# Regression for Prediction

Balancing Complexity and Simplicity
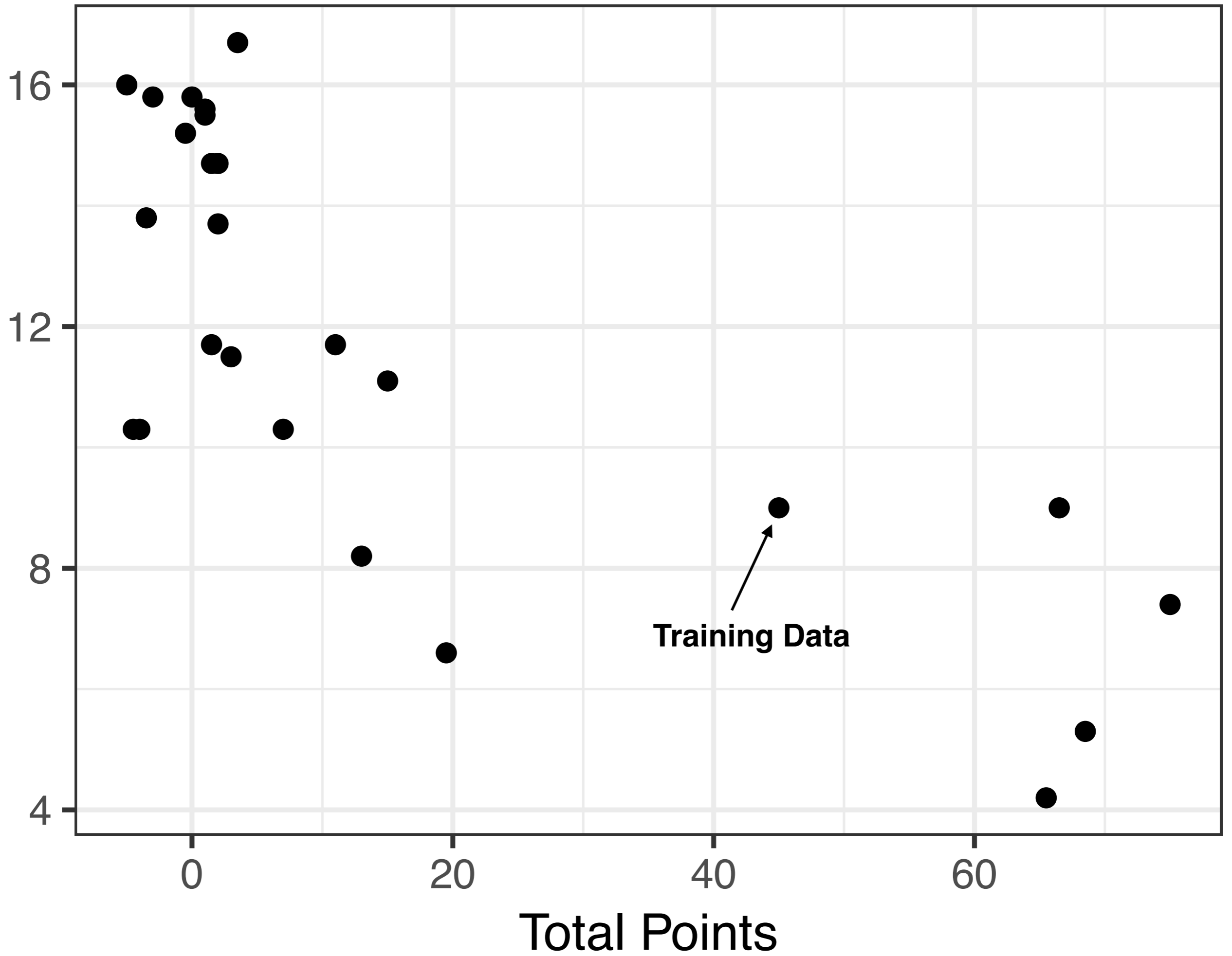
# Goldilocks
## and the
# Three Bears

# How do I balance simplicity and complexity?
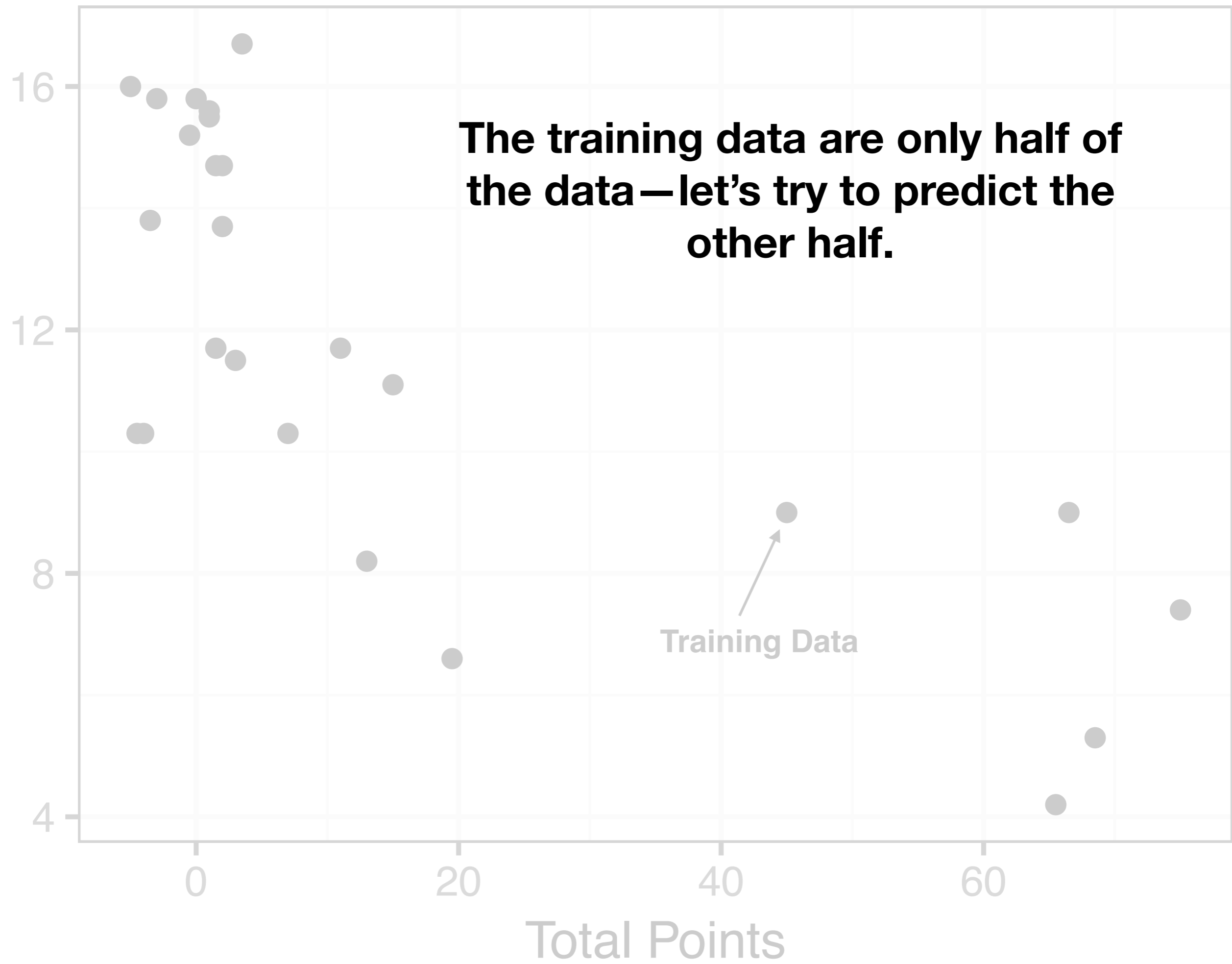
# Three Types of Data Sets

- **training set**: cases and variables used to fit the models

- **prediction set**: cases to be predicted—includes same explanatory variables as training set, but missing the outcome of interest

- **test set**: has same cases as prediction set, but also includes the outcome of interest.

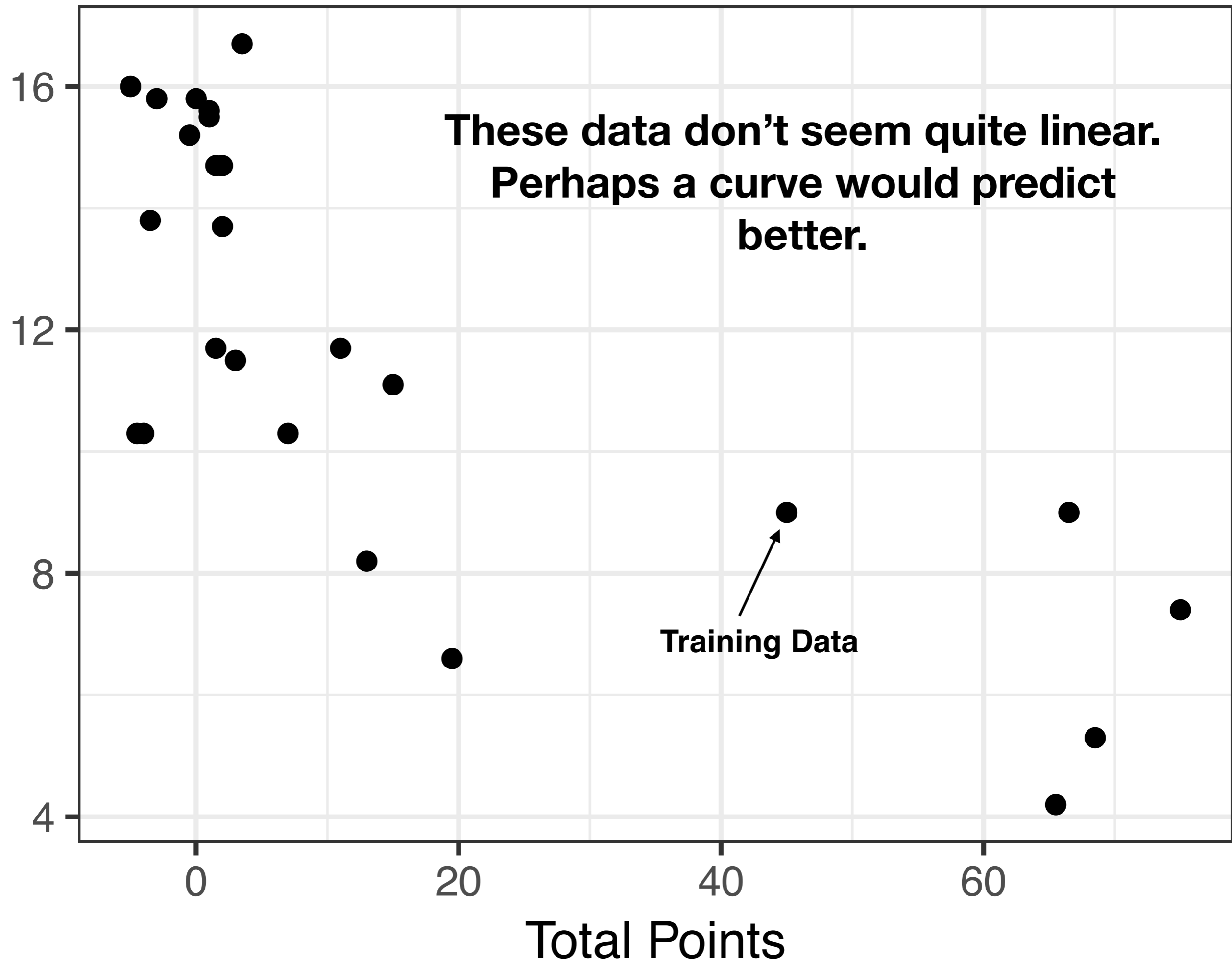The training data are only half of the data—let's try to predict the other half.

Training Data

These data don't seem quite linear. Perhaps a curve would predict better.

Training Data

# How curvy should the line be?

$$y = \beta_0 + \beta_1 x + u$$

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + u$$

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + u$$

$$\vdots$$

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \cdots + \beta_k x^k + u$$

```
# 1st-order polynomial (i.e., linear)
> m <- lm(firearm_death_rate ~ total_points,
+         data = train)
> e <- residuals(m)
> sqrt(mean(e^2))
[1] 2.332054
```
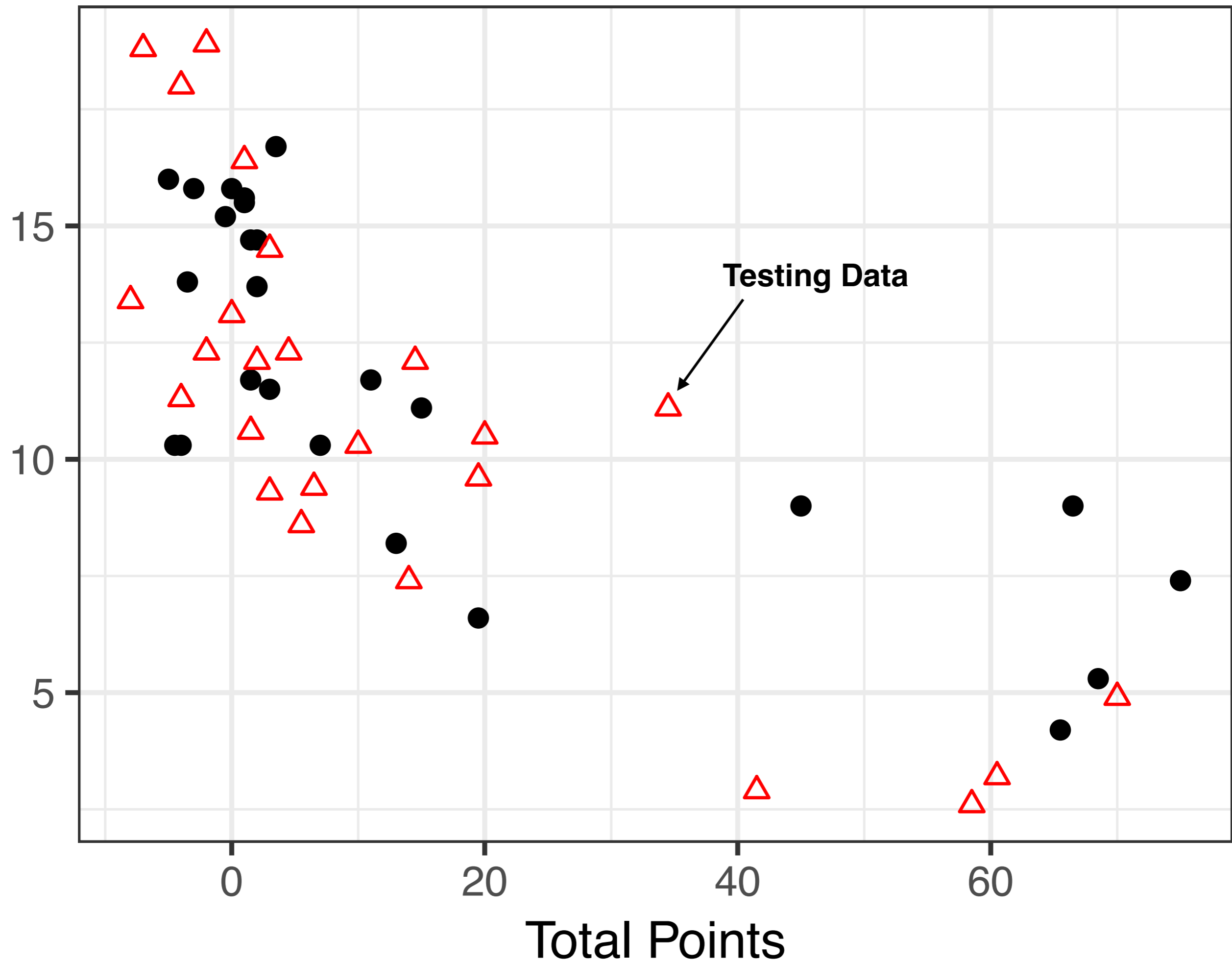
```
# 2nd-order polynomial (i.e.,quadratic)
> m <- lm(firearm_death_rate ~ poly(total_points, 2),
+          data = train)
> e <- residuals(m)
> sqrt(mean(e^2))
[1] 2.198679
```

```
# 3rd-order polynomial (i.e., cubic)
> m <- lm(firearm_death_rate ~ poly(total_points, 3),
+         data = train)
> e <- residuals(m)
> sqrt(mean(e^2))
[1] 2.197419
```
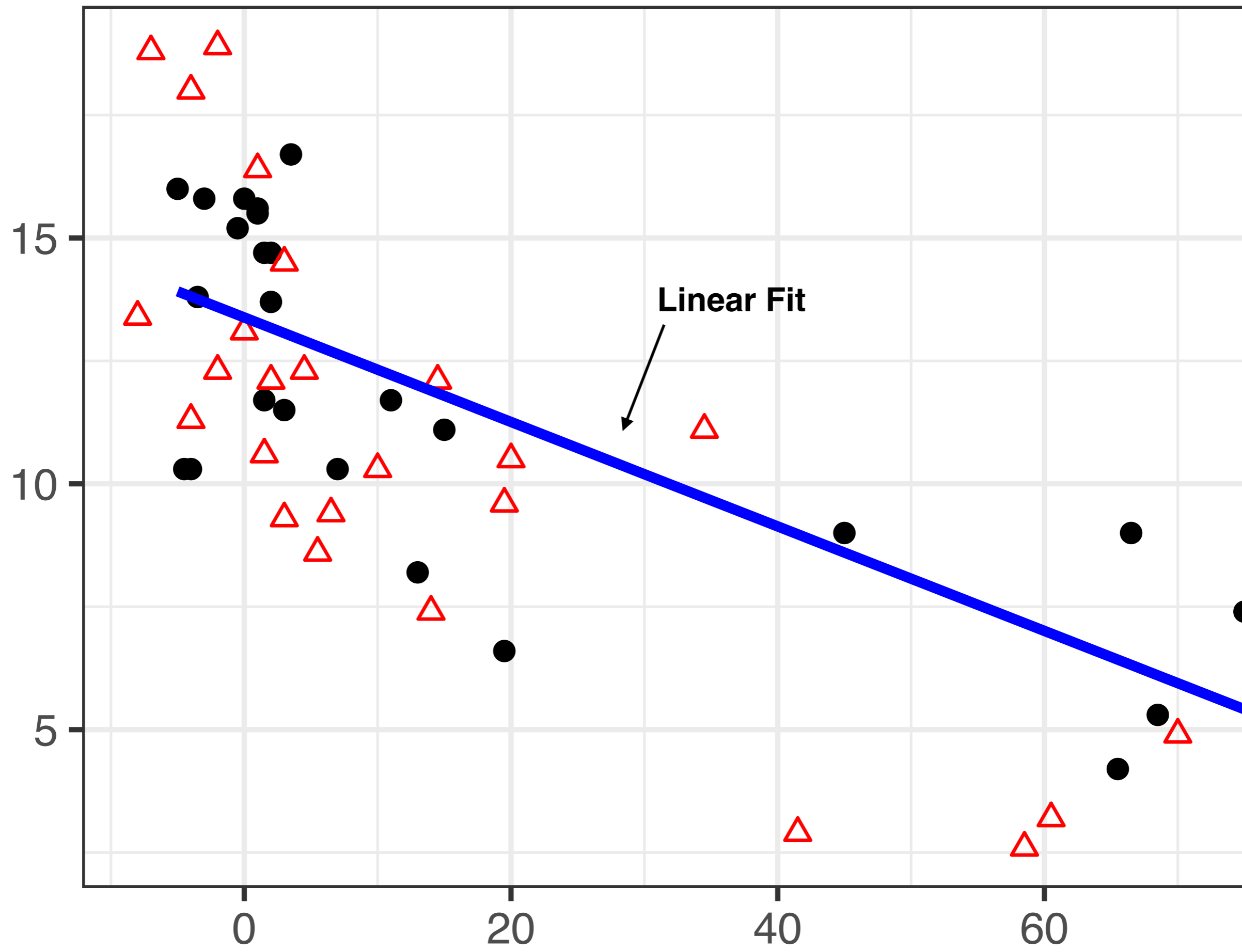
```
# 10th-order polynomial
> m <- lm(firearm_death_rate ~ poly(total_points, 10),
+         data = train)
> e <- residuals(m)
> sqrt(mean(e^2))
[1] 1.569003
```

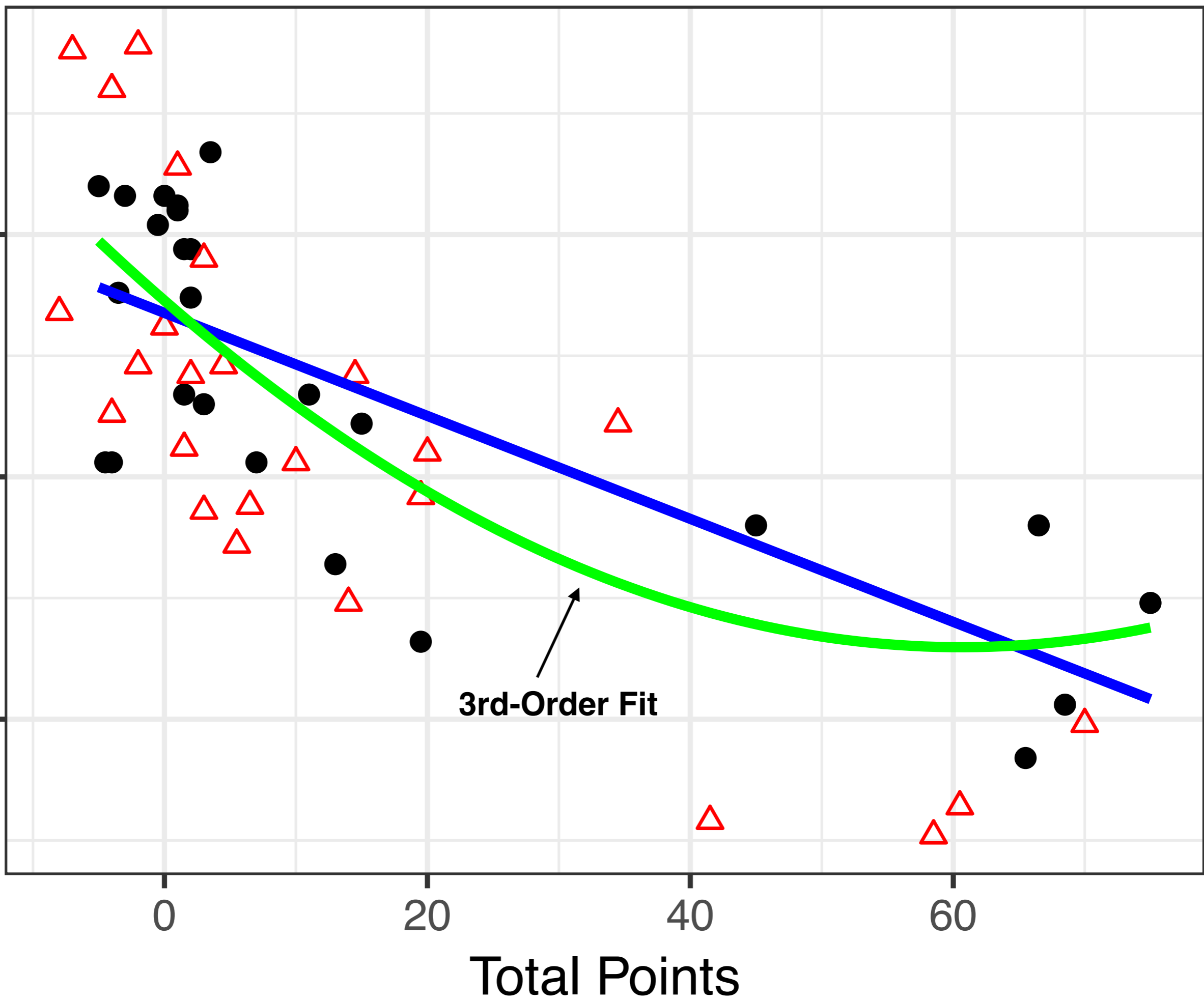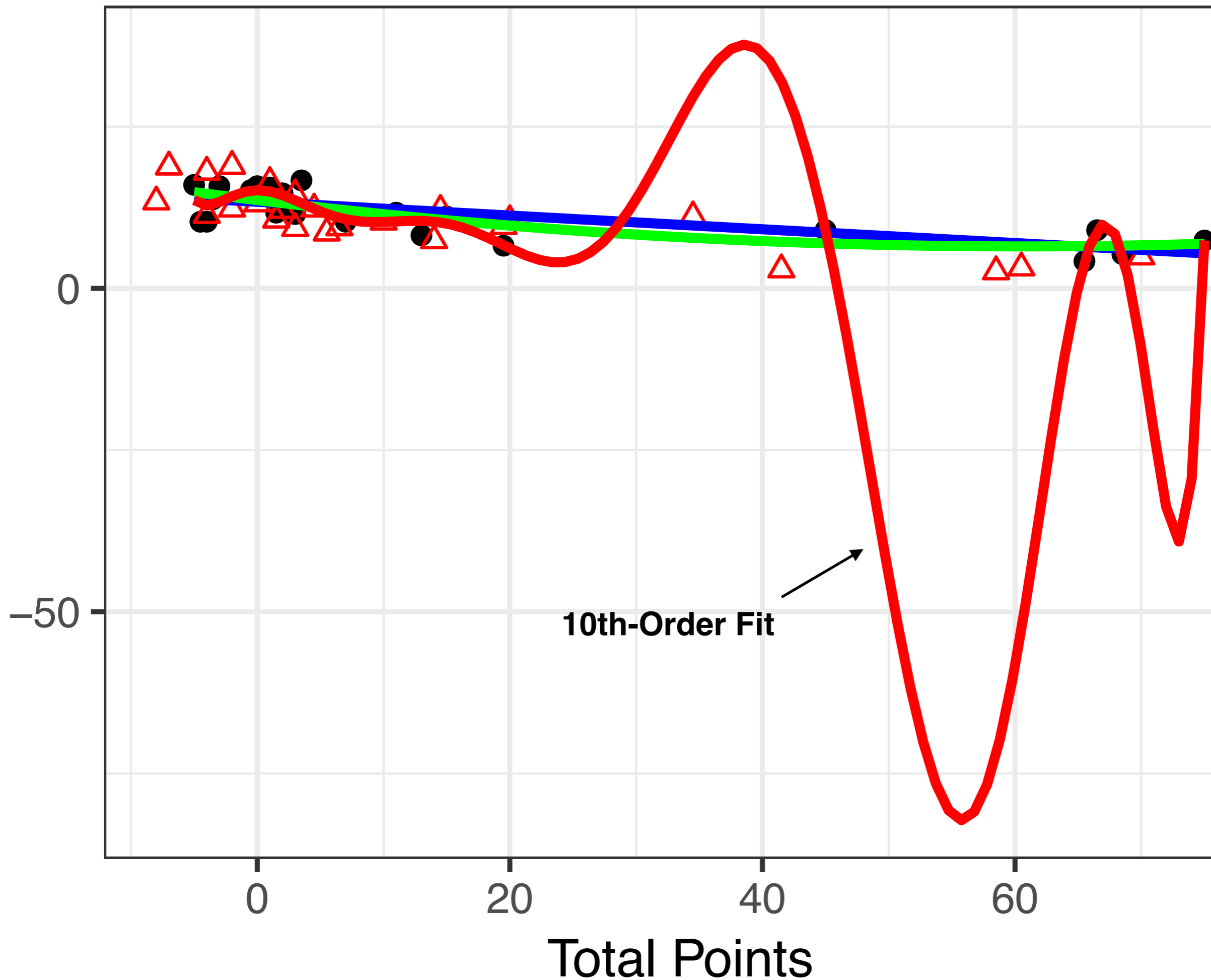Will the regression with the lowest r.m.s. error have the most predictive power?

3rd-Order Fit

10th-Order Fit

In-sample error always
decreases with complexity.

# How do I balance simplicity and complexity?

How do balance simplicity and complexity?

How do I balance simplicity and complexity?

INTUITION