

Lesson: Linear regression diagnostic methods

Module: Regression Diagnostics

Photo by [Janine Robinson](#) on [Unsplash](#)





The Linear Regression Model

Definition/Assumptions of the linear regression model:

1. *Linearity*
2. *Independence*
3. *Homoskedasticity (constant variance)*
4. *Normality*



Motivation

So far, we've mainly taken for granted that the MLR assumptions have been met. But what if they aren't?

Diagnostic methods:

1. Graphical techniques
2. Numerical techniques



Motivation

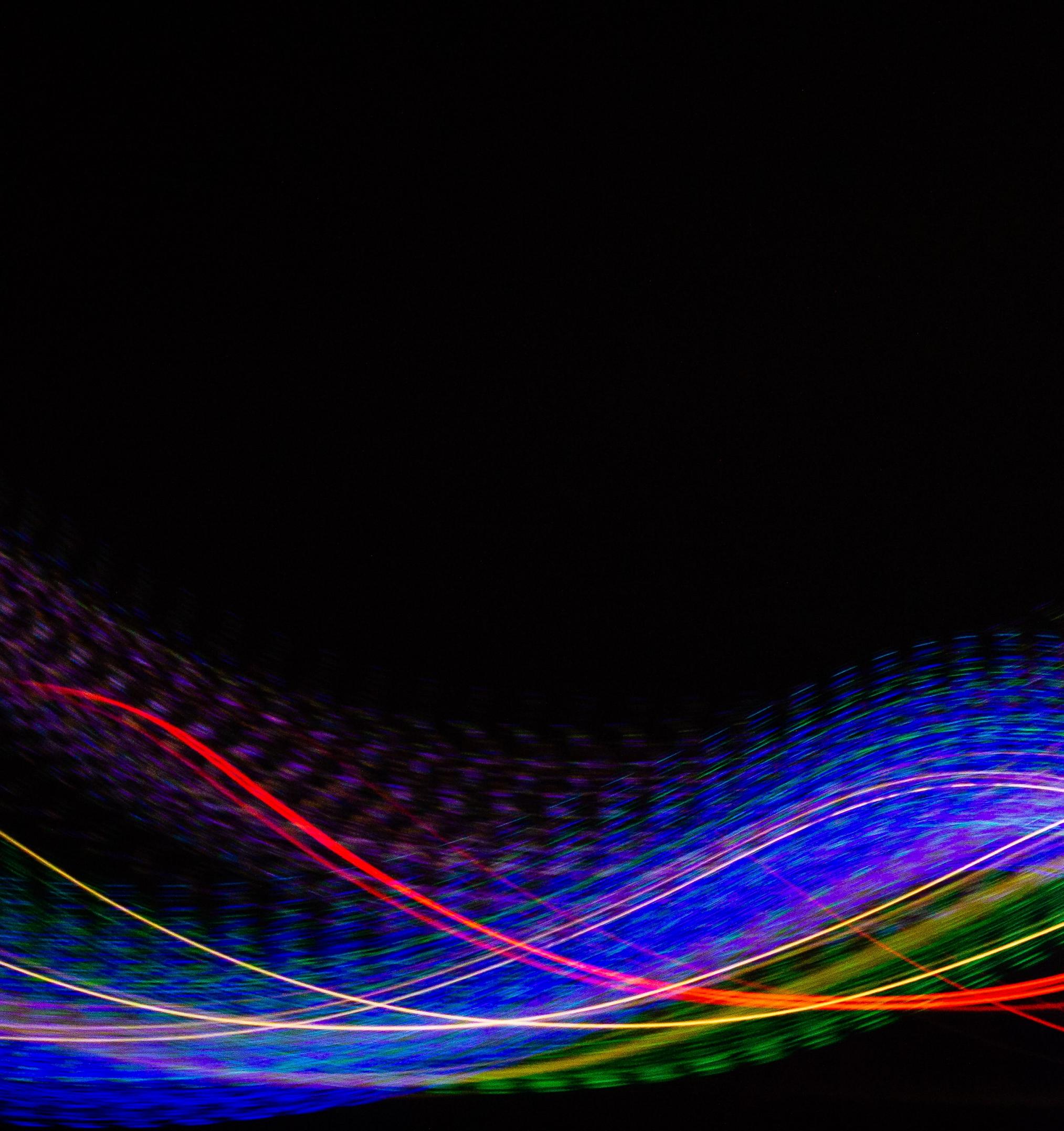
A note on the validity assumption

Lesson: Violations of the linearity assumption

Module: Regression Diagnostics

Photo by [Janine Robinson](#) on [Unsplash](#)

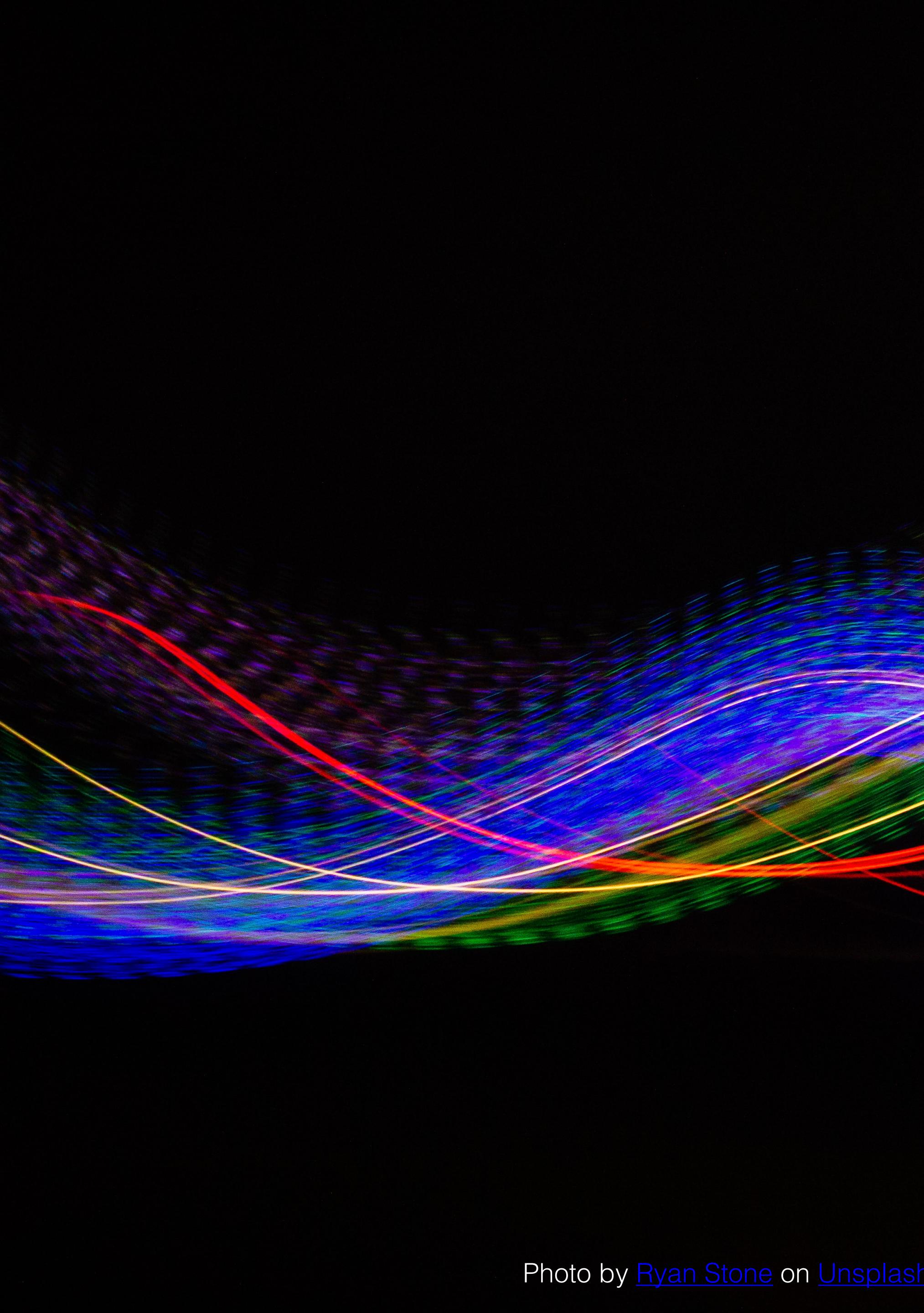




$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \text{ where } E(\boldsymbol{\varepsilon}) = 0.$$

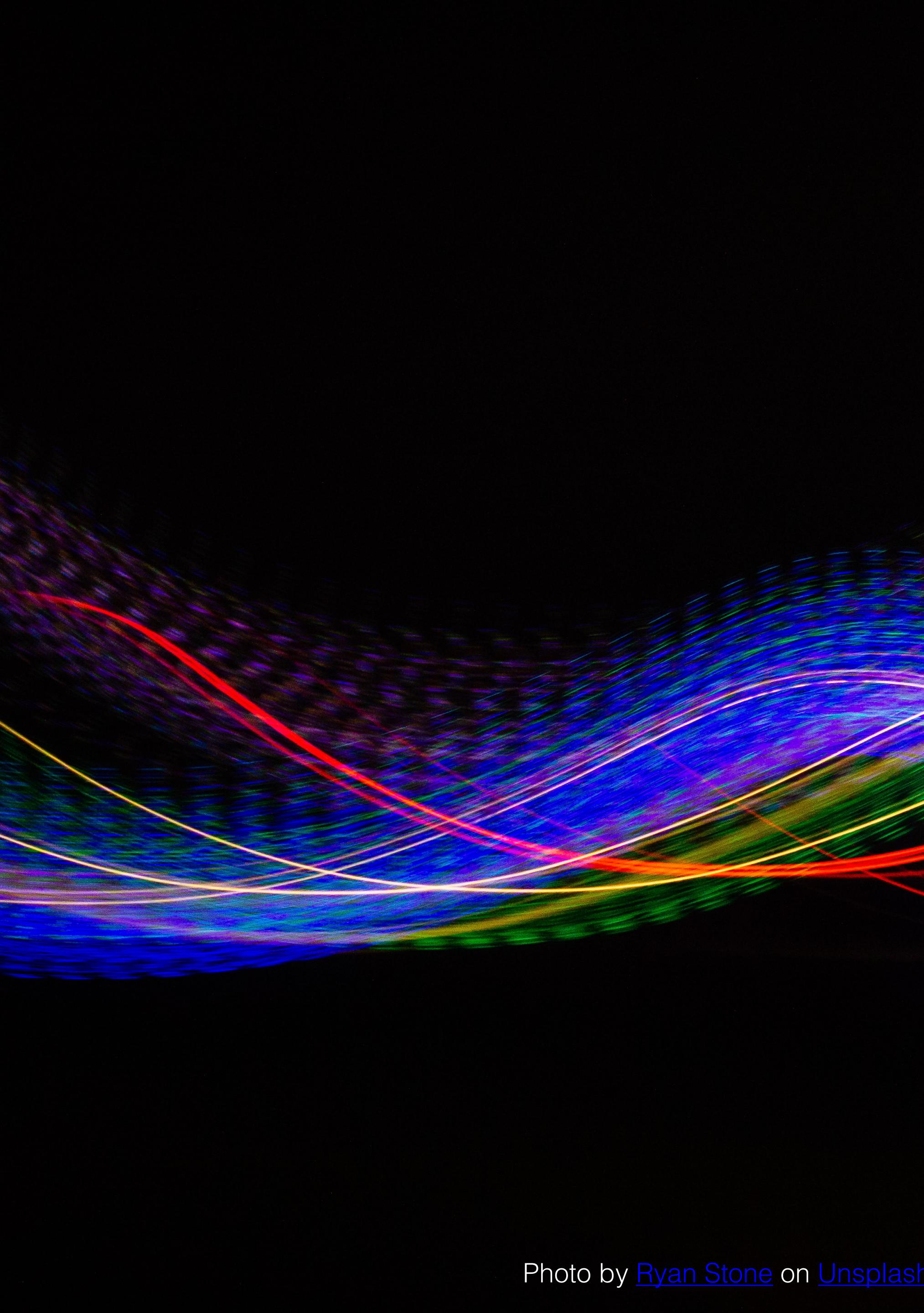
Or equivalently, that for $i = 1, \dots, n$,

$$Y_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_p x_{i,p} + \varepsilon_i, \text{ where } E(\varepsilon_i) = 0.$$



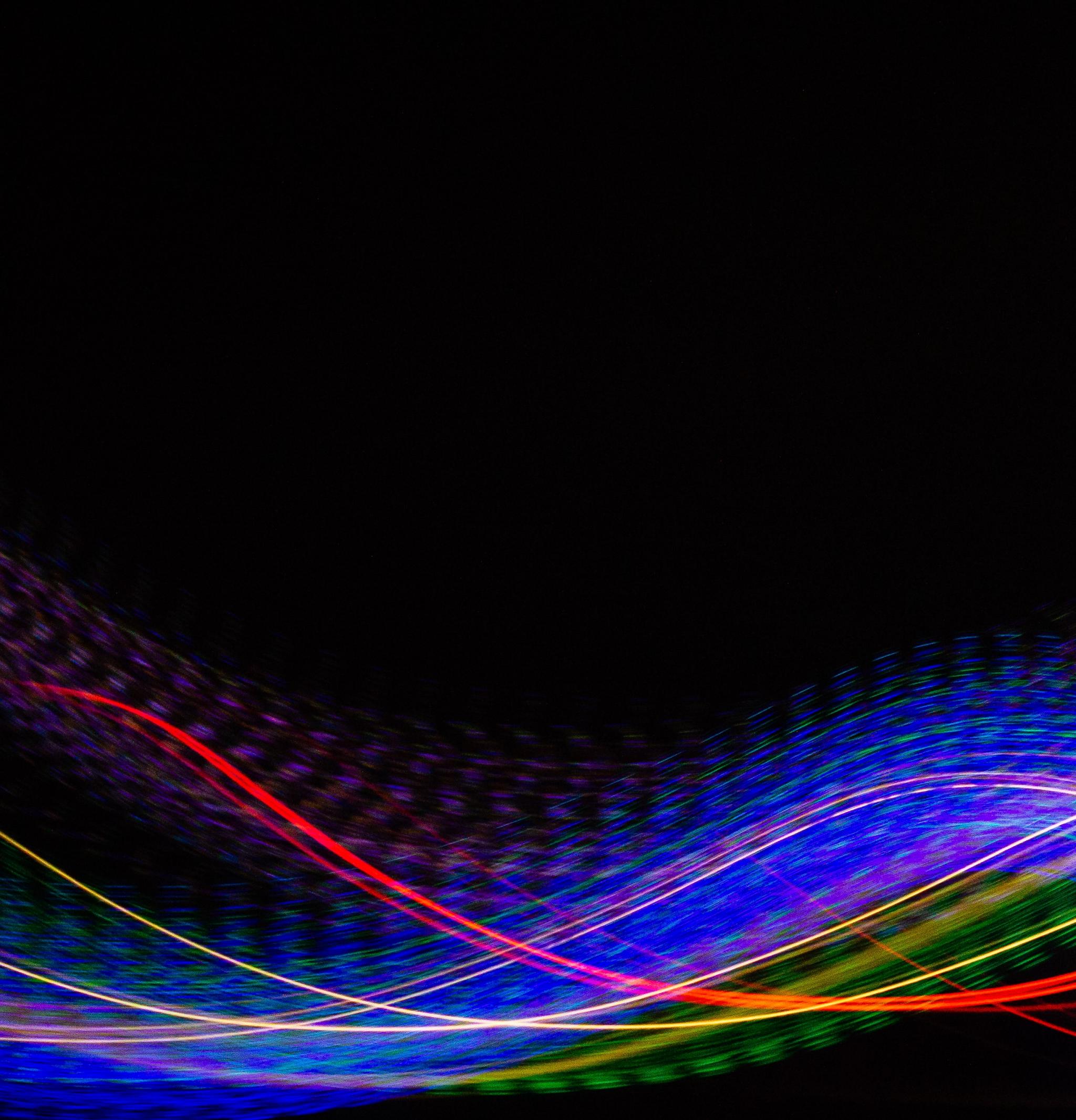
This specification means that, for each $i = 1, \dots, n$ and $j = 1, \dots, p$:

1. $E(Y_i)$ is a linear function of $x_{i,j}$, holding all other $x_{i,k}$, $k \neq j$ fixed.
2. The slope of the line that relates $E(Y_i)$ to $x_{i,j}$ does not depend on the values of any other $x_{i,k}$, $k \neq j$.
3. The effects of each $x_{i,j}$ on the



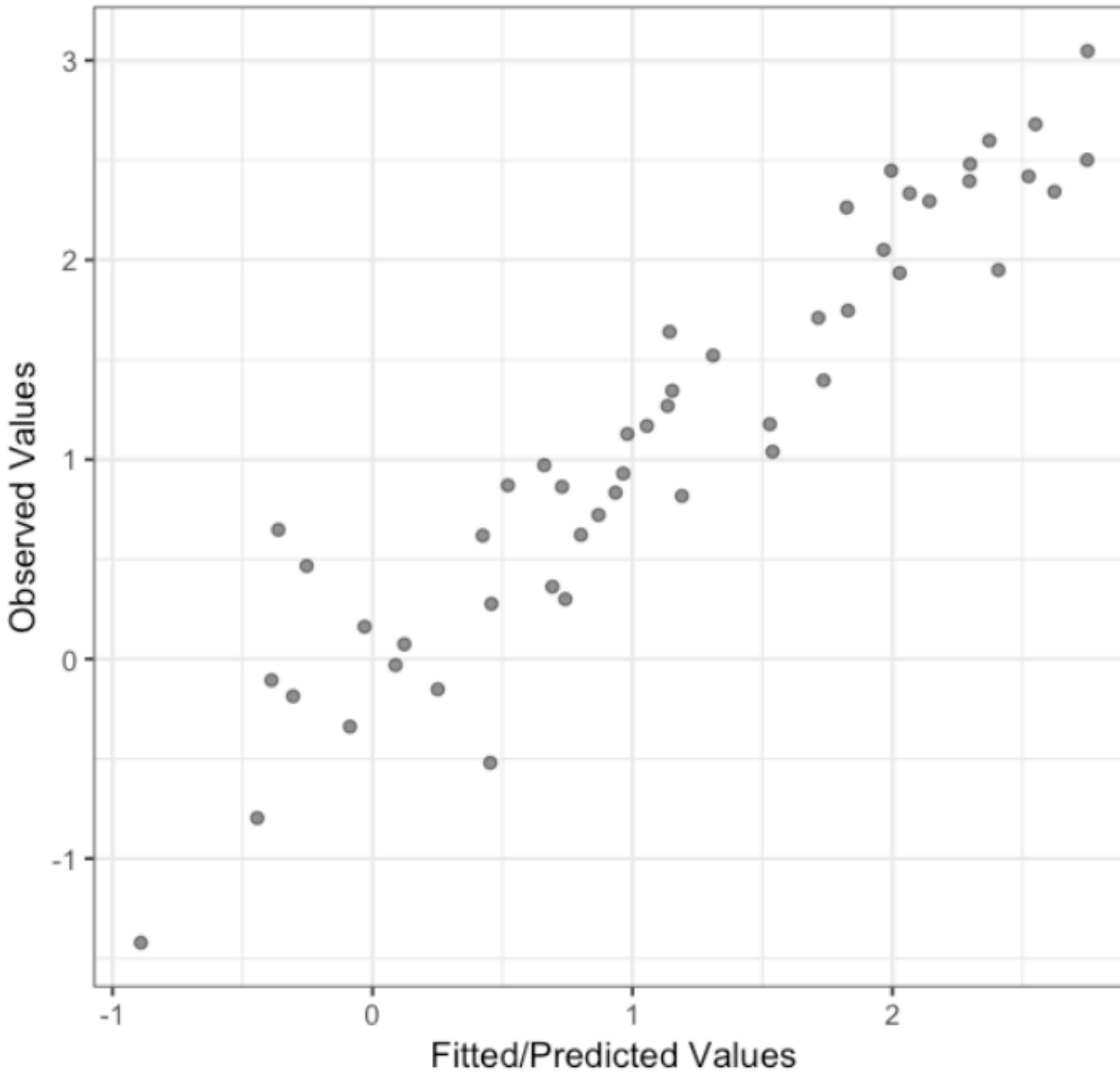
If this assumption is not met:

1. The LS estimator is not unbiased.
2. The model will have low explanatory power.
3. The inferences (e.g., t-tests, F-test) based on the misspecified regression model will be biased and misleading.
4. The predictions from the model are less likely to be accurate.

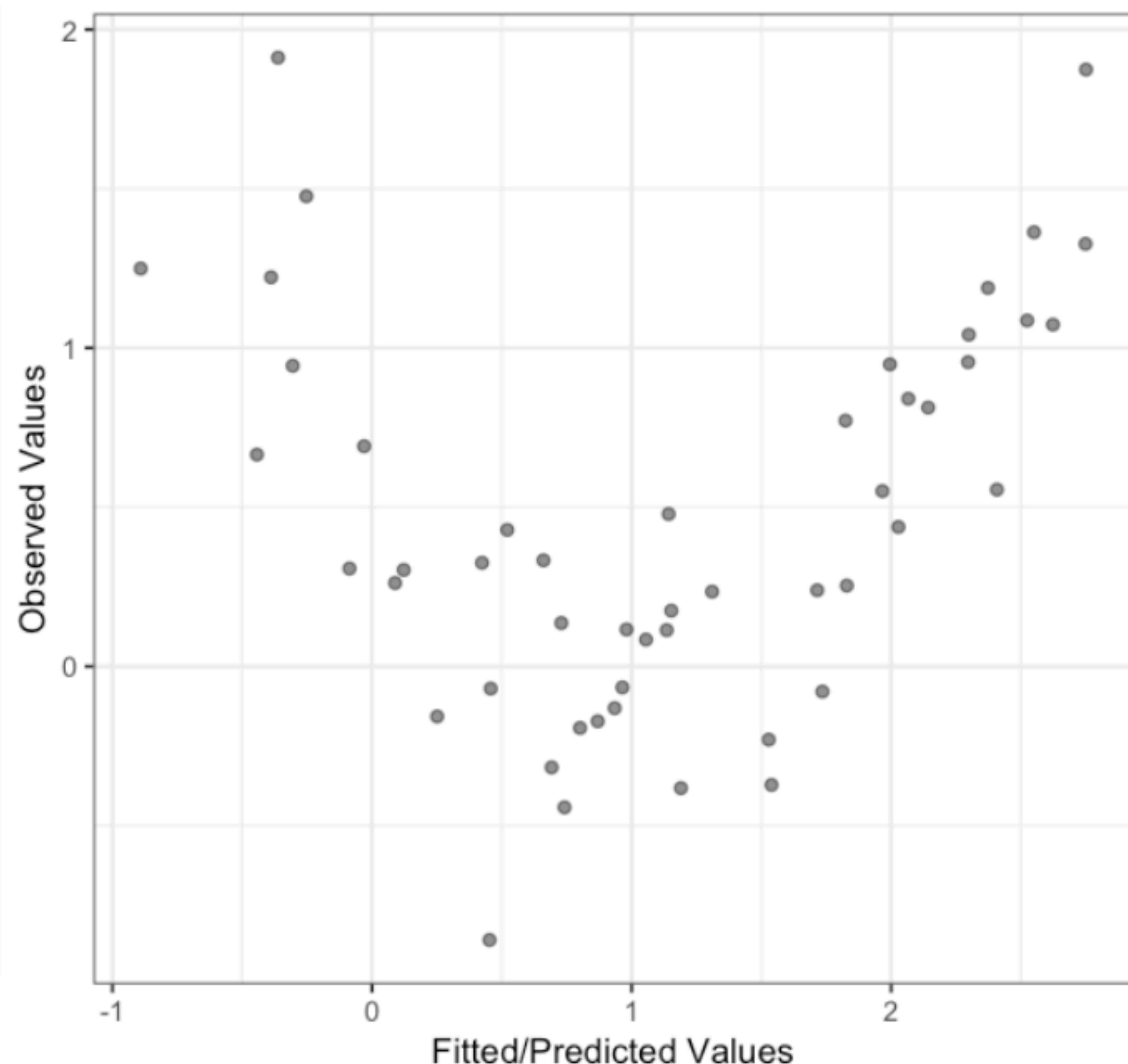


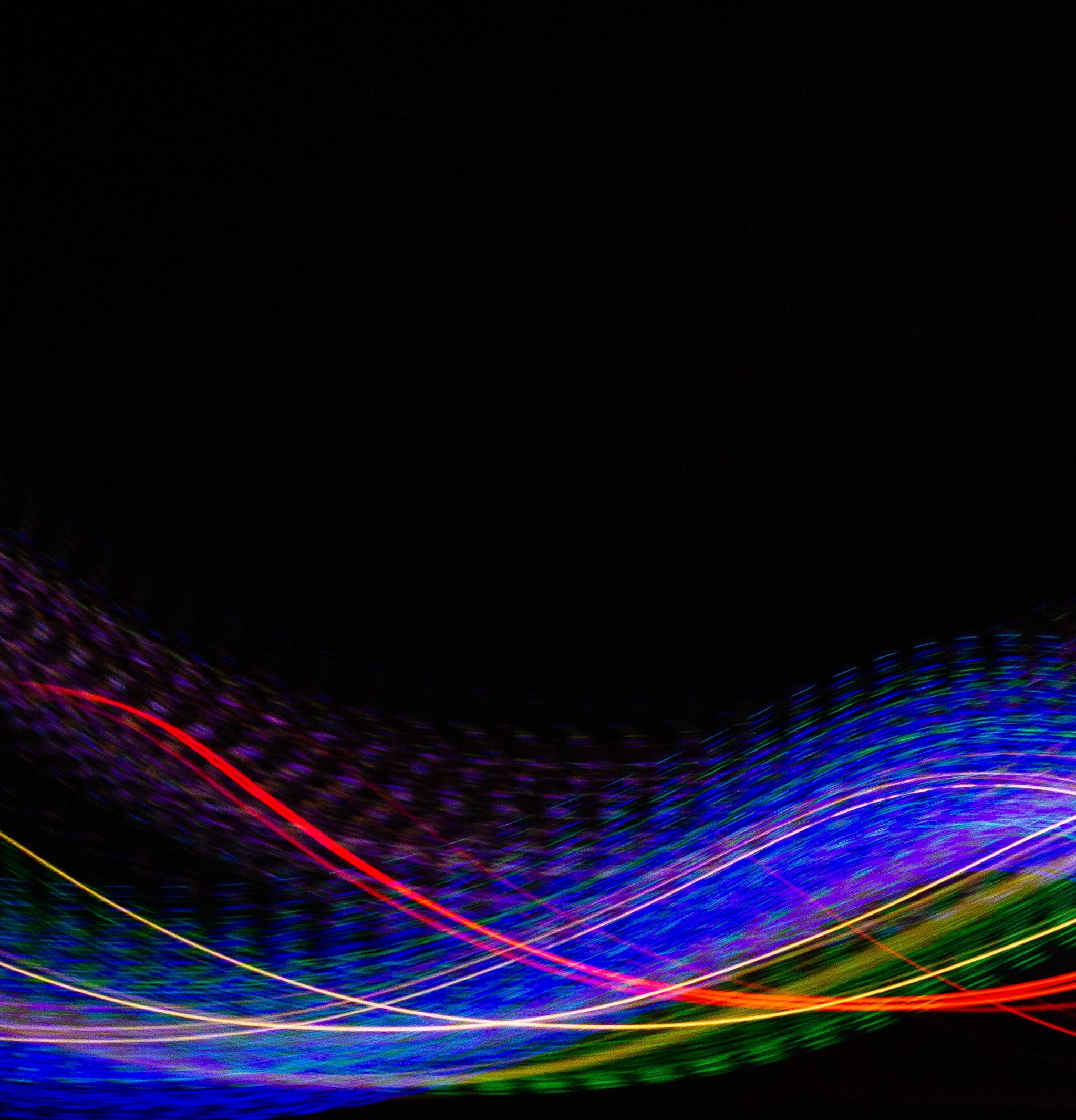
Fitted/Predicted vs. observed value plots

Correct Model



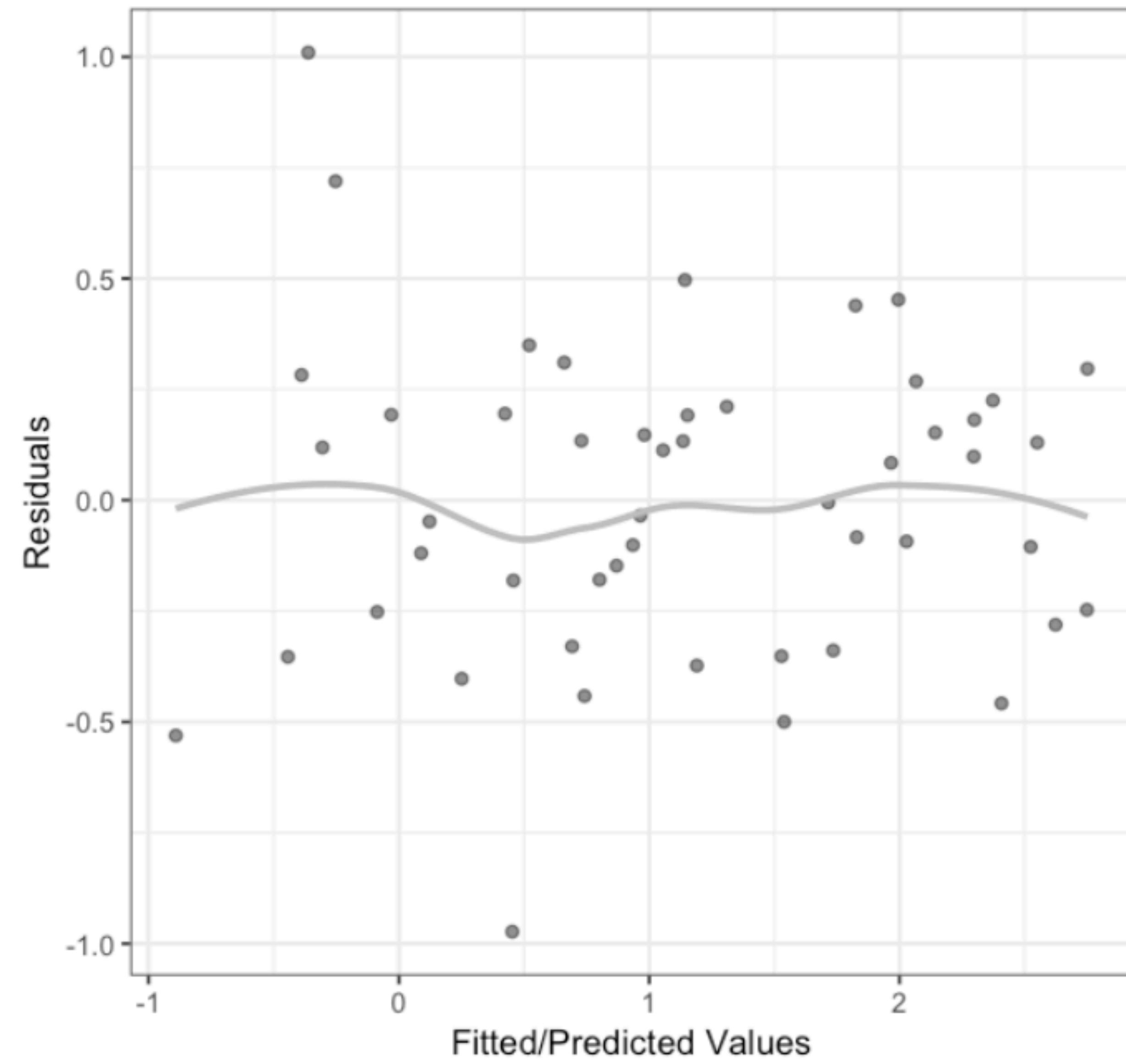
Incorrect Model



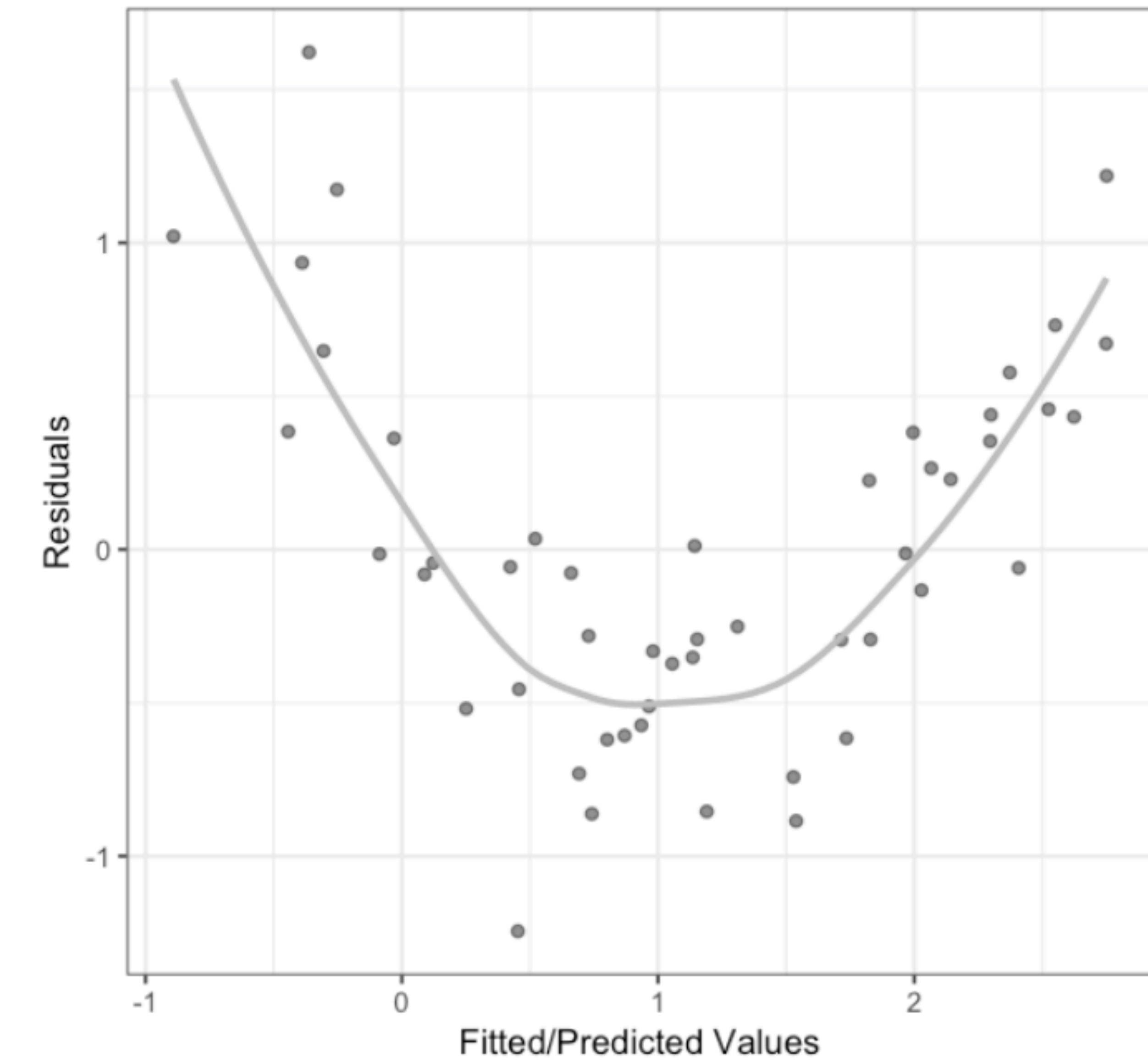


Residual vs. fitted/predicted value plots

Correct Model



Incorrect Model

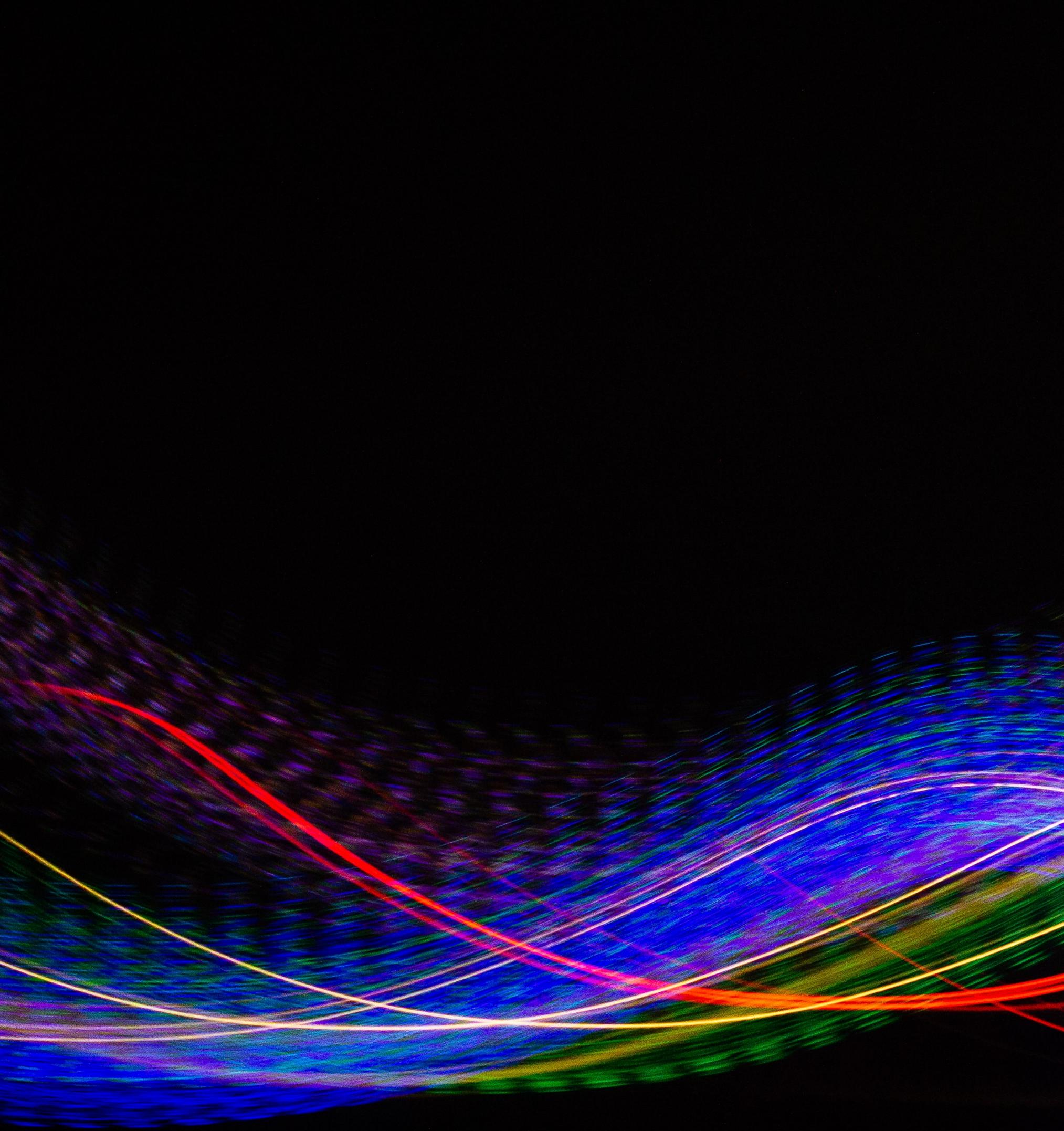


Lesson: Violations of the independence assumption

Module: Regression Diagnostics

Photo by [Janine Robinson](#) on [Unsplash](#)





Uncorrelated/Independent Errors. This violation might occur, and be serious, when measuring *spatial* or *temporal* data.

1. Residuals vs time/index.
2. Successive residual plot.
3. Durbin-Watson test.

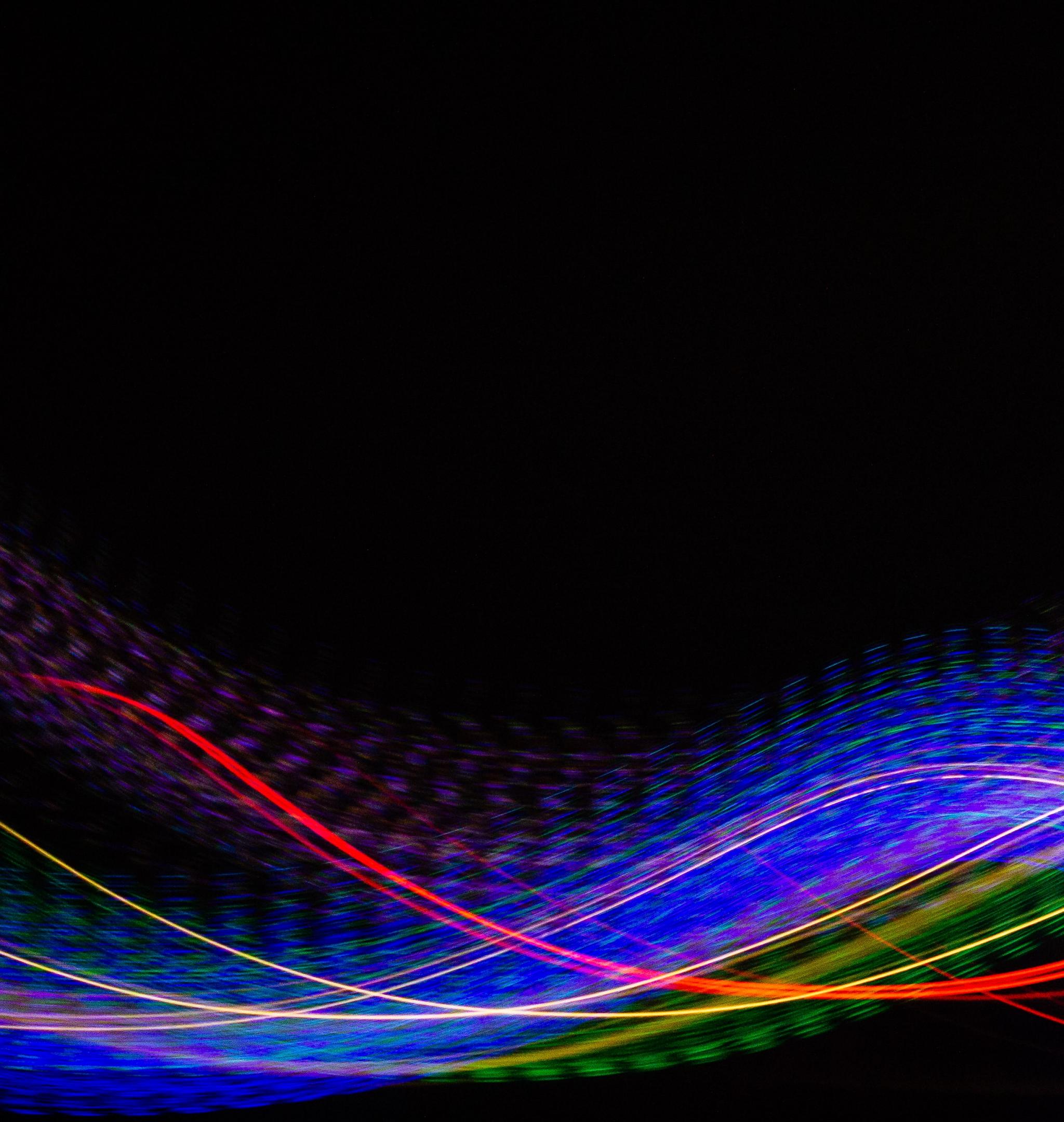
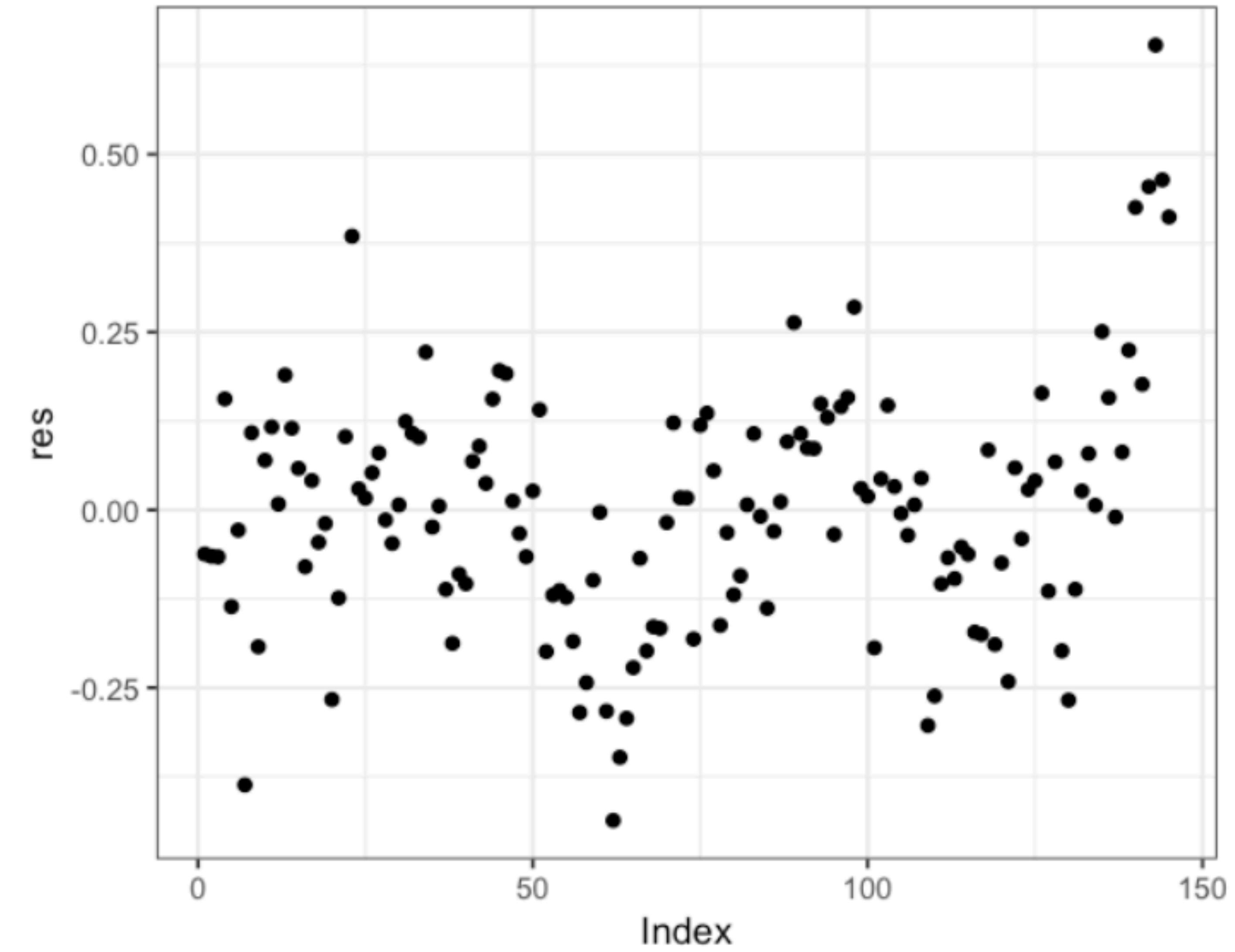
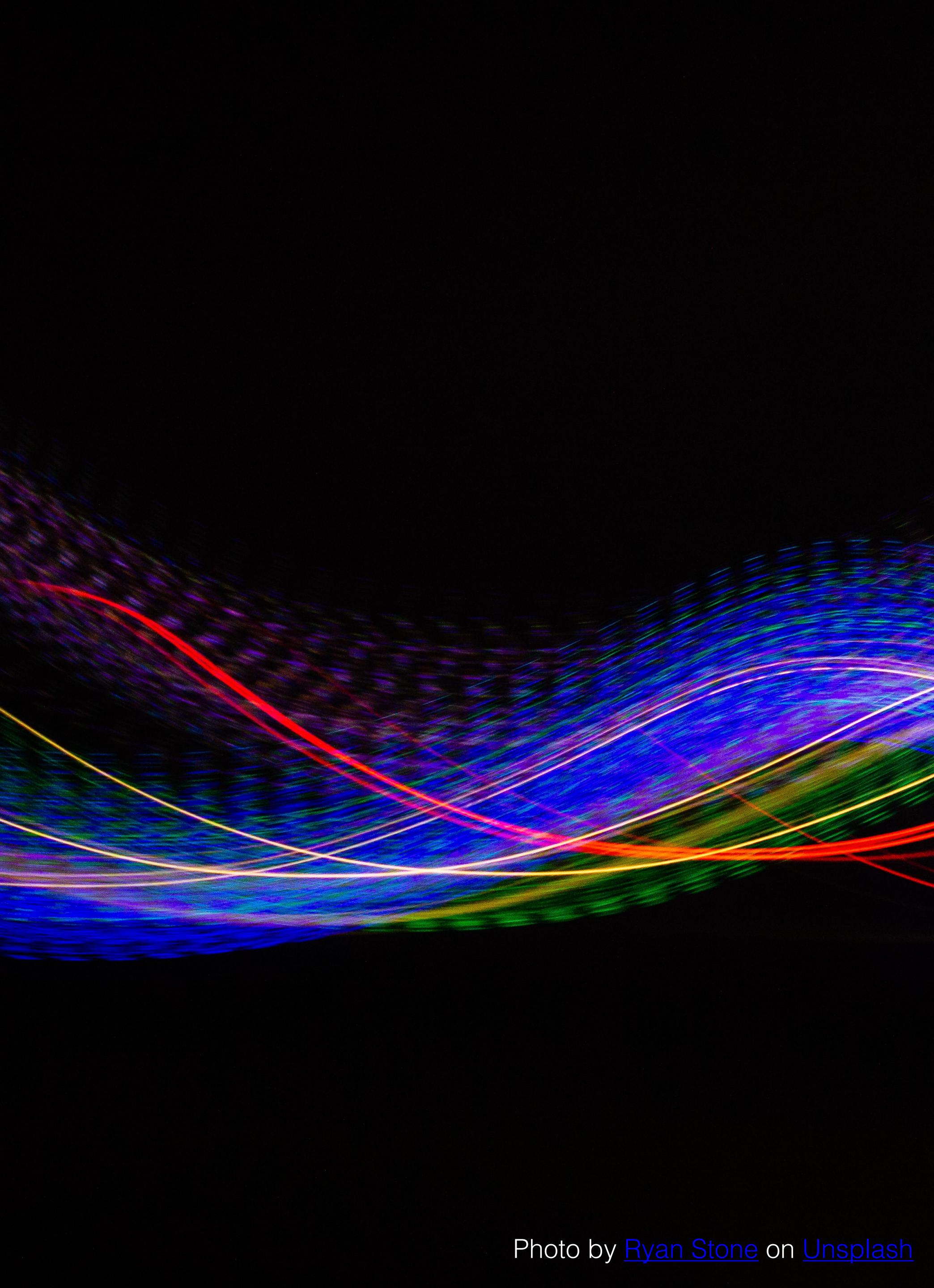


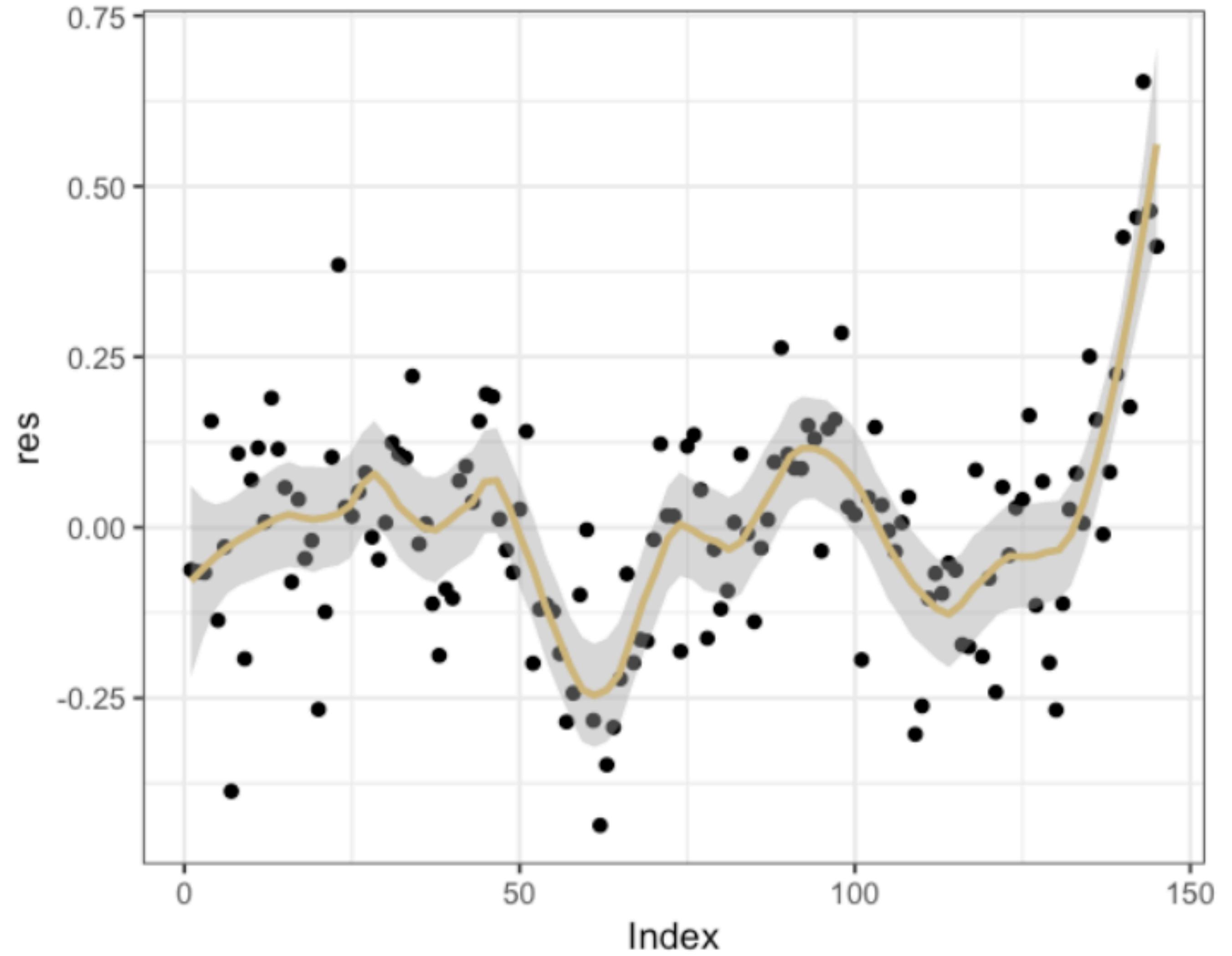
Photo by [Ryan Stone](#) on [Unsplash](#)

Residual vs index plot





Residual vs index plot



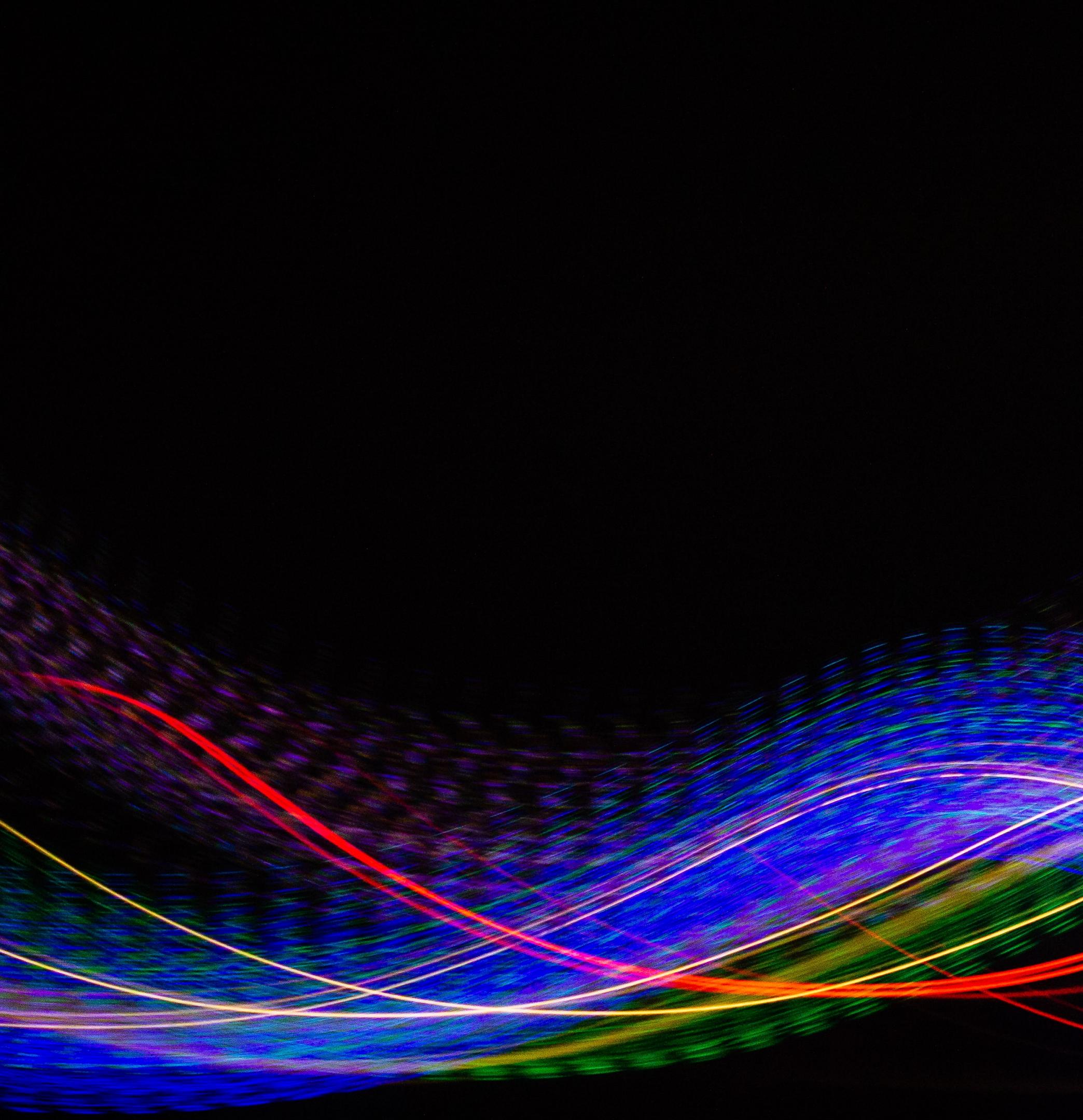
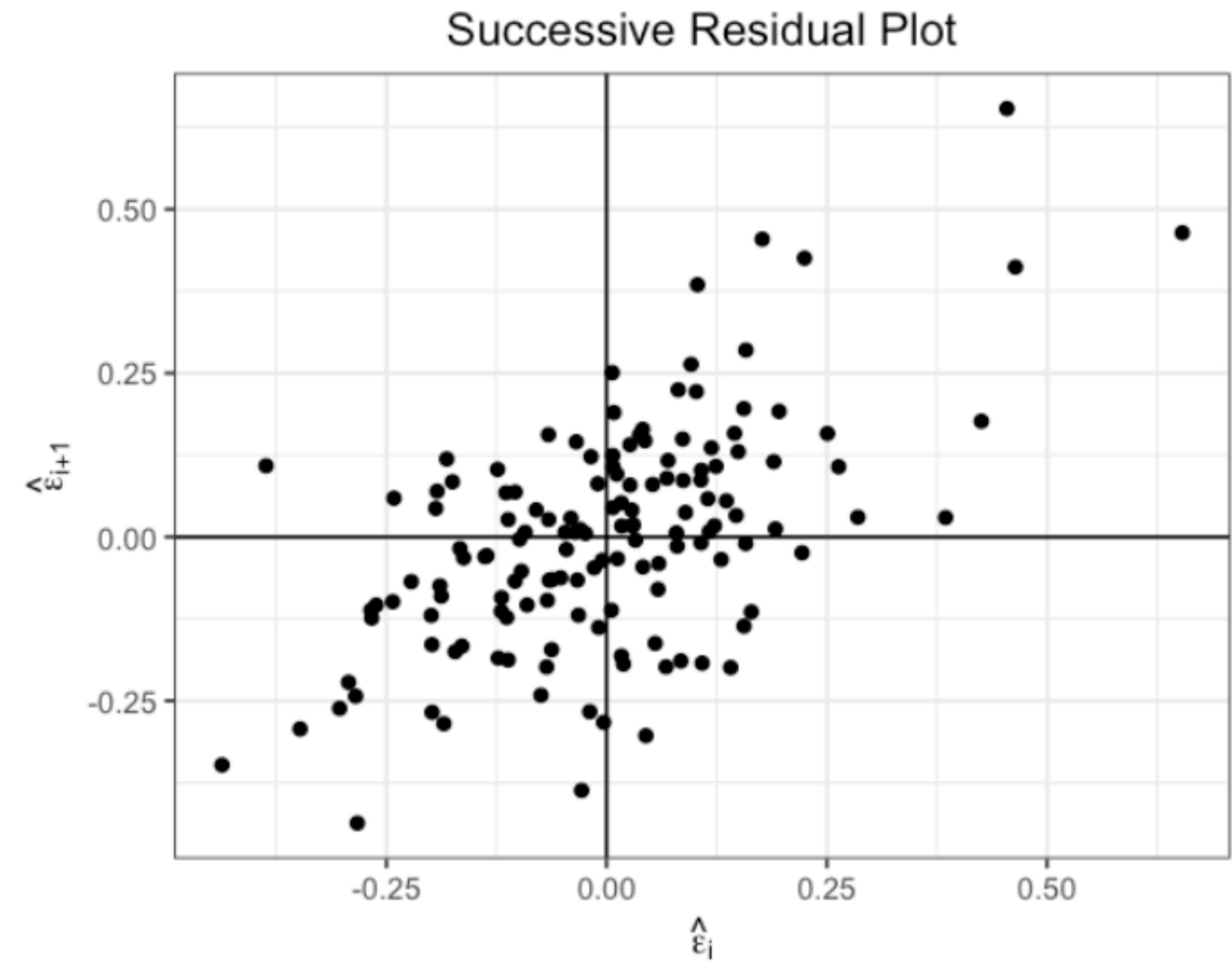
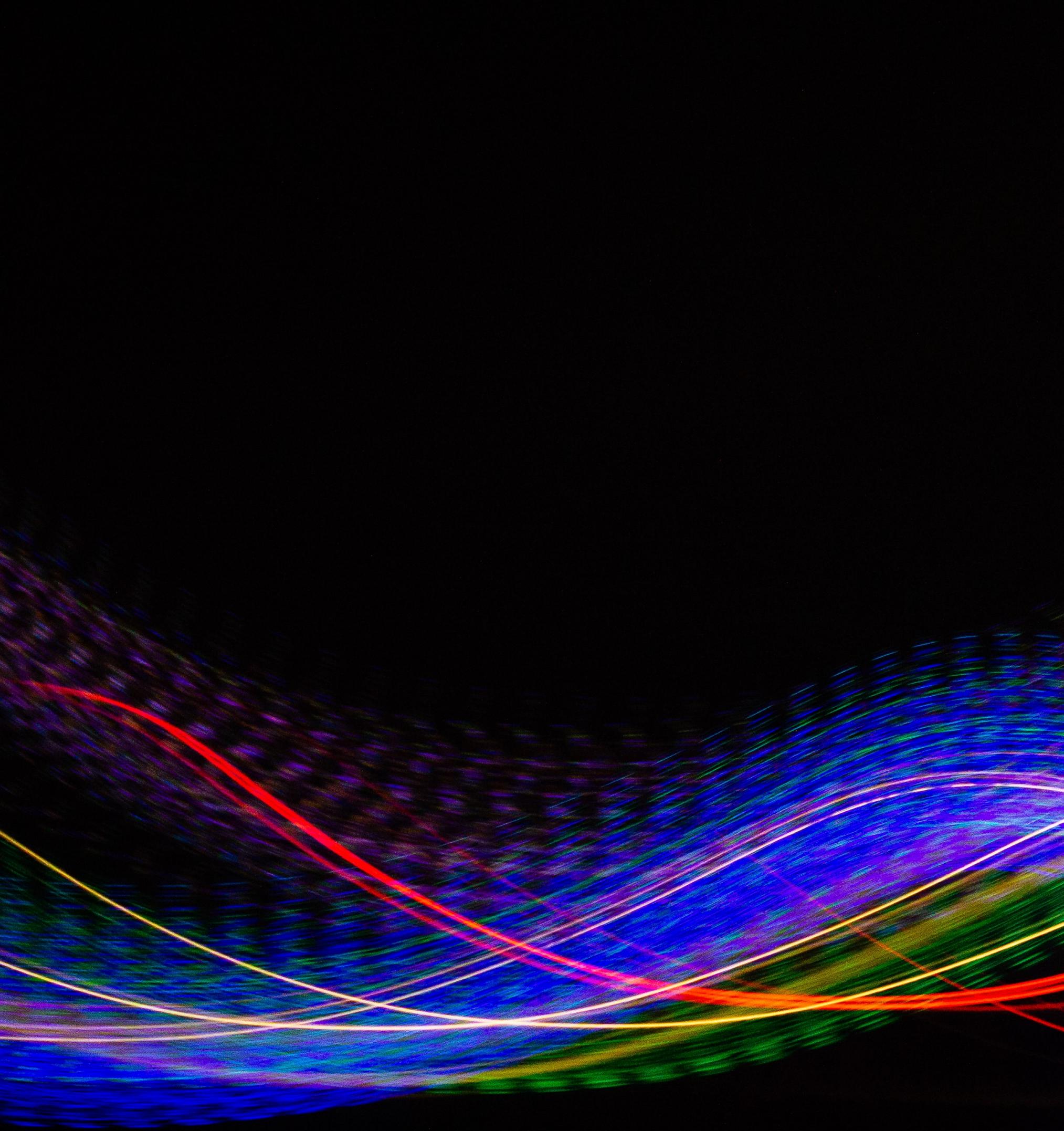


Photo by [Ryan Stone](#) on [Unsplash](#)





What solutions are available for violations
of this assumption?

1. Generalized least squares
2. Time Series/Spatial Statistics

Lesson: Violations of the constant variance assumption

Module: Regression Diagnostics

Photo by [Janine Robinson](#) on [Unsplash](#)





Photo by [Serenity Mitchell](#) on [Unsplash](#)

Constant Variance Assumption. Recall that one of our assumptions is:

Even under this assumption, the residuals—our estimator of the error—are not constant variance:



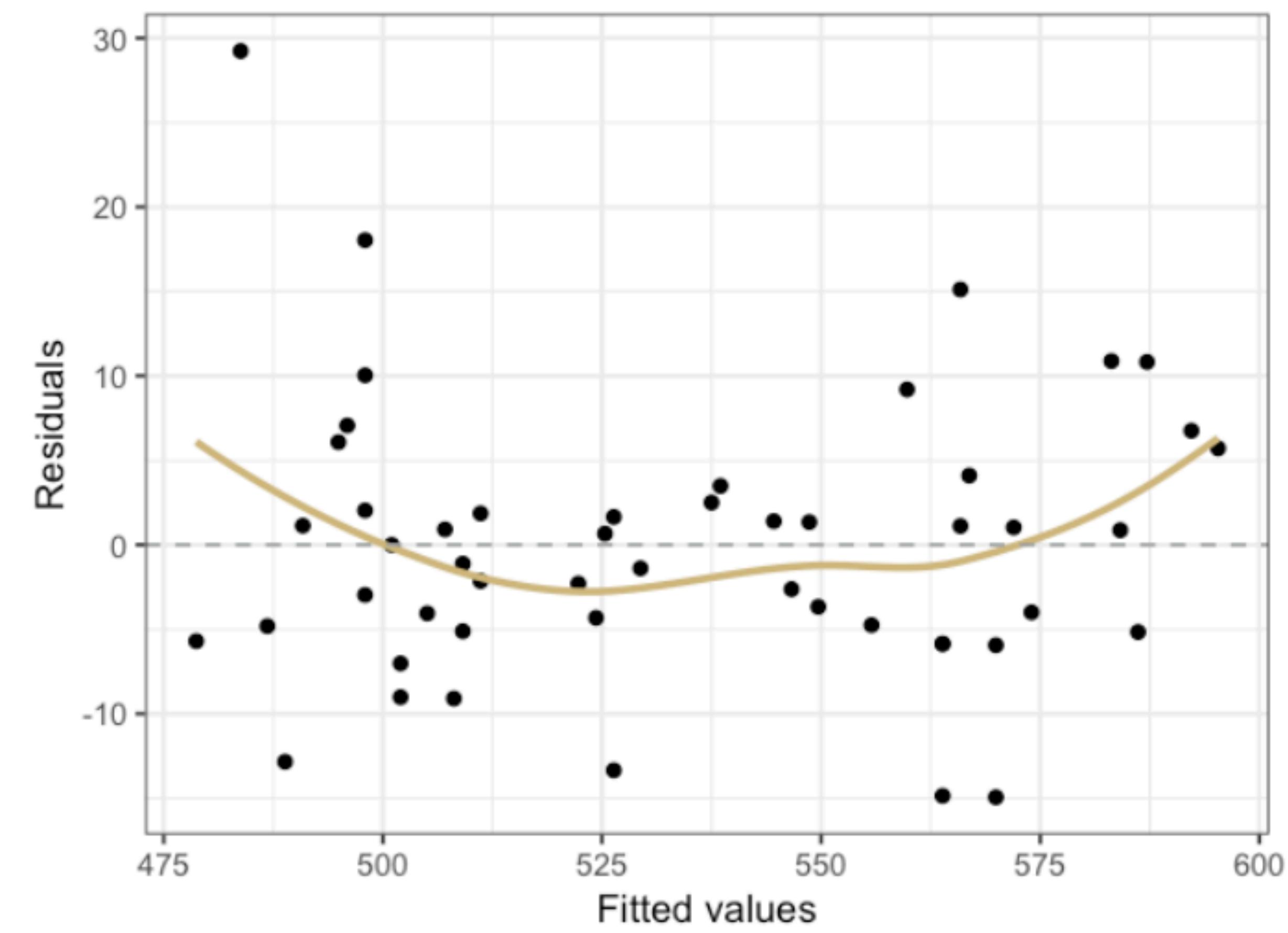
Photo by [Serenity Mitchell](#) on [Unsplash](#)

Constant Variance Diagnostics:

1. Residual vs. Fitted plot.
2. Residual vs. potential predictor plot.

Checking Error Assumptions

Residual vs Fitted Plot



Residual vs Fitted Plot

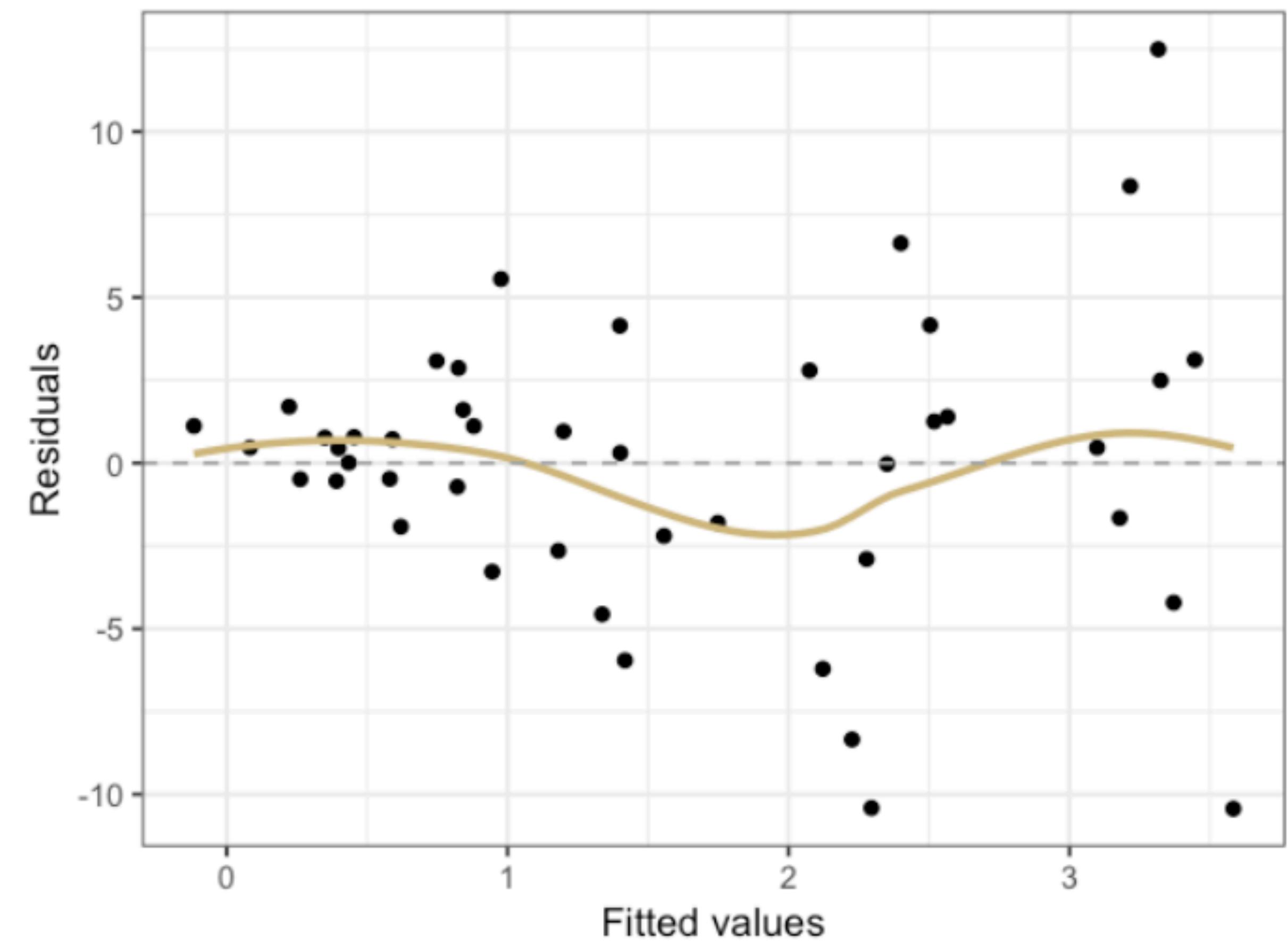




Photo by [Serenity Mitchell](#) on [Unsplash](#)

Solutions to non constant variance:

1. Transformation of the response.
2. Weighted/Generalized Least Squares.

Lesson: Violations of the normality assumption

Module: Regression Diagnostics

Photo by [Janine Robinson](#) on [Unsplash](#)





Normality Assumption.

To diagnose deviations from normality, we can look at:

1. QQ plots
2. Residual vs fitted plots
3. Shapiro-Wilk test for normality

(Height of Hallgrímskirkja in Reykjavík, Iceland)



Correct Model

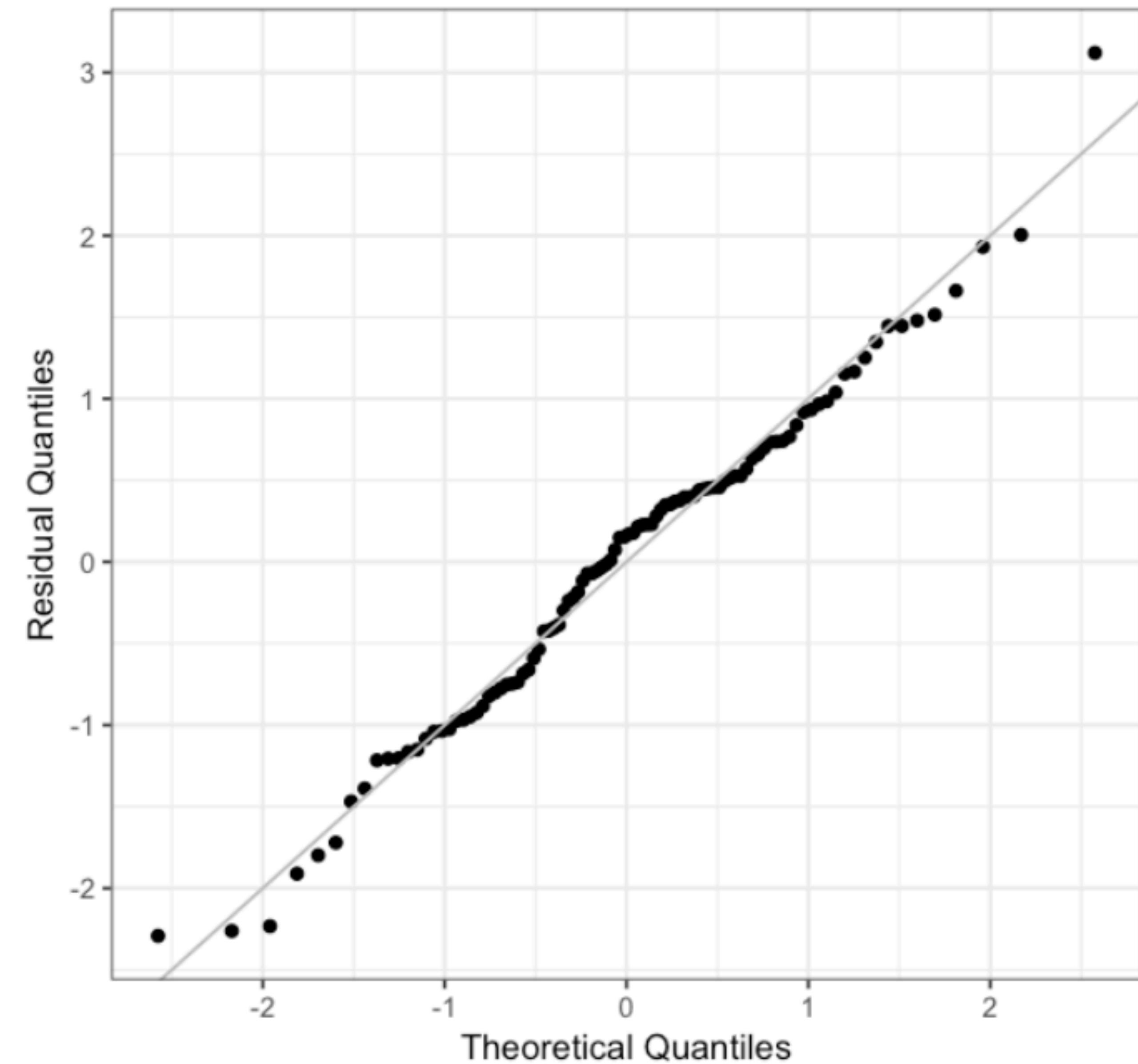
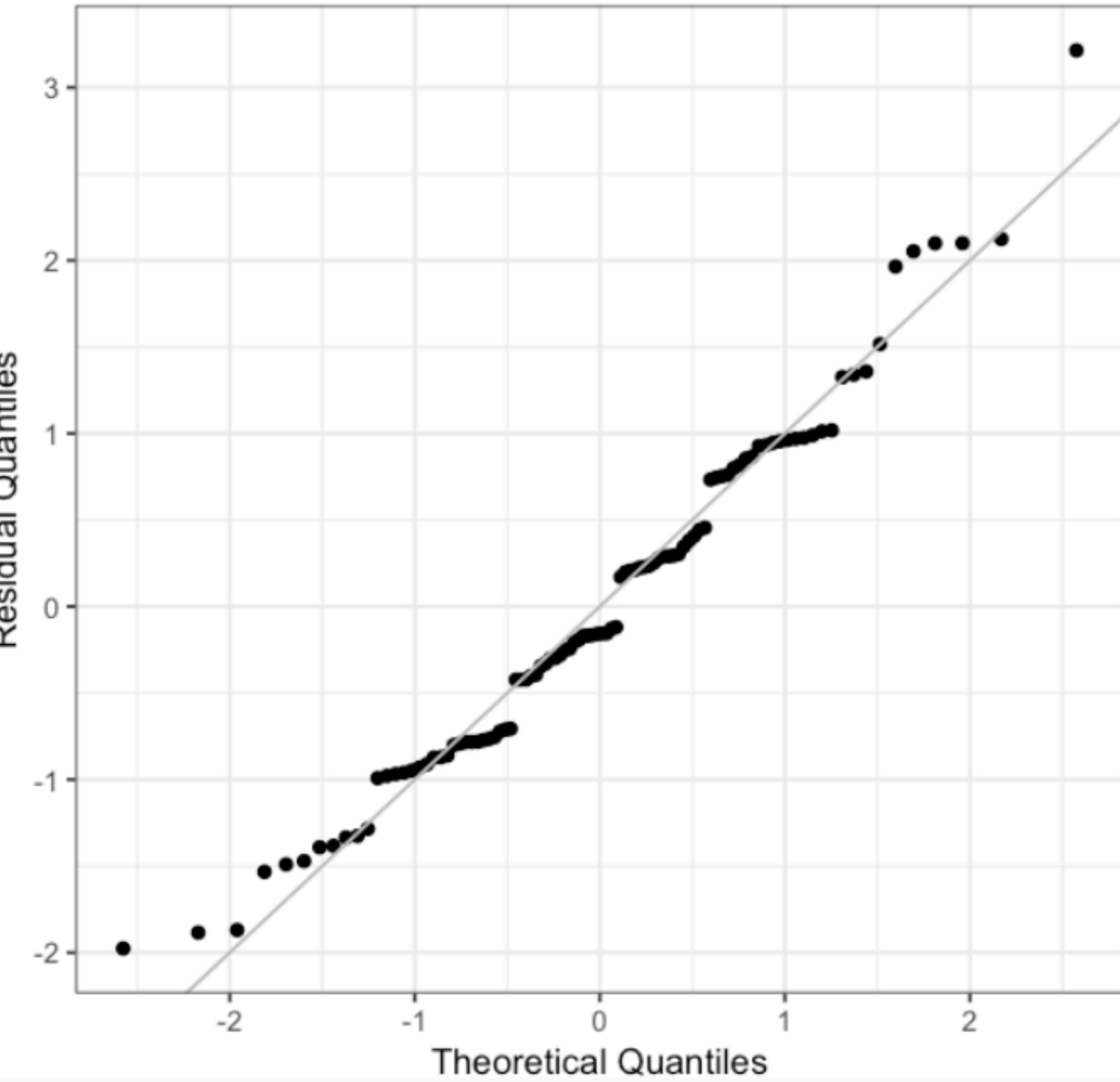


Photo by [Cassidy Dickens](#) on [Unsplash](#)

(Height of Hallgrímskirkja in Reykjavík, Iceland)



Incorrect Model



(Height of Hallgrímskirkja in Reykjavík, Iceland)



Correct Model

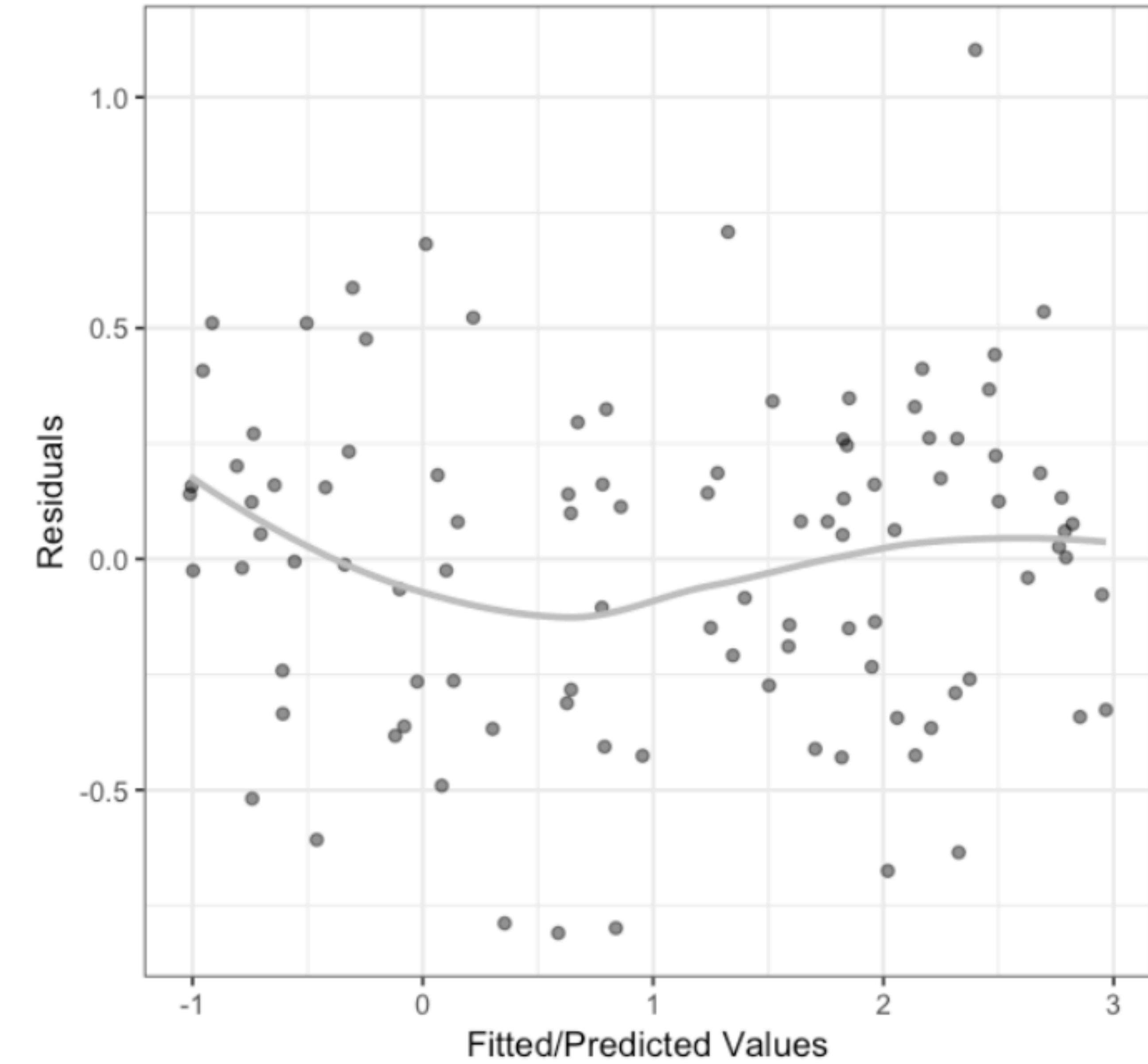
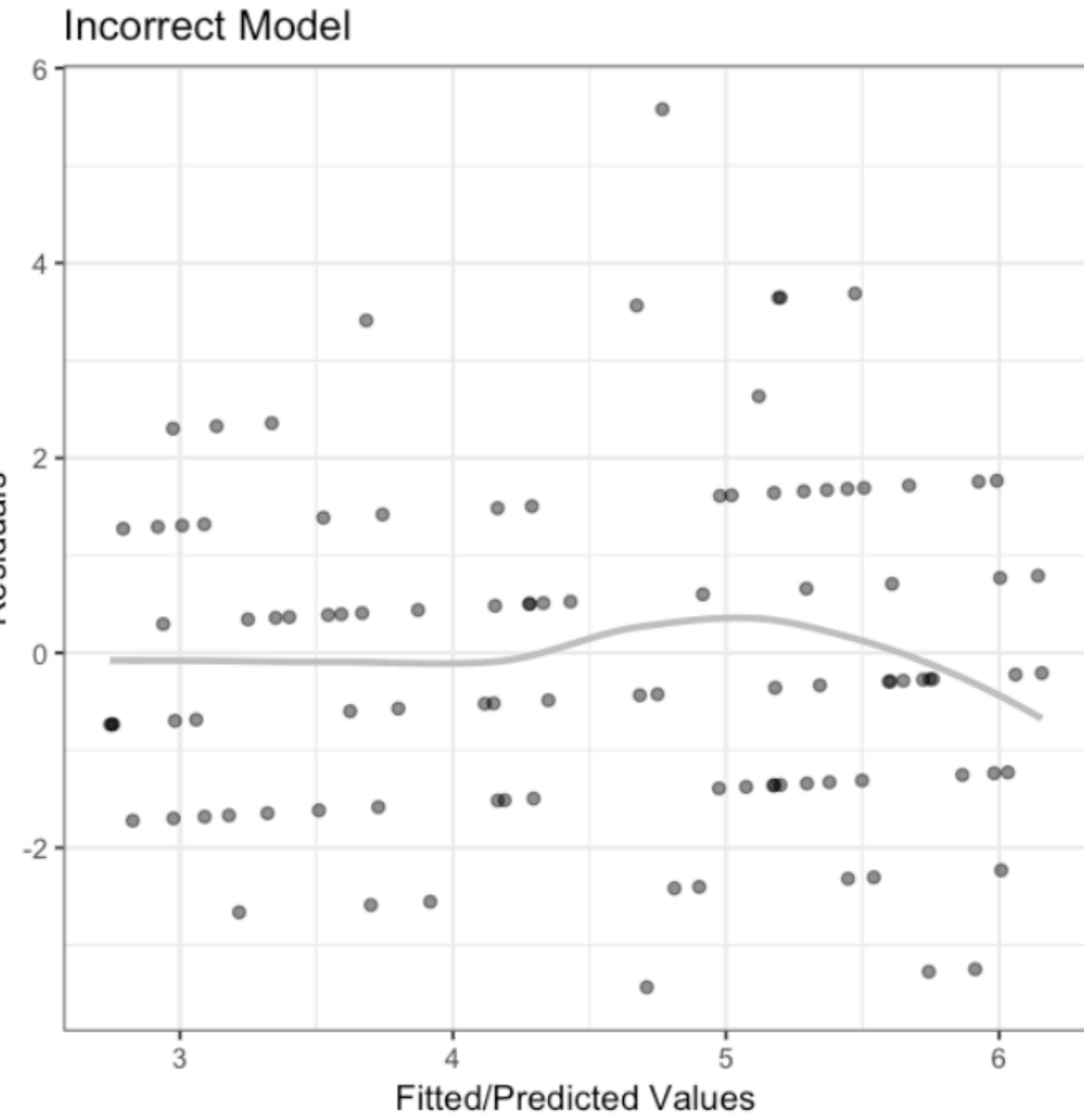


Photo by [Cassidy Dickens](#) on [Unsplash](#)

(Height of Hallgrímskirkja in Reykjavík, Iceland)





H_0 : The sample (residuals) came from a normal distribution

H_1 : The sample (residuals) came from a non-normal distribution

```
shapiro.test(df.diagnostics$r_correct)
```

Shapiro-Wilk normality test

```
data: df.diagnostics$r_correct  
W = 0.98711, p-value = 0.4449
```

(Height of Hallgrímskirkja in Reykjavík, Iceland)



Solutions:

1. Transformations.
2. *Generalized linear models (GLMs)*