

# Public Health Risk in Bear Creek

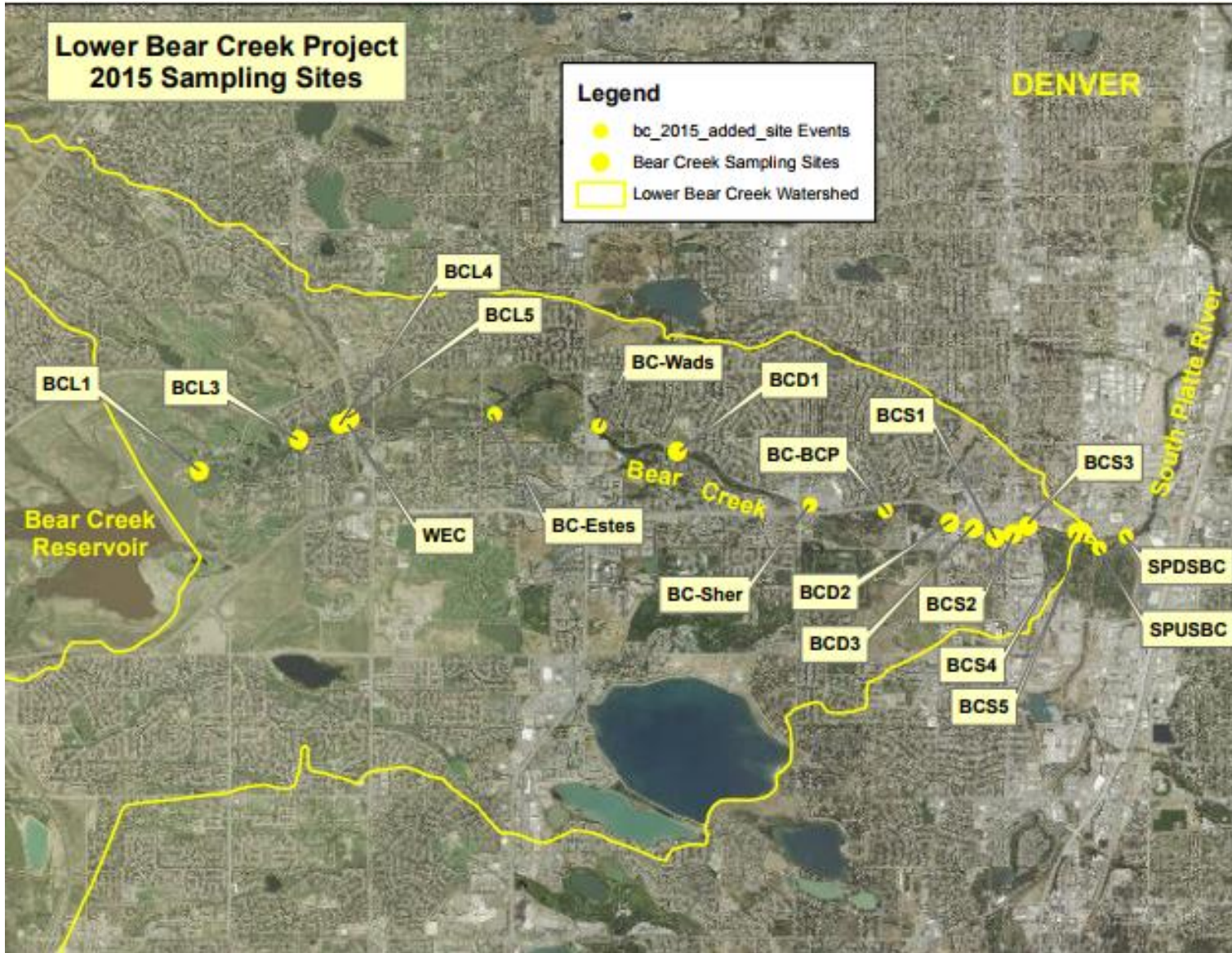
*A Time Series Analysis Approach to Estimating E. coli Levels in Bear Creek Watershed: Implications for Further Study.*

**Christopher Campbell, Ahern Nelson, and Keenan O'Brian,  
with Faculty Advisor Dr. Elizabeth Ribble**





# Introduction to Project



Sampling  
Sites

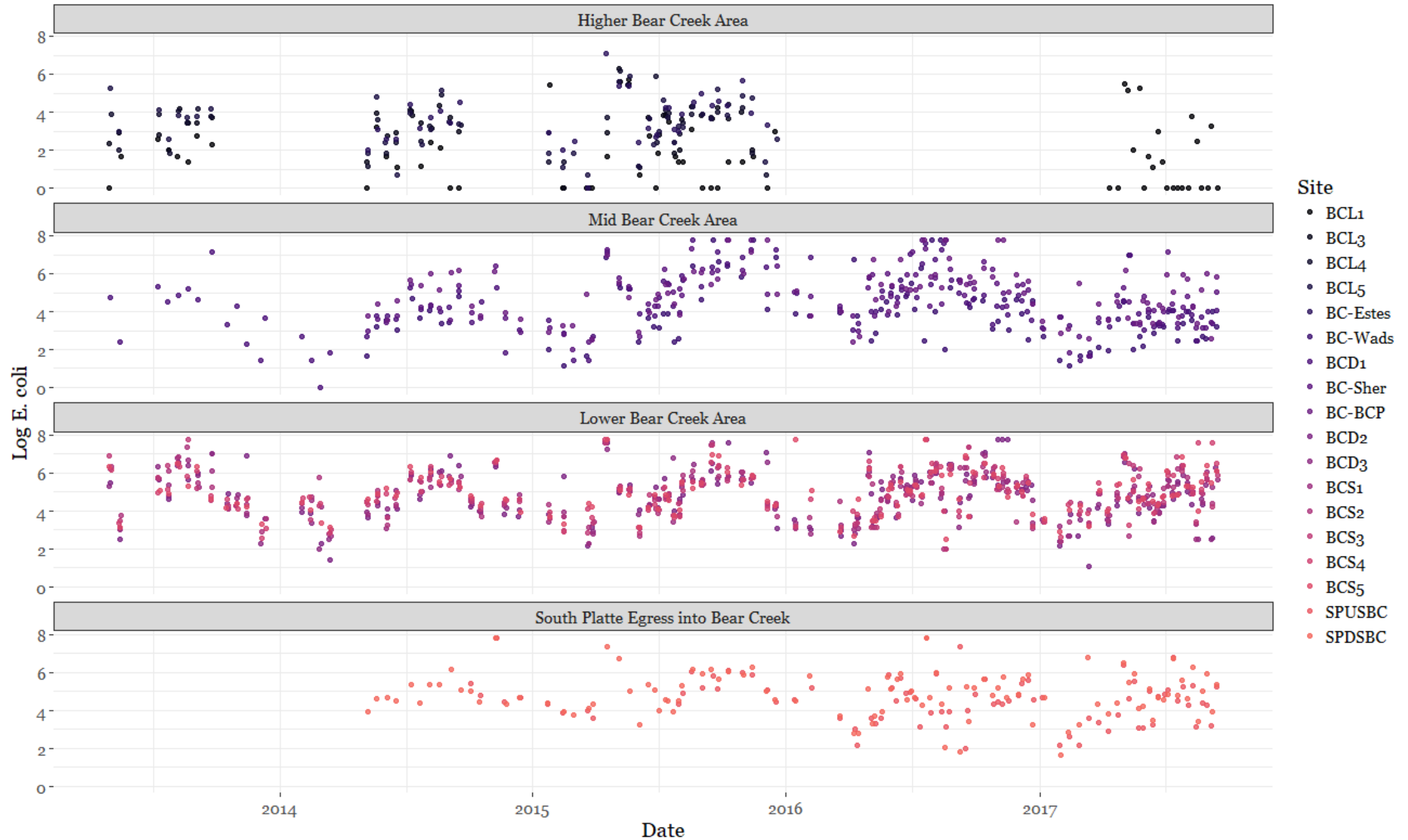
# Aims of Project

1. Understanding risk posed to humans by bacteria similar to E. coli
2. Determine possible predictors of E. coli
3. Understand variation we see in the data

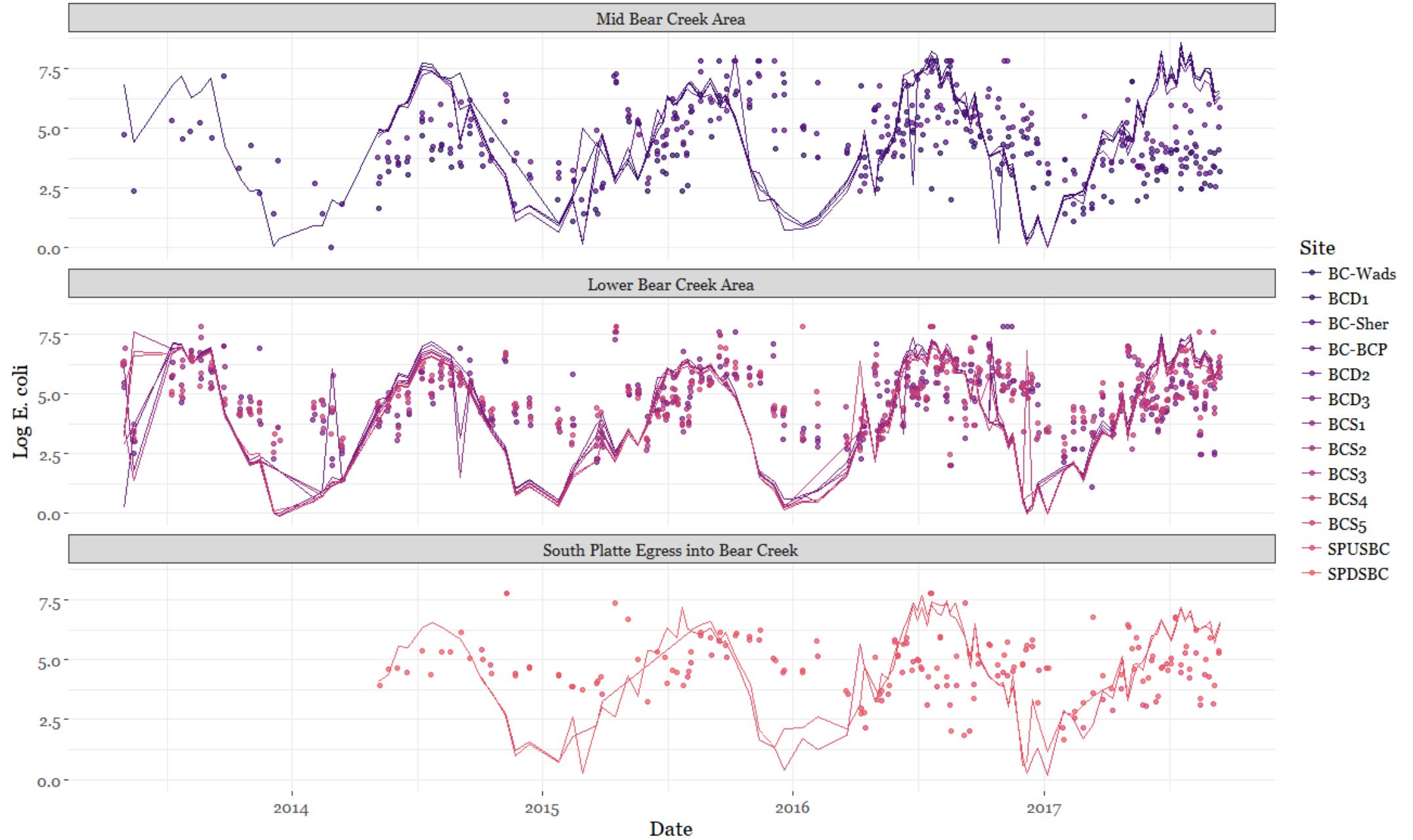
A thin vertical black line is positioned to the left of the text, extending from the top of the first line of text to the bottom of the second line.

# Exploratory Data Analysis

# E. coli Over Time Per Geographic Area



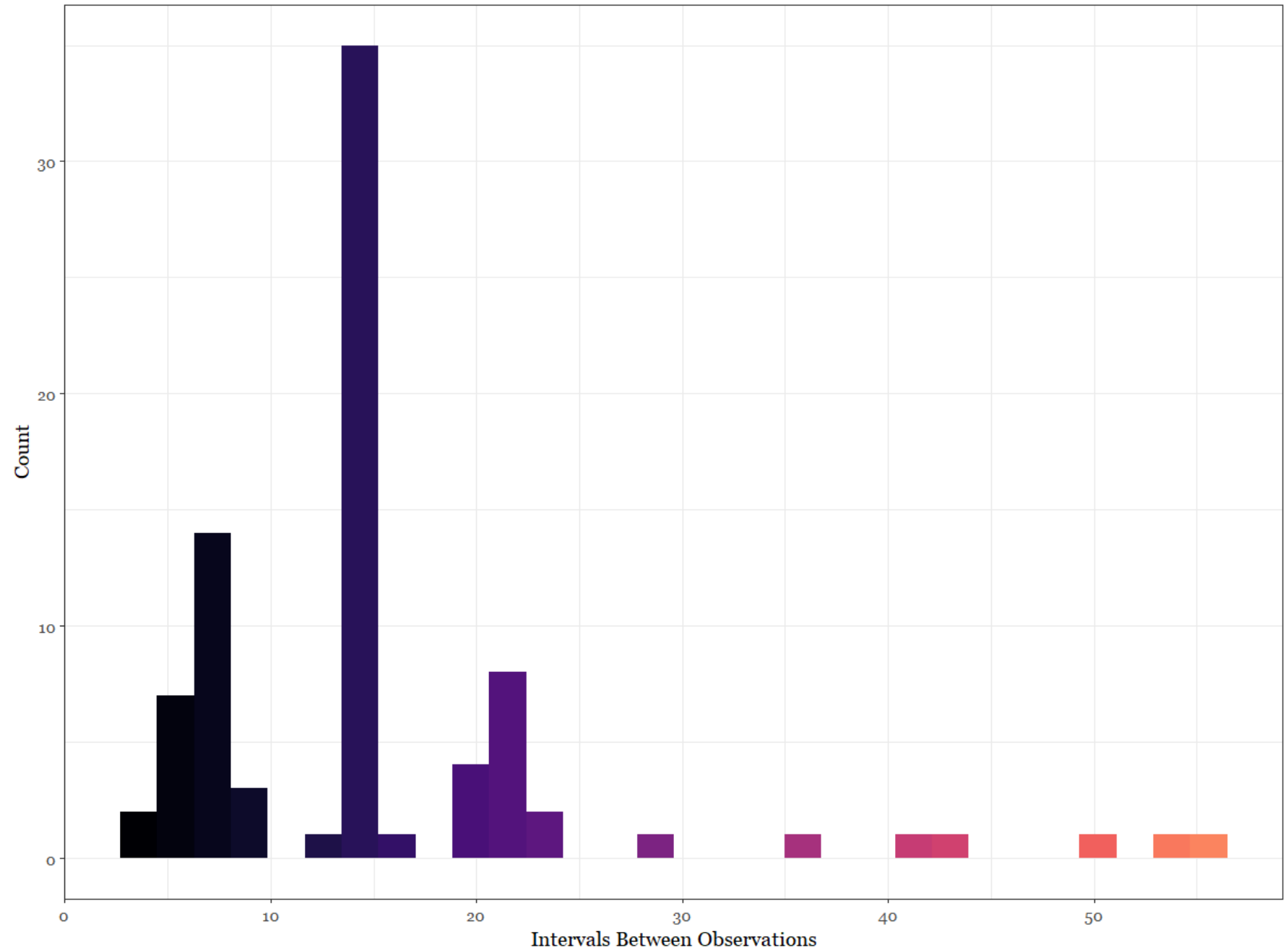
# E. coli Over Time Per Geographic Area



# Some Control Issues

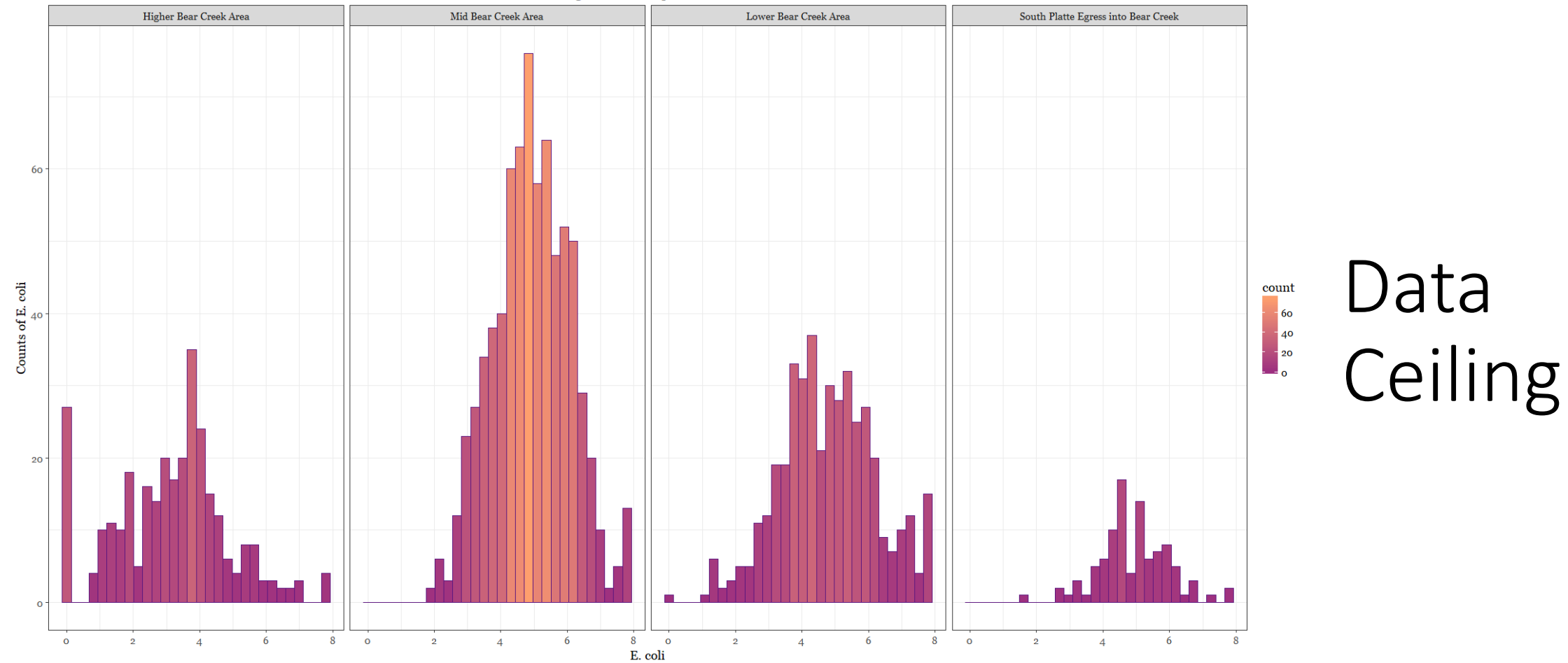
- Data collected over sporadic intervals...

- ...Can bin almost all data bi-weekly.





Histogram of E. coli per Binned Site



# Past Modeling Attempts and Independence

- Statistical tests rely on the assumption of independence of the data. This includes ANOVA and T-Tests.
- Also, regressions will not produce accurate t-values when the data has autocorrelation.

# We may want to try regression...

---

```
Call:
lm(formula = logEColi ~ tempC + pH + turbidity)

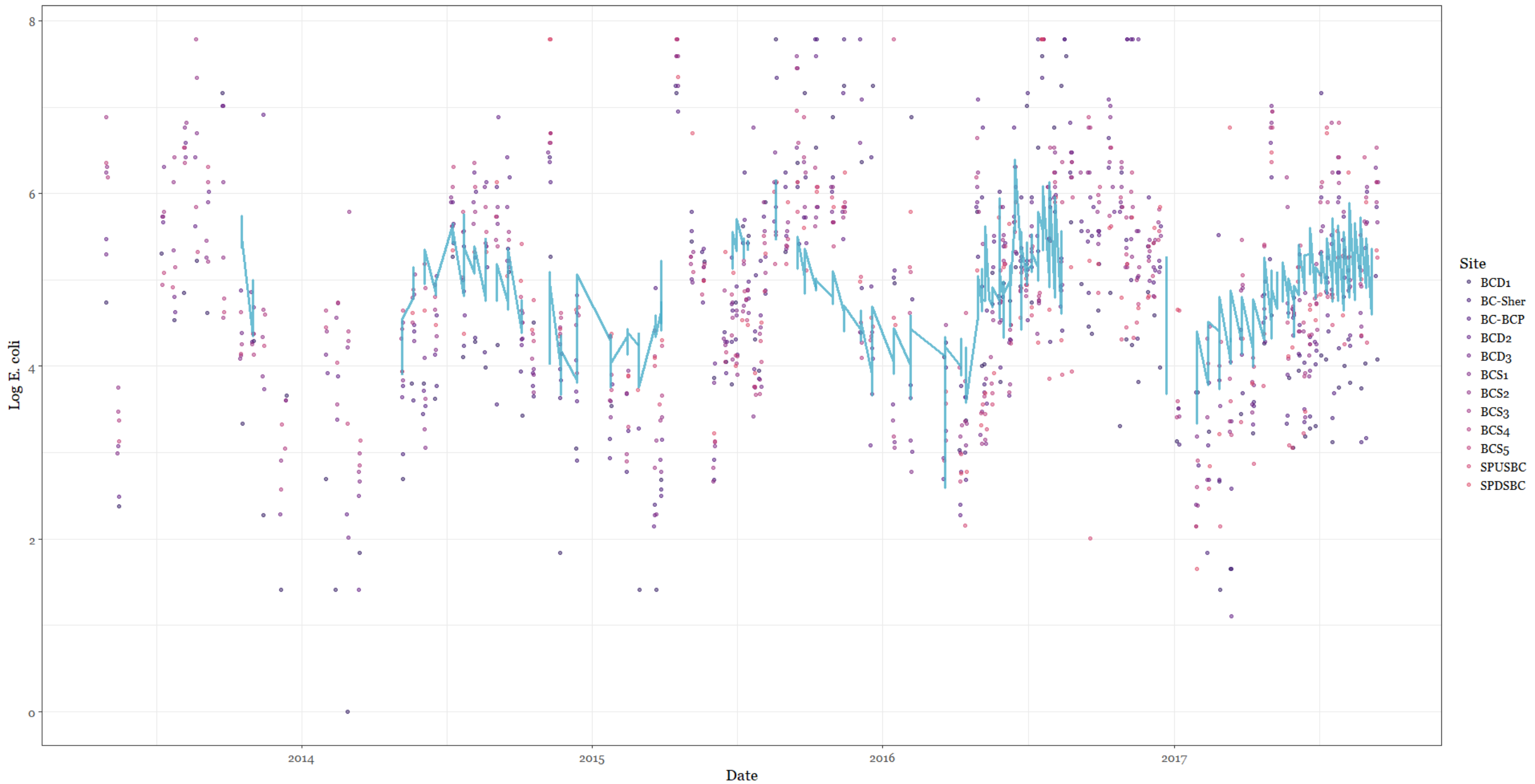
Residuals:
    Min       1Q   Median       3Q      Max
-3.0745 -0.7365 -0.0902  0.6455  3.8253

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 10.285571  0.826594  12.443 < 2e-16 ***
tempC        0.065491  0.006478  10.109 < 2e-16 ***
pH          -0.837798  0.105711  -7.925 7.21e-15 ***
turbidity    0.027899  0.006466   4.315 1.79e-05 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

This regression has highly significant coefficients,

**however....**

Linear Regression Predictions for E. coli



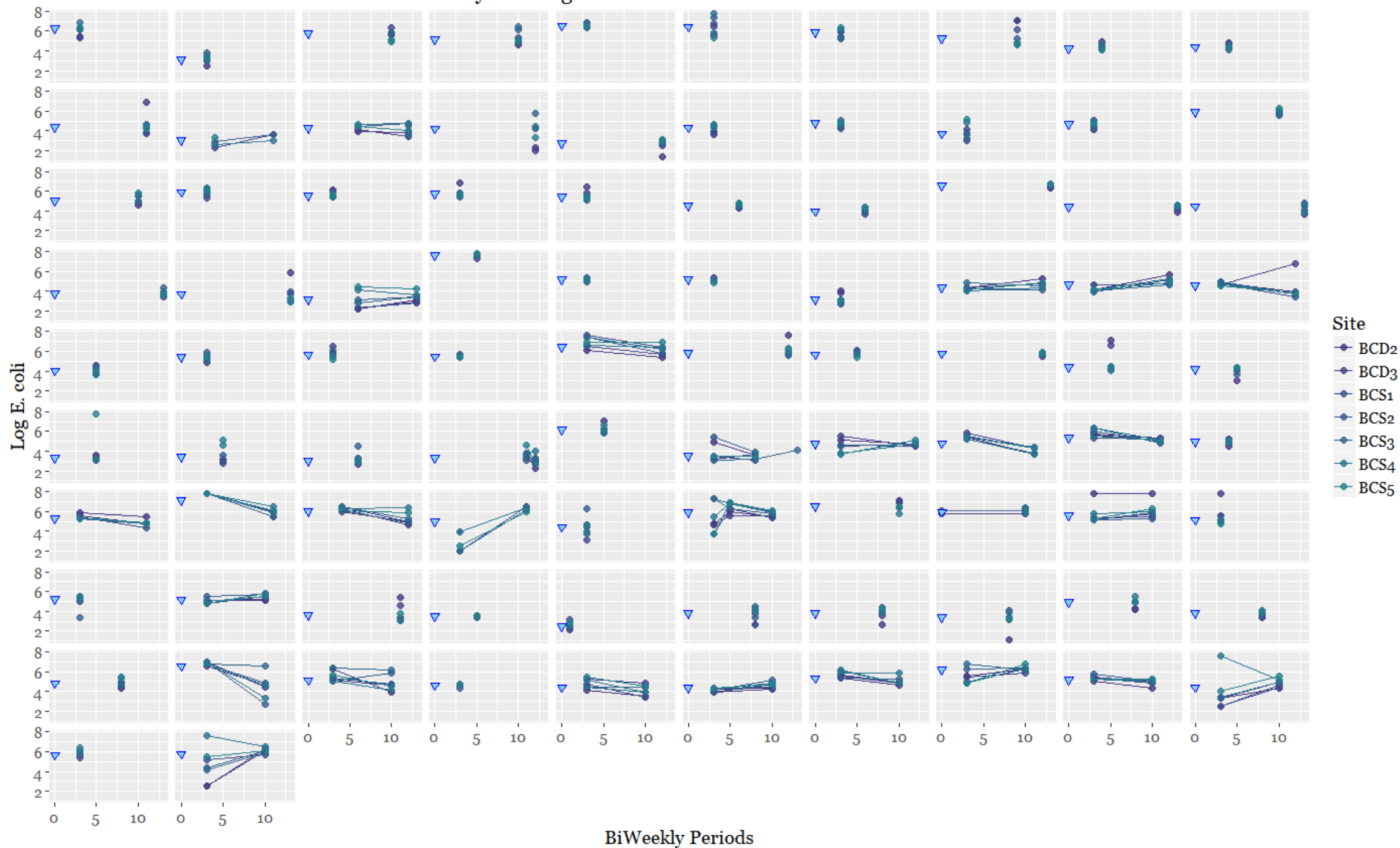
A thin vertical black line is positioned to the left of the title text.

# Time Series Analysis

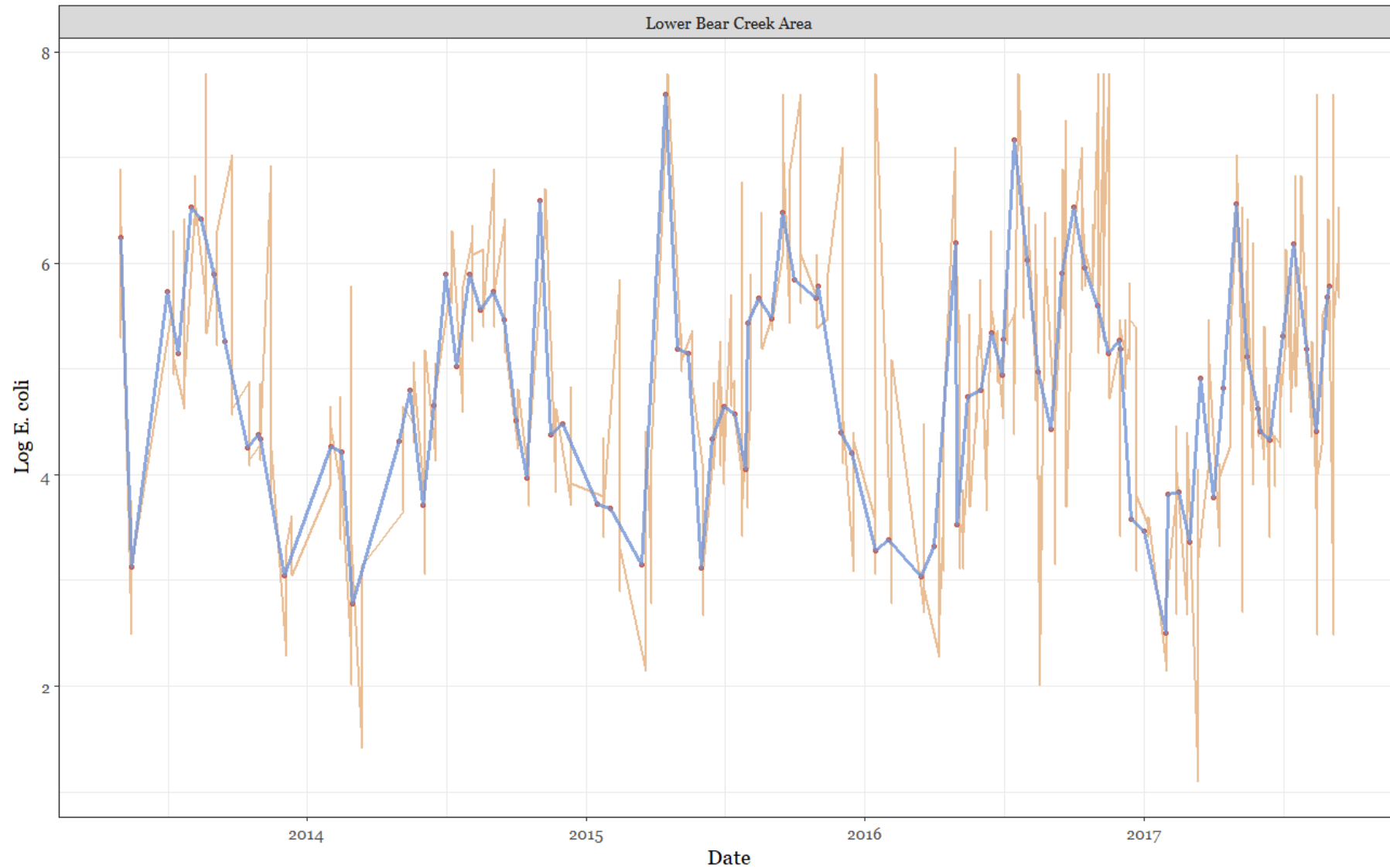
# E. coli Over Time Per Geographic Area



BiWeekly Readings of *E. coli* in the Lower Bear Creek Area



# Median as a Binning Statistic

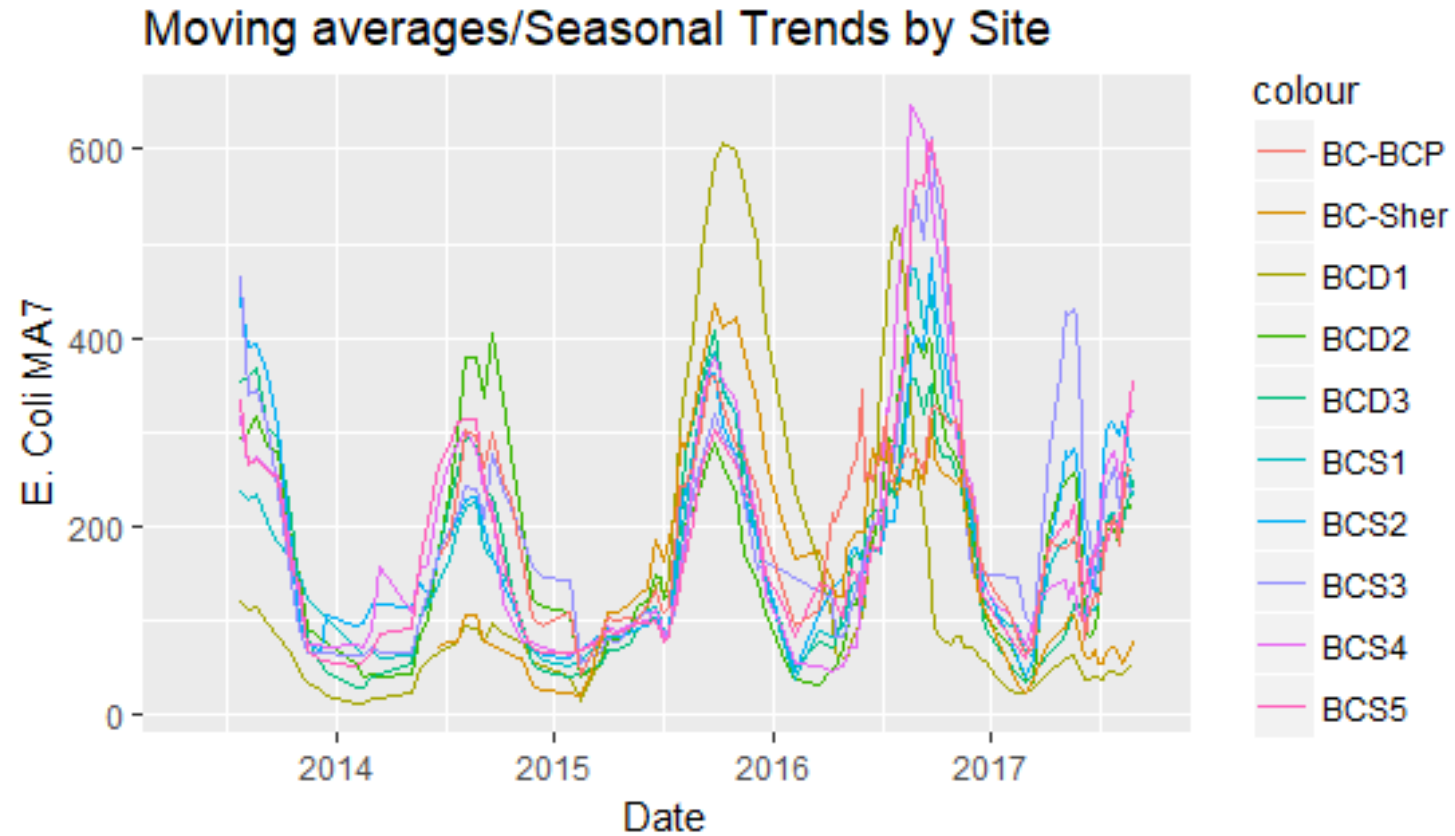




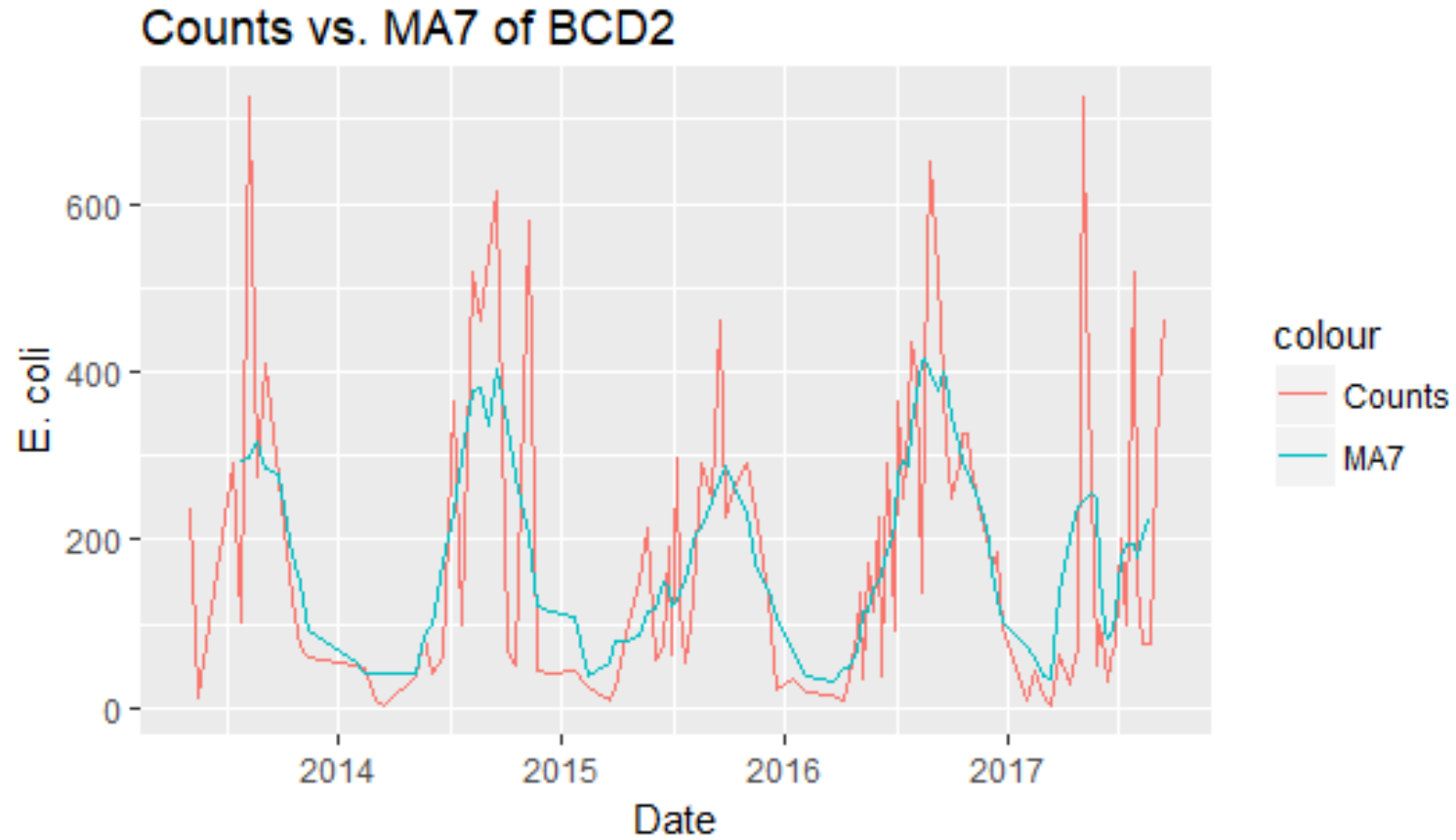


# Time Series Model Specification

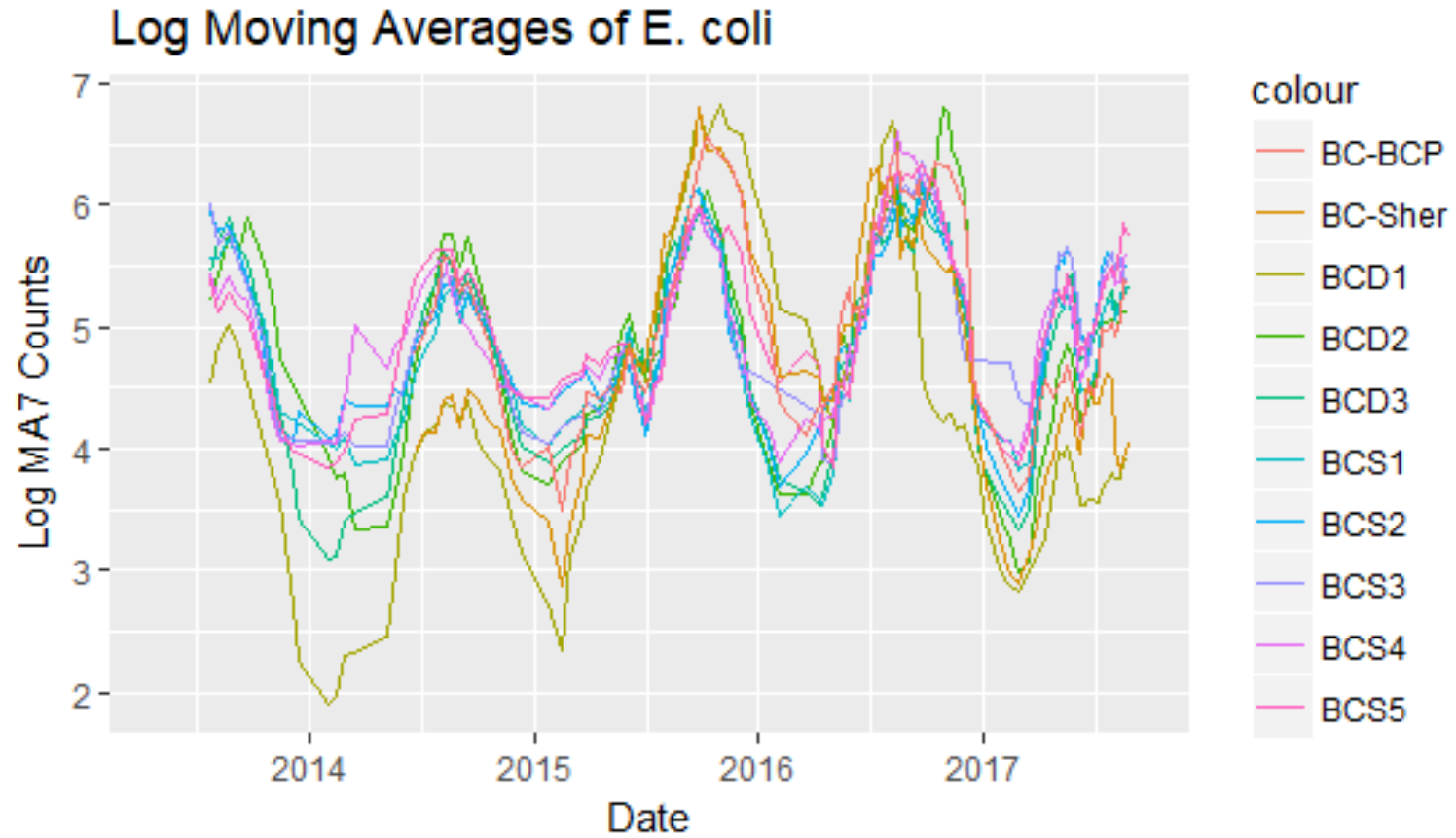
# Moving Averages of Individual Sites demonstrates clear seasonal trends



# Example: Moving Average smooths trends



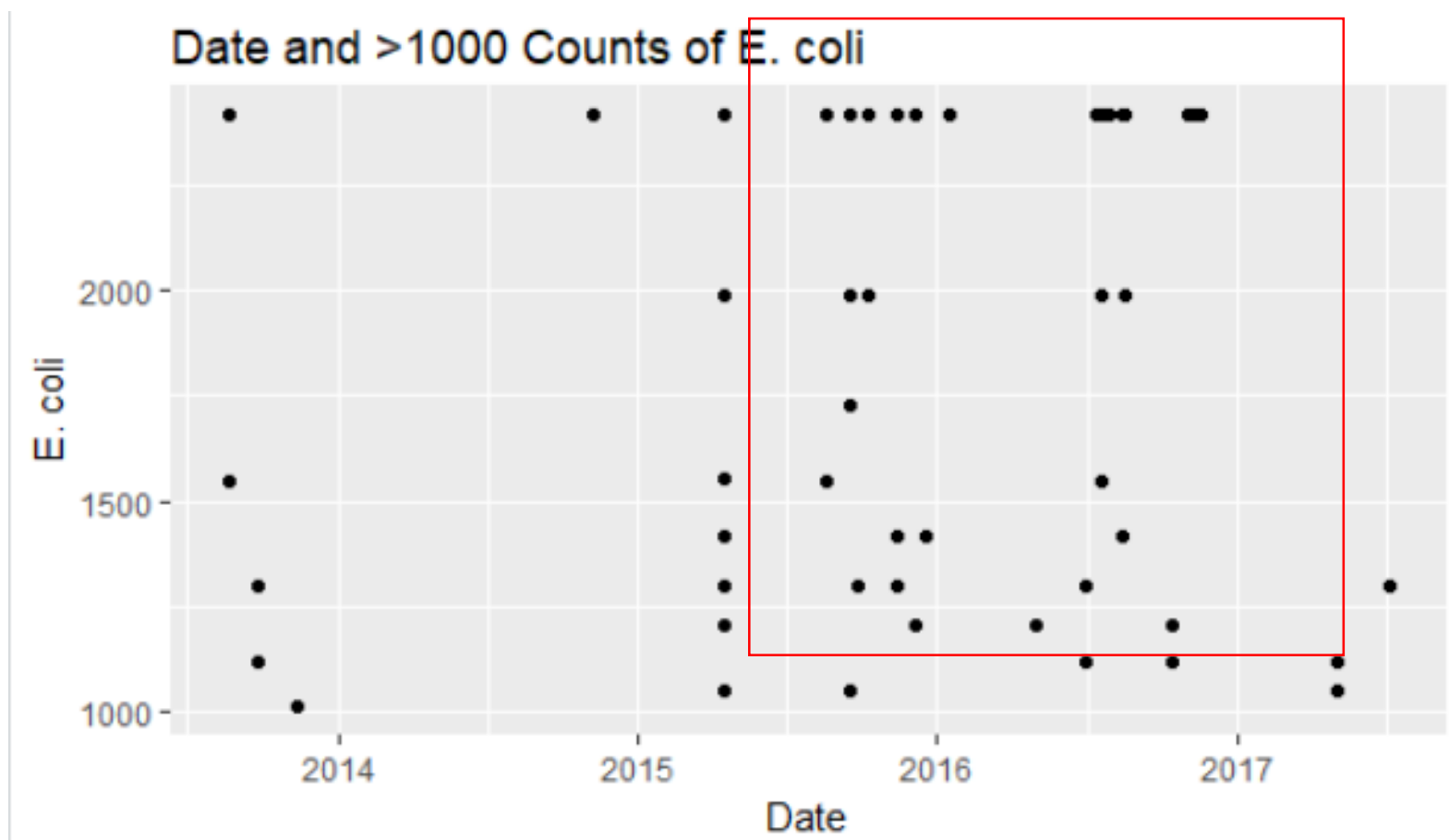
# Log Transformation: Possible upward trend?



# Possible Upward Trend in Data?

- Observed trend for 2014-2016, data should be monitored to see if trend persists in future years
- Trend wasn't apparent from 2013 to 2014 or 2016 to 2017
- Upon inspection, it appears that 2015 and 2016 had a uniquely high number of high counts compared to other years
- Sewage Break incidents?

# Clusters of high counts for 2015 and 2016



# Activation Temperature effects

- Activation temperature range of E. coli falls between 15°C and 45°C
- T-test analysis (with log transformation) reveals a significant difference of e. coli counts for observations in activation range vs. not, though not a particularly large difference...

Welch Two Sample t-test

data: activation\_log\$logE.coli and nonactivation\_log\$logE.coli

t = 5.9528, df = 1466.8, p-value = 3.292e-09

alternative hypothesis: true difference in means is not equal to 0

95 percent confidence interval:

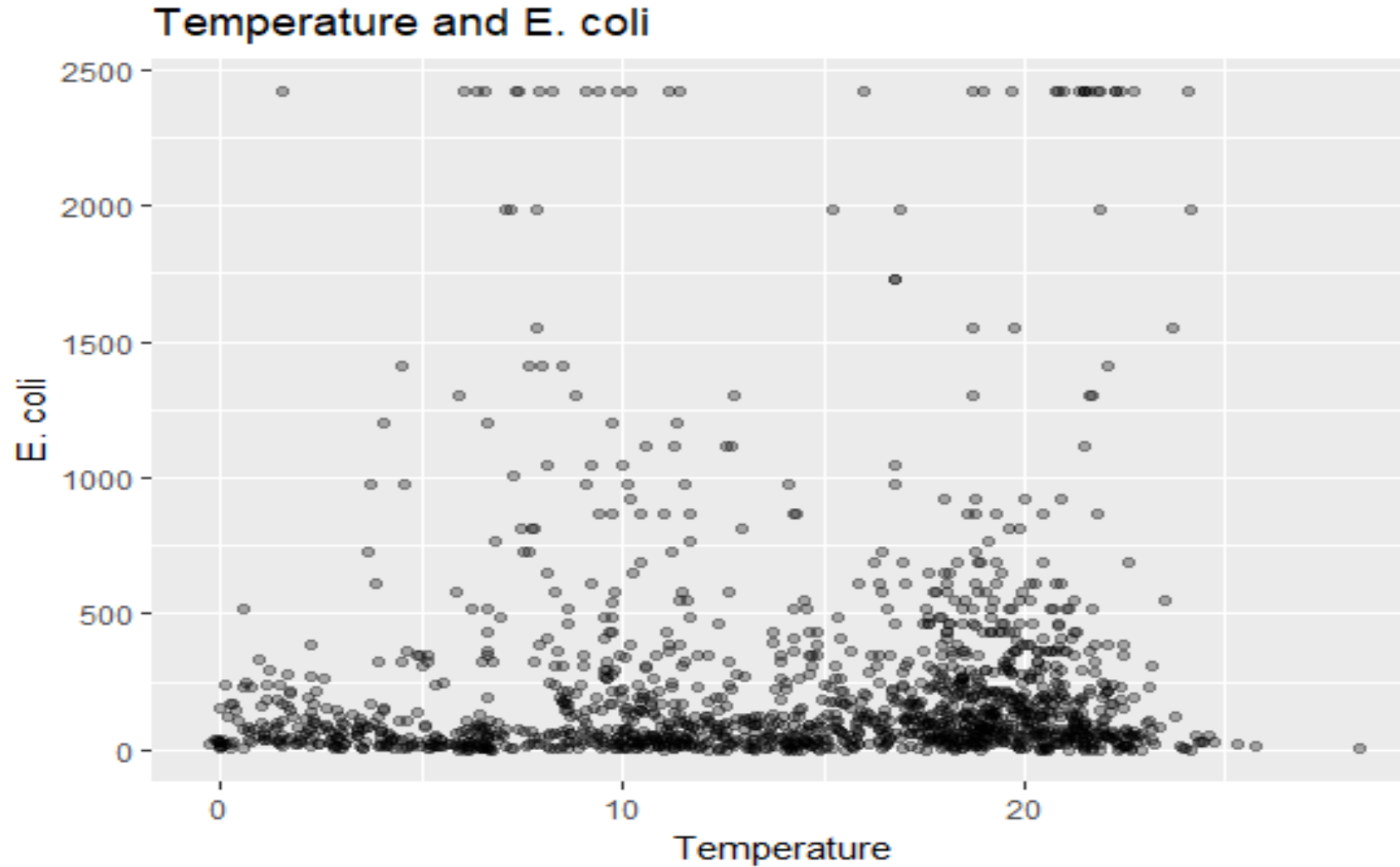
0.3060827 0.6069460

sample estimates:

mean of x mean of y

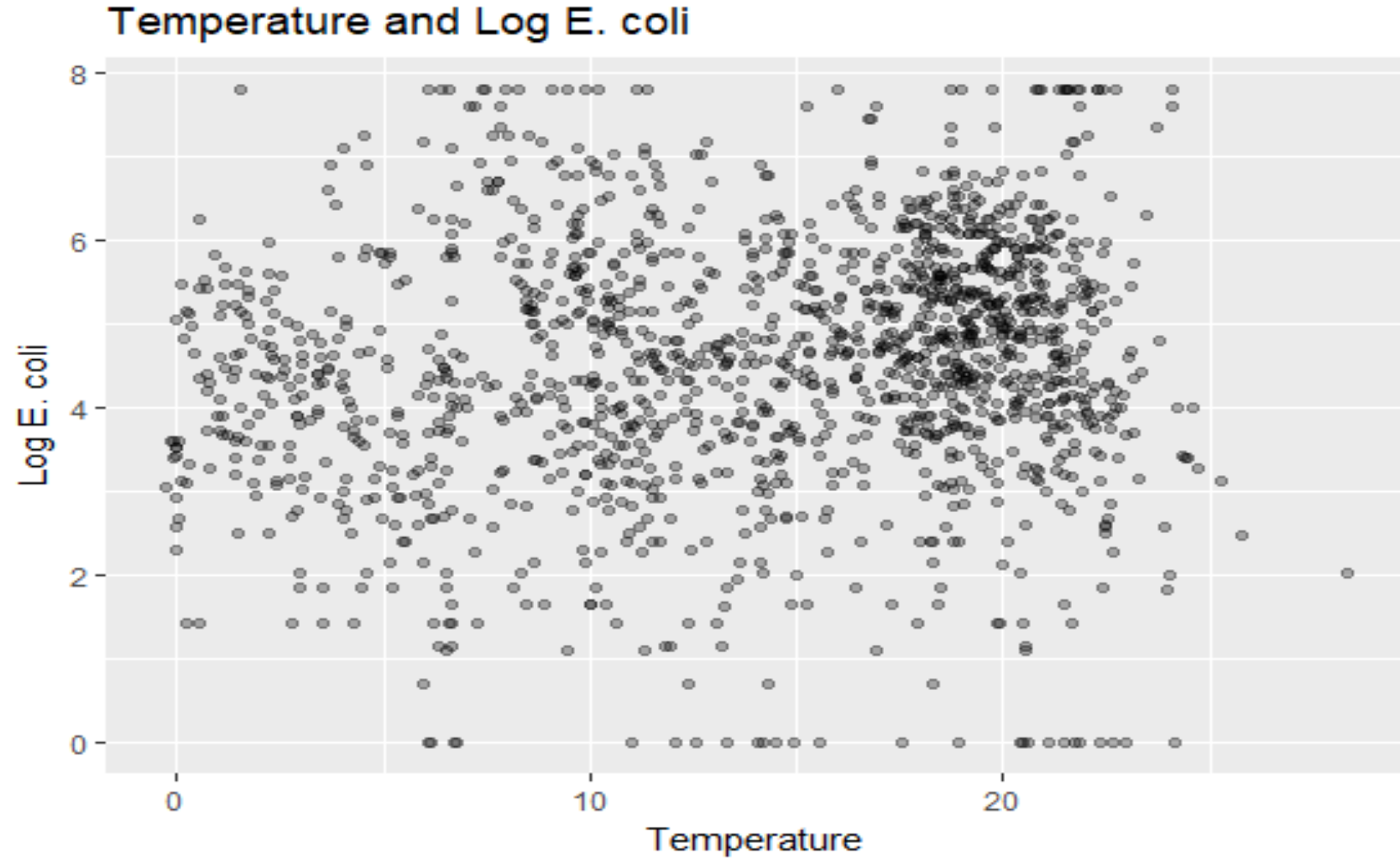
4.775455 4.318941

# Activation Temperature effects

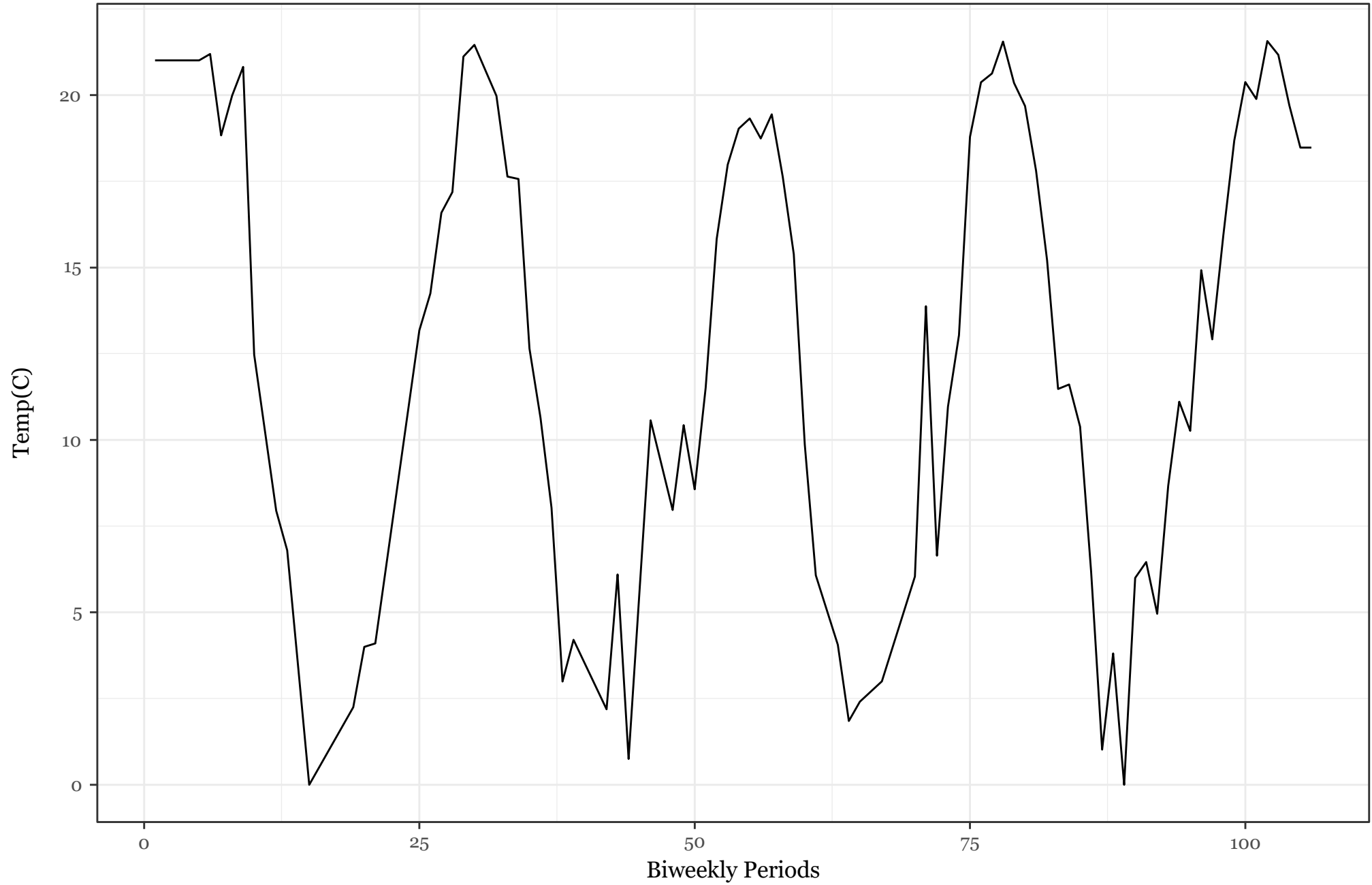




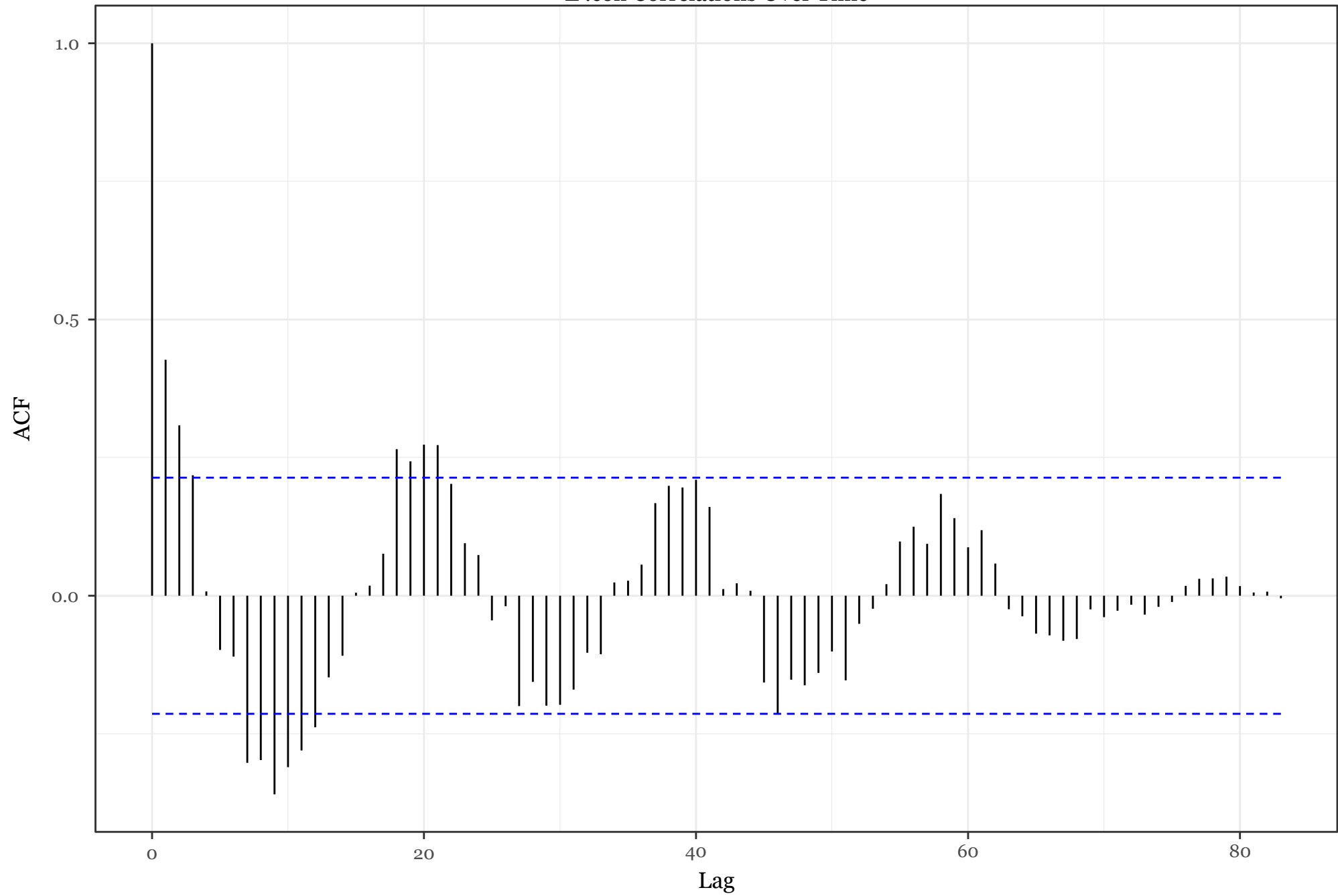
# Activation Temperature effects

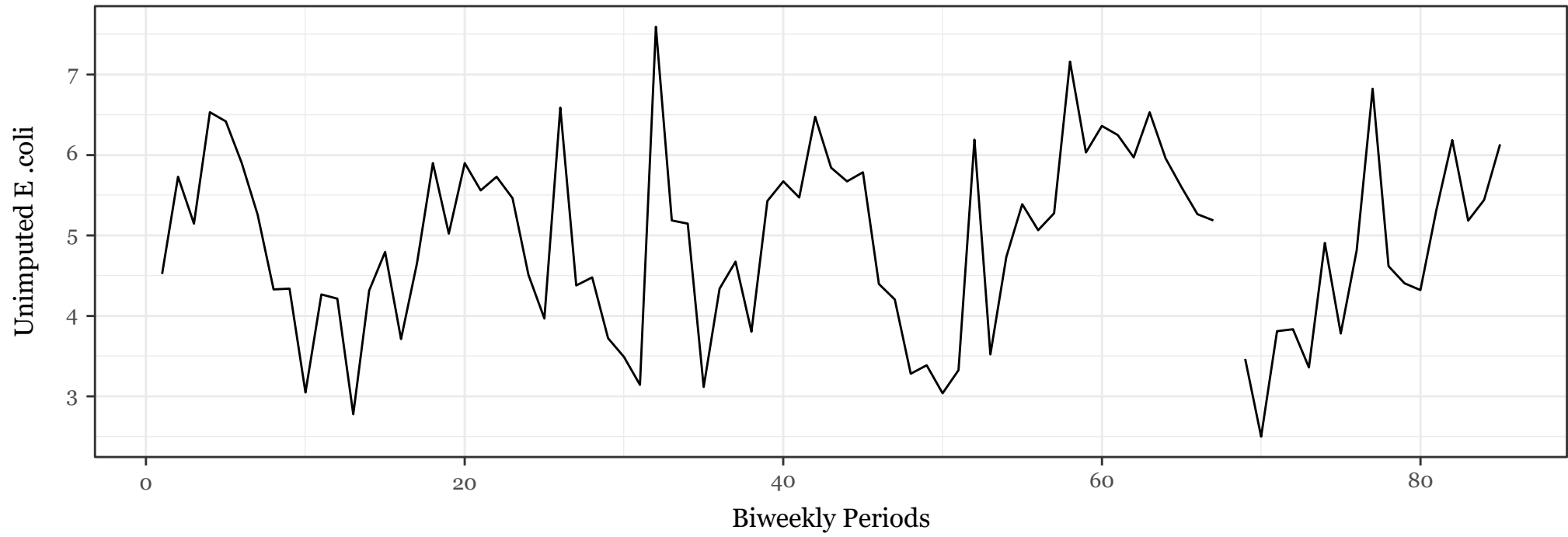
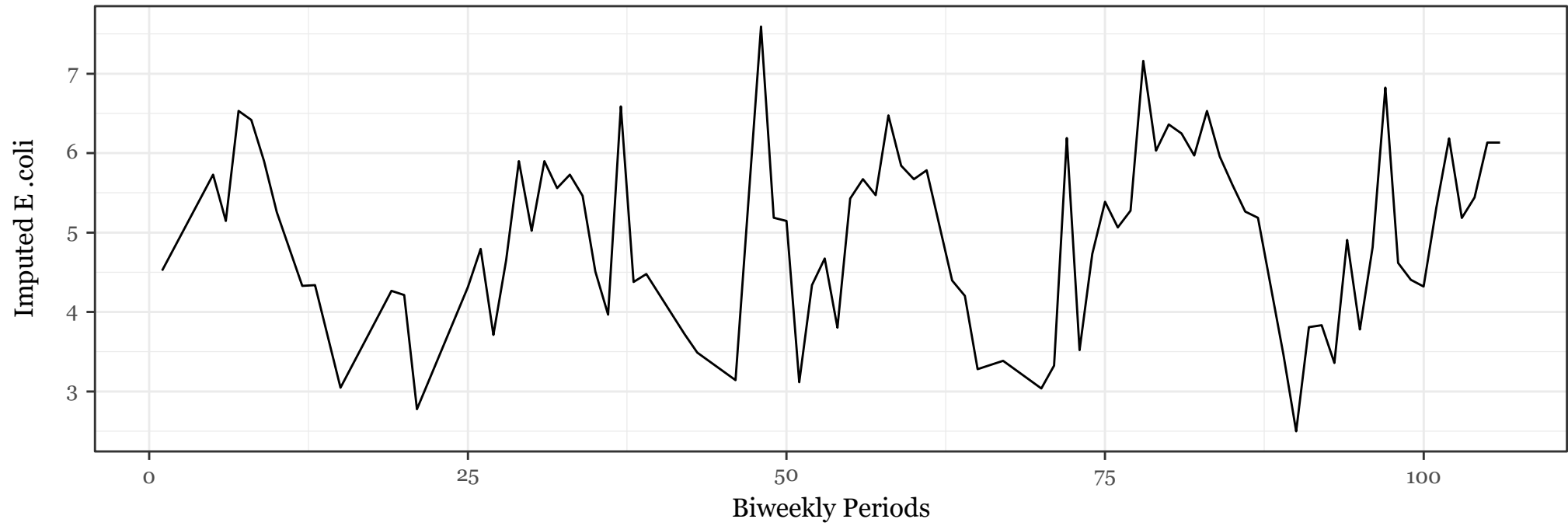


Temperature(C) vs. Time

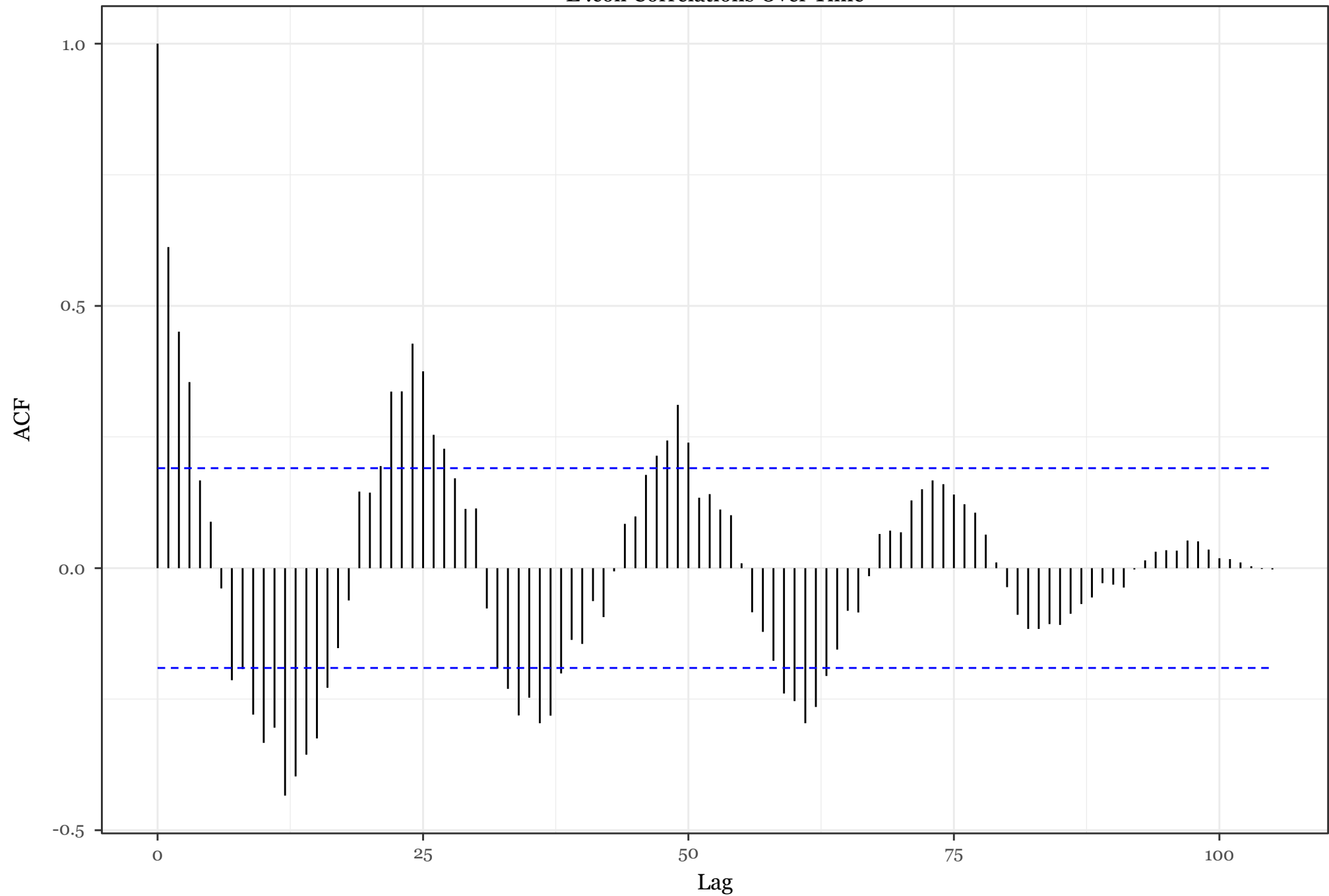


E .coli Correlations Over Time

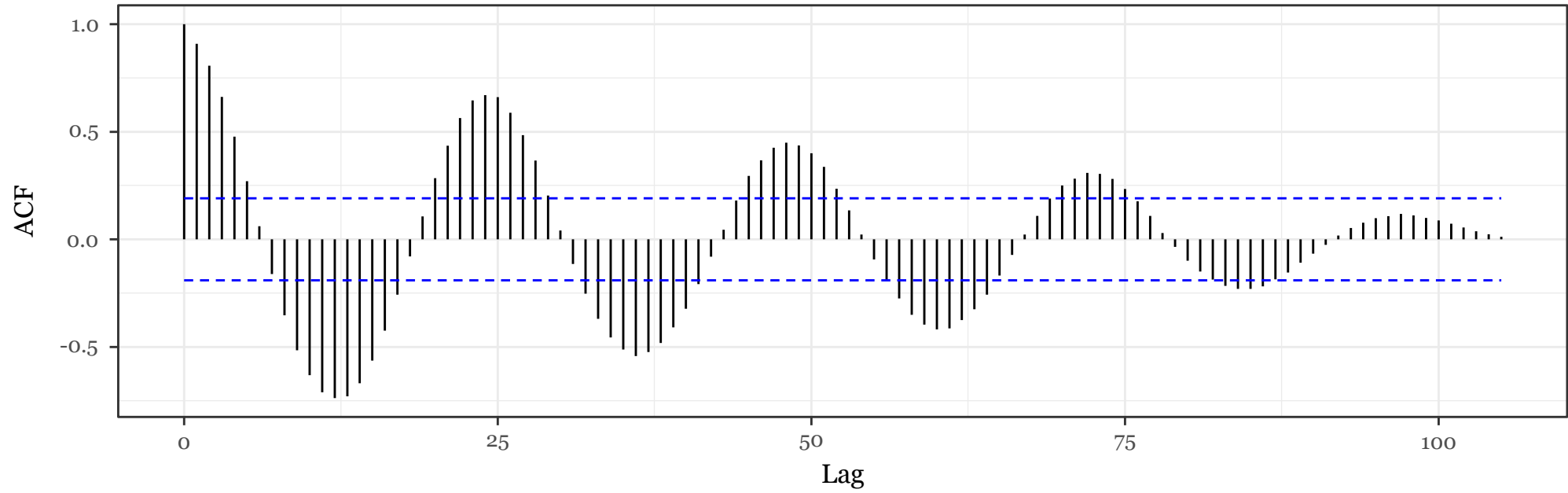




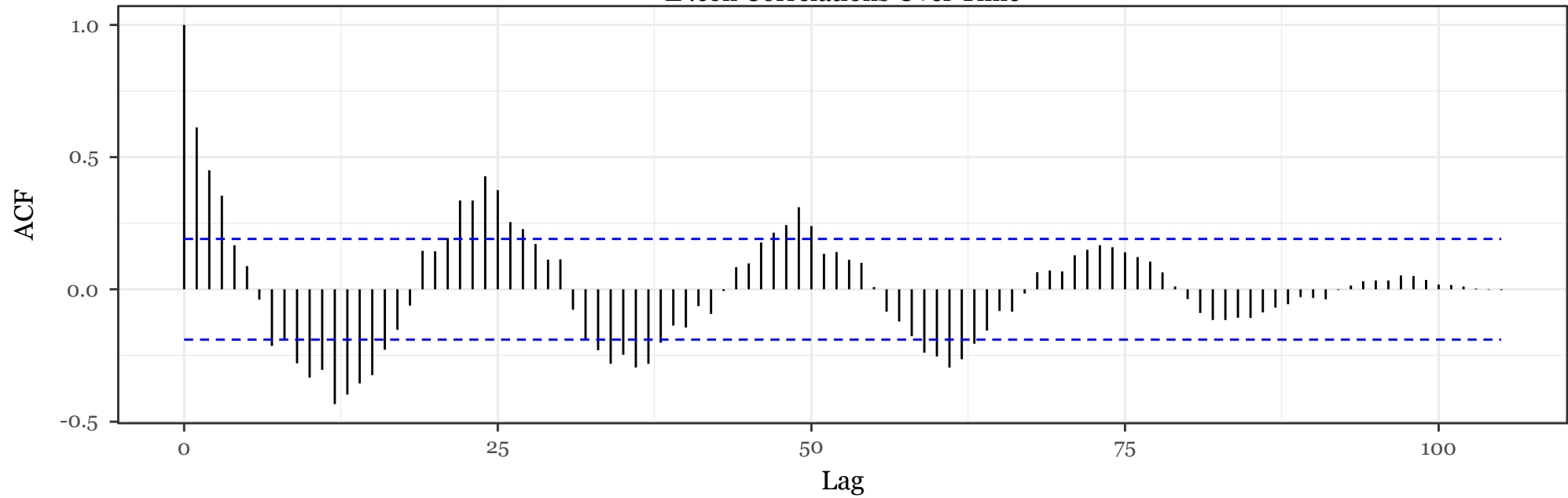
E .coli Correlations Over Time



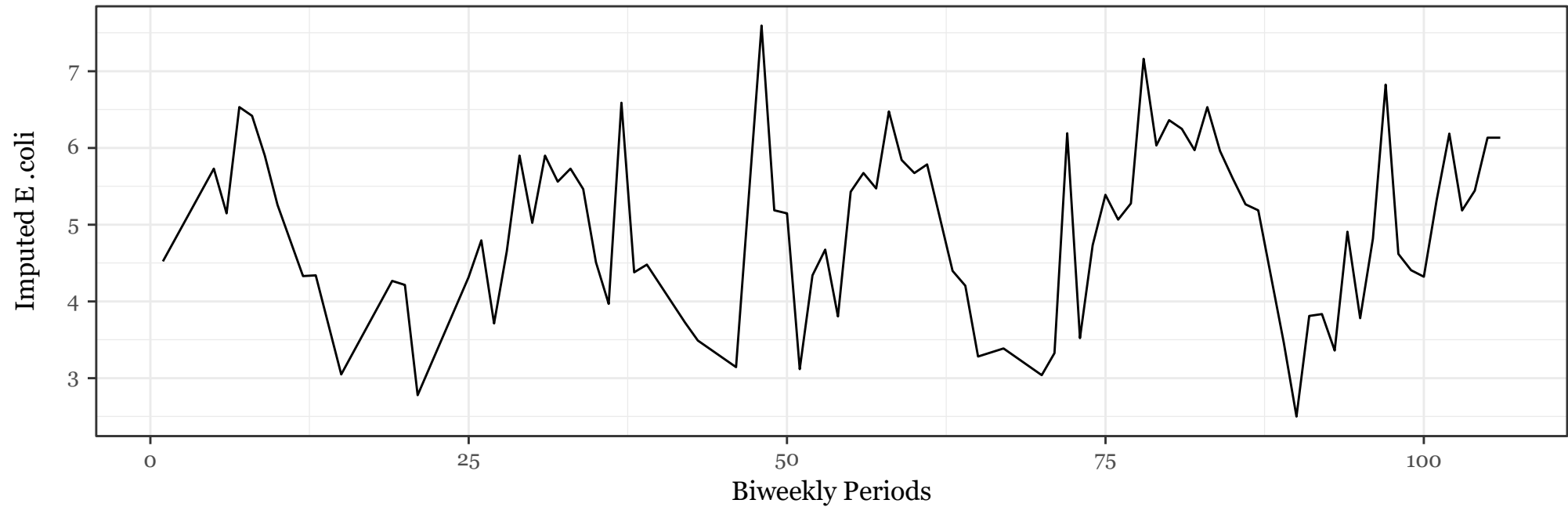
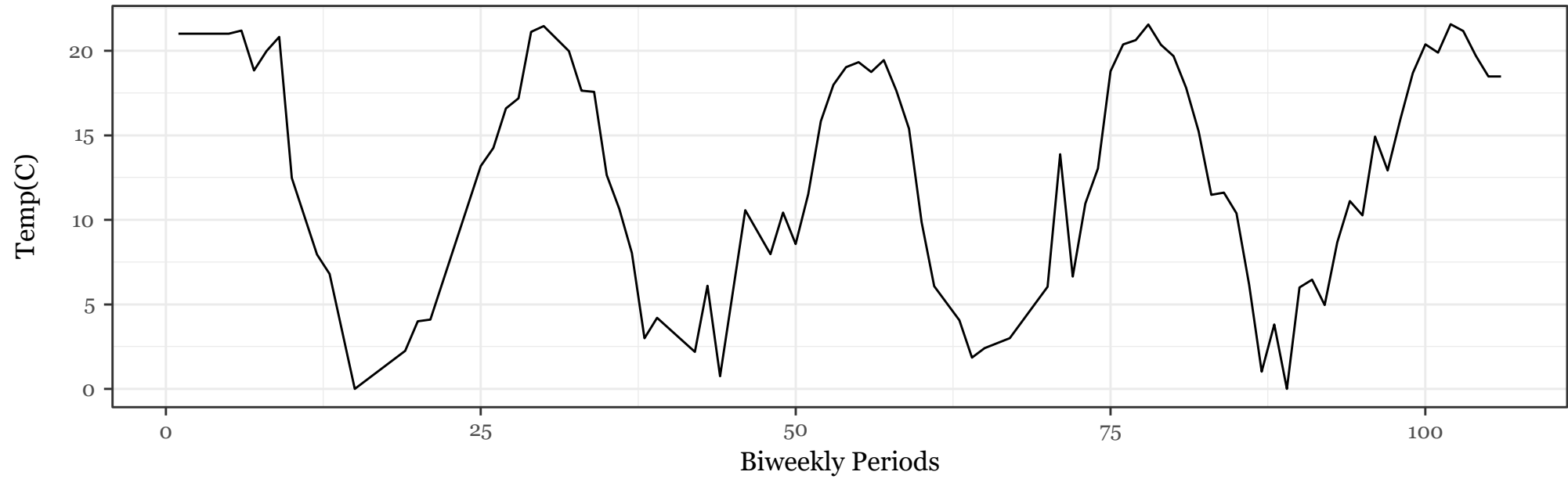
Temperature Correlations Over Time



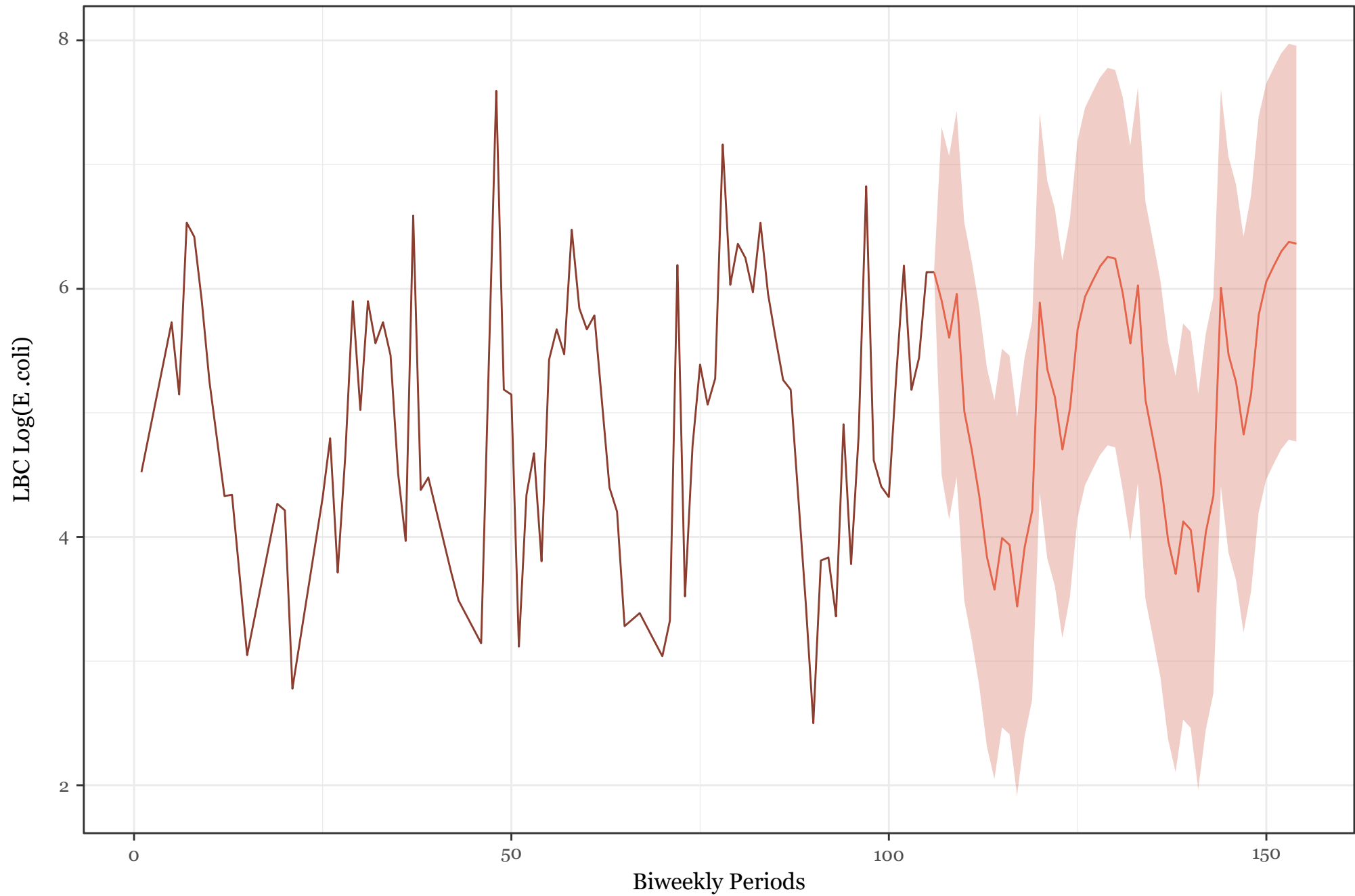
E .coli Correlations Over Time



Temperature(C) vs. Time

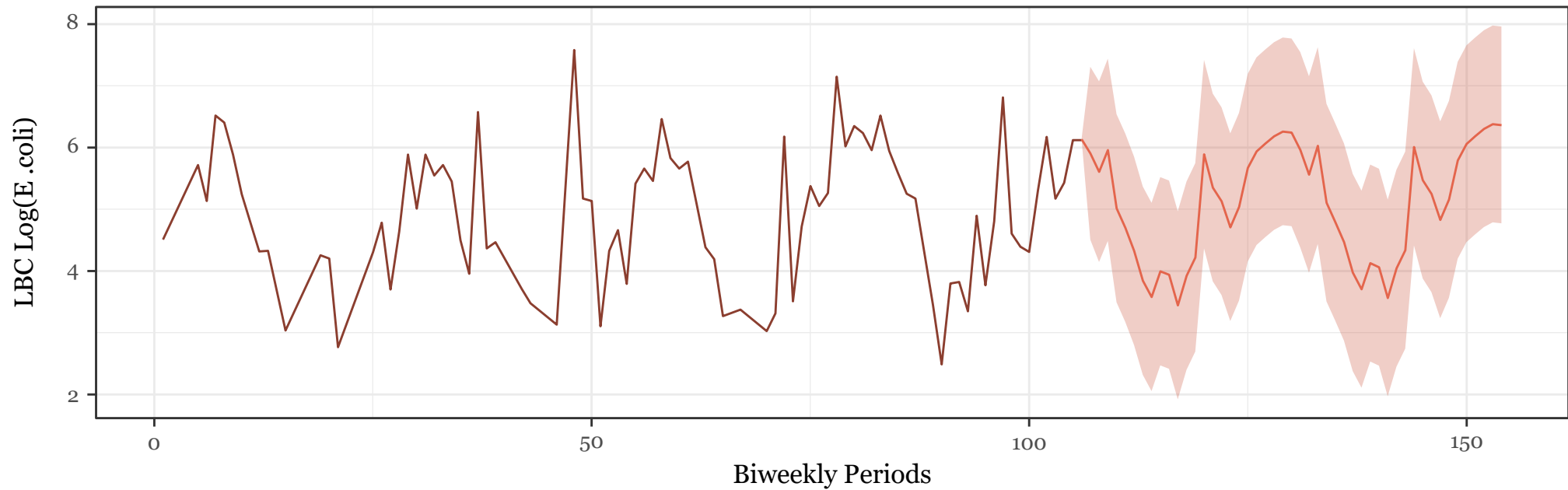


Forecast from ARIMA(4,1,1)(0,1,3)[24]

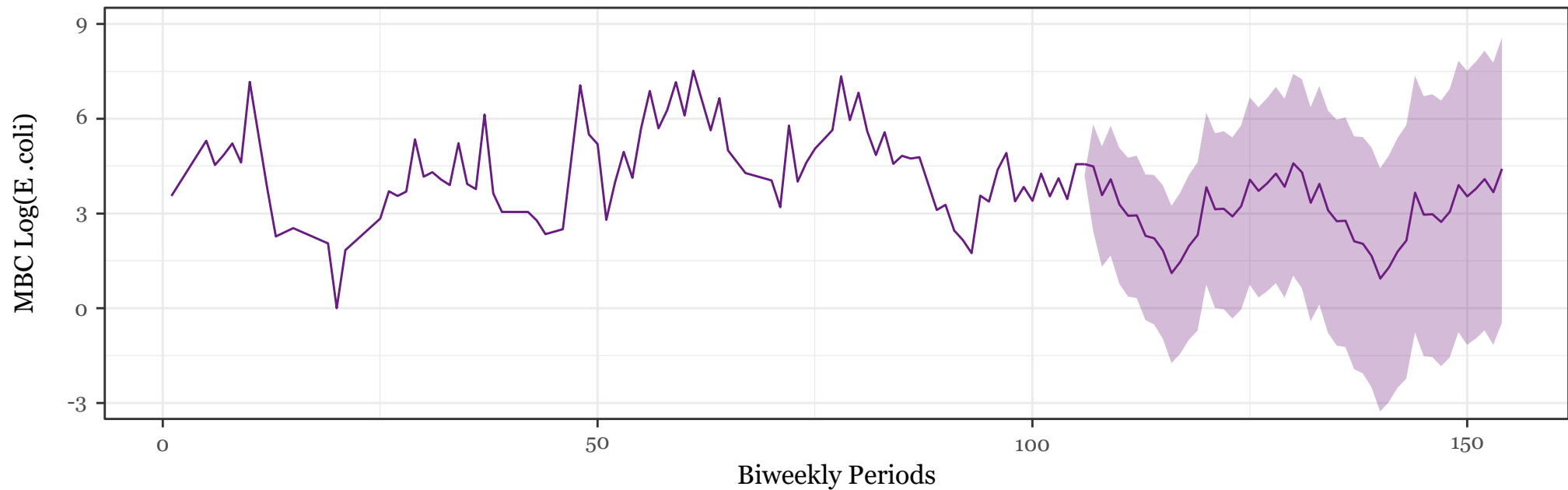




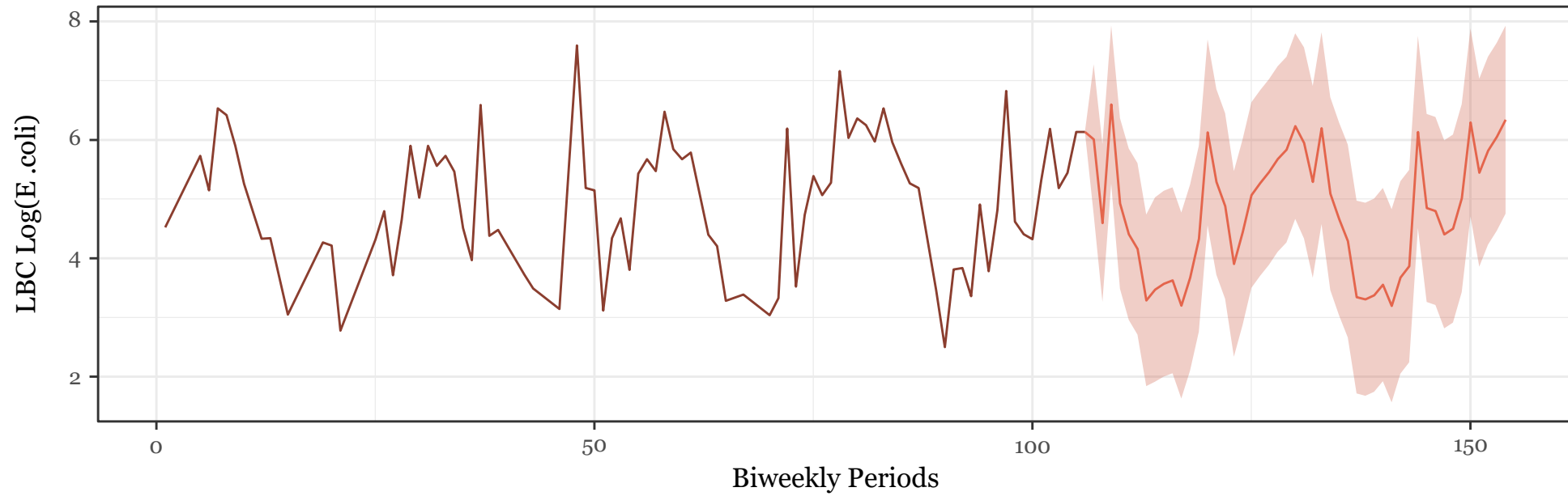
Forecast from ARIMA(4,1,1)(0,1,3)[24]



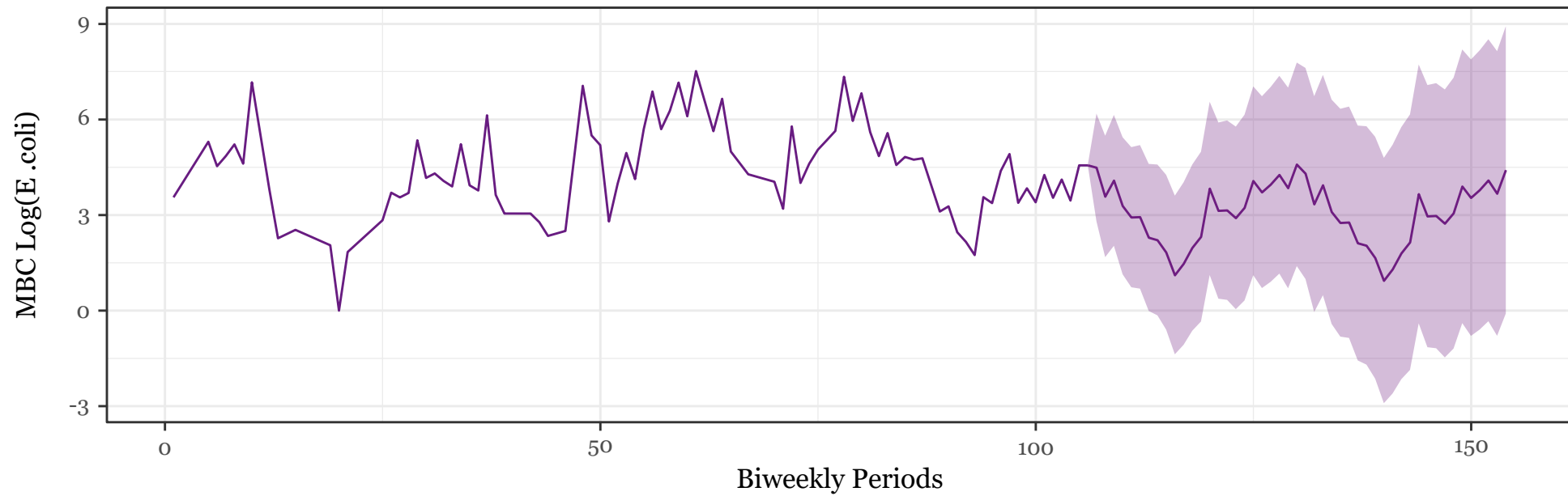
Forecast from ARIMA(4,1,1)(0,1,3)[24]



Forecast from ARIMA(1,0,7)(0,1,3)[24]



Forecast from ARIMA(4,1,1)(0,1,3)[24]





# Some Implications

## Final Thoughts/Inferences

- Seasonality appears to be largest predictor of E. coli levels
- Weather data should be monitored in closer time intervals in future to observe runoff effects
- Confounds: What caused occasional random high counts of E. coli in cold temperatures, particularly in 2015 and 2016?

Can we  
answer the  
questions we  
wanted to  
answer?

- Dataset may not be amenable to finding causes of E. coli in stream
- E. coli may or may not be related to human activity

# Data Swamped by Seasonality?

- Difficulties with inferences about:
  - E. coli and human sources
  - Non-seasonal predictors of E. coli
- Focus could be on local effects
  - Monitor local behavior using control site with little human presence

# References

- Hogg, R. V., & Craig, A. T. (1995). *Introduction to mathematical statistics.* (7th edition) (pp. 269-278). Upper Saddle River, New Jersey: Prentice Hall.