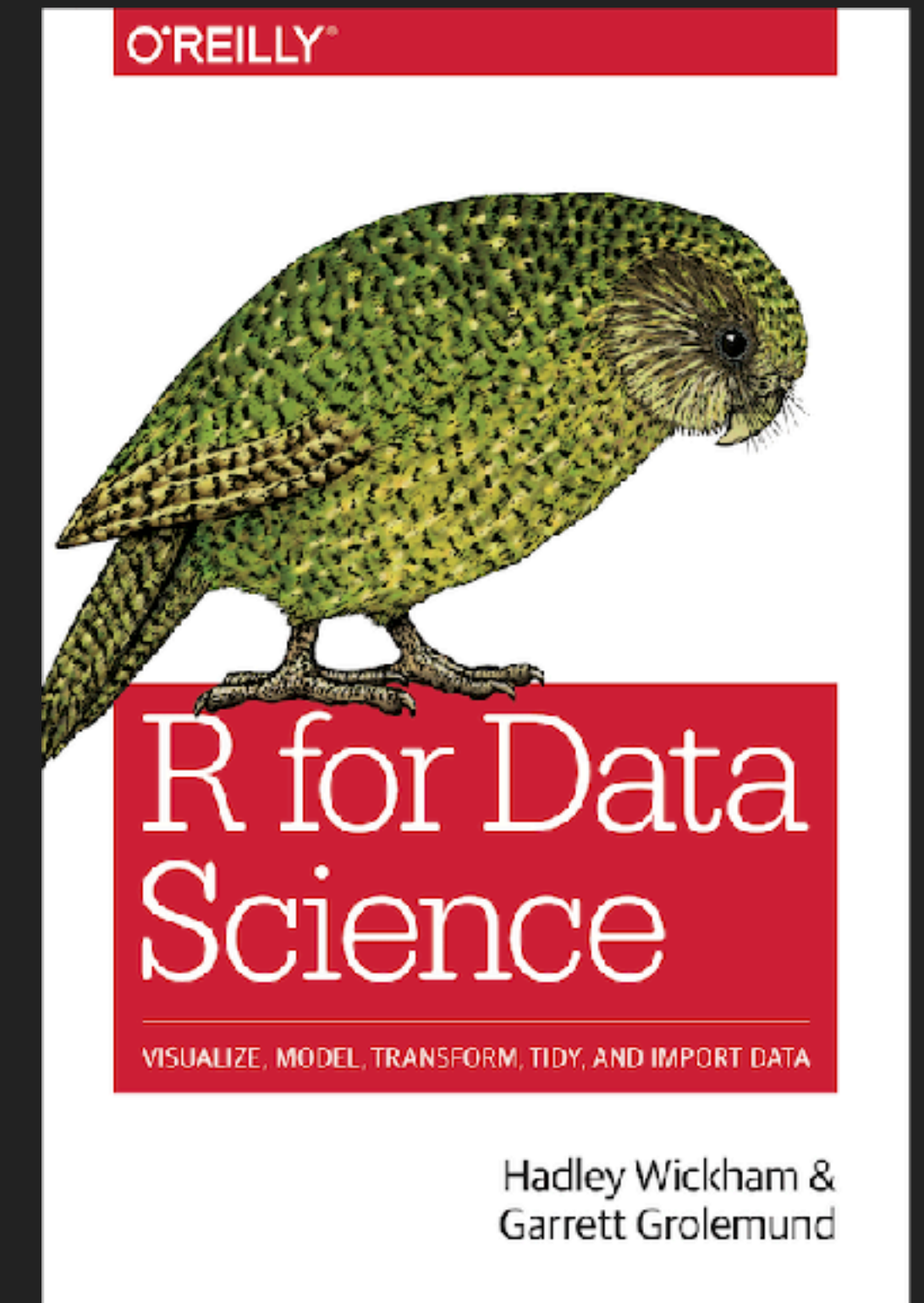
Thursday, October 27, 2022

Outline for Today

- ▶ Intro to linear regression
- ▶ Common pitfalls
- ▶ Diagnostics
- ▶ Activity



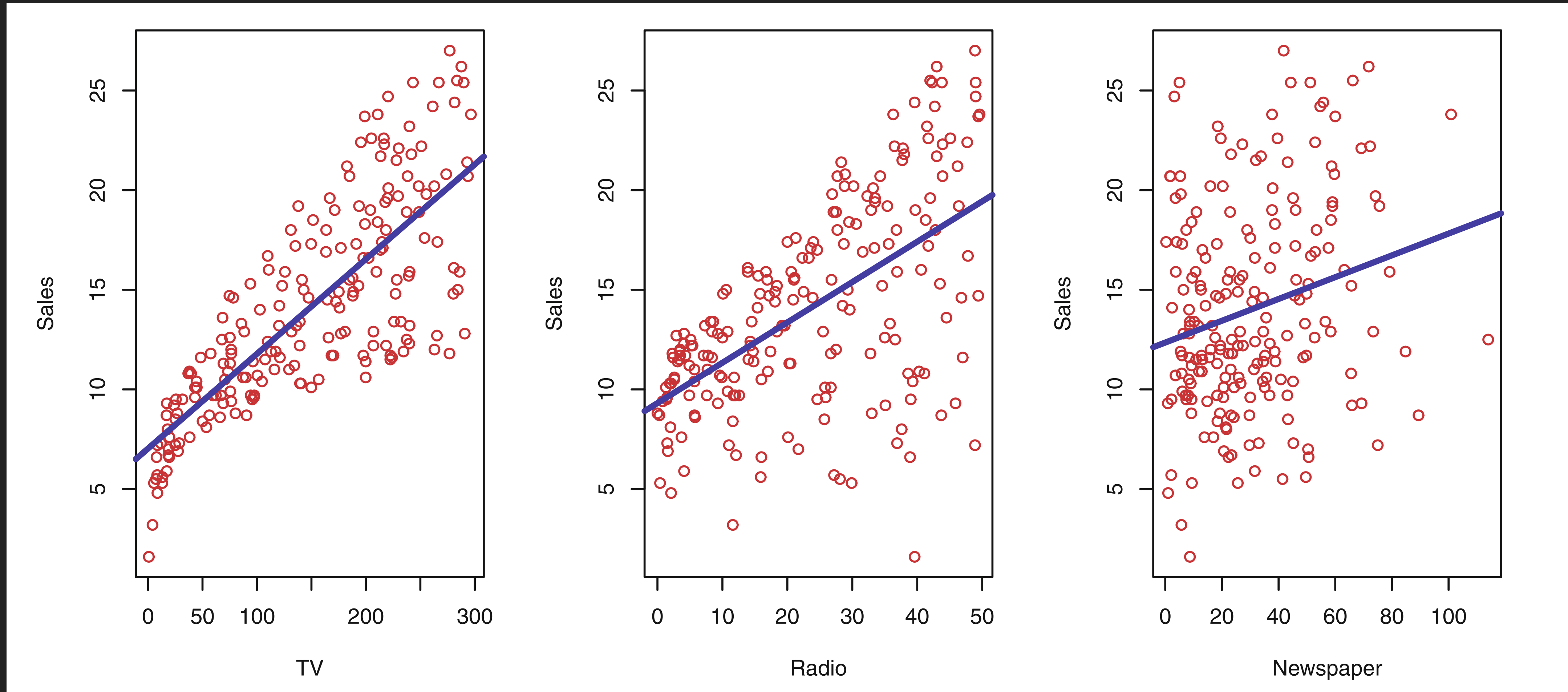
Ch 3 (Linear regression)



Ch 22-24 (Modeling)

Let's start with a case study.

Sales (in thousands of dollars) as a function of TV, radio, and newspaper advertising budgets (in thousands of dollars), for 200 cities.

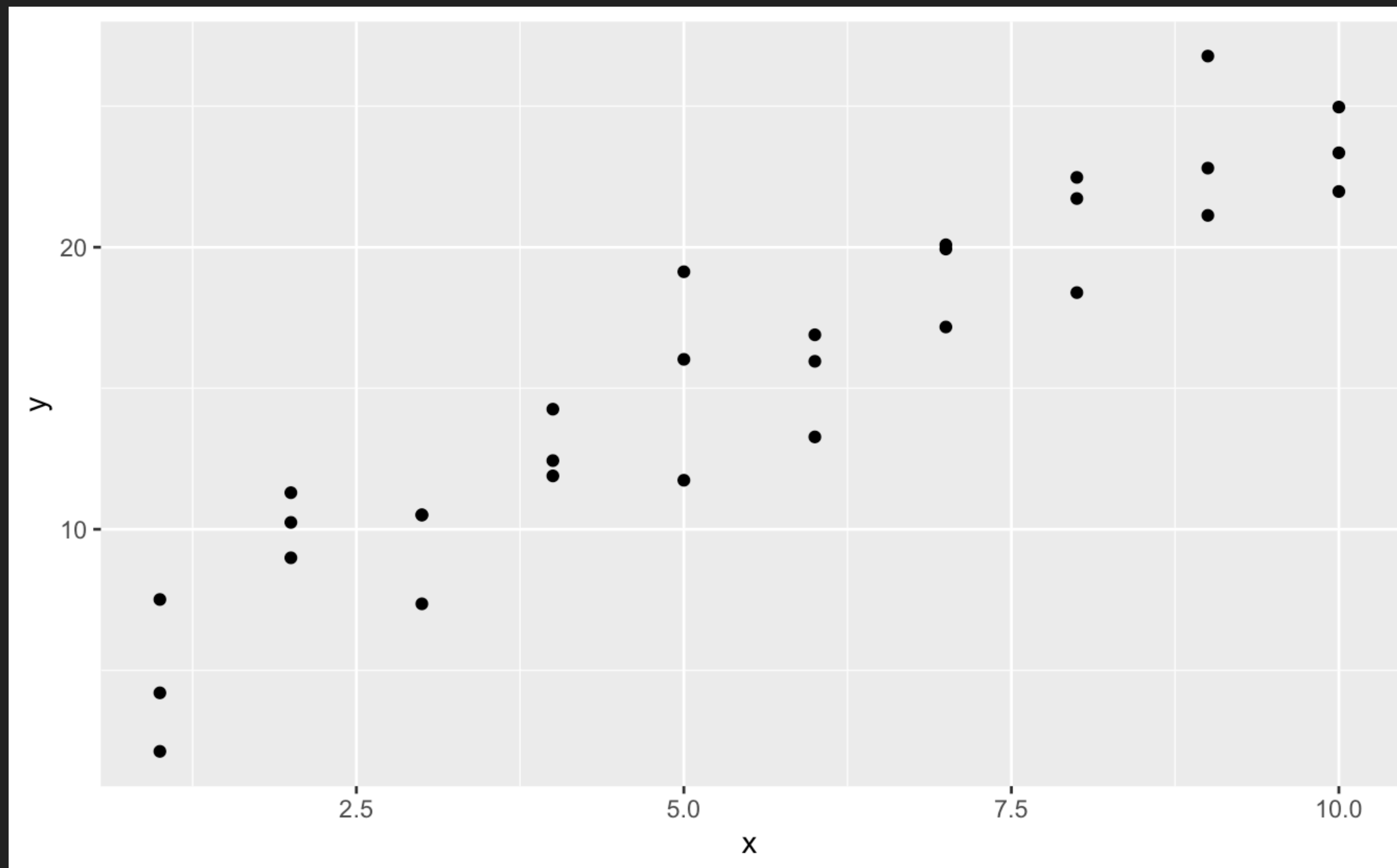


A Few Important Questions That We Might Seek To Address

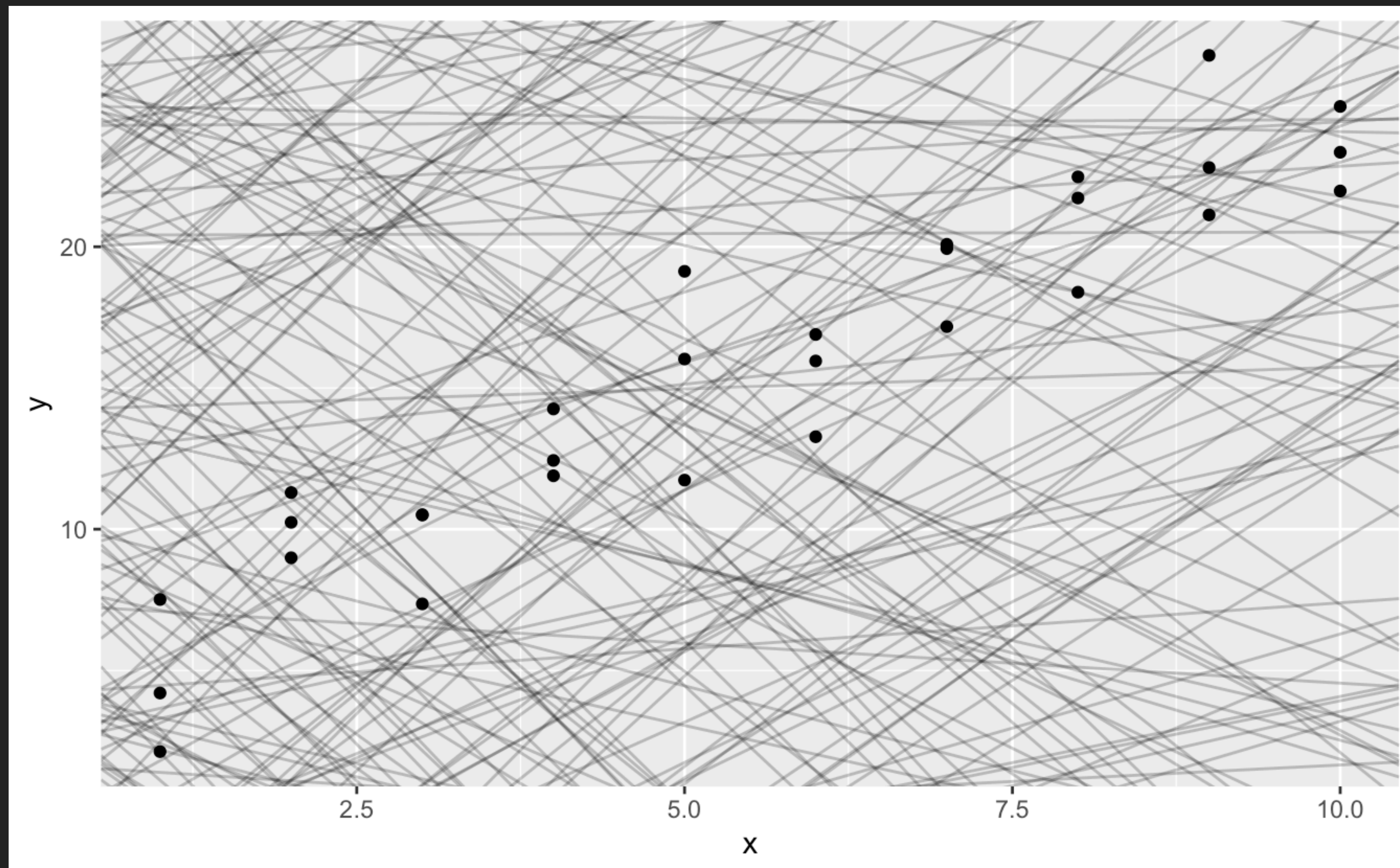
- ▶ Is there a relationship between advertising budget and sales?
- ▶ How strong is the relationship between advertising budget and sales?
- ▶ Which media contribute to sales?
- ▶ How accurately can we estimate the effect of each medium on sales?
- ▶ How accurately can we predict future sales?
- ▶ Is the relationship linear?
- ▶ Is there synergy among the advertising media?

Simple Linear Regression

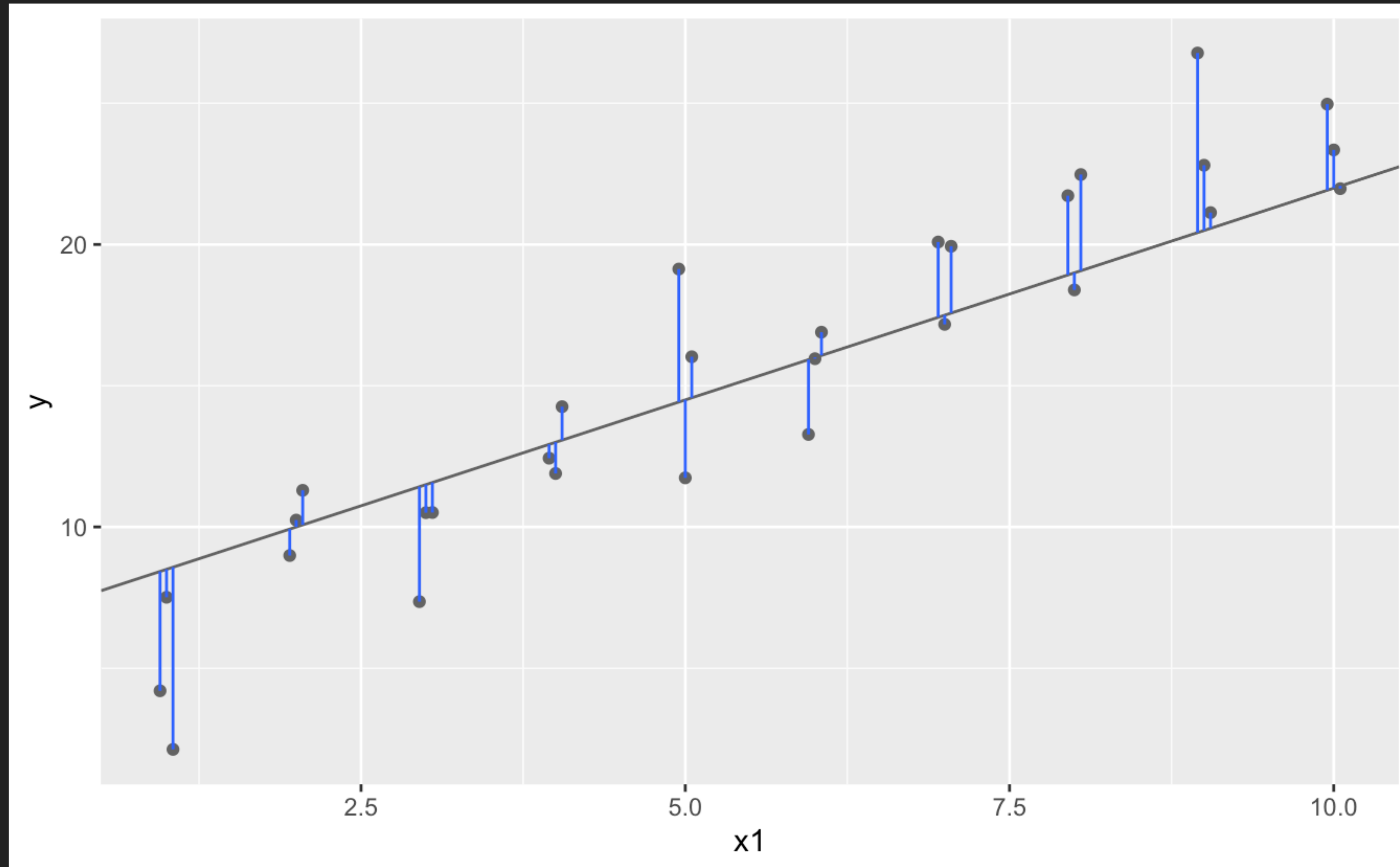
$$Y \approx \beta_0 + \beta_1 X.$$



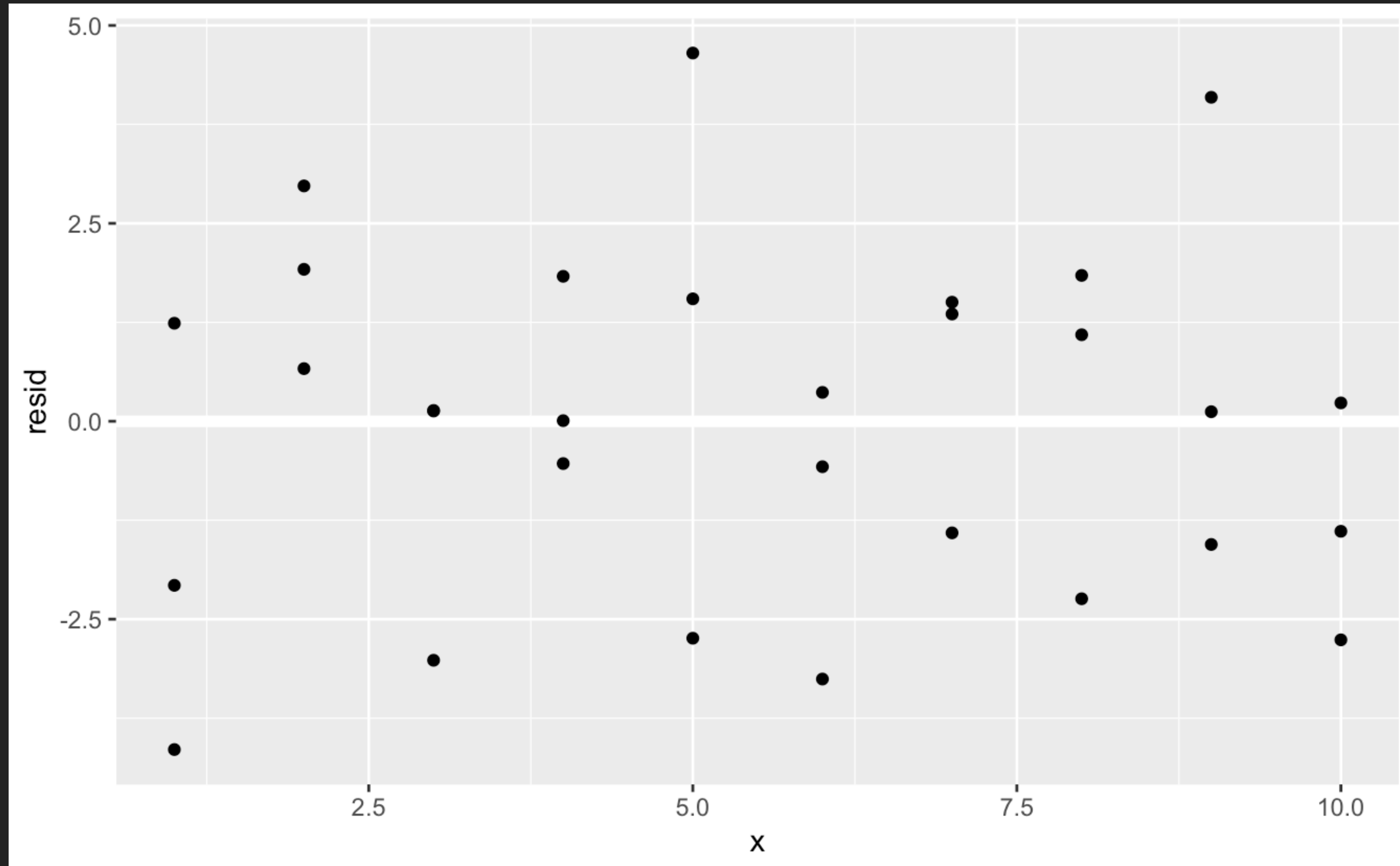
Many Possible Linear Models



Best Model? Minimize Error



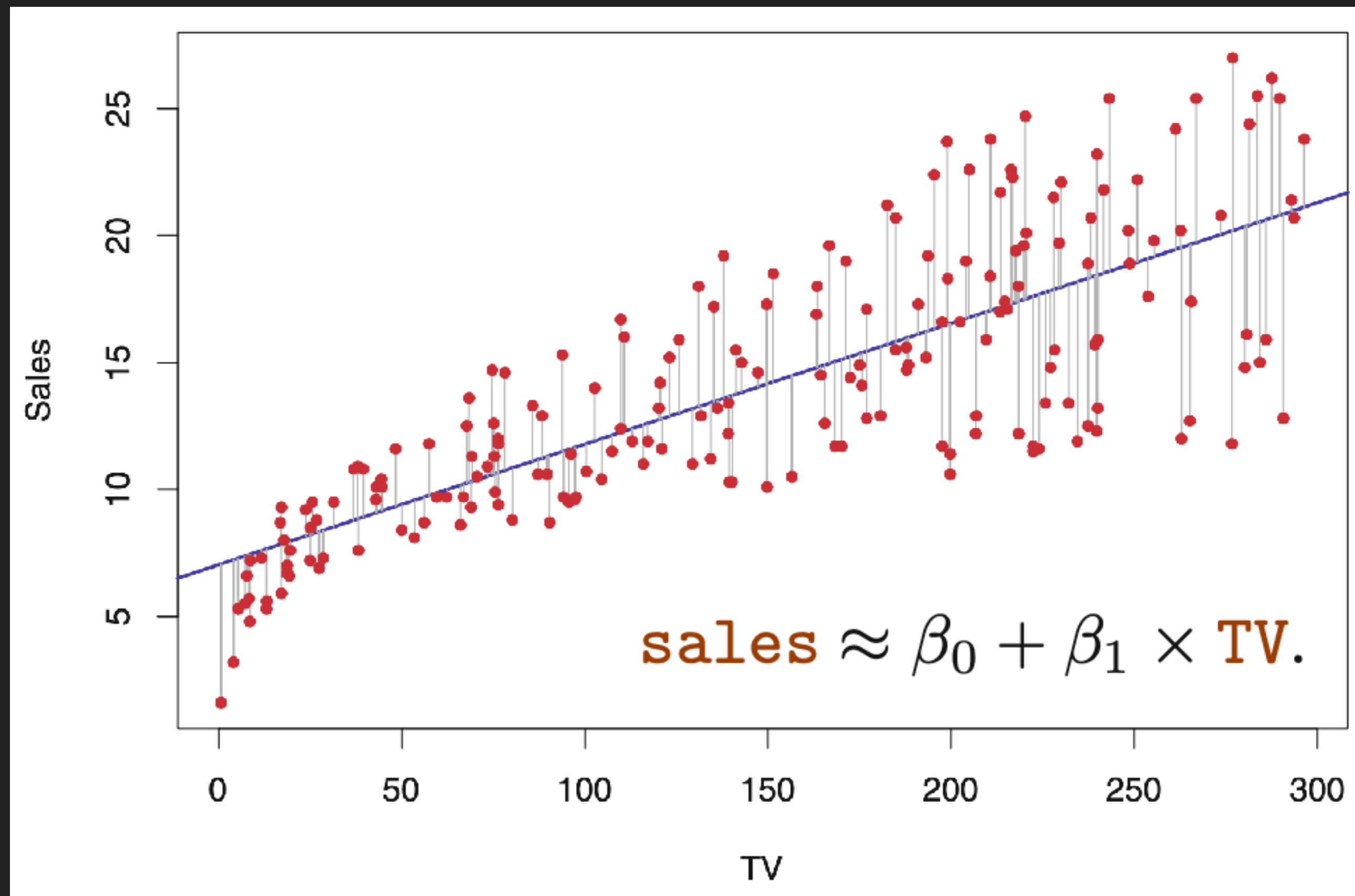
Residuals



Simple Linear Regression

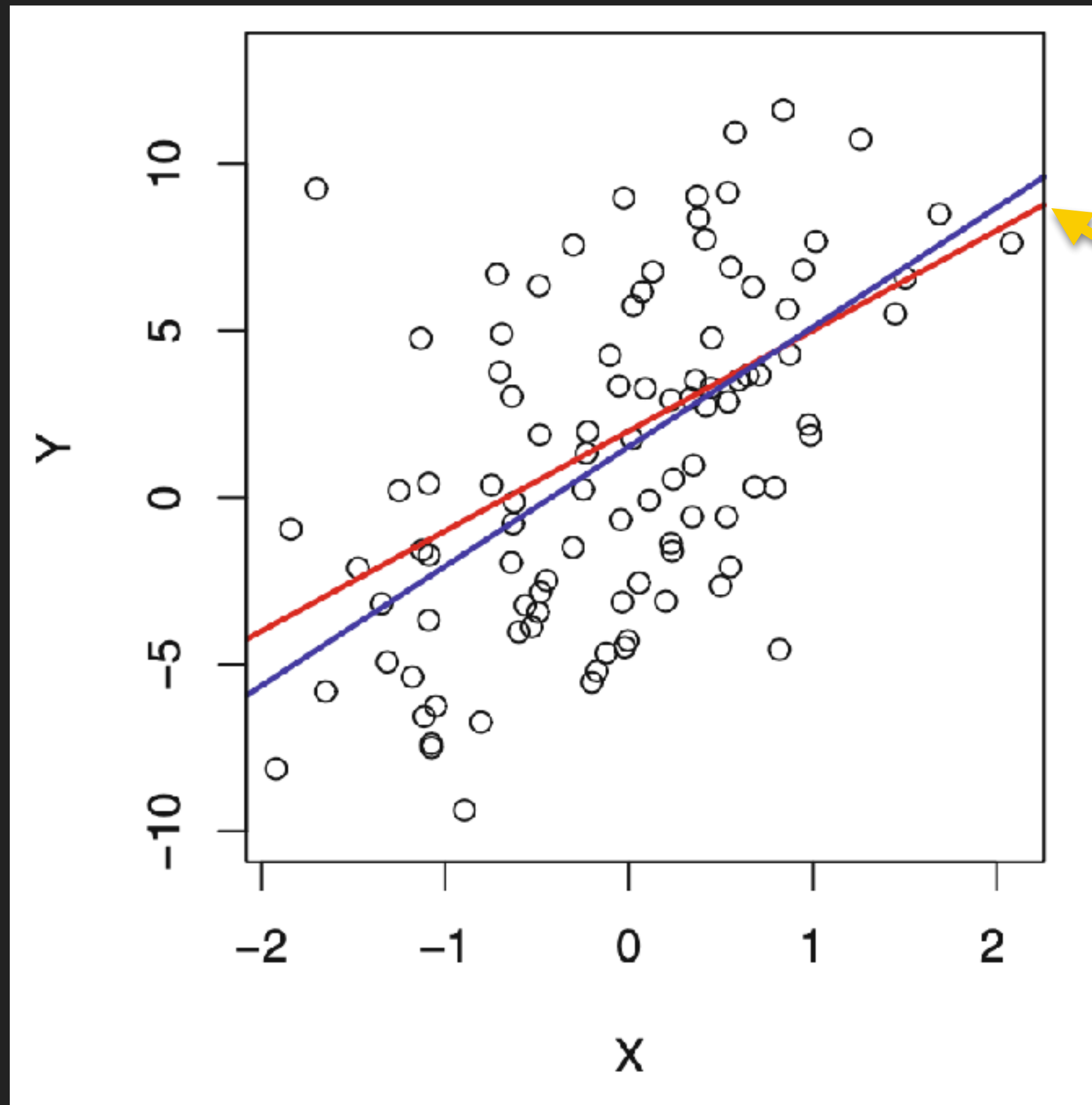
$$Y \approx \beta_0 + \beta_1 X.$$

The least squares fit for the regression of sales onto TV



- ▶ The least squares fit for the regression of sales onto TV is found by minimizing the sum of squared errors.

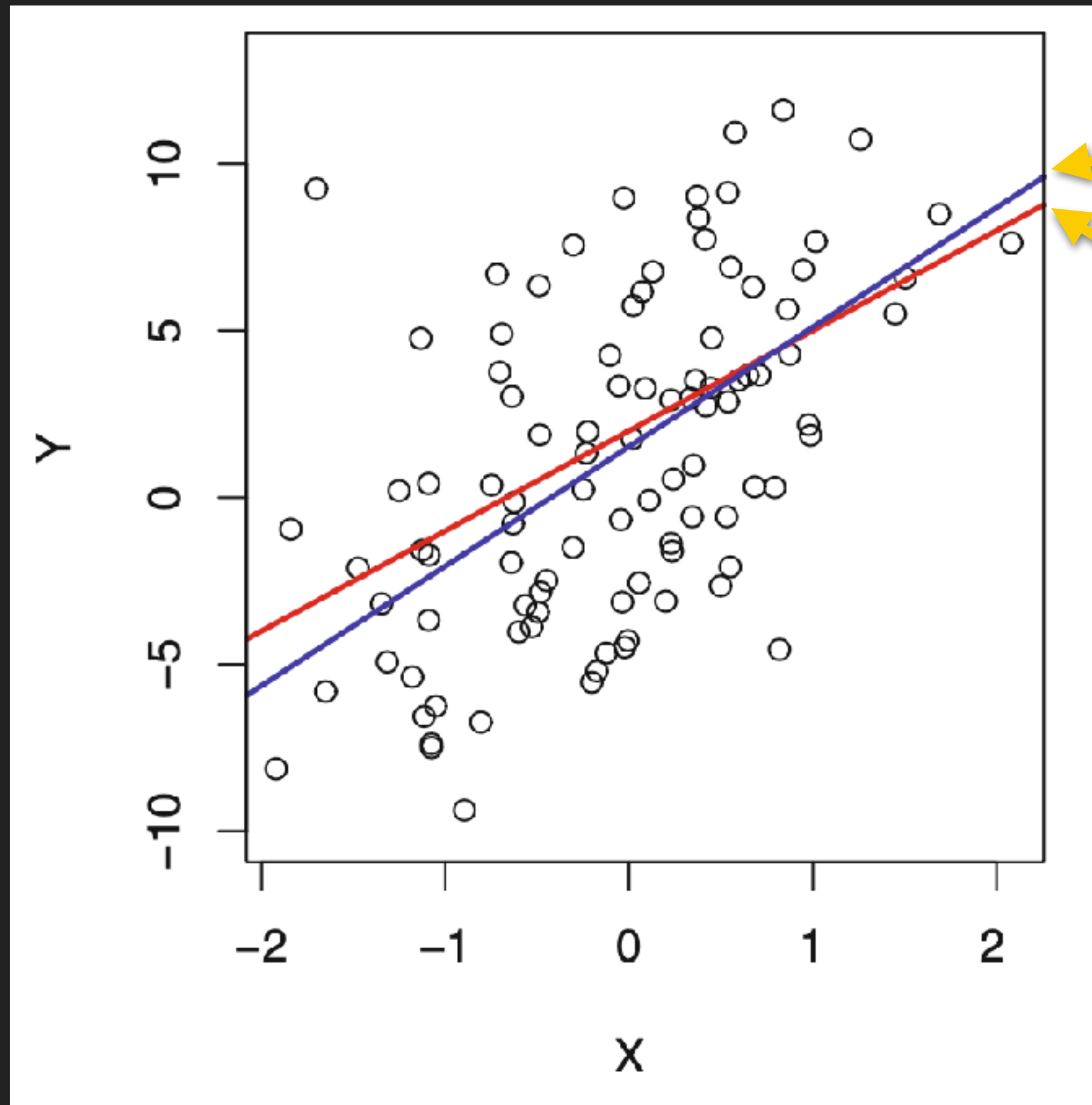
Assessing the Accuracy of the Coefficient Estimates



The true relationship:

$$f(X) = 2 + 3X$$

Assessing the Accuracy of the Coefficient Estimates

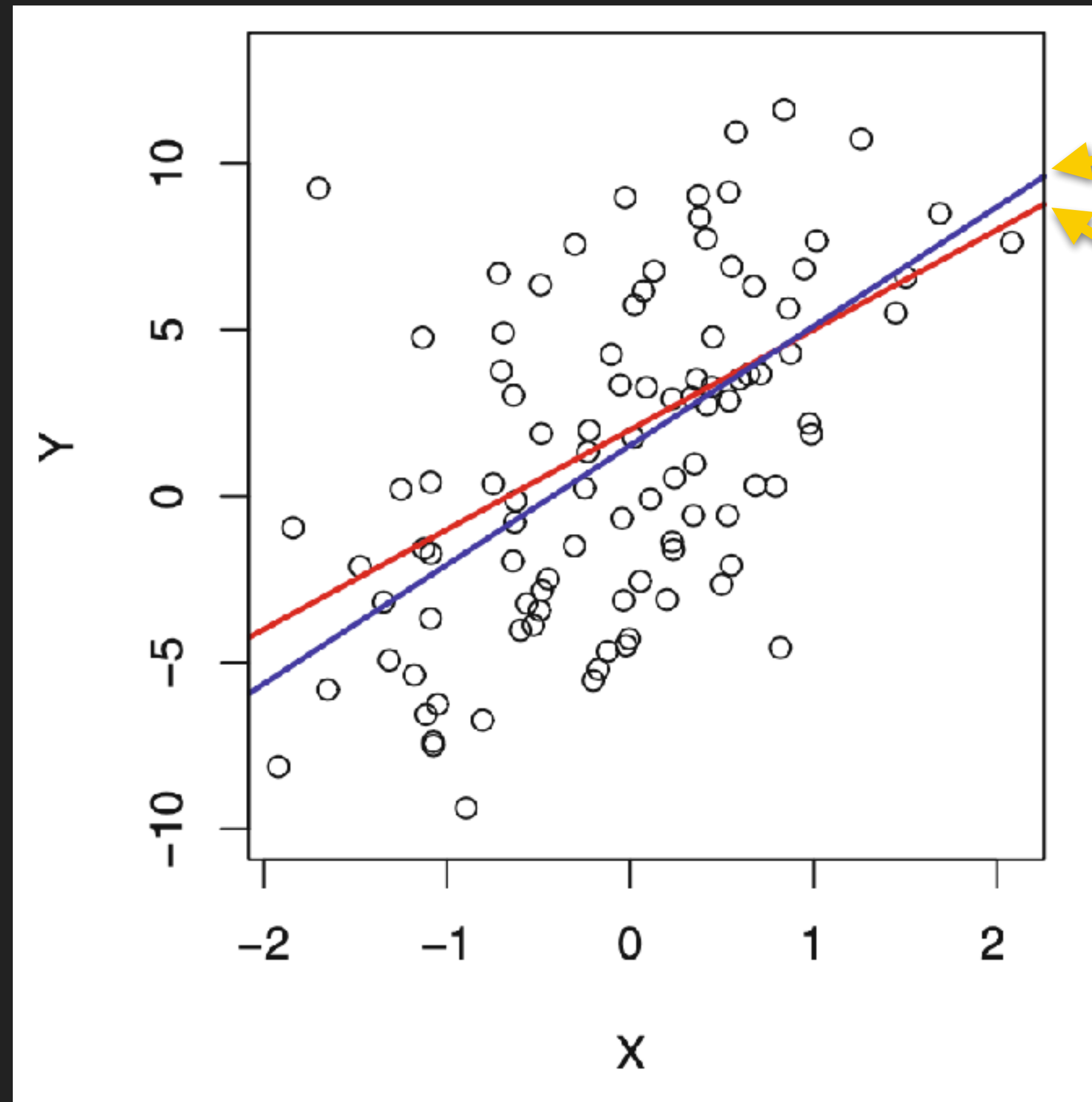


The least squares estimate for $f(X)$ based on the observed data.

The true relationship:

$$f(X) = 2 + 3X$$

Assessing the Accuracy of the Coefficient Estimates



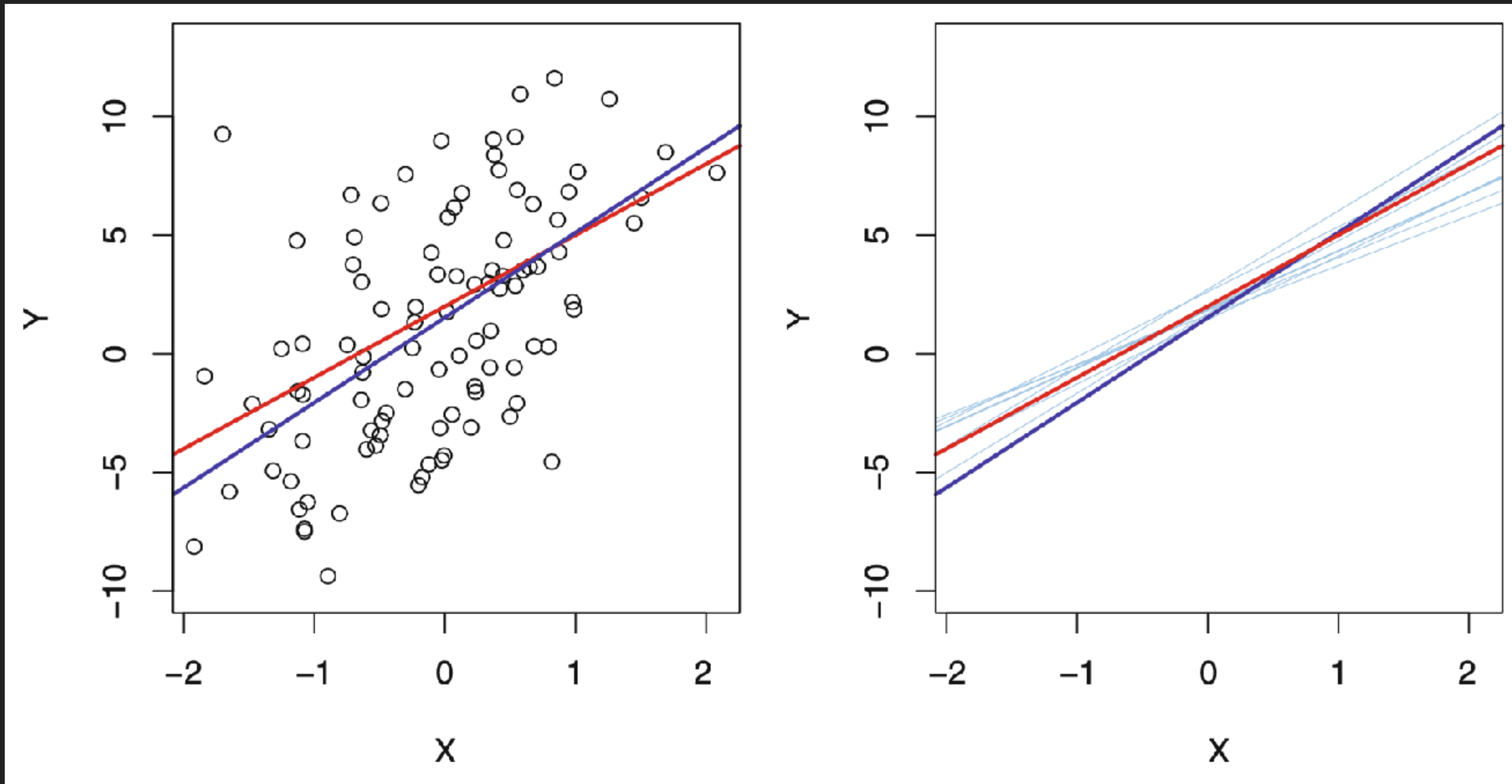
The least squares estimate for $f(X)$ based on the observed data.

The true relationship: $f(X) = 2 + 3X$

In real applications, the population regression line is unobserved.

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

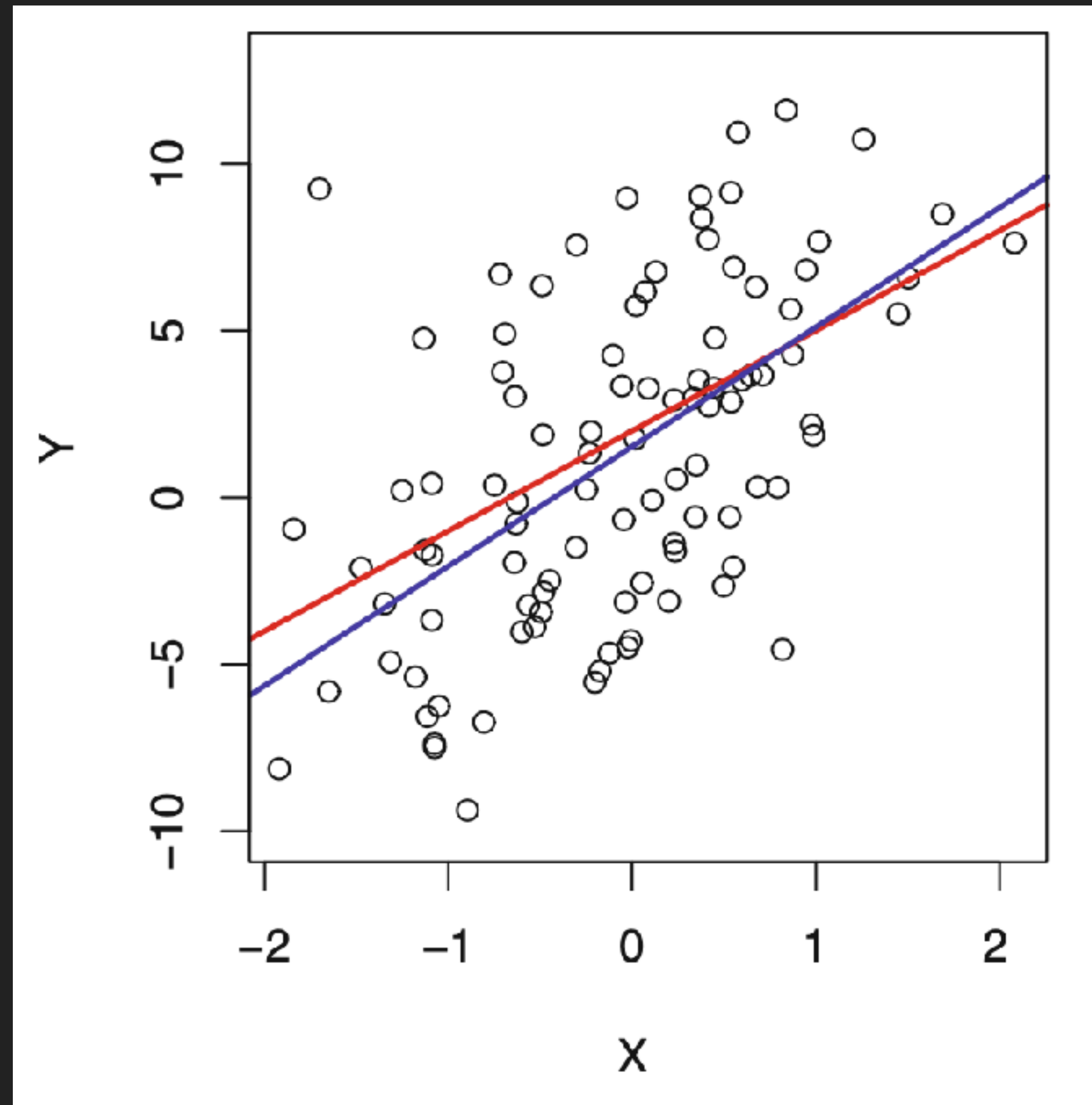
Assessing the Accuracy of the Coefficient Estimates



Ten least squares lines, each computed on the basis of a separate random set of observations.

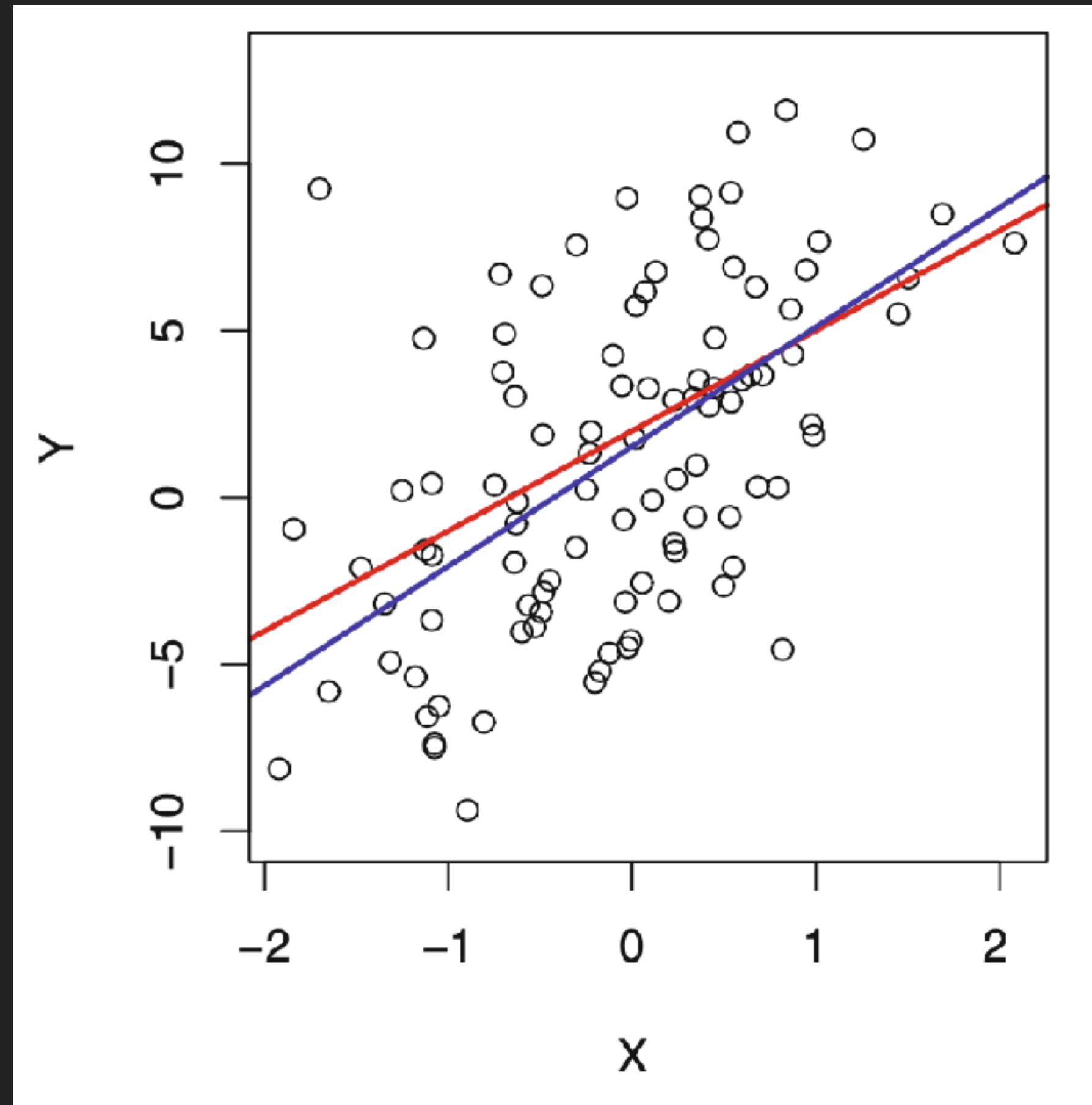
The average of many least squares lines is pretty close to the true population regression line.

Analogy with the estimation of the population mean μ of a random variable Y



- ▶ A natural question: how accurate is the sample mean $\hat{\mu}$ as an estimate of μ ?
 - ▶ Standard error
- ▶ Standard errors can be used to compute confidence intervals.
 - ▶ A 95% confidence interval is defined as a range of values such that with 95% probability, the range will contain the true unknown value of the parameter.

Analogy with the estimation of the population mean μ of a random variable Y



- ▶ For linear regression, the 95% confidence interval for β_1 approximately takes the form

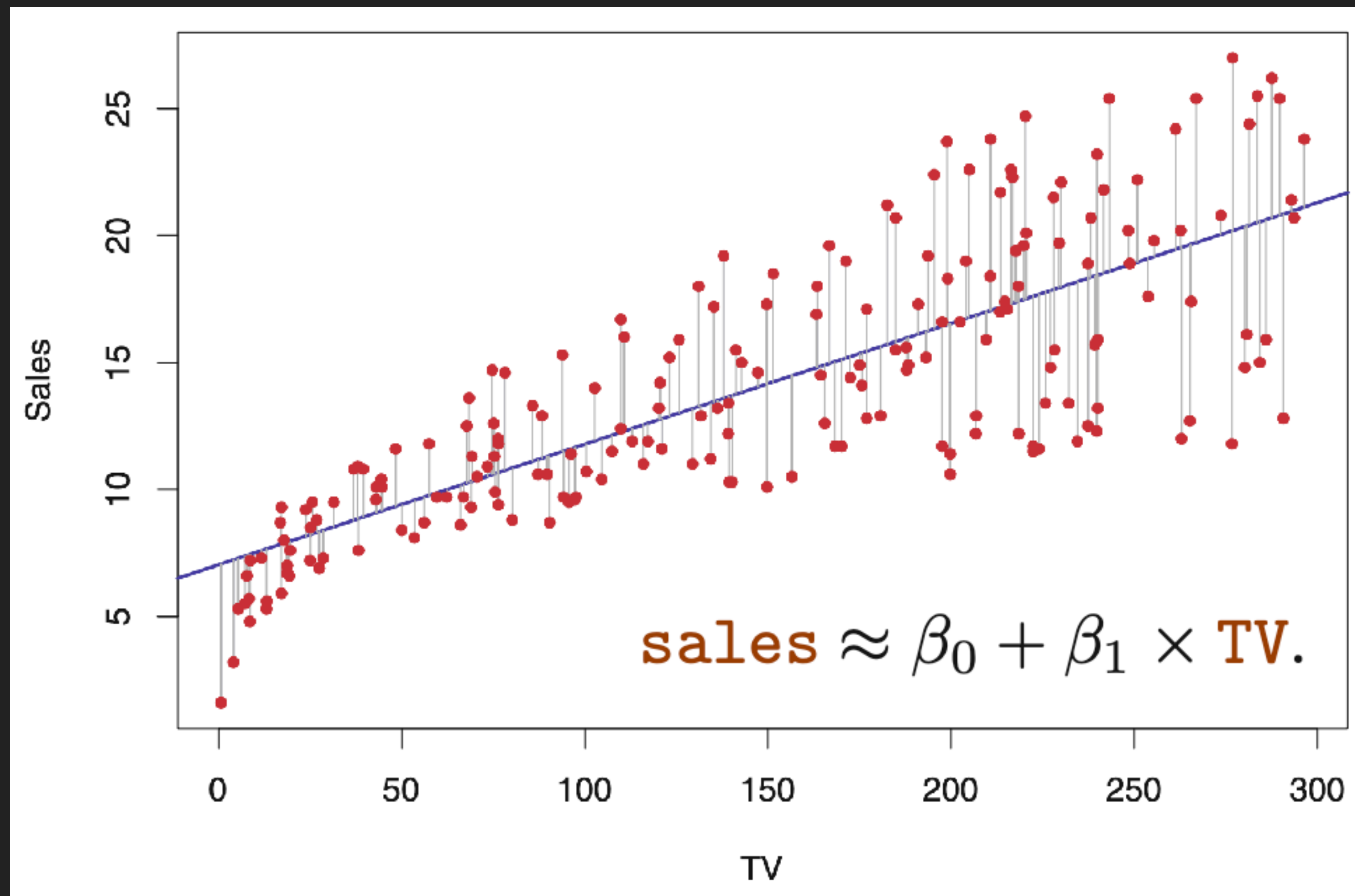
$$\hat{\beta}_1 \pm 2 \cdot \text{SE}(\hat{\beta}_1).$$

- ▶ Similarly, a confidence interval for β_0 approximately takes the form

$$\hat{\beta}_0 \pm 2 \cdot \text{SE}(\hat{\beta}_0).$$

Back to our example

The least squares fit for the regression of sales onto TV



- ▶ The 95 % CI for β_0 : [6.130, 7.935]
The 95 % CI for β_1 : [0.042, 0.053]
- ▶ In the absence of any advertising, sales will, on average, fall somewhere between 6,130 and 7,935 units.
- ▶ For each \$1,000 increase in TV advertising, there will be an average increase in sales of between 42 and 53 units.

Key idea for empirical research

Standard Errors Can Also Be Used To Perform Hypothesis Tests on the Coefficients.

- ▶ Testing the null hypothesis:
 - ▶ H_0 : There is no relationship between X and Y
- ▶ vs the alternative hypothesis
 - ▶ H_a : There is some relationship between X and Y

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

Standard Errors Can Also Be Used To Perform Hypothesis Tests on the Coefficients.

- ▶ Testing the null hypothesis:

- ▶ H_0 : There is no relationship between X and Y

- ▶ Corresponds to testing

$$H_0 : \beta_1 = 0$$

- ▶ vs the alternative hypothesis

- ▶ H_a : There is some relationship between X and Y

- ▶ vs

$$H_a : \beta_1 \neq 0,$$

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

=> Compute a t-statistic and associated p-value

Standard Errors Can Also Be Used To Perform Hypothesis Tests on the Coefficients.

- ▶ Testing the null hypothesis:

- ▶ H_0 : There is no relationship between X and Y

- ▶ Corresponds to testing

$$H_0 : \beta_1 = 0$$

- ▶ vs the alternative hypothesis

- ▶ H_a : There is some relationship between X and Y

- ▶ vs

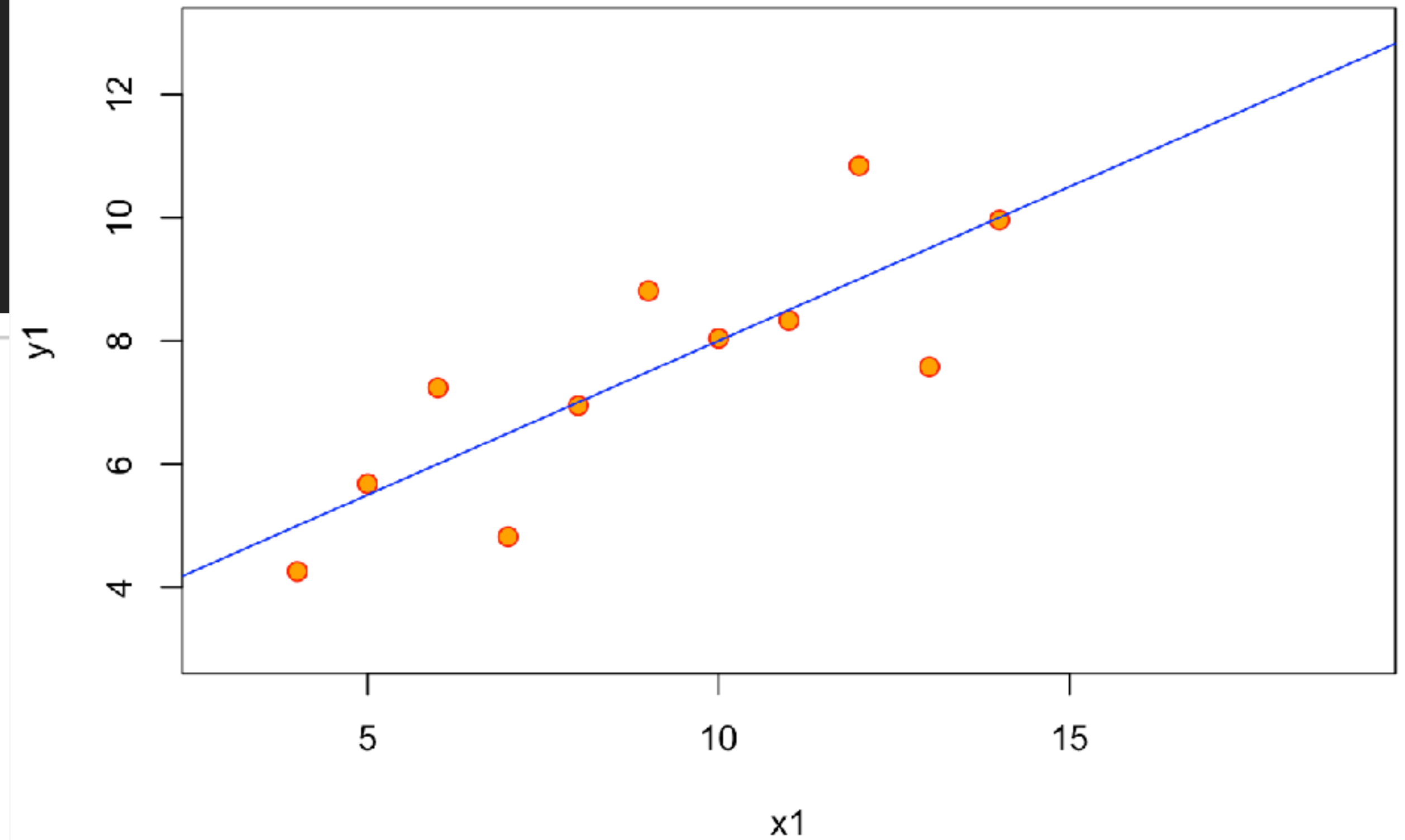
$$H_a : \beta_1 \neq 0,$$

	Coefficient	Std. error	t-statistic	p-value
Intercept	7.0325	0.4578	15.36	< 0.0001
TV	0.0475	0.0027	17.67	< 0.0001

An increase of \$1,000 in the TV advertising budget is associated with an increase in sales by around 50 units.

Another Example

```
##  
## Call:  
## lm(formula = y1 ~ x1, data = anscombe)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -1.92127 -0.45577 -0.04136  0.70941  1.83882   
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)      
## (Intercept)   3.0001     1.1247   2.667  0.02573 *      
## x1            0.5001     0.1179   4.241  0.00217 **     
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 1.237 on 9 degrees of freedom  
## Multiple R-squared:  0.6665, Adjusted R-squared:  0.6295   
## F-statistic: 17.99 on 1 and 9 DF,  p-value: 0.00217
```



Let's make it more realistic

How To Extend our Analysis To Accommodate all Predictors?

- ▶ One option is to run three separate simple linear regressions.

	Coefficient	Std. error	t-statistic	p-value
Intercept	7.0325	0.4578	15.36	< 0.0001
TV	0.0475	0.0027	17.67	< 0.0001

	Coefficient	Std. error	t-statistic	p-value
Intercept	9.312	0.563	16.54	< 0.0001
radio	0.203	0.020	9.92	< 0.0001

	Coefficient	Std. error	t-statistic	p-value
Intercept	12.351	0.621	19.88	< 0.0001
newspaper	0.055	0.017	3.30	0.00115

How To Extend our Analysis To Accommodate all Predictors?

- ▶ A better option is to give each predictor a separate slope coefficient in a single model:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p + \epsilon,$$

$$\text{sales} = \beta_0 + \beta_1 \times \text{TV} + \beta_2 \times \text{radio} + \beta_3 \times \text{newspaper} + \epsilon.$$

- ▶ We interpret β_j as the average effect on Y of a one unit increase in X_j , *holding all other predictors fixed*.

Aside: Ingredients for Establishing a Causal Relationship

The cause preceded the effect

The cause was related to the effect

We can find no plausible alternative explanation for the effect other than the cause

Back to our Advertising Example

	Coefficient	Std. error	t-statistic	p-value
Intercept	2.939	0.3119	9.42	< 0.0001
TV	0.046	0.0014	32.81	< 0.0001
radio	0.189	0.0086	21.89	< 0.0001
newspaper	-0.001	0.0059	-0.18	0.8599

	Coefficient	Std. error	t-statistic	p-value
Intercept	7.0325	0.4578	15.36	< 0.0001
TV	0.0475	0.0027	17.67	< 0.0001

	Coefficient	Std. error	t-statistic	p-value
Intercept	9.312	0.563	16.54	< 0.0001
radio	0.203	0.020	9.92	< 0.0001

	Coefficient	Std. error	t-statistic	p-value
Intercept	12.351	0.621	19.88	< 0.0001
newspaper	0.055	0.017	3.30	0.00115

Interaction Effects

- ▶ Consider the standard linear regression model with two variables

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon.$$

- ▶ According to this model, if we increase X_1 by one unit, then Y will increase by an average of β_1 units

Interaction Effects

- ▶ Extending this model with an interaction term gives:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \epsilon.$$

Interaction Effects

- ▶ Extending this model with an interaction term gives:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \epsilon.$$

$$= \beta_0 + (\beta_1 + \beta_3 X_2) X_1 + \beta_2 X_2 + \epsilon$$

$$= \beta_0 + \tilde{\beta}_1 X_1 + \beta_2 X_2 + \epsilon$$

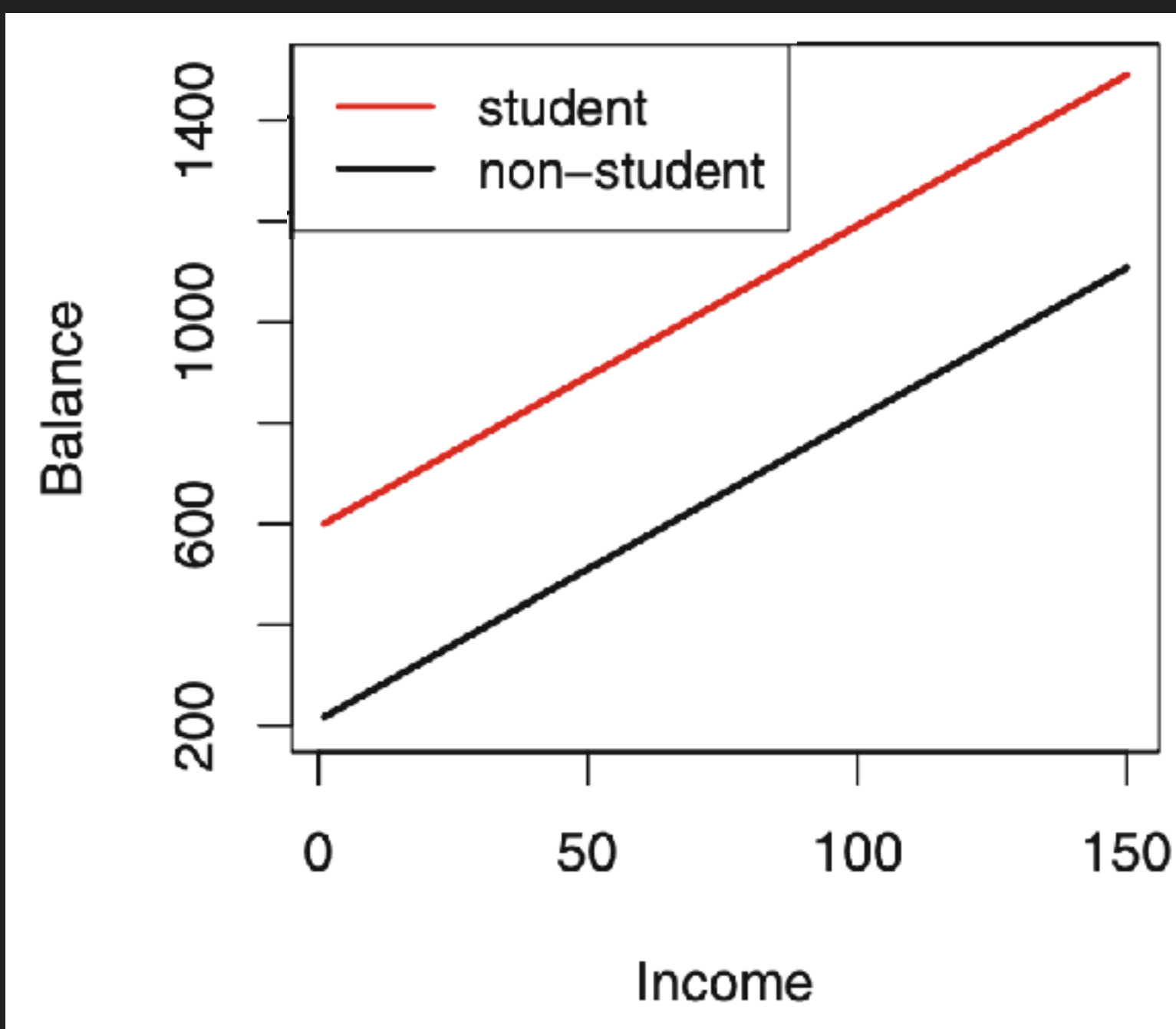
- ▶ According to this model, adjusting X_2 will change the impact of X_1 on Y

Example: Model Credit Card Balance Using Income (Numerical) and Student (Categorical)

$$\begin{aligned} \text{balance}_i &\approx \beta_0 + \beta_1 \times \text{income}_i + \begin{cases} \beta_2 & \text{if } i\text{th person is a student} \\ 0 & \text{if } i\text{th person is not a student} \end{cases} \\ &= \beta_1 \times \text{income}_i + \begin{cases} \beta_0 + \beta_2 & \text{if } i\text{th person is a student} \\ \beta_0 & \text{if } i\text{th person is not a student.} \end{cases} \end{aligned}$$

Example: Model Credit Card Balance Using Income (Numerical) and Student (Categorical)

$$\begin{aligned} \text{balance}_i &\approx \beta_0 + \beta_1 \times \text{income}_i + \begin{cases} \beta_2 & \text{if } i\text{th person is a student} \\ 0 & \text{if } i\text{th person is not a student} \end{cases} \\ &= \beta_1 \times \text{income}_i + \begin{cases} \beta_0 + \beta_2 & \text{if } i\text{th person is a student} \\ \beta_0 & \text{if } i\text{th person is not a student.} \end{cases} \end{aligned}$$



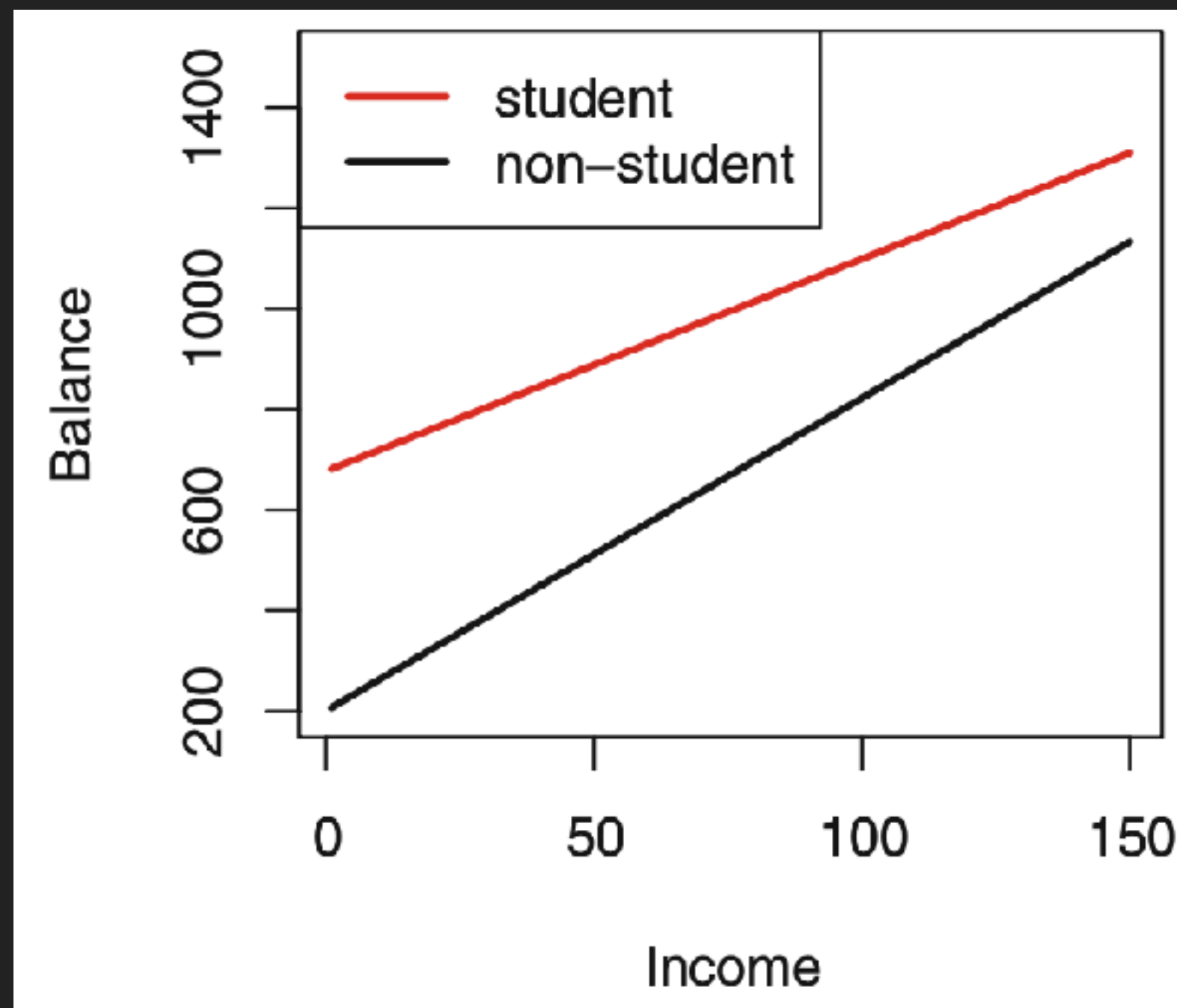
- ▶ Without an interaction term: fitting two parallel lines to the data, one for students and one for non-students.
- ▶ The lines for students and non-students have **different intercepts**, $\beta_0 + \beta_2$ versus β_0 , but the **same slope**, β_1 .

Example: Model Credit Card Balance Using Income (Numerical) and Student (Categorical)

$$\begin{aligned} \text{balance}_i &\approx \beta_0 + \beta_1 \times \text{income}_i + \begin{cases} \beta_2 + \beta_3 \times \text{income}_i & \text{if student} \\ 0 & \text{if not student} \end{cases} \\ &= \begin{cases} (\beta_0 + \beta_2) + (\beta_1 + \beta_3) \times \text{income}_i & \text{if student} \\ \beta_0 + \beta_1 \times \text{income}_i & \text{if not student} \end{cases} \end{aligned}$$

Example: Model Credit Card Balance Using Income (Numerical) and Student (Categorical)

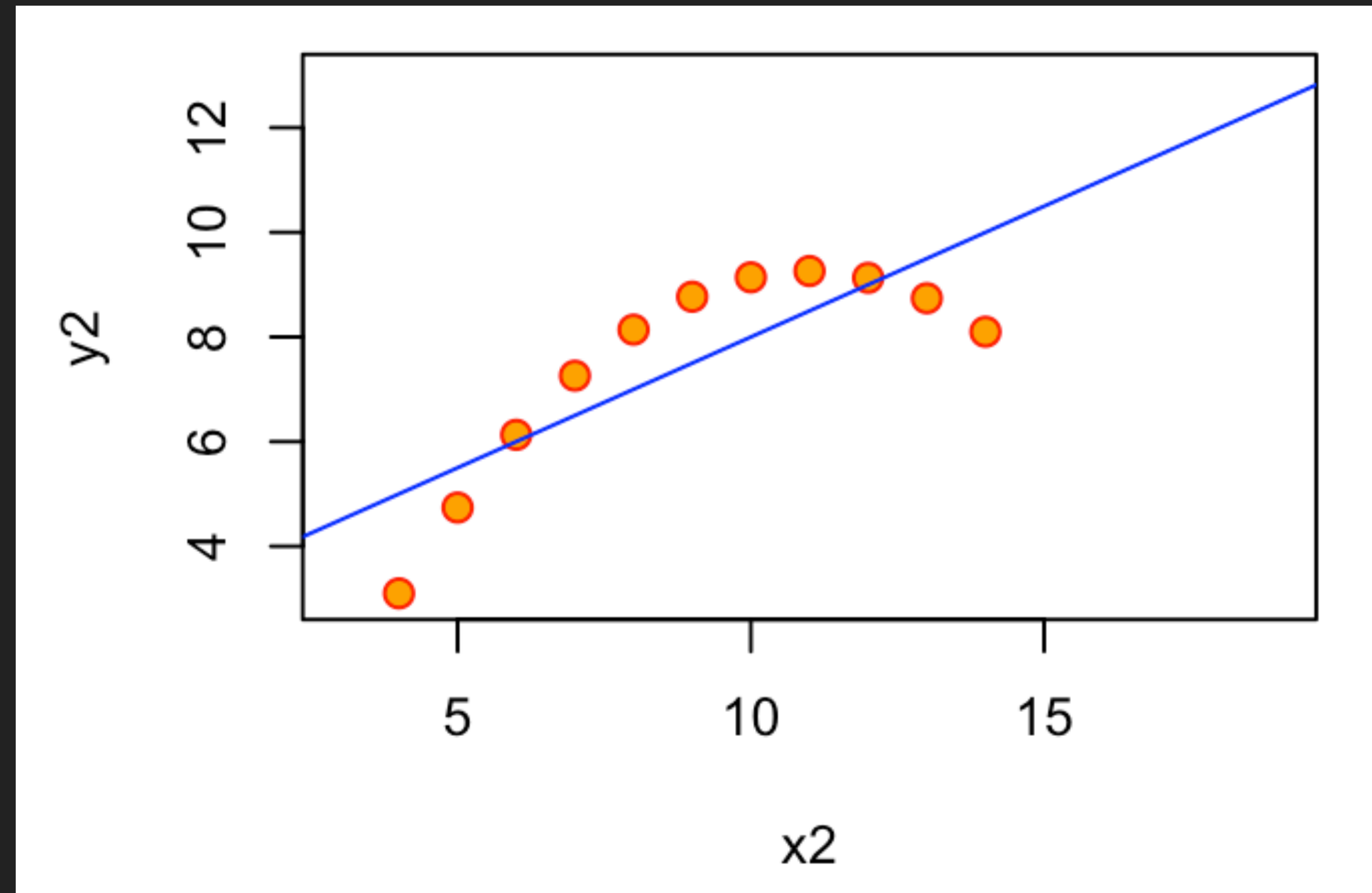
$$\begin{aligned} \text{balance}_i &\approx \beta_0 + \beta_1 \times \text{income}_i + \begin{cases} \beta_2 + \beta_3 \times \text{income}_i & \text{if student} \\ 0 & \text{if not student} \end{cases} \\ &= \begin{cases} (\beta_0 + \beta_2) + (\beta_1 + \beta_3) \times \text{income}_i & \text{if student} \\ \beta_0 + \beta_1 \times \text{income}_i & \text{if not student} \end{cases} \end{aligned}$$



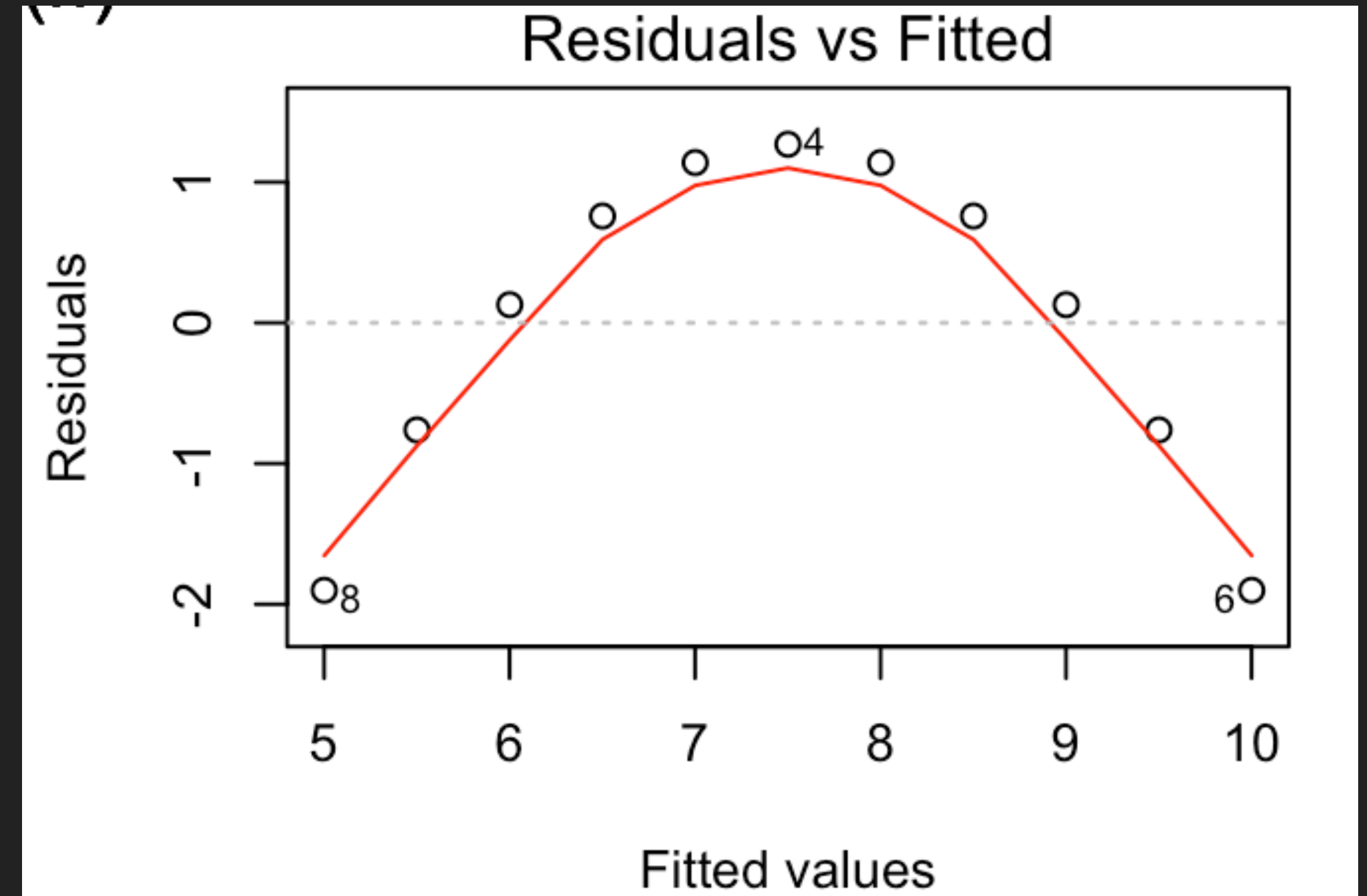
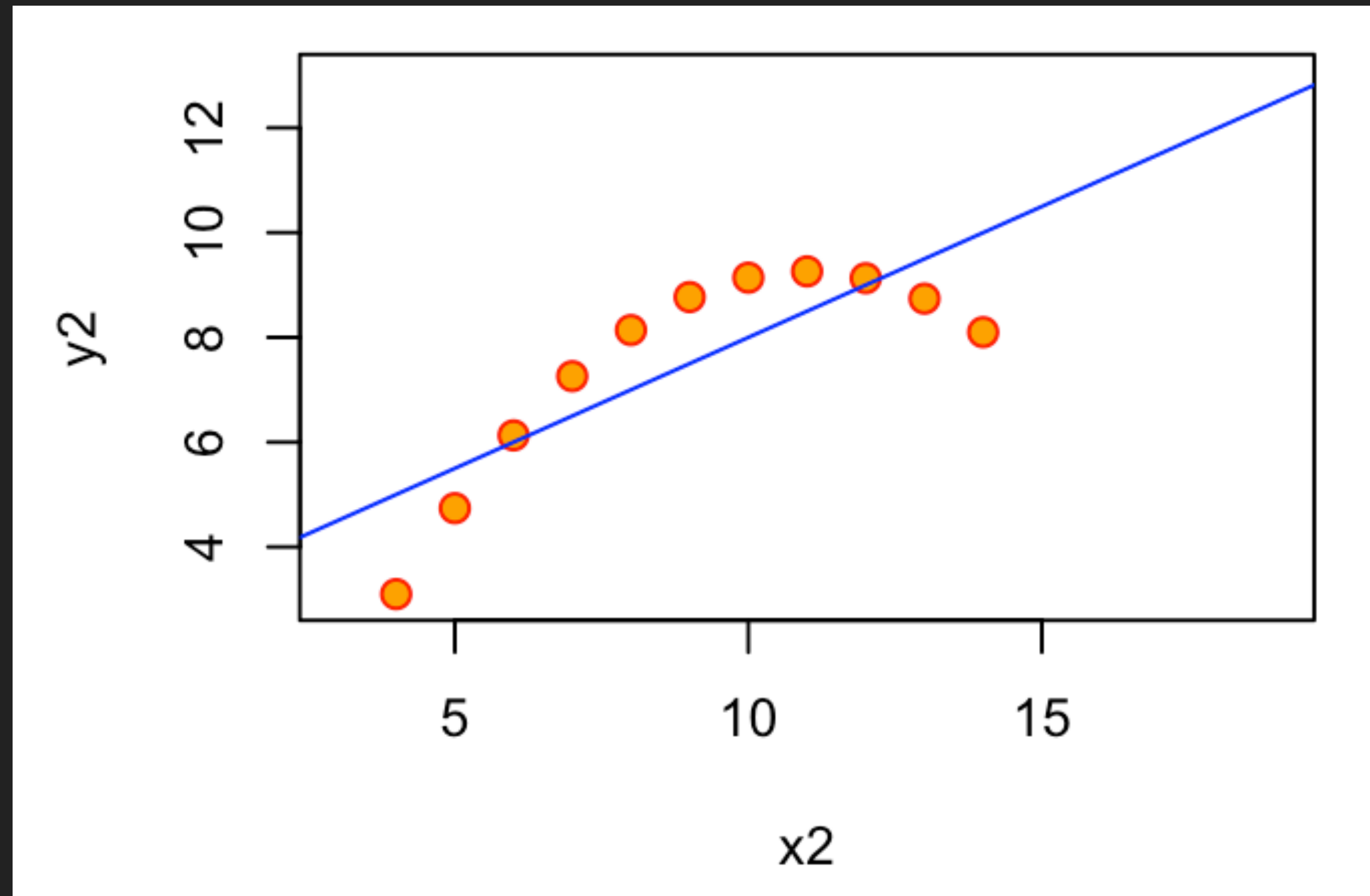
- ▶ With an interaction term: the regression lines for the students and the non-students have **different intercepts**, $\beta_0 + \beta_2$ versus β_0 , and **different slopes**, $\beta_1 + \beta_3$ versus β_1 .
- ▶ Allows for the possibility that changes in income may affect the credit card balances of students and non-students differently.

It's complicated.

Potential Problem: Non-Linearity of the Data

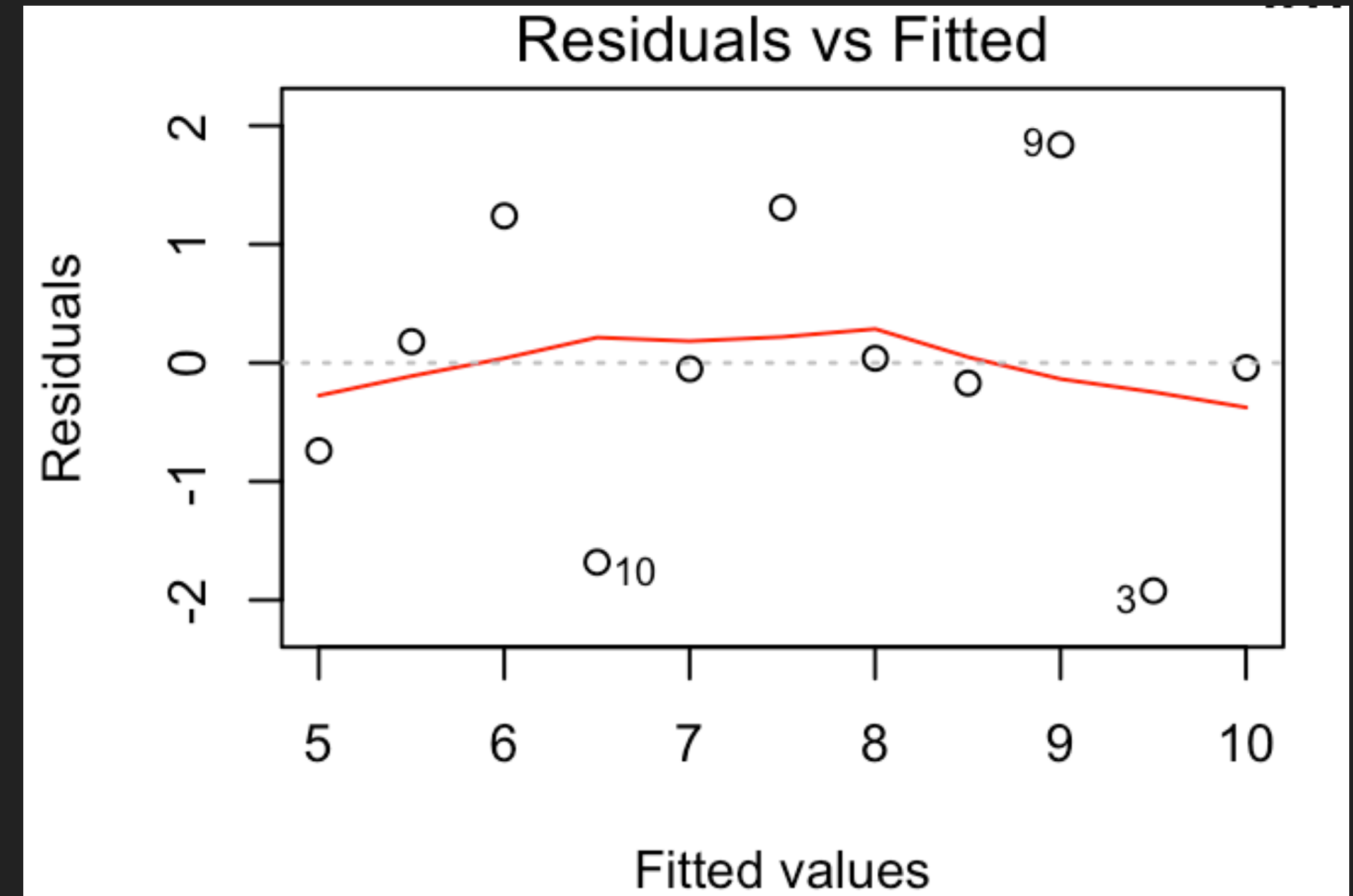
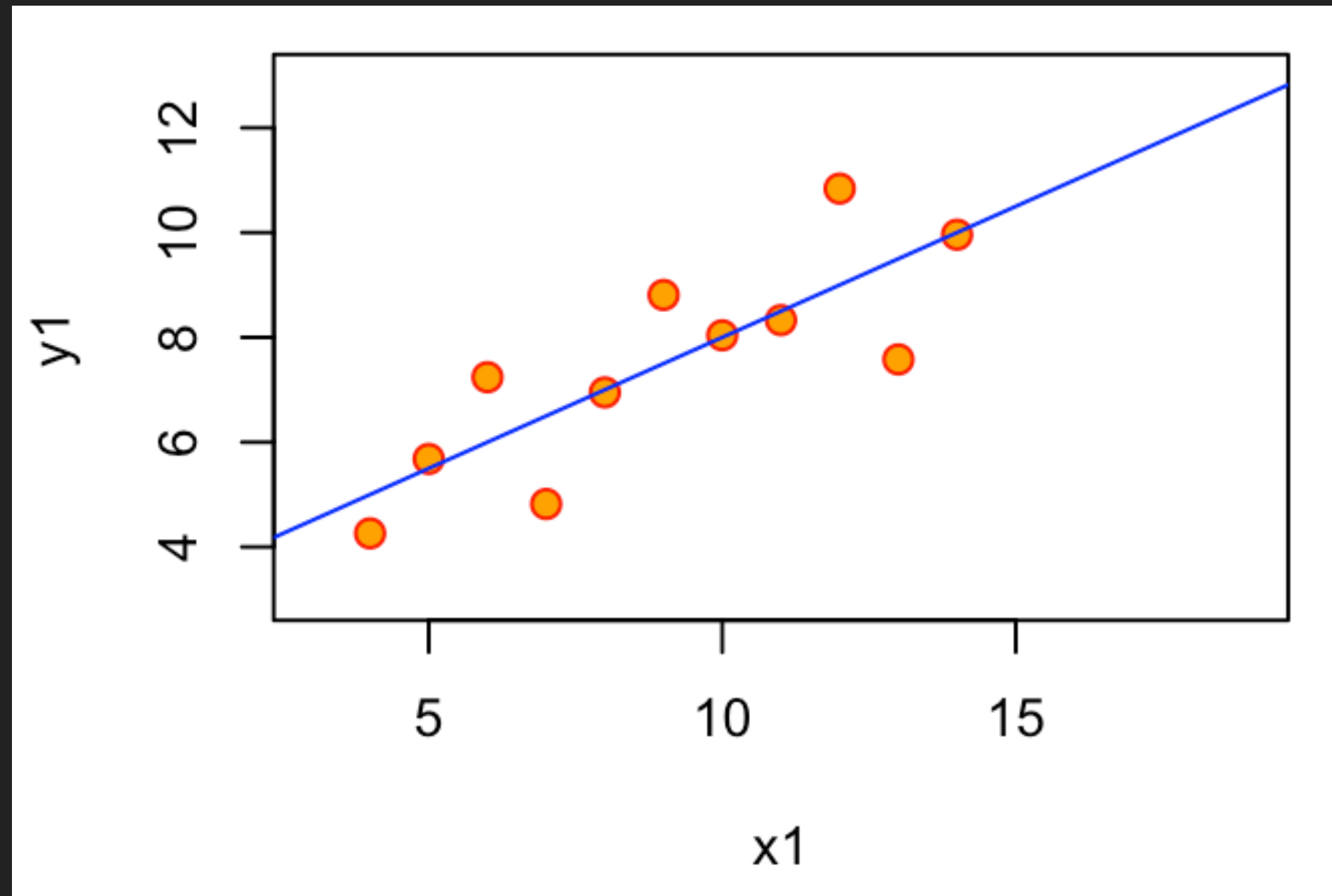


Potential Problem: Non-Linearity of the Data

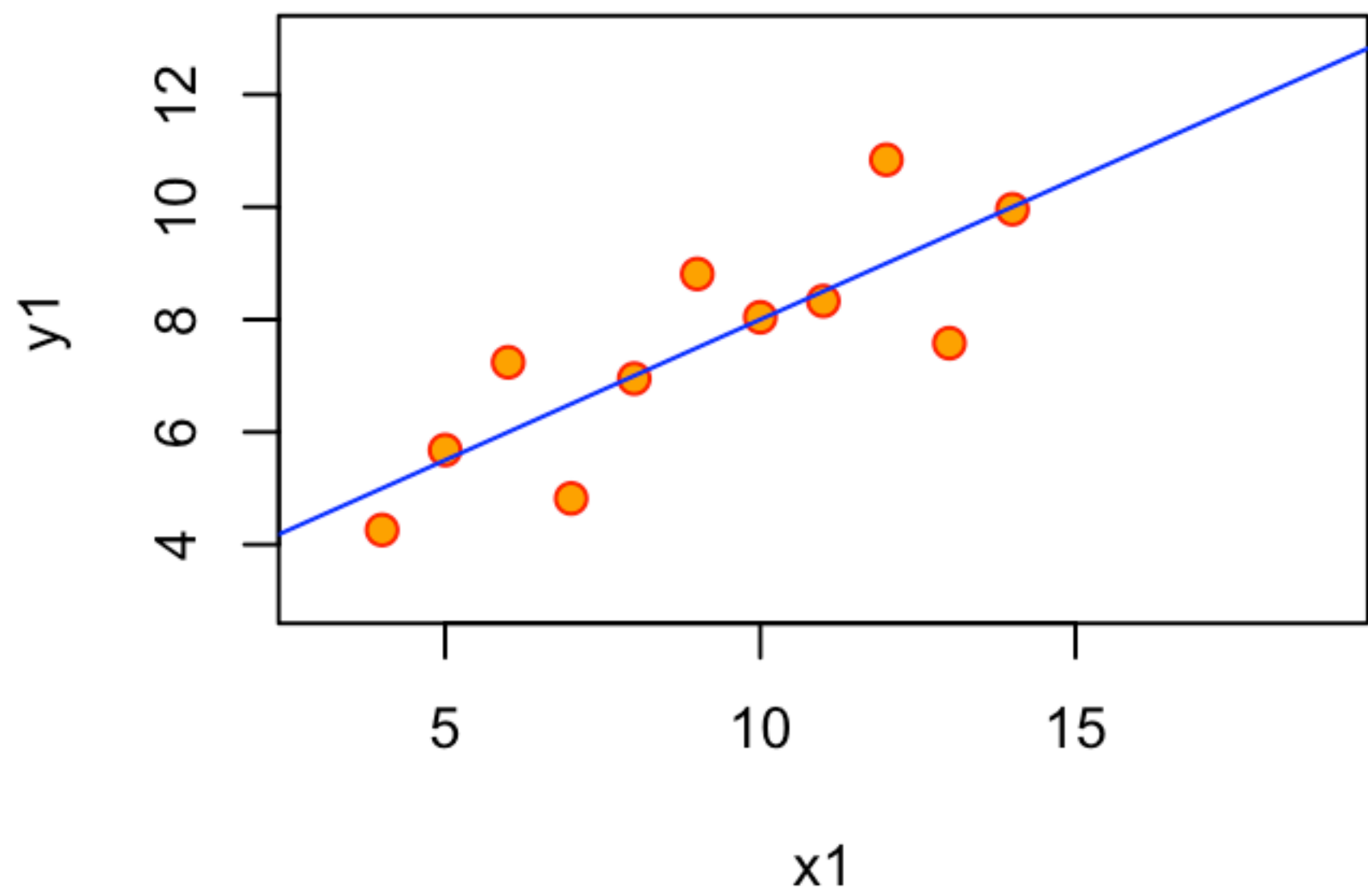


- ▶ Ideally, the residual plot will show no discernible pattern.
 - ▶ Otherwise, indicates nonlinear relationship in the data.

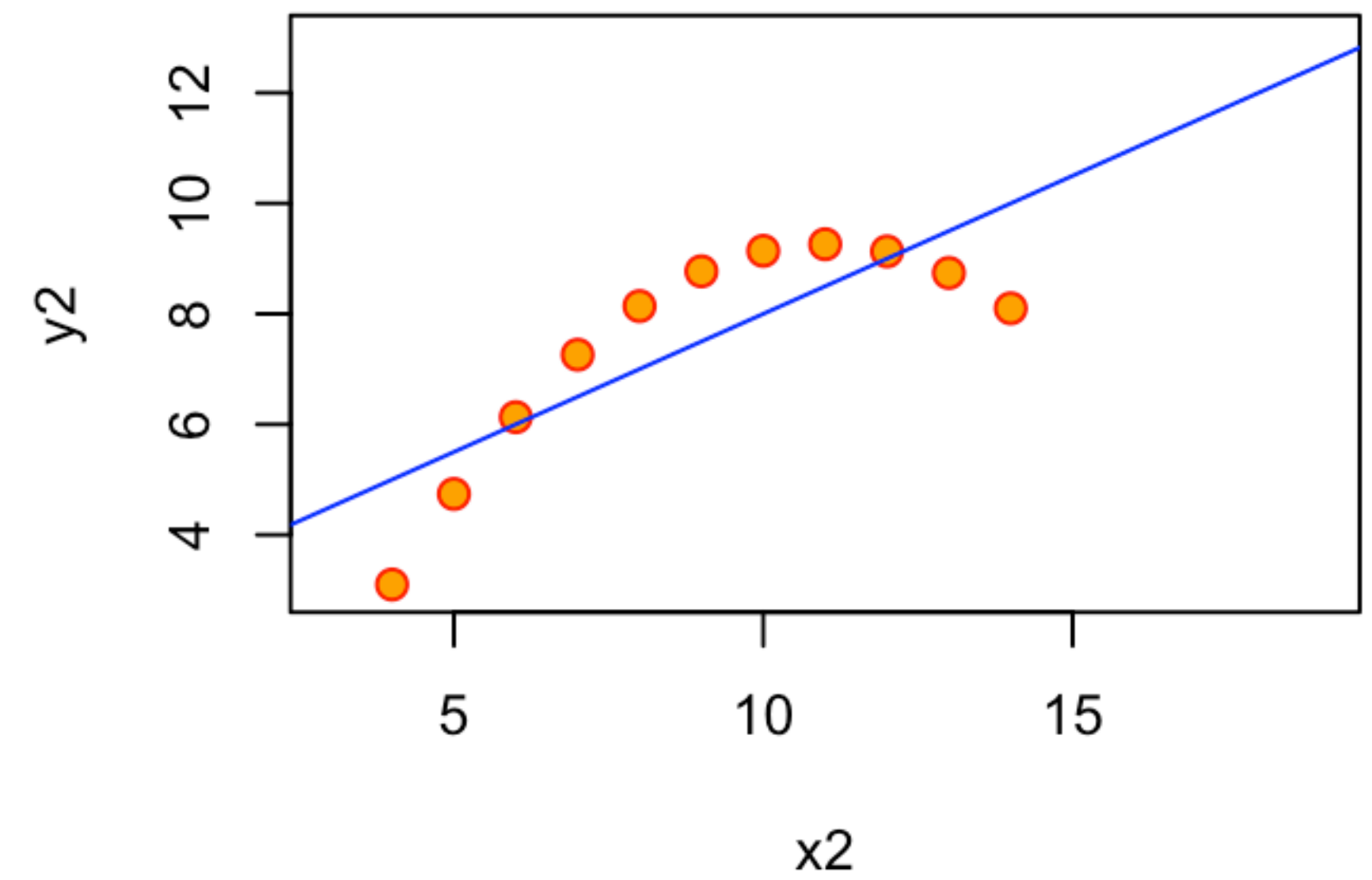
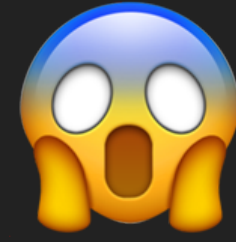
Potential Problem: Non-Linearity of the Data



- ▶ Contrast the example on the previous slide to the one we had earlier.



Ouch!

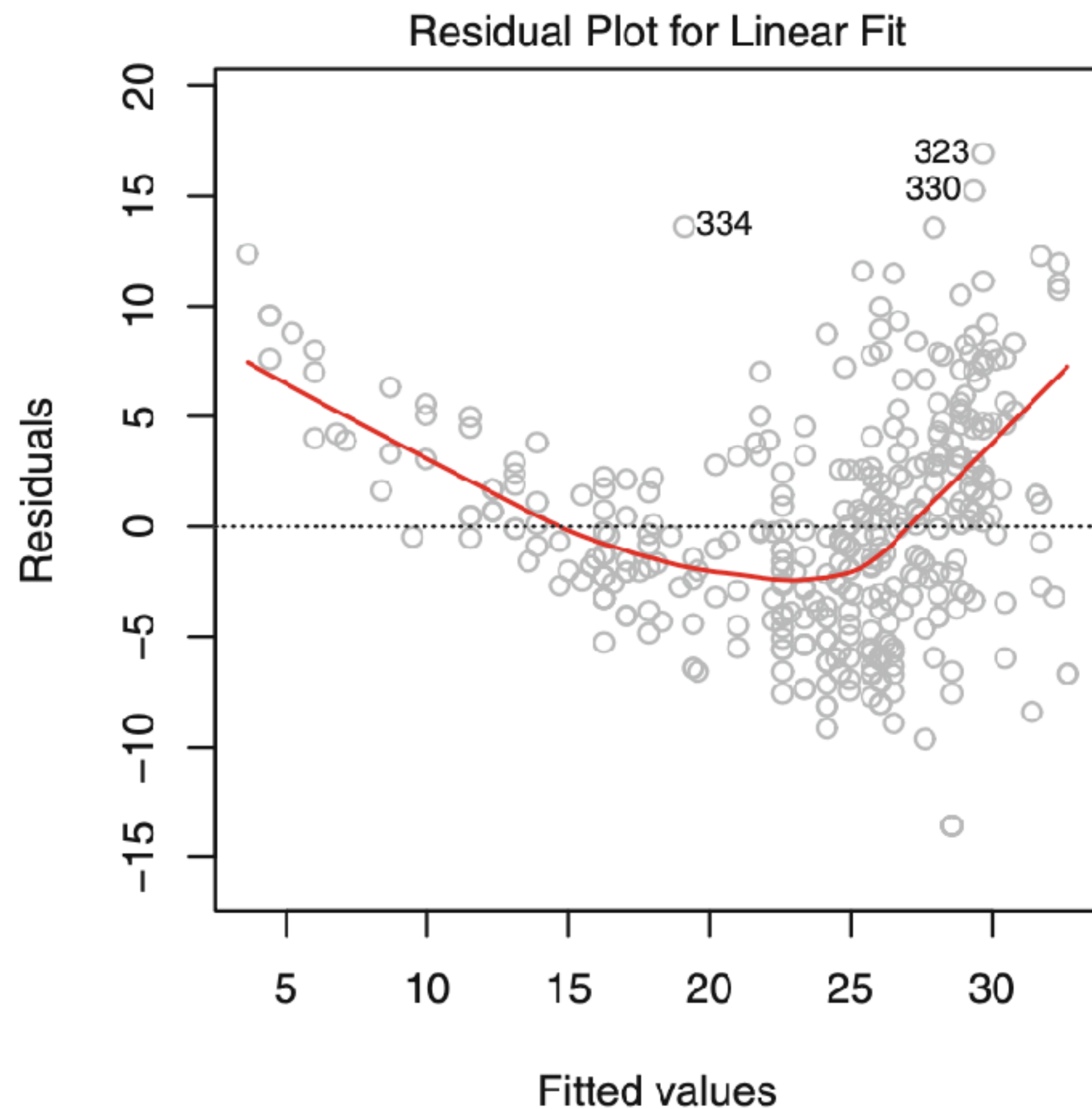


```
##
## Call:
## lm(formula = y1 ~ x1, data = anscombe)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.92127 -0.45577 -0.04136  0.70941  1.83882
##
## Coefficients:
##      Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3.0001     1.1247   2.667  0.02573 *
## x1            0.5001     0.1179   4.241  0.00217 **
##
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.237 on 9 degrees of freedom
## Multiple R-squared:  0.6665, Adjusted R-squared:  0.6295
## F-statistic: 17.99 on 1 and 9 DF, p-value: 0.00217
```

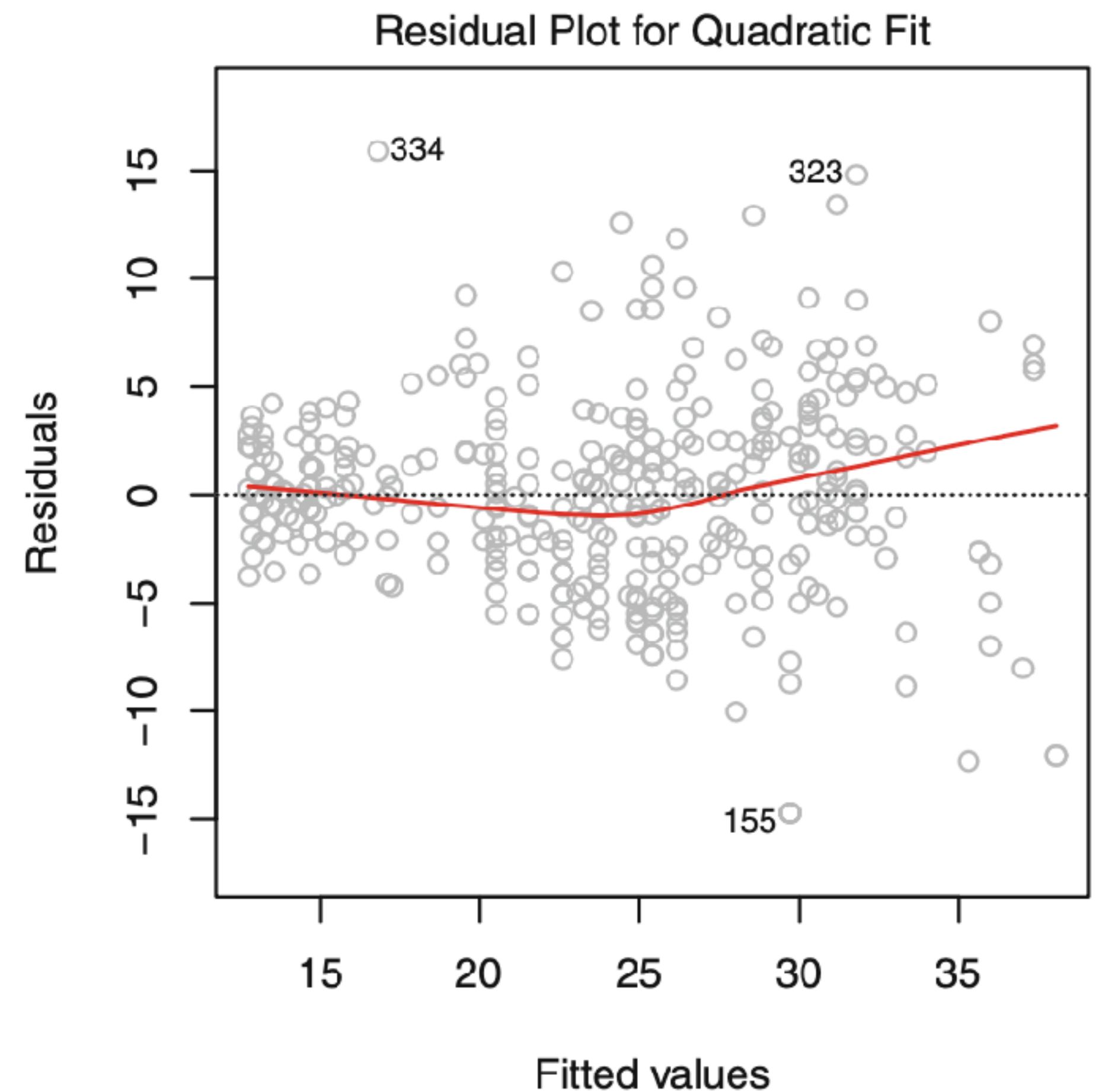
```
##
## Call:
## lm(formula = y2 ~ x2, data = anscombe)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.9009 -0.7609  0.1291  0.9491  1.2691
##
## Coefficients:
##      Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3.001     1.125   2.667  0.02576 *
## x2            0.500     0.118   4.239  0.00218 **
##
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.237 on 9 degrees of freedom
## Multiple R-squared:  0.6662, Adjusted R-squared:  0.6292
## F-statistic: 17.97 on 1 and 9 DF, p-value: 0.002179
```

Another Example: Dealing With Non-Linearity

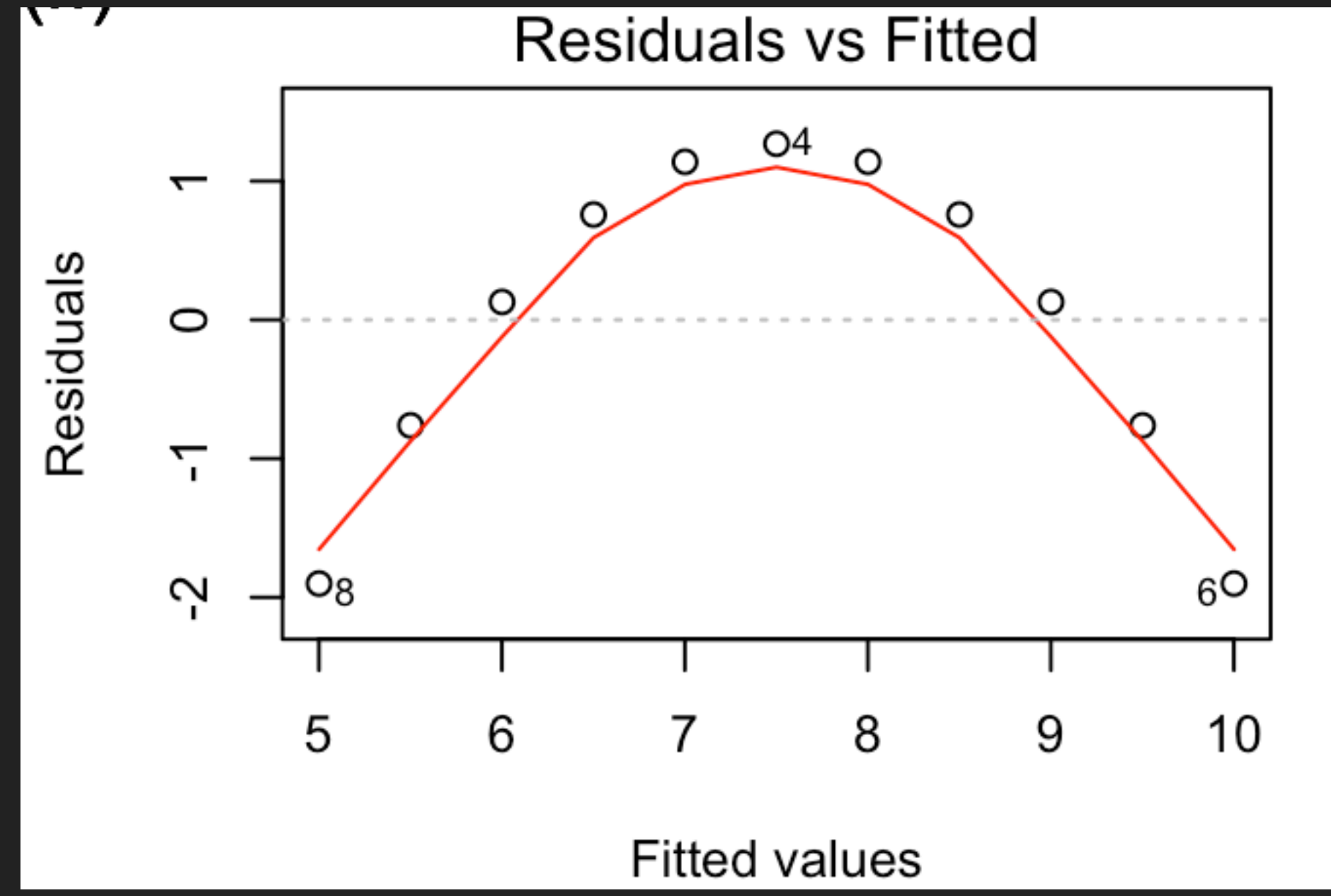
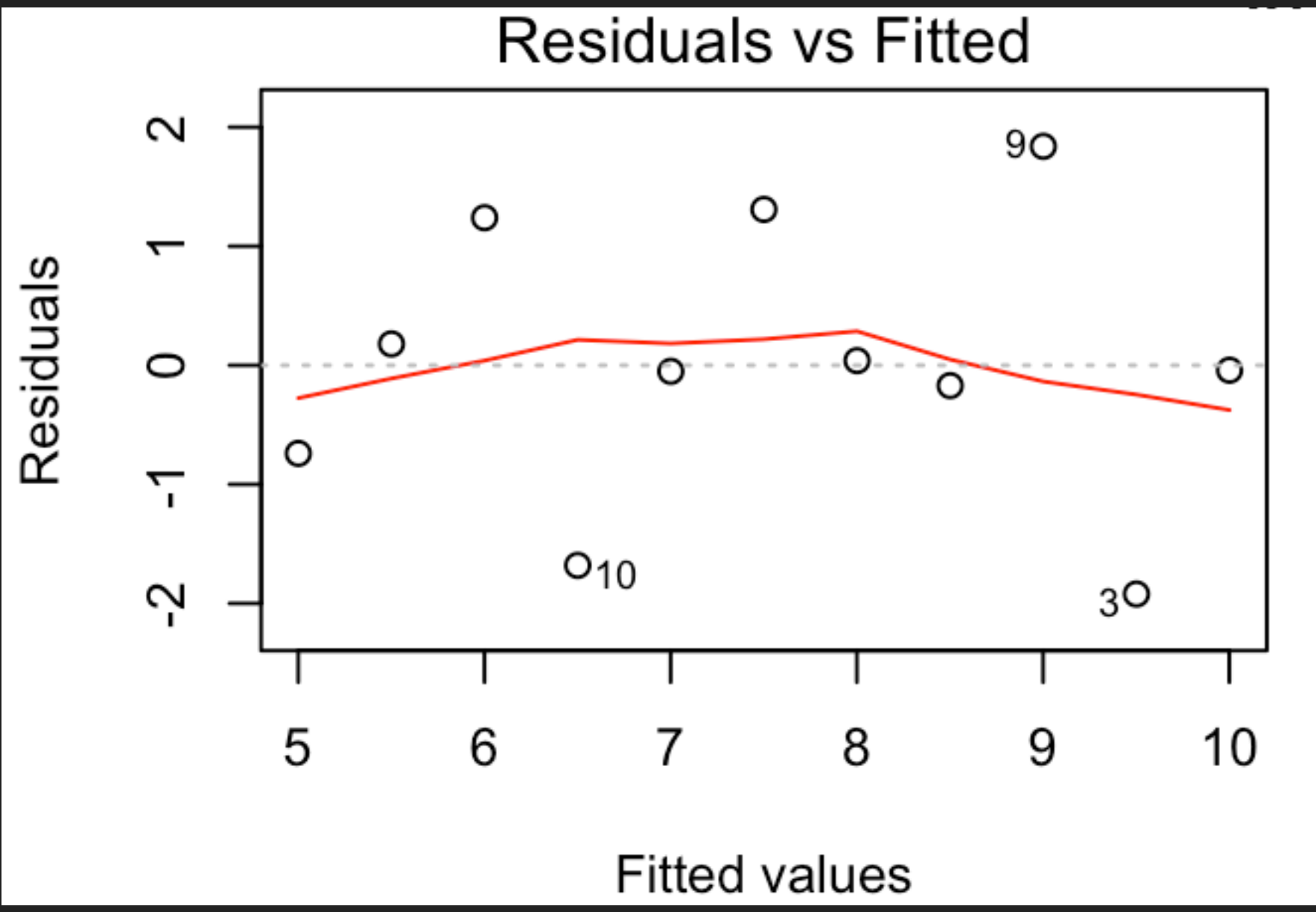
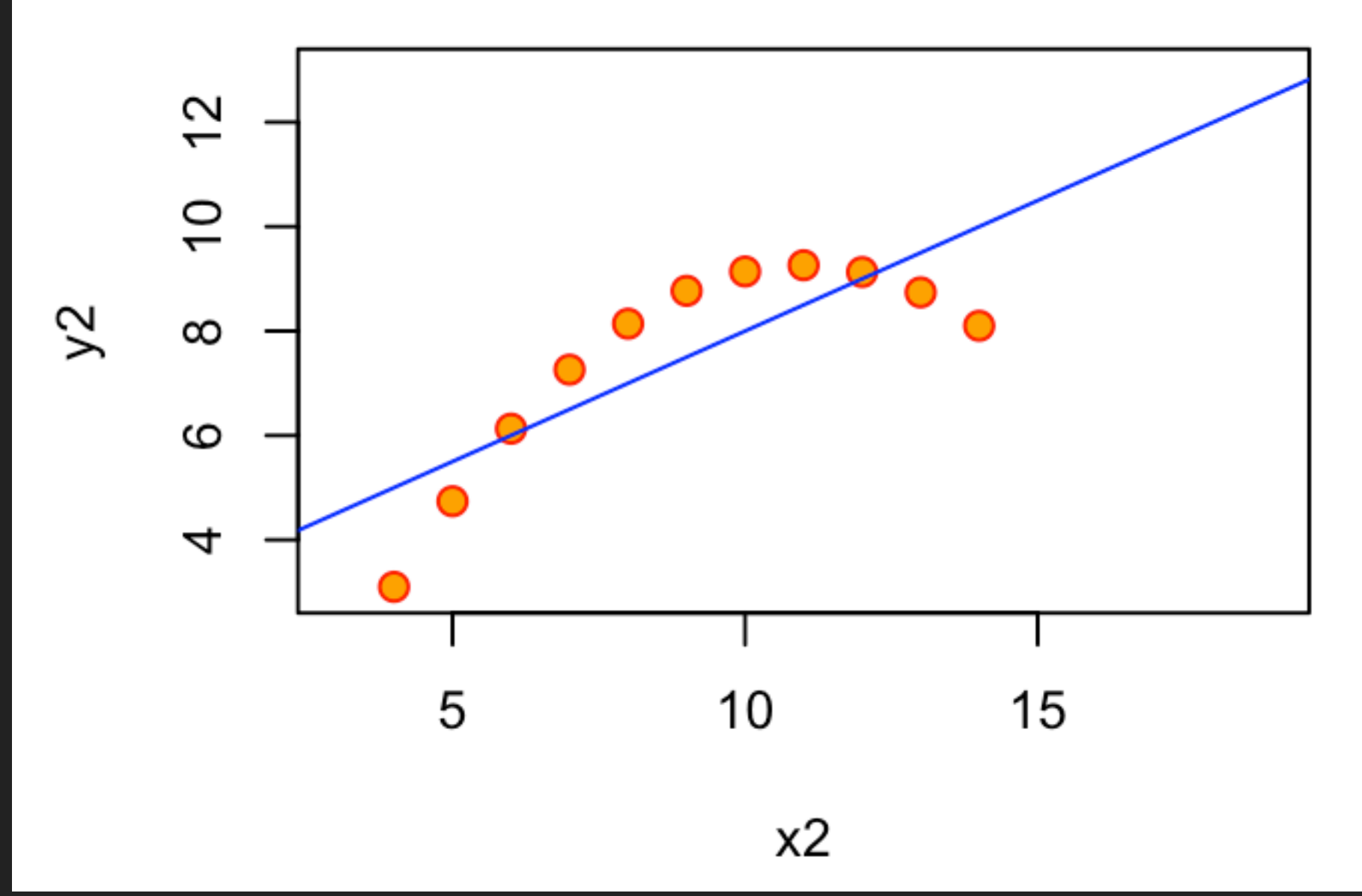
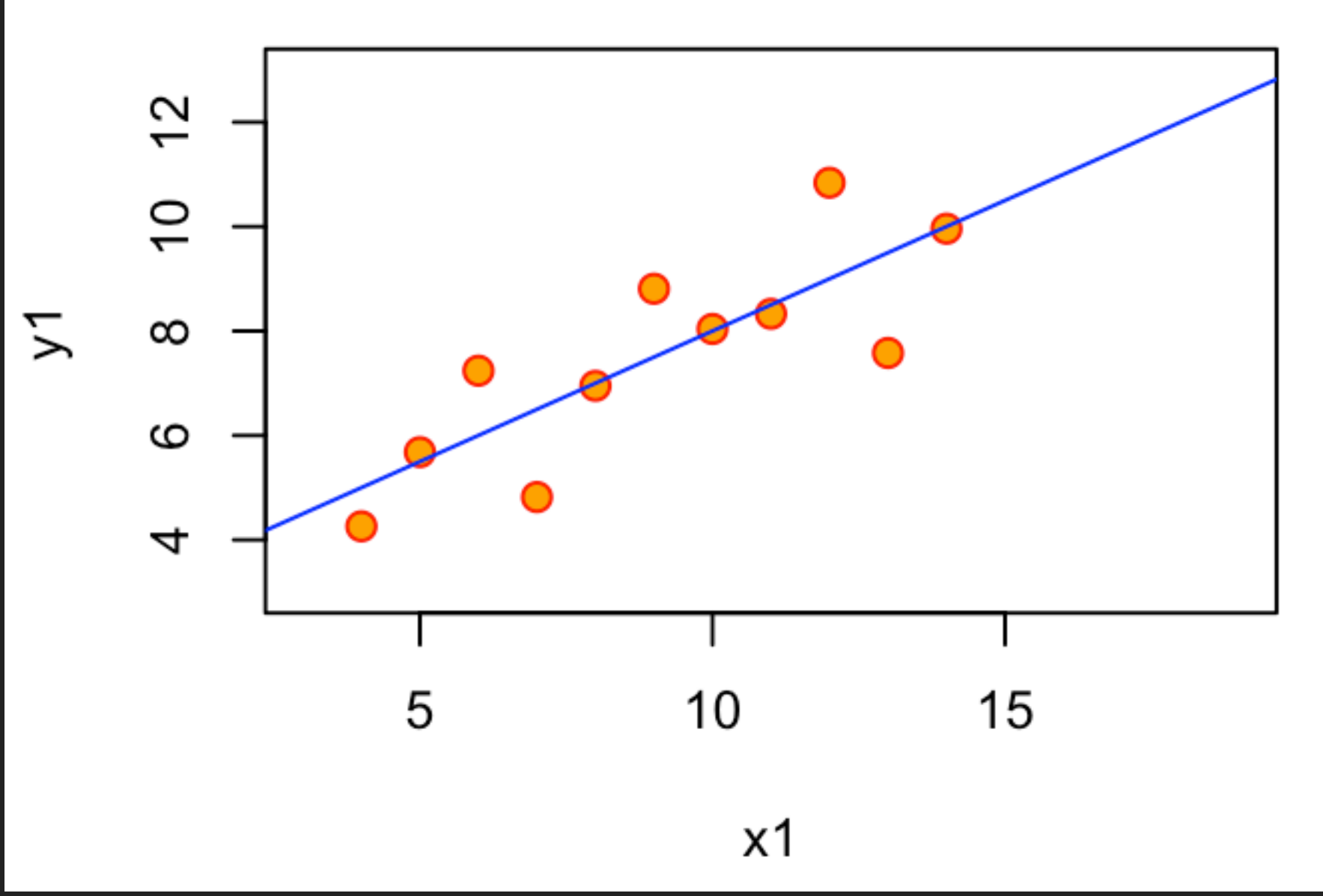
linear regression of mpg on horsepower



linear regression of mpg on horsepower and horsepower²



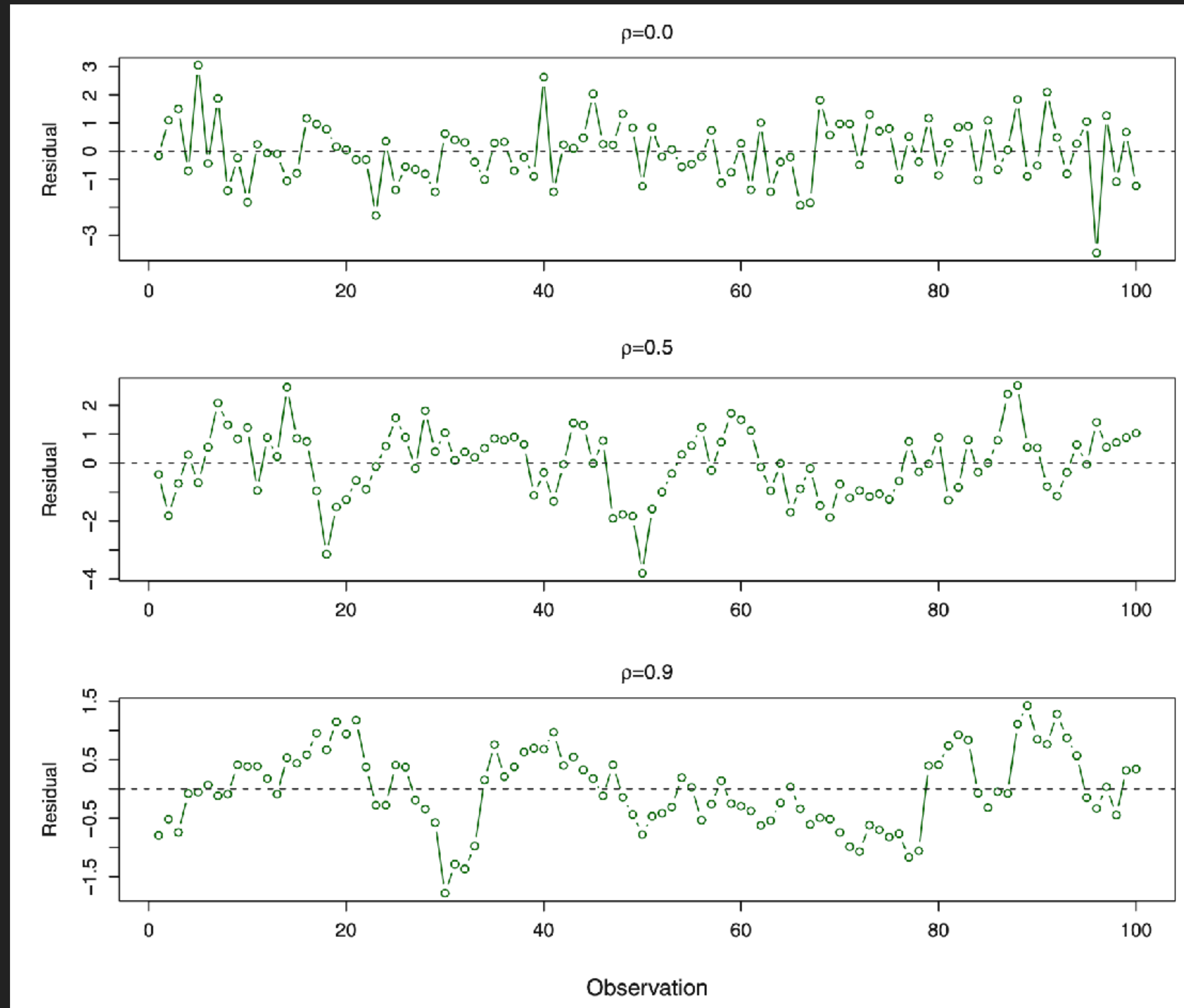
Remember This Example?



Potential Problem: Correlation of Error Terms

- ▶ Some causes:
 - ▶ Time series: observations at adjacent time points will have positively correlated errors
 - ▶ Also non time-series causes
- ▶ Effect:
 - ▶ The estimated standard errors will tend to underestimate the true standard errors.

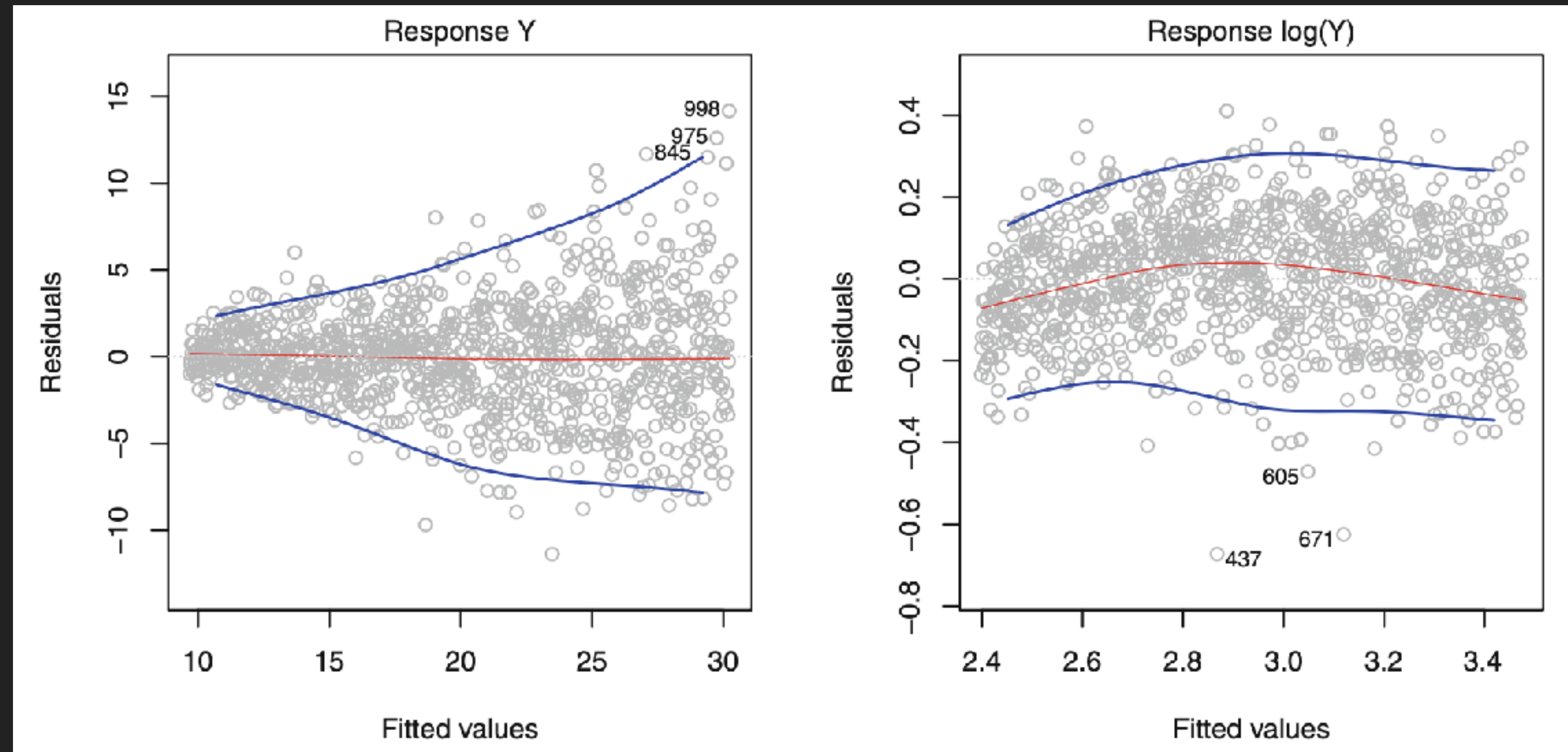
Plots of residuals from simulated time series data sets generated with differing levels of correlation between error terms for adjacent time points.



Potential Problem: Non-Constant Variance of Error Terms (“Heteroscedasticity”)

- ▶ Symptom: the variances of the error terms may increase with the value of the response.

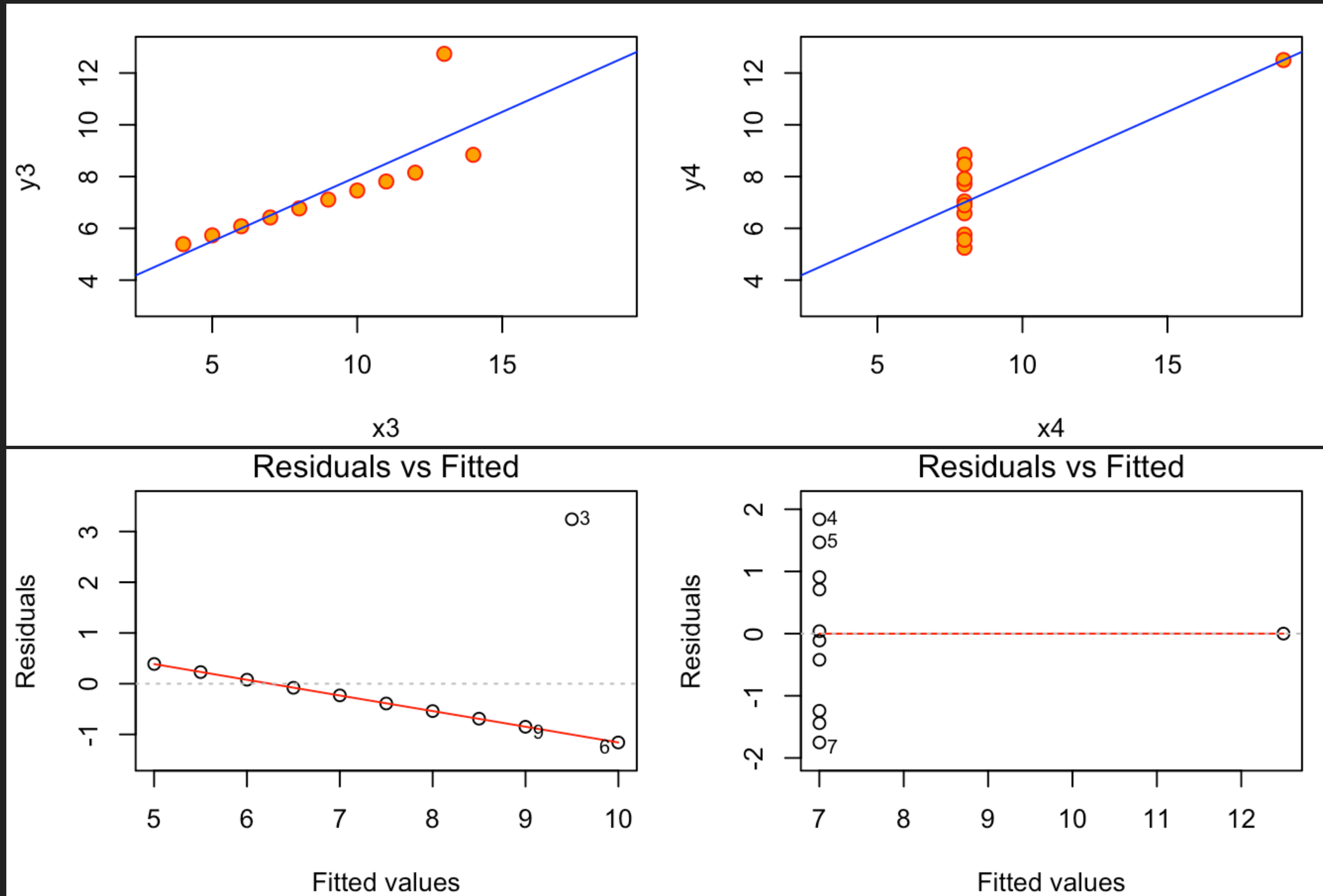
The funnel shape indicates heteroscedasticity.



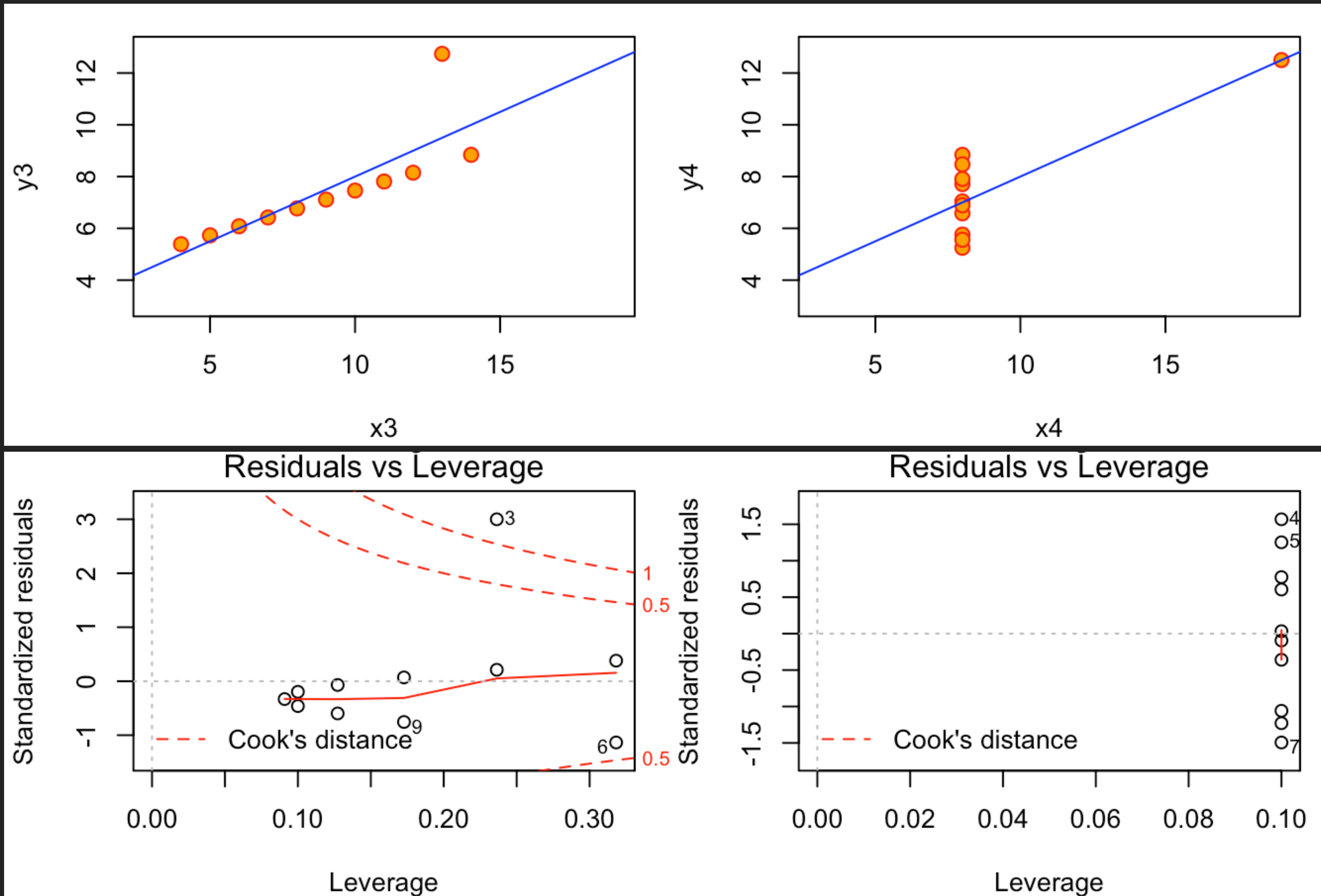
The response has been log transformed, and there is now no evidence of heteroscedasticity.

Heteroscedasticity tends to produce p-values that are smaller than they should be.

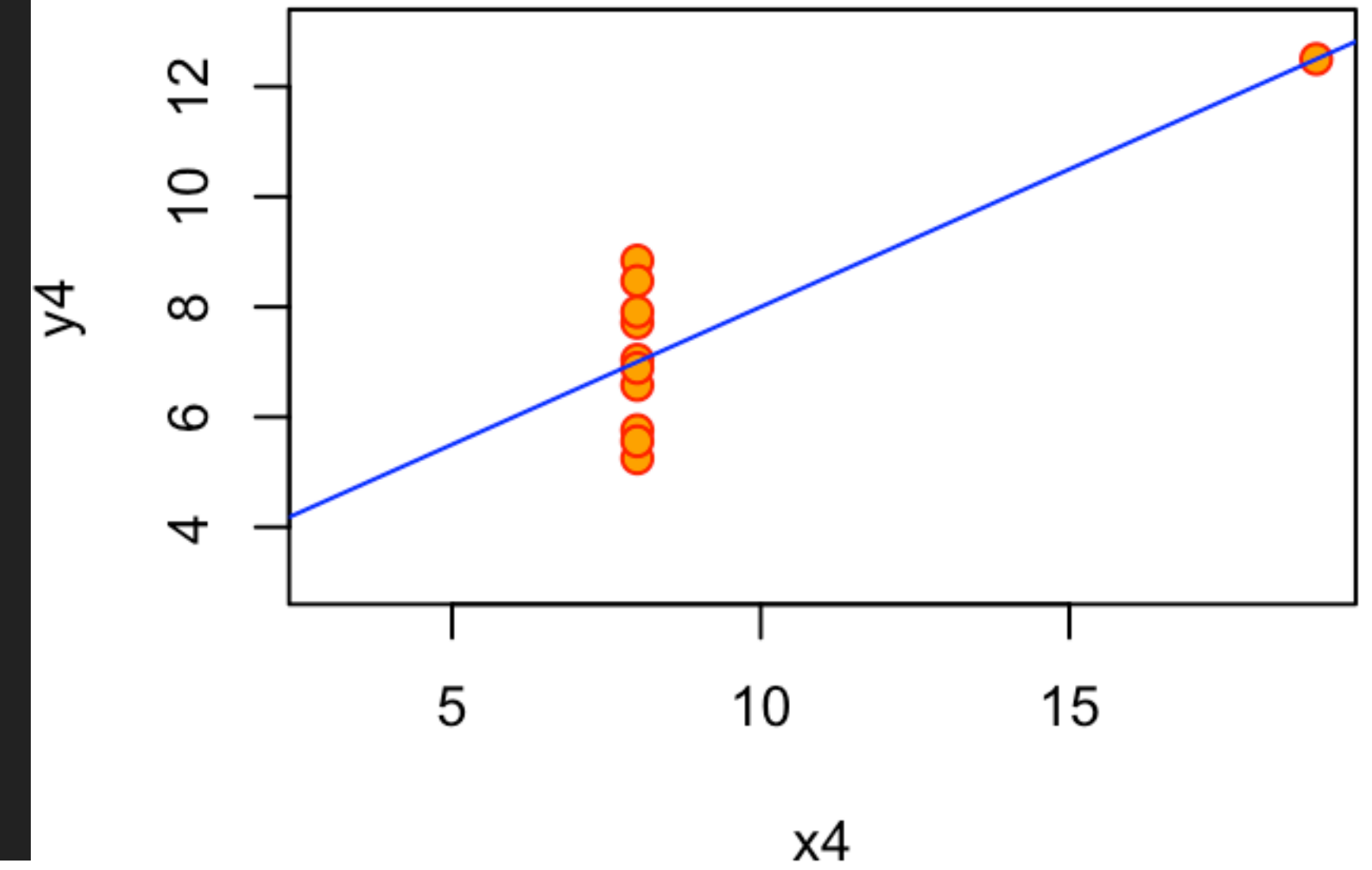
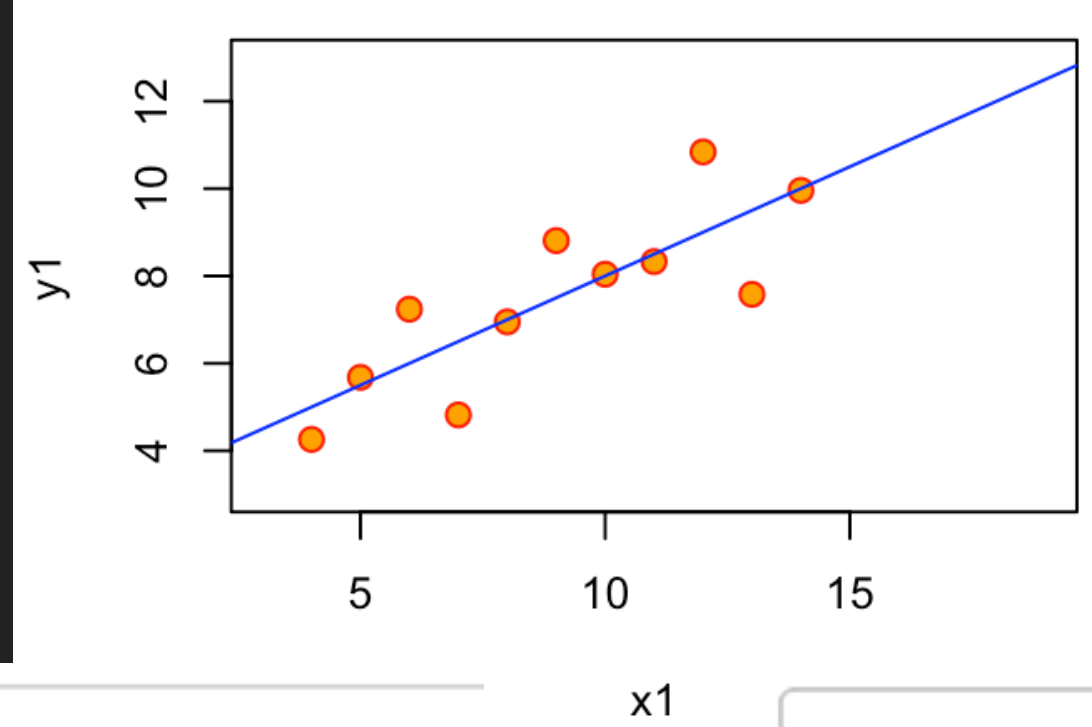
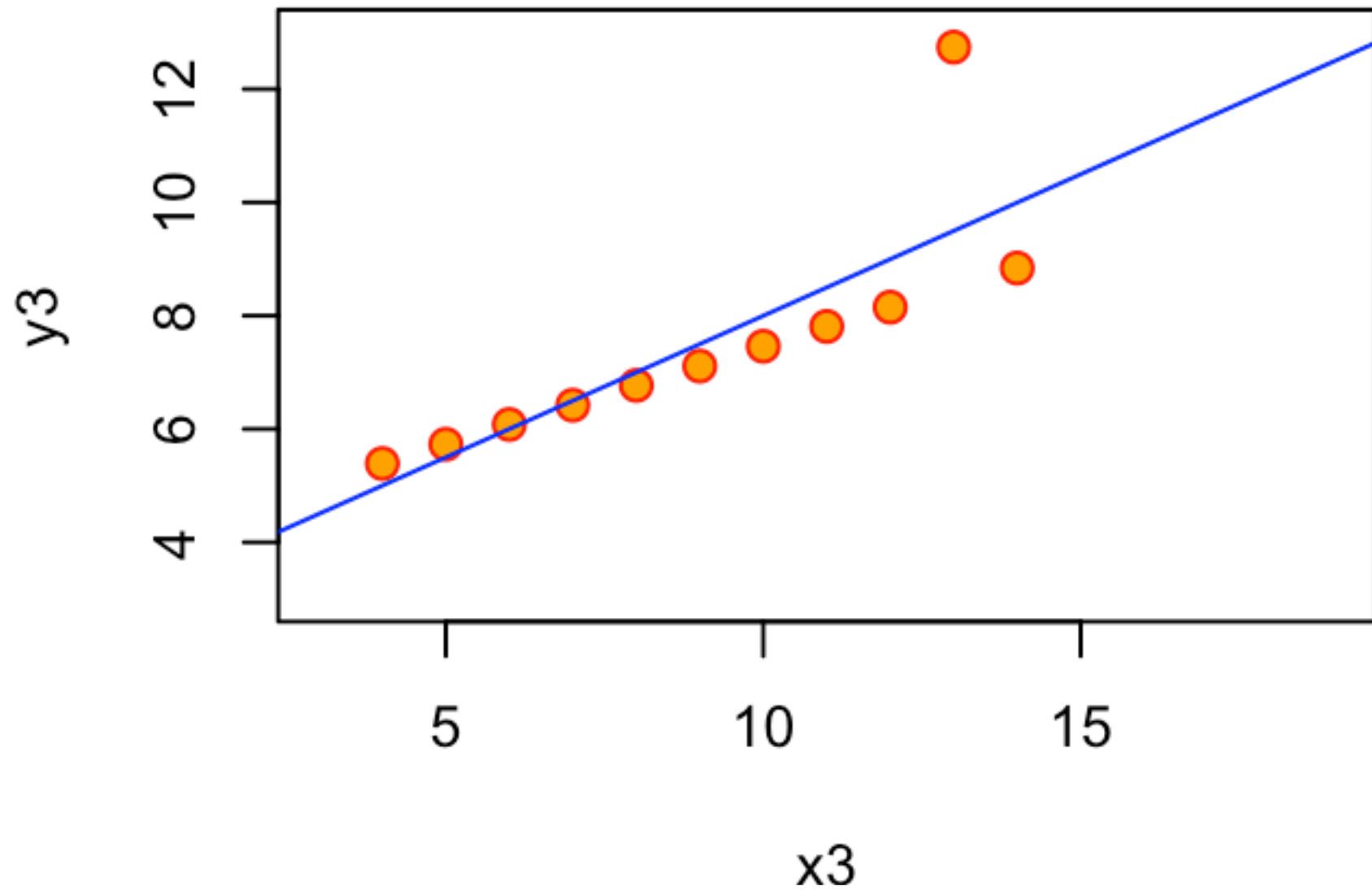
Potential Problems: Outliers and High Leverage Points



Potential Problems: Outliers and High Leverage Points



Remember the Earlier Example?



```
##
## Call:
## lm(formula = y3 ~ x3, data = anscombe)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.1586 -0.6146 -0.2303  0.1540  3.2411
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3.0025     1.1245   2.670  0.02562 *
## x3            0.4997     0.1179   4.239  0.00218 **
##
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.236 on 9 degrees of freedom
## Multiple R-squared:  0.6663, Adjusted R-squared:  0.6292
## F-statistic: 17.97 on 1 and 9 DF, p-value: 0.002176
```

```
##
## Call:
## lm(formula = y4 ~ x4, data = anscombe)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.751 -0.831  0.000  0.809  1.839
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   3.0017     1.1239   2.671  0.02559 *
## x4            0.4999     0.1178   4.243  0.00216 **
##
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.236 on 9 degrees of freedom
## Multiple R-squared:  0.6667, Adjusted R-squared:  0.6297
## F-statistic: 18 on 1 and 9 DF, p-value: 0.002165
```

Potential Problem: Collinearity

- ▶ Here's an extreme example of perfectly collinear data.
- ▶ By construction, x_1 and x_2 are exactly the same variable, and the outcome y is perfectly modeled as $y = x_1 + x_2$

```
my.data <- data.frame(y = c(12, 13, 10, 5, 7, 12, 15),  
                     x1 = c(6, 6.5, 5, 2.5, 3.5, 6, 7.5),  
                     x2 = c(6, 6.5, 5, 2.5, 3.5, 6, 7.5))  
  
my.data
```

Potential Problem: Collinearity

- ▶ Here's an extreme example of perfectly collinear data.
- ▶ By construction, x_1 and x_2 are exactly the same variable, and the outcome y is perfectly modeled as $y = x_1 + x_2$
- ▶ But there's a problem... because the following are also true

```
my.data <- data.frame(y = c(12, 13, 10, 5, 7, 12, 15),  
                      x1 = c(6, 6.5, 5, 2.5, 3.5, 6, 7.5),  
                      x2 = c(6, 6.5, 5, 2.5, 3.5, 6, 7.5))  
  
my.data
```

$$y = 2x_1$$

$$y = 3x_1 - x_2$$

$$y = -400x_1 + 402x_2$$

Potential Problem: Collinearity

- ▶ Here's an extreme example of perfectly collinear data.
- ▶ By construction, x_1 and x_2 are exactly the same variable, and the outcome y is perfectly modeled as $y = x_1 + x_2$
- ▶ But there's a problem... because the following are also true

$$y = 2x_1$$

$$y = 3x_1 - x_2$$

$$y = -400x_1 + 402x_2$$

```
my.data <- data.frame(y = c(12, 13, 10, 5, 7, 12, 15),  
                      x1 = c(6, 6.5, 5, 2.5, 3.5, 6, 7.5),  
                      x2 = c(6, 6.5, 5, 2.5, 3.5, 6, 7.5))  
  
my.data
```

Effects:

- ▶ The model is unable to accurately distinguish between many nearly equally plausible linear combinations of collinear variables.
- ▶ This can lead to large standard errors on coefficients, and even coefficient signs that don't make sense.

Potential Problem: Collinearity

- ▶ Here's an extreme example of perfectly collinear data.
- ▶ By construction, x_1 and x_2 are exactly the same variable, and the outcome y is perfectly modeled as $y = x_1 + x_2$
- ▶ But there's a problem... because the following are also true

$$y = 2x_1$$

$$y = 3x_1 - x_2$$

$$y = -400x_1 + 402x_2$$

```
my.data <- data.frame(y = c(12, 13, 10, 5, 7, 12, 15),  
                     x1 = c(6, 6.5, 5, 2.5, 3.5, 6, 7.5),  
                     x2 = c(6, 6.5, 5, 2.5, 3.5, 6, 7.5))  
  
my.data
```

```
# Evaluate Collinearity  
library(car)  
vif(fit) # variance inflation factors  
sqrt(vif(fit)) > 2 # problem?
```

Activity: How To Address These Questions?

- ▶ Is there a relationship between advertising budget and sales?
- ▶ How strong is the relationship between advertising budget and sales?
- ▶ Which media contribute to sales?
- ▶ How accurately can we estimate the effect of each medium on sales?
- ▶ How accurately can we predict future sales?
- ▶ Is the relationship linear?
- ▶ Is there synergy among the advertising media?

... to be continued

Credits

- ▶ Graphics: Dave DiCello photography (cover)
- ▶ Bruce, P., Bruce, A., & Gedeck, P. (2020). Practical Statistics for Data Scientists: 50+ Essential Concepts Using R and Python. O'Reilly Media.
- ▶ Goodman, S. (2008). A dirty dozen: Twelve p-value misconceptions. In Seminars in Hematology (Vol. 45, No. 3, pp. 135-140). WB Saunders.
- ▶ James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). An introduction to statistical learning (Vol. 112, p. 18). New York: springer.
- ▶ Grolemund, G., & Wickham, H. (2018). R for data science.