

1 **Inclusion of environmentally themed search terms improves Elastic Net regression**  
2 **nowcasts of regional Lyme disease rates**

3 Eric Kontowicz<sup>1,2</sup>, Grant Brown<sup>3</sup>, Jim Torner<sup>1</sup>, Margaret Carrel<sup>4</sup>, Kelly Baker<sup>5</sup>, Christine A.  
4 Petersen<sup>1,2,6\*</sup>

5 <sup>1</sup>Department of Epidemiology, College of Public Health, University of Iowa, Iowa City, Iowa  
6 52242, USA

7 <sup>2</sup>Center for Emerging Infectious Diseases, University of Iowa Research Park, Coralville, Iowa,  
8 52241, USA

9 <sup>3</sup>Department of Biostatistics, College of Public Health, University of Iowa, Iowa City, Iowa,  
10 52242, USA

11 <sup>4</sup>Department of Geographical and Sustainability Sciences, College of Liberal Arts & Sciences,  
12 University of Iowa, Iowa City, Iowa, 52242, USA

13 <sup>2,5</sup>Department of Occupational and Environmental Health, College of Public Health, University of  
14 Iowa, Iowa City, 52242, USA

15 <sup>6</sup>Immunology Program, Carver College of Medicine, University of Iowa, Iowa City, Iowa, 52242,  
16 USA

17 \*Correspondence: [christine-petersen@uiowa.edu](mailto:christine-petersen@uiowa.edu)

18 EK: [eric-kontowicz@uiowa.edu](mailto:eric-kontowicz@uiowa.edu)

19 GB: [grant-brown@uiowa.edu](mailto:grant-brown@uiowa.edu)

20 JT: [james-torner@uiowa.edu](mailto:james-torner@uiowa.edu)

21 MC: [margaret-carrel@uiowa.edu](mailto:margaret-carrel@uiowa.edu)

22 KB: [kelly-k-baker@uiowa.edu](mailto:kelly-k-baker@uiowa.edu)

23 CAP: [christine-petersen@uiowa.edu](mailto:christine-petersen@uiowa.edu)

24

25

26

27

28

29

30

31

32

33

34 **Abstract**

35 Lyme disease is the most widely reported vector-borne disease in the United States. 95% of  
36 human cases are reported in the Northeast and upper Midwest. Human cases typically occur in  
37 the spring and summer months when an infected nymph *Ixodid* tick takes a blood meal. Current  
38 federal surveillance strategies report data on an annual basis, leading to nearly a year lag in  
39 national data reporting. These lags in reporting make it difficult for public health agencies to  
40 assess and plan for the current burden of Lyme disease. Implementation of a nowcasting  
41 model, using historical data to predict current trends, provides a means for public health  
42 agencies to evaluate current Lyme disease burden and make timely priority-based budgeting  
43 decisions. The objective of this study was to develop and compare the performance of  
44 nowcasting models using free data from Google Trends and Centers of Disease Control and  
45 Prevention surveillance reports for Lyme Disease. We developed two sets of elastic net models  
46 for five regions of the United States first using monthly proportional hit data from 21 disease  
47 symptoms and tick related terms and second using monthly proportional hit data from all terms  
48 identified via Google correlate plus 21 disease symptom and vector terms. Elastic net models  
49 using the larger term list were highly accurate (Root Mean Square Error: 0.74, Mean Absolute  
50 Error: 0.52,  $R^2$ : 0.97) for four of the five regions of the United States. Including these more  
51 environmental terms improved accuracy 1.33-fold while reducing error 0.5-fold compared to  
52 predictions from models using disease symptom and vector terms alone. Models using Google  
53 data similar to this could help local and state public health agencies accurately monitor Lyme  
54 disease burden during times of reporting lag from federal public health reporting agencies.

55

56

57

58 **Introduction**

59 Lyme disease is the most widely reported vector-borne disease in the United States (1),  
60 with 95% of human cases occurring in the Northeast and upper Midwest (2). *Borrelia burgdorferi*  
61 *sensu lato* (including *Borrelia mayonii*, hereafter *B. burgdorferi*) is the causative agent of Lyme  
62 disease. It is transmitted to people predominantly when nymph or, to a lesser extent, adult ticks  
63 infected with *B. burgdorferi* take a blood meal (3, 4). Hard to detect nymphal *Ixodes* ticks quest  
64 for blood meals during spring and early summer months. People are at greatest risk of  
65 contracting Lyme disease during and immediately following this time (5-9) when spending time  
66 in the environment for either work or recreation (10). Areas with sandy soil and wooded  
67 vegetation are environmental factors associated with higher tick densities (11). With increased  
68 geographic spread of Lyme disease, there has been increased incidence since 2000 (12). Lyme  
69 disease has a large economic burden on patients and their surrounding communities (13, 14).

70 Surveillance of Lyme disease in the United States requires participation from many  
71 different areas of the health care system (15). This surveillance relies on case reports from  
72 physicians, lab reports from diagnostic labs and collation of this data as cases by local and state  
73 health departments. These case reports are forwarded to the CDC, which then aggregates the  
74 data and produces summary reports on national Lyme disease incidence. Due to differences in  
75 reporting from states and localities, compilation of data at the federal level can take several  
76 years, resulting in a time lag for release of nationwide surveillance and summary reports. This  
77 lag in federal reporting has been problematic for local health departments (LHDs), as they must  
78 predict current and emerging public health needs based on Federal data that is several years  
79 old (16). LHDs not only play a vital role in surveillance of Lyme disease, but also help mitigate  
80 disease incidence through the implementation of local interventions. Funded prevention  
81 efforts/campaigns by LHDs can have a positive effect on health in communities (17).  
82 Unfortunately, there are often many important competing health priorities in communities. As  
83 such, LHDs must make critical decisions to allocate their limited fund to areas of highest need.

84 Modeling methods that accurately nowcast, or predict the present, Lyme disease incidence in a  
85 region would allow for better planning on the part of LHDs to allocate their efforts. Using  
86 statistical learning methods for nowcasting can also discover, or highlight, patterns that are  
87 associated with disease and can be used to generate future hypotheses.

88 Usage of non-traditional indicators of disease spread, like Google search traffic history,  
89 has gained credibility from public health audiences (18, 19). Google search data has been used  
90 with a variety of mathematical and statistical models to predict obesity rates, unemployment  
91 rates, and infectious diseases with varying levels of accuracy (20-23). The principal insight of  
92 these approaches is that search data is available at a wide temporal and geographical scale,  
93 and such queries may be correlated with a phenomena or disease process of interest or human  
94 behaviors (24). This correlation can be leveraged to make predictions of current or future health  
95 outcome rates. In addition, relative frequencies of search terms may generate interesting  
96 hypotheses concerning human behaviors and their relationship with disease outcomes.

97 Given the complex and potentially high dimensional nature of search data, statistical and  
98 machine learning tools are a natural fit for model development. There are a variety of parametric  
99 and non-parametric statistical learning approaches used in the literature for infectious disease  
100 prediction as discussed in a recent review (25). In this work, we do not provide a comprehensive  
101 review of such options, but rather seek to demonstrate that nowcasting is a promising  
102 opportunity for Lyme disease specifically. For this reason, we employ Elastic net regression.  
103 Elastic net regression provides a flexible parametric approach which strikes a compromise  
104 between the L1 and L2 penalties of LASSO and ridge regression, respectively. It is also  
105 computationally straightforward, being easily employed on modest hardware. An additional  
106 advantage of elastic net regression is the grouping effect, where strongly correlated features  
107 tend to remain or be excluded from the model together (26).

108           In this study, we built elastic net regression models capable of nowcasting Lyme disease  
109 rates in five different regions of the United States. We developed two models for each region, 1.  
110 Using search traffic data from only disease name, symptom and vector related terms and 2.  
111 Using search traffic from terms identified via Google Correlate™ in addition to disease name,  
112 symptom and vector related terms to identify trends using information recently sought by the  
113 general public on the disease, it's symptoms, and correlated terms (27, 28). We hypothesized  
114 that nowcasting models would have better predictive accuracy and lower error when using a full  
115 list of search terms that the average person would search compared to models that only use  
116 terms related to disease name, symptom and vectors of Lyme disease. Further, the three most  
117 important terms from accurate models will be potential exposure/location themed and their  
118 search patterns will align temporally with the timing of Lyme disease incidence in endemic  
119 areas, the Northeast and Midwest, and less well in non-endemic areas, the Southwest and  
120 West.

## 121 **Materials and Methods**

### 122 *Outcome Data*

123           All Lyme disease incidence data for this study was provided by the United States  
124 Centers for Disease Control and Prevention (CDC) (<https://www.cdc.gov/lyme/stats/tables.html>).  
125 In 2008, the CDC switched to a Suspected, Probable, or Confirmed case reporting approach.  
126 Cases were considered confirmed if an individual presents with erythema migrans and with a  
127 known exposure, a case of erythema migrans with laboratory evidence of infection and without  
128 known exposure, or a case with at least one late manifestation that has laboratory evidence of  
129 infection. Any other case of physician-diagnosed Lyme disease that has laboratory evidence of  
130 infection were considered probable cases. Both confirmed and probable case definitions were  
131 included to provide a more sensitive and inclusive criterion. Laboratory evidence of infection in  
132 both definitions allowed for strong confidence in a Lyme disease diagnosis. Even so,

133 heterogeneity remains in reporting strategy; between 2015 and 2016, Massachusetts changed  
134 their reporting strategy to only report laboratory confirmed cases to the CDC. Only reporting  
135 laboratory confirmed cases is likely to lead to underreporting of the true burden of disease (29).  
136 Lyme disease incidence is reported by the CDC on a per county of diagnosis for each US state.  
137 For the purposes of this study, we aggregated these counts to state and month based on date  
138 of diagnosis. Next, regional incidence rates were calculated for five different regions: Northeast,  
139 Midwest, Southeast, Southwest and West (**Figure 1**). Regions were developed as a hybrid of  
140 known high incidence regions and the US Census regions (15). Regional monthly Lyme disease  
141 incidence rates were calculated using combined state level population data from the 2010 US  
142 Census. Data was split into training and hold-out sets; models were fit on observations between  
143 February 2004 and December 2014 and validated on the hold-out observations which had  
144 available surveillance data from January 2015 to December 2017.

145 **Figure 1 Regions of the United States.** United States divided into 5 different regions  
146 (Northeast, Midwest, Southeast, Southwest, and West) used to calculate regional Lyme  
147 incidence, and regional search term data. Map created using ArcGIS software.

148 *Google Search Term data*

149 Regional Lyme disease incidence trends from the training period were used with Google  
150 Correlate™ to identify the top 100 correlated search terms on which monthly proportional search  
151 hit data was collected (30). Google Correlate™ was not able to identify terms at state levels.  
152 These correlations can only be made on a nationwide basis for a submitted time series. Thus,  
153 we were not able to limit our search term identification by region. However, using regional Lyme  
154 disease time series data, provided many regionally specific terms in the top 100 correlated  
155 terms for each region (Supplemental Table 1). High correlation was determined when the  
156 correlation value was greater than 0.8, moderate if correlation value was between 0.5 and 0.8,  
157 and poor when less than 0.5.

158 Google Correlate™ implements an Approximate Nearest Neighbor (ANN) system to  
159 identify candidate search terms that matched similar temporal trends from supplied data. This  
160 system implemented a two-pass hash-base system. The first pass computed the approximate  
161 distance from the supplied time series to a hash of each series in Google's database (30). The  
162 second pass computed the exact distance function using the top results supplied from the first  
163 pass (30). For each region, the 100 terms identified from Google Correlate™ and the 21 Lyme  
164 disease symptom and *Ixodid*- vector related terms were entered into *gtrendsR* (31) to collect  
165 proportional monthly search hit data for each term per region (31, 32). Search hit data was  
166 collected at the state level for each term and averaged to regional aggregates. This was then  
167 used as feature data for nowcasting Lyme disease incidence trends (33, 34). Search hit data  
168 from the relevant time periods (2004-2018) was collected between September 18, 2019 and  
169 September 26, 2019.

## 170 *Modeling*

171 For each region, two groups of elastic net regression models were fit for comparison: 1.  
172 a model using only monthly proportional hit data from the 21 disease symptoms and tick related  
173 terms list, and 2. a model using monthly proportional hit data from terms identified via Google  
174 Correlate™ in combination with the disease symptom and tick term list (this will be referred to as  
175 the full-term list for the remainder of the paper). The training data was from February 2004  
176 through December 2014. To help prevent overfitting we implemented a rolling training window  
177 for the statistical learning process with a twelve-month learning window and one month  
178 validation window. To further address the potential for overfitting, we excluded data between  
179 January 2015 until December 2017 from the model training process. The hold-out data set was  
180 not used in any model training or in-sample validation and was only used to determine how  
181 models would respond to new data and to determine if the models overfitted to the training data.

182 We collected all search data in September of 2019 therefore all nowcasting done by  
183 developed models presented in this article will not exceed September 2019. All elastic net  
184 models were built and run in R version 3.6.2 using the *caret* and *glmnet* packages (35, 36).  
185 Model fit was determined using Root Mean Square Error (RMSE), Mean Absolute Error (MAE),  
186 and R<sup>2</sup>. All graphics of model fit, and search term correlation were created using *ggplot2* in R  
187 version 3.6.2, and search terms are presented as directly provided by Google Correlate™.

188 Elastic net regression is a penalized form of ordinary least squares regression and  
189 contains a hybrid of ridge and Least Absolute Shrinkage and Selection Operator (LASSO)  
190 regression penalties (26). Elastic Net regression was implemented to both reduce the impact or  
191 outright eliminate non-essential feature data as it compromises the L1 and L2 penalties of  
192 LASSO and ridge regression respectively. Alpha and lambda hyper parameters are used in  
193 Elastic net regression to balance the tuning of the L1 (LASSO) and L2 (Ridge) norm penalty  
194 parameters (Equation 1A).

Equation 1A: Combined Penalty

$$\sum_{i=1}^n (Y - X\hat{\beta})^2 + \lambda \sum_{j=1}^p [(1 - \alpha)(\hat{\beta}_j)^2 + \alpha|\hat{\beta}_j|]$$

195 Alpha determines the relative weights of the two penalty parameters and lambda determines the  
196 overall weight of the summation of the individual penalties. For each region and elastic net  
197 model group (disease symptom and vector terms alone vs. full-term list), we tested a  
198 combination of 50 and 150 different automatically generated values of alpha and lambda to  
199 select optimal values.

200 Regional monthly Lyme disease incidence rates, as calculated from CDC surveillance  
201 data, was the outcome of the nowcasting models. Feature data was regional monthly search hit  
202 data from each region. We only used data from search terms where *gtrendsR* was able to  
203 appropriately return proportional monthly hit data. Despite terms having a correlation at the



204 national level and therefore identified via Google Correlate™, some terms held non-variable  
205 values of zero for their monthly proportional hit data at the region level. These terms with their  
206 zero variance would cause model failure and thus were excluded from the modeling process.

### 207 *Variable Importance*

208 Elastic net regression can reduce or outright eliminate feature data from final models.  
209 We wanted to determine which search terms had the greatest influence in the final, best tuned  
210 models. To determine search terms influence, the *varImp* function from the caret package was  
211 used to calculate the scaled importance of each term in the final models. The *varImp* function  
212 takes the absolute value of each coefficient and ranks these coefficients and stores them as  
213 variable importance from zero to one hundred. Put simply, larger coefficients have greater  
214 influence and thus are associated with increased importance.

### 215 **Results**

216 Between 2004 and 2017, the Northeast consistently had the highest counts of Lyme  
217 disease followed by the Midwest. The lowest incidence areas were consistently the West and  
218 Southwest regions (**Figure 2**). All regions showed seasonal oscillation of Lyme disease  
219 incidence with typical peaks in summer months (July, August, and September) and falling in  
220 winter (**Figure 3A**). Seasonal oscillation occurs at a lower incidence in the West and southwest  
221 regions compared to the high-incidence regions of the Northeast and Midwest (**Figure 3B and**  
222 **3C**). These regional temporal trends were used with Google Correlate™ to identify 100 terms  
223 with correlated search patterns. Across all regions, there were environmental themes of outdoor  
224 activities that included concerts, camping, and water parks; places where people are likely to be  
225 exposed to Ixodes ticks during the late spring, summer and early fall (11). (**Table 1**, complete  
226 list of candidate search terms provided in supplemental Table 1). Gtrends was used to collect  
227 regional monthly proportional search data for each term identified from Google Correlate along

228 with the symptom and vector terms (120 total terms for each region). Some terms identified with  
 229 Google Correlate™ at the national level were identified as having no search traffic at the  
 230 regional level and were removed from the regional list (**Table 2**).

231 **Table 1 Candidate search terms identified via Google Correlate™ by region with**  
 232 **Symptom/Vector terms**

<b>Northeast Search Terms Identified by Google Correlate™</b>	<b>Midwest Search Terms Identified by Google Correlate™</b>
free concerts, july calendar, necbl, little league all stars, alive at five, movies under the stars, prospect park bandshell, summer recipe, harwich mariners, freezer jam	festivals milwaukee, beaches in michigan, kings island discount, easy summer recipes, lake beaches, motel wisconsin dells, movies in the park, summer desserts, dorm bedding, drive in ohio
<b>Southeast Search Terms Identified by Goggle Correlate™</b>	<b>Southwest Search Terms Identified by Google Correlate™</b>
intex, cloudy pool, summer things, alabama water park, blue bayou in baton rouge, cloudy pool water, baking soda pool, summer things to do, green pool, springtails	loans for, how to make string bracelets, pigeon forge hotels, recipes on the grill, sandstone amphitheater, cheap bmx bikes, cataratas del niagara, world rv, cave of the winds colorado springs, produce stand
<b>West Search Terms Identified by Google Correlate™</b>	<b>Symptom and Ixodid-vector Terms added for Each Region</b>
concert in the park, berry picking, movies in park, concert in park, blueberry picking, outdoor movies, soak city, lake water park, blueberry farm, broomfield bay	tick, black tick, lyme, lyme disease, rash ,bullseye rash, bell's palsy, facial paralysis, side of face paralyzed, knee pain, swollen knees, swollen joint, swollen joints, joint pain, fever, tired, deer tick, black-legged tick, black legged tick, black leg tick

233  
 234 **Table 2 Number of search terms that had monthly proportional hit data available from**  
 235 **Gtrends™**

<b>Region</b>	<b>Terms Into Gtrends™</b>	<b>Terms From Gtrends™</b>
Northeast	120	87
Midwest	120	86
Southeast	120	80
Southwest	120	42
West	120	83

236  
 237 **Figure 2 Regional Lyme disease incidence count from CDC surveillance.** Incidence counts  
 238 calculated by summing monthly state incidence form CDC surveillance in each region.  
 239 Calculations and graphs made suing RStudio version 3.6.2.

240 **Figure 3 Regional Lyme Disease Incidence. (A)** All regions relative to Northeast incidence  
 241 rates, **(B)** Southwest, **(C)** West. Incidence rates calculated by summing monthly state incidence  
 242 from CDC surveillance in each region. Denominator values calculated from 2010 US Census  
 243 state populations and aggregated to region. Calculations and graphs made using RStudio  
 244 version 3.6.2.

245 For accurate modeling predictions, or nowcasts, it is important to use feature data that is  
246 correlated to the outcome data of interest. Pearson's correlation was performed for each term's  
247 proportional monthly search traffic and regional Lyme disease rates within the training  
248 timeframe. Individual term correlation with Lyme disease incidence had a large range for each  
249 region of the US with moderate mean and median correlation values (**Table 3**, complete results  
250 provided in supplemental Table 2). Each region, except the Southwest, had sixteen terms with a  
251 correlation greater than 0.7 (complete results provided in supplemental Table 2). Over the  
252 regions that have suitable *Ixodes* climate and habitat (Northeast, Midwest, Southeast, and West  
253 regions), we found high maximum correlation values (0.893, 0.898, 0.840, and 0.836,  
254 respectively) for the top correlated search terms. Many of the 21 terms based on known Lyme  
255 disease symptoms or vectors had poor bivariate correlation with regional Lyme disease  
256 incidence. For example, fever, which is more often searched in winter months (37), was  
257 negatively correlated with Lyme disease incidence in every region for the entire timeframe of the  
258 study (**Figure 4**).

259 **Table 3 Summary values of bivariate correlation of full-term list search terms to regional**  
260 **Lyme disease rates of model training data**

Region	Range	Mean Correlation	Median Correlation
Northeast	-0.279, 0.893	0.560	0.663
Midwest	-0.245, 0.898	0.602	0.691
Southeast	-0.137, 0.840	0.524	0.590
Southwest	-0.065, 0.612	0.229	0.231
West	-0.165, 0.836	0.421	0.416

261

262 **Figure 4 Negative bivariate correlation of fever to Lyme Disease Incidence for all regions**  
263 **of the United States.** Correlation calculated using Pearson method with independent variable  
264 as proportional Google hits for each term and dependent variable Lyme Incidence per 100,000  
265 for each region.

266 The variance of feature data is also important for making accurate predictions. Features  
267 that have little to no variance overtime make for poor predictors. The variability of each term  
268 was assessed per region. It was found that two terms in the Northeast, one term in the Midwest,

269 and ten terms in the Southeast had zero variance. These terms were excluded from the  
 270 nowcasting process.

271 To evaluate the hypothesis that nowcast predictions would be more accurate when  
 272 including the full list of candidate search terms as compared to a list of Lyme disease specific  
 273 terms, two sets of elastic net regression models were constructed: 1. models with only Lyme  
 274 disease symptoms and vector terms as features and 2. models with the full list of non-zero  
 275 variance terms identified from Google Correlate coupled with symptom and vector terms  
 276 (supplemental table 1). Predictions from regression models developed using data from  
 277 symptom and vector terms exclusively, produced accurate nowcasting models (assessed via  
 278  $R^2$ ) with low error (assessed via RMSE and MAE) in four of the five US regions (**Table 4**, results  
 279 for both models provided in supplemental table 3). The predictions from these models provide  
 280 accurate estimations of the timing of the seasonal pattern of Lyme disease (**Figure 5**).

281 **Table 4 Predictions from symptoms and vector terms only models produce accurate**  
 282 **predictions with low error**

	Northeast	Midwest	Southeast	Southwest	West
$\alpha, \lambda$	0.47, 0.60	0.33, 0.20	0.29, 0.07	0.11, 0.01	0.1, 0.01
<b>Training</b>					
<b>RMSE</b>	1.32	0.36	0.11	0.01	0.01
<b>MAE</b>	0.89	0.21	0.07	0.01	0.01
<b>R<sup>2</sup></b>	0.77	0.65	0.67	0.32	0.50
<b>In-sample Validation</b>					
<b>RMSE</b>	1.50	0.38	0.11	0.01	0.01
<b>MAE</b>	1.01	0.25	0.07	0.01	0.01
<b>R<sup>2</sup></b>	0.71	0.59	0.69	0.38	0.29
<b>Out of Sample</b>					
<b>RMSE</b>	1.65	0.43	0.14	0.01	0.01
<b>MAE</b>	1.38	0.34	0.10	0.01	0.01
<b>R<sup>2</sup></b>	0.79	0.76	0.82	0.37	0.63

283  
 284 **Figure 5 Elastic Net modeling using disease symptom and vector terms only produces**  
 285 **accurate nowcasting model for Lyme disease. (A) Northeast, (B) Midwest, (C) Southeast,**  
 286 **(D) Southwest, and (E) West. Same one-year period from Northeast region with accurate**  
 287 **nowcast model. Two elastic net models were developed for each region. Elastic net models**  
 288 **trained using CDC surveillance data and search term data from February 2004 through**  
 289 **December 2014. Vertical dashed line starts at January 2015 and indicates the start of the hold-**

290 out data set. Nowcasting performed using search term data from January 2018 until September  
291 2019.

292 Nowcasting models developed using the full list of search terms produced predictions  
293 that had a 1.33-fold improvement in accuracy and a 0.5-fold reduction in error compared to the  
294 symptom and vector only models (**Table 5**, results for both models provided in supplement table  
295 4). For each region it was found that using the full-term lists, which often included  
296 environmentally themed terms, increased the accuracy and reduced the error of predictions. On  
297 average, model accuracy ( $R^2$ ) improved by 0.2 when using the full list of search terms. The  
298 greatest improvement in accuracy when using the full-term list models was seen in the West ( $R^2$   
299 difference was 0.31). The Southeast had the least improvement (0.12) in accuracy. RMSE was  
300 reduced by 0.18 on average across all regions and MAE was reduced by 0.14 when comparing  
301 predictions between the full-term list models and the symptom and vector only models. The  
302 greatest reduction in error was seen in the Northeast region. It was found that predictions from  
303 the full-term list compared to the symptom and vector only models reduced RMSE by 0.69 and  
304 MAE by 0.56 cases per 100,000 population in the Northeast region. Reduction in error for the  
305 Southwest and West were found to be approximately 0 (RMSE = 0.001 and 0.004 respectively;  
306 MAE = 0.001 and 0.002 respectively). Predictions from the full list models also produced  
307 accurate timing of seasonal patterns of Lyme disease, but with improved mimicking of peaks  
308 and recessions (**Figure 6**). Compared to the symptom and vector term only models, predictions  
309 from the full-term list model showed more accurate variation in the spring and summer peaks of  
310 Lyme disease across all regions. In both modeling efforts, the Southwest consistently had the  
311 poorest predictive accuracy.

312 **Table 5 Predictions form full-term list models produce highly accurate predictions with**  
313 **low error**

	Northeast	Midwest	Southeast	Southwest	West
$\alpha, \lambda$	0.1, 0.85	0.93, 0.00	0.1, 0.07	0.1, 0.01	0.1, 0.00
<b>Training</b>					
<b>RMSE</b>	0.66	0.12	0.06	0.01	0.01
<b>MAE</b>	0.46	0.09	0.04	0.01	0.00

<b>R<sup>2</sup></b>	0.94	0.95	0.91	0.56	0.84
<b>In-sample Validation</b>					
<b>RMSE</b>	0.99	0.23	0.08	0.01	0.01
<b>MAE</b>	0.62	0.14	0.05	0.01	0.01
<b>R<sup>2</sup></b>	0.87	0.85	0.84	0.44	0.70
<b>Out of Sample</b>					
<b>RMSE</b>	0.74	0.29	0.14	0.01	0.01
<b>MAE</b>	0.52	0.17	0.09	0.01	0.01
<b>R<sup>2</sup></b>	0.97	0.94	0.91	0.45	0.82

314  
 315 **Figure 6 Elastic Net modeling using the full-term list produces predictions with greater**  
 316 **accuracy and less error. (A) Northeast, (B) Midwest, (C) Southeast, (D) Southwest, and (E)**  
 317 **West. Same one-year period from Northeast region with accurate nowcast model. Two elastic**  
 318 **net models were developed for each region. Elastic net models trained using CDC surveillance**  
 319 **data and search term data from February 2004 through December 2014. Vertical dashed line**  
 320 **starts at January 2015 and indicates the start of the hold-out data set. Nowcasting performed**  
 321 **using search term data from January 2018 until September 2019.**

322  
 323 In some years, Lyme disease incidence in the Northeast and Midwest showed  
 324 secondary peaks or plateaus in the post-summer spike of incident cases. These secondary  
 325 spikes or plateaus typically occur in late summer and early fall months as infected adult ticks  
 326 take blood meals transmitting Lyme disease to people. Predictions from models using only  
 327 symptoms and vector terms did not have sufficient sensitivity to detect to these changes (**Figure**  
 328 **7A and 7C**). Alternatively, predictions from the full-term list models had sufficient sensitivity to  
 329 detect these secondary spikes or plateaus of decreasing incidence at the regional level (**Figure**  
 330 **7B and 7D**).

331 **Figure 7 Elastic net modeling using full-term list is sensitive to secondary spikes of Lyme**  
 332 **disease incidence in Northeast and Midwest regions. (A) Northeast Lyme disease incidence**  
 333 **and disease symptom and vector terms only model predictions, (B) Northeast Lyme disease**  
 334 **incidence and full-term list model predictions, (C) Midwest Lyme disease incidence and disease**  
 335 **symptom and vector terms only model predictions, and (D) Midwest Lyme disease incidence**  
 336 **and full-term list model predictions. Elastic net models trained using CDC surveillance data and**  
 337 **search term data from February 2004 through December 2014 and hold-out data from January**  
 338 **2015 and December 2017.**

339 Statistical learning techniques can help highlight specific areas in which future  
 340 hypothesis or interventions could be generated. We identified the three most important terms  
 341 from the accurate full-term list nowcasting models. (**Table 6**). As hypothesized, many of the top

342 three most important terms for producing accurate nowcasts were regionally specific and  
 343 environmentally themed. The Northeast and Southeast were the only regions that had a  
 344 potential symptom term (bull's-eye rash, rash) identified in the top three important terms. We  
 345 further hypothesized that due to the importance of these environmentally related themes, the  
 346 time series of these search terms trends would mimic the same general trends for Lyme  
 347 disease. These patterns are particular evident in areas with higher incidence of Lyme disease;  
 348 the Northeast, Midwest and Southeast (**Figure 8**). It was found that the search traffic for these  
 349 top three terms aligns with the peaks and recessions of Lyme disease on the same monthly  
 350 scale.

351 **Table 6 Three most important terms for each model often environmentally themed**

Northeast			
Elastic Net 1		Elastic Net 2	
Search Term	Scaled Importance	Search Term	Scaled Importance
July Calendar	100.00	July Calendar	100.00
Fresh Cherry Pie	82.12	Fresh Cherry Pie	83.29
Bullseye Rash	75.51	Bullseye Rash	75.47
Midwest			
Elastic Net 1		Elastic Net 2	
Search Term	Scaled Importance	Search Term	Scaled Importance
Festivals Milwaukee	100.00	Festivals Milwaukee	100.00
Lake Beaches	97.35	Kings Island Discount	99.16
Kings Island Discount	96.35	Lake Beaches	97.40
Southeast			
Elastic Net 1		Elastic Net 2	
Search Term	Scaled Importance	Search Term	Scaled Importance
Intex Pool Cover	100.00	Intex Pool Cover	100.00
Rash	87.07	Rash	88.06
Swampdogs	85.64	Swampdogs	85.45
Southwest			
Elastic Net 1		Elastic Net 2	
Search Term	Scaled Importance	Search Term	Scaled Importance
Loans for	100.00	Loans for	100.00
CA Water	67.20	CA Water	66.82
Hotels CA	61.00	Hotels CA	60.14
West			

Elastic Net 1		Elastic Net 2	
Search Term	Scaled Importance	Search Term	Scaled Importance
Movies in the Park	100.00	Movies in the Park	100.00
Concert in the Park	69.18	Concert in the Park	69.65
Waterworld Denver	62.13	Waterworld Denver	62.44

352  
353 **Figure 8 Time series of regional candidate search terms for simple Lyme disease**  
354 **tracking. (A).** Northeast, **(B).** Midwest, and **(C).** Southeast. The top three most important terms  
355 from each region model identified by *varImp* function in R. (a-c). Candidate terms scaled to align  
356 with regional Lyme disease incidence. Terms presented directly as provided by Google  
357 Correlate™.

## 358 Discussion

359 With the growing incidence of Lyme disease in the United states, novel methods that  
360 help health departments to prepare for years of increased Lyme disease exposure are critical.  
361 We found that when using google search history data in nowcasting, accurate predictions of  
362 Lyme disease can be generated. Importantly, the search traffic for the top three search terms  
363 generally followed the same temporal nature of regional Lyme disease incidence. These terms  
364 and nowcasting methods could help Health Departments determine approximate trends of Lyme  
365 disease in their area by monitoring the search traffic trends of the terms via the free tool of  
366 Google Trends™. Additionally, many of the terms that remained in these accurate models were  
367 environmentally themed and can be used to generate future hypotheses for intervention and  
368 prevention activities.

369 Overall, each elastic net model performed well and provided accurate estimations of the  
370 of regional Lyme disease incidence provided by surveillance data from the CDC (**Table 4** and  
371 **5**). Results showed that predictions were more accurate from models using a full list of  
372 colloquial search terms the average person is likely to search compared to models that only  
373 used symptom, disease or vector terms. It was also found that predictions from models that  
374 included the full-term list were more sensitive to detecting secondary spikes and recession  
375 plateaus in the fall months of the Northeast and Midwest (**Figure 7**). Moreover, many of the  
376 search terms identified via Google Correlate which had high levels of bivariate correlation and



377 remained important throughout the elastic net modeling process were environmentally related.  
378 While not all these terms directly relate to an activity that have obvious risk of tick exposure and  
379 transmission of Lyme disease, environmentally related terms can serve as a proxy for an  
380 intention for people to spend time outdoors. Increased time spent outdoors has been shown to  
381 increase exposure to ticks in the environment (38-40). Causal inference cannot be directly  
382 drawn from these results, however given the common pattern of environmental terms and many  
383 of their high correlations a pattern has emerged. These terms can help LHDs generate  
384 hypothesis on where to perform future tick surveillance, implement intervention measures, or  
385 spread tick awareness. These findings suggest the importance of including colloquial search  
386 terms over symptom or vector related terms alone for current and future prediction efforts. Our  
387 models can be implemented by LHDs as they currently are, or terms that more specific the local  
388 populations search habits can be substituted to further improve performance.

389         The Southwest, a non-endemic region for Lyme disease (15), continually had the  
390 poorest performing predictions. *Ixodes* ticks in the Southwest are more suspected to feed on  
391 lizards and other non-reservoir hosts (41), thus it is not surprising that Lyme disease incidence  
392 was low. The CDC also classifies county of residence and not county of acquisition in  
393 surveillance reports therefore it is likely that those diagnosed in this region were exposed  
394 elsewhere. The Southwest also had the lowest number of feature data compared to all other  
395 regions. These all likely led to the low performance of predictions in this region. On the other  
396 hand, the West region, which also had a low number of incident cases, but had a greater  
397 number of feature data had better performing model predictions. The West also has suitable  
398 habitat for *Ixodes pacificus*, a known vector of *B. burgdorferi* (15). These results indicate that in  
399 addition to having an appropriate number of feature data and outcomes, regions also need to  
400 have a suitable environment for the tick vectors in order to produce accurate nowcasts. These  
401 findings continue to show the importance of inducing environment related feature data for

402 current or future prediction efforts in areas that are either endemic with Lyme disease or have  
403 suitable *Ixodid* tick habitats.

404 To our knowledge, two prior studies have been performed using google search data to  
405 try and improve model performance (42, 43). One study concluded that using a single term,  
406 “Borreliose”, was not helpful in improving model accuracy (42). While “Borreliose” is a medically  
407 accurate term for Lyme disease, we found that colloquial disease terms had moderate to high  
408 levels of correlation. Our findings found that the bivariate correlation for disease symptoms and  
409 colloquial disease terms ranged from -0.33 to 0.85 across five U.S. regions. Terms often  
410 moderately (correlation value > 0.5), or highly correlated (correlation value > 0.8), with regional  
411 monthly Lyme disease incidence included: “lyme disease”, “lyme”, “rash”, and “tick”. Further,  
412 environmentally related terms often had the highest levels of correlation across all regions.  
413 Another study developed a tool, Lymelight, which monitored the incidence of Lyme disease in  
414 real time using Lyme disease symptom web searches in a two-year period to predict future  
415 Lyme disease burden and treatment impacts (43). Despite producing accurate models, this  
416 method only used symptom terms which may not predict true patterns of Lyme disease or risky  
417 behaviors. Our findings show using symptom, disease and vector terms in combination with  
418 terms that focus on environments in which one may have the risk of being exposed can greatly  
419 improve model performance over symptom and vector terms alone. These findings continue to  
420 suggest the importance of direct or proxy measures for time spent outdoors when predicting  
421 vector-borne diseases.

422 An advantage of using data from Google search history, R studio as a modeling  
423 software, and elastic net regression is that accurate predictions can be made quickly  
424 (approximately 24 hours from start to finish) and free. This can allow LHDs to have more up to  
425 date estimations of regional Lyme disease incidence beyond federal report schedules without  
426 additional financial burden. We found when graphing the search traffic for three most important

427 terms from regional models, in endemic areas of the Northeast and Midwest, as hypothesized  
428 they provide a very good broad scale of timing. Following these terms, or more locally specific  
429 environmental terms could provide even quicker tracking of general temporal trends of Lyme  
430 disease for LHDs. Most of the top three important terms were environmentally related. This  
431 further suggests the importance of including terms or variables that focus on the environment for  
432 current and future prediction efforts.

433         While there are strengths of statistical learning approaches, there are limitations to our  
434 approach as well. These models were developed at the regional level and are subject to less  
435 accurate predictions at the state or local level without refitting the model. Additionally, grouping  
436 states into different regions will alter results of these findings as both regional rate and search  
437 term identification using Google Correlate™ were performed regional aggregation strategy.  
438 These models are not generalizable to other vector-borne diseases in their current form. Similar  
439 approaches could be used for other vector-borne diseases such as Anaplasmosis, as this is  
440 also vectored by *Ixodid* ticks and therefore will have similar temporal trends and environmental  
441 risk factors. Additionally, these models are not generalizable to other countries. All the Lyme  
442 disease and search data were based on US disease and Google habits, it is unlikely that our  
443 developed models would produce accurate results in other countries. However, a similar  
444 approach could be used in other countries that have strong surveillance data and a free access  
445 database of the countries' most utilized search engine. Moreover, other sources of data on  
446 human behavior (i.e. data from social networks like Twitter) present additional opportunities for  
447 such models, potentially at greater spatial and temporal granularity. Greater consideration or  
448 different modeling techniques may need to be implemented for communicable diseases.  
449 However, these models can be incorporated to get a general idea of surrounding areas for  
450 those LHDs that are vastly underfunded. Local or regionally specific terms could easily be  
451 substituted into these models which could help improve model fit on a case-by-case basis.

452 These findings highlight the importance of strong disease surveillance and computational  
453 modeling efforts working together. Predictions over time are likely to improve not only due to  
454 increases in statistical and computing power, but in the maintenance and enhancement of  
455 strong disease surveillance efforts performed nationwide.

## 456 **References**

- 457 1. Wormser GP, Dattwyler RJ, Shapiro ED, Halperin JJ, Steere AC, Klemperer MS, et al. The clinical  
458 assessment, treatment, and prevention of Lyme disease, human granulocytic anaplasmosis, and  
459 babesiosis: clinical practice guidelines by the Infectious Diseases Society of America. *Clin Infect Dis*.  
460 2006;43(9):1089-134.
- 461 2. Mead PS. Epidemiology of Lyme disease. *Infect Dis Clin North Am*. 2015;29(2):187-210.
- 462 3. Piesman J. Transmission of Lyme disease spirochetes (*Borrelia burgdorferi*). *Experimental &*  
463 *applied acarology*. 1989;7(1):71-80.
- 464 4. Wood CL, Lafferty KD. Biodiversity and disease: a synthesis of ecological perspectives on Lyme  
465 disease transmission. *Trends in ecology & evolution*. 2013;28(4):239-47.
- 466 5. Gilmore Jr RD, Mbow ML, Stevenson B. Analysis of *Borrelia burgdorferi* gene expression during  
467 life cycle phases of the tick vector *Ixodes scapularis*. *Microbes and infection*. 2001;3(10):799-808.
- 468 6. Mather TN, Nicholson MC, Donnelly EF, Matyas BT. Entomologic index for human risk of Lyme  
469 disease. *American Journal of Epidemiology*. 1996;144(11):1066-9.
- 470 7. Stafford III KC. Survival of immature *Ixodes scapularis* (Acari: Ixodidae) at different relative  
471 humidities. *Journal of medical entomology*. 1994;31(2):310-4.
- 472 8. Simmons T, Shea J, Myers-Claypole M, Kruse R, Hutchinson M. Seasonal activity, density, and  
473 collection efficiency of the blacklegged tick (*Ixodes scapularis*)(Acari: Ixodidae) in mid-western  
474 Pennsylvania. *Journal of medical entomology*. 2015;52(6):1260-9.
- 475 9. Berger KA, Ginsberg HS, Dugas KD, Hamel LH, Mather TN. Adverse moisture events predict  
476 seasonal abundance of Lyme disease vector ticks (*Ixodes scapularis*). *Parasites & vectors*. 2014;7(1):181.
- 477 10. Mead PS. Epidemiology of Lyme disease. *Infectious Disease Clinics*. 2015;29(2):187-210.
- 478 11. Guerra M, Walker E, Jones C, Paskewitz S, Cortinas MR, Stancil A, et al. Predicting the risk of  
479 Lyme disease: habitat suitability for *Ixodes scapularis* in the north central United States. *Emerging*  
480 *infectious diseases*. 2002;8(3):289.
- 481 12. Steere AC, Coburn J, Glickstein L. The emergence of Lyme disease. *The Journal of clinical*  
482 *investigation*. 2004;113(8):1093-101.
- 483 13. Adrion ER, Aucott J, Lemke KW, Weiner JP. Health care costs, utilization and patterns of care  
484 following Lyme disease. *PloS one*. 2015;10(2):e0116767.
- 485 14. Maes E, Lecomte P, Ray N. A cost-of-illness study of Lyme disease in the United States. *Clinical*  
486 *Therapeutics*. 1998;20(5):993-1008.
- 487 15. Schwartz AM, Hinckley AF, Mead PS, Hook SA, Kugeler KJ. Surveillance for Lyme disease—United  
488 States, 2008–2015. *MMWR Surveillance Summaries*. 2017;66(22):1.
- 489 16. CDC. Lyme disease Surveillance and available data 2020 [updated November 22, 2019. Available  
490 from: Lyme disease surveillance and available data.
- 491 17. Boeke MC, Zahner SJ, Booske BC, Remington PL. Local public health department funding: trends  
492 over time and relationship to health outcomes. *Wisconsin Medical Journal (WMJ)*. 2008;107(1):25.

- 493 18. Tang SL, Subramanian P. Review on nowcasting using least absolute shrinkage selector operator  
494 (LASSO) to predict dengue occurrence in San Juan and Iquitos as part of disease surveillance system.  
495 *Periodicals of Engineering and Natural Sciences*. 2019;7(2):608-17.
- 496 19. Schmidt CW. Trending now: using social media to predict and track disease outbreaks. *National*  
497 *Institute of Environmental Health Sciences*; 2012.
- 498 20. Generous N, Fairchild G, Deshpande A, Del Valle SY, Priedhorsky R. Global disease monitoring  
499 and forecasting with Wikipedia. *PLoS computational biology*. 2014;10(11):e1003892.
- 500 21. Ogden N, AbdelMalik P, Pulliam J. Emerging Infections: Emerging infectious diseases: prediction  
501 and detection. *Canada Communicable Disease Report*. 2017;43(10):206.
- 502 22. Pavlicek J, Kristoufek L. Nowcasting unemployment rates with google searches: Evidence from  
503 the visegrad group countries. *PloS one*. 2015;10(5):e0127084.
- 504 23. Sarigul S, Rui H. Nowcasting obesity in the US using Google search volume data. 2014.
- 505 24. Scharkow M, Vogelgesang J. Measuring the public agenda using search engine queries.  
506 *International Journal of Public Opinion Research*. 2011;23(1):104-13.
- 507 25. Siettos CI, Russo L. Mathematical modeling of infectious disease dynamics. *Virulence*.  
508 2013;4(4):295-306.
- 509 26. Zou H, Hastie T. Regularization and variable selection via the elastic net. *Journal of the royal*  
510 *statistical society: series B (statistical methodology)*. 2005;67(2):301-20.
- 511 27. Ginsberg J, Mohebbi MH, Patel RS, Brammer L, Smolinski MS, Brilliant L. Detecting influenza  
512 epidemics using search engine query data. 2009;457(7232):1012-4.
- 513 28. Polgreen PM, Chen Y, Pennock DM, Nelson FD, Weinstein RA. Using internet searches for  
514 influenza surveillance. 2008;47(11):1443-8.
- 515 29. Schiffman E, McLaughlin C, Ray J, Kemperman M, Hinckley A, Friedlander H, et al.  
516 Underreporting of Lyme and Other Tick-Borne Diseases in Residents of a High-Incidence County,  
517 Minnesota, 2009. *Zoonoses and public health*. 2018;65(2):230-7.
- 518 30. Mohebbi M, Vanderkam D, Kodysh J, Schonberger R, Choi H, Kumar S. Google correlate  
519 whitepaper. 2011.
- 520 31. Massicotte P, Edelbuettel D, Massicotte MP. Package 'gtrendsR'. 2016.
- 521 32. Google Trends [Available from: <https://trends.google.com/trends/?geo=US>].
- 522 33. Preis T, Moat HS. Adaptive nowcasting of influenza outbreaks using Google searches.  
523 2014;1(2):140095.
- 524 34. Lamos V, Cristianini NJ. Technology. Nowcasting events from the social web with  
525 statistical learning. 2012;3(4):1-22.
- 526 35. Kuhn M. The caret package. R Foundation for Statistical Computing, Vienna, Austria URL  
527 [https://cran](https://cran.r-project.org/package=caret) r-project org/package= caret. 2012.
- 528 36. Hastie T, Qian J. Glmnet vignette. Retrieve from [http://www](http://www.web.stanford.edu/~hastie/Papers/Glmnet_Vignette.pdf) web stanford edu/~  
529 hastie/Papers/Glmnet\_Vignette.pdf Accessed September. 2014;20:2016.
- 530 37. Kang M, Zhong H, He J, Rutherford S, Yang F. Using google trends for influenza surveillance in  
531 South China. 2013;8(1).
- 532 38. Mead P, Hook S, Niesobecki S, Ray J, Meek J, Delorey M, et al. Risk factors for tick exposure in  
533 suburban settings in the Northeastern United States. *Ticks and tick-borne diseases*. 2018;9(2):319-24.
- 534 39. Salkeld DJ, Porter WT, Loh SM, Nieto NC. Time of year and outdoor recreation affect human  
535 exposure to ticks in California, United States. *Ticks and tick-borne diseases*. 2019;10(5):1113-7.
- 536 40. Schwartz BS, Goldstein MD. Lyme disease in outdoor workers: risk factors, preventive measures,  
537 and tick removal methods. *American journal of epidemiology*. 1990;131(5):877-85.
- 538 41. Tietjen M, Esteve-Gasent MD, Li AY, Medina RF. A comparative evaluation of northern and  
539 southern *Ixodes scapularis* questing height and hiding behaviour in the USA. *Parasitology*.  
540 2020;147(13):1569-76.

541 42. Kapitány-Fövény M, Ferenci T, Sulyok Z, Kegele J, Richter H, Vályi-Nagy I, et al. Can Google  
542 Trends data improve forecasting of Lyme disease incidence? Zoonoses and public health.  
543 2019;66(1):101-7.

544 43. Sadilek A, Hswen Y, Bavadekar S, Shekel T, Brownstein JS, Gabrilovich E. Lymelight: forecasting  
545 Lyme disease risk using web search data. npj Digital Medicine. 2020;3(1):1-12.

546

547

548 **S1 Table Complete list of search terms identified by Google Correlate from Each Region.**

549 Terms for each region were identified via Google Correlate using region specific Lyme disease  
550 rates from training period data.

551 **S2 Table Bivariate correlations of each search term to the regional Lyme disease rates.**

552 Pearson Correlations values were calculated between each term monthly proportional search  
553 data and corresponding Lyme disease rates for each term and region.

554 **S3 Table Predictions from symptoms and vector terms only models produce accurate  
555 predictions with low error**

556

557 **S4 Table Predictions form full list models produce highly accurate predictions with low  
558 error**

559

560

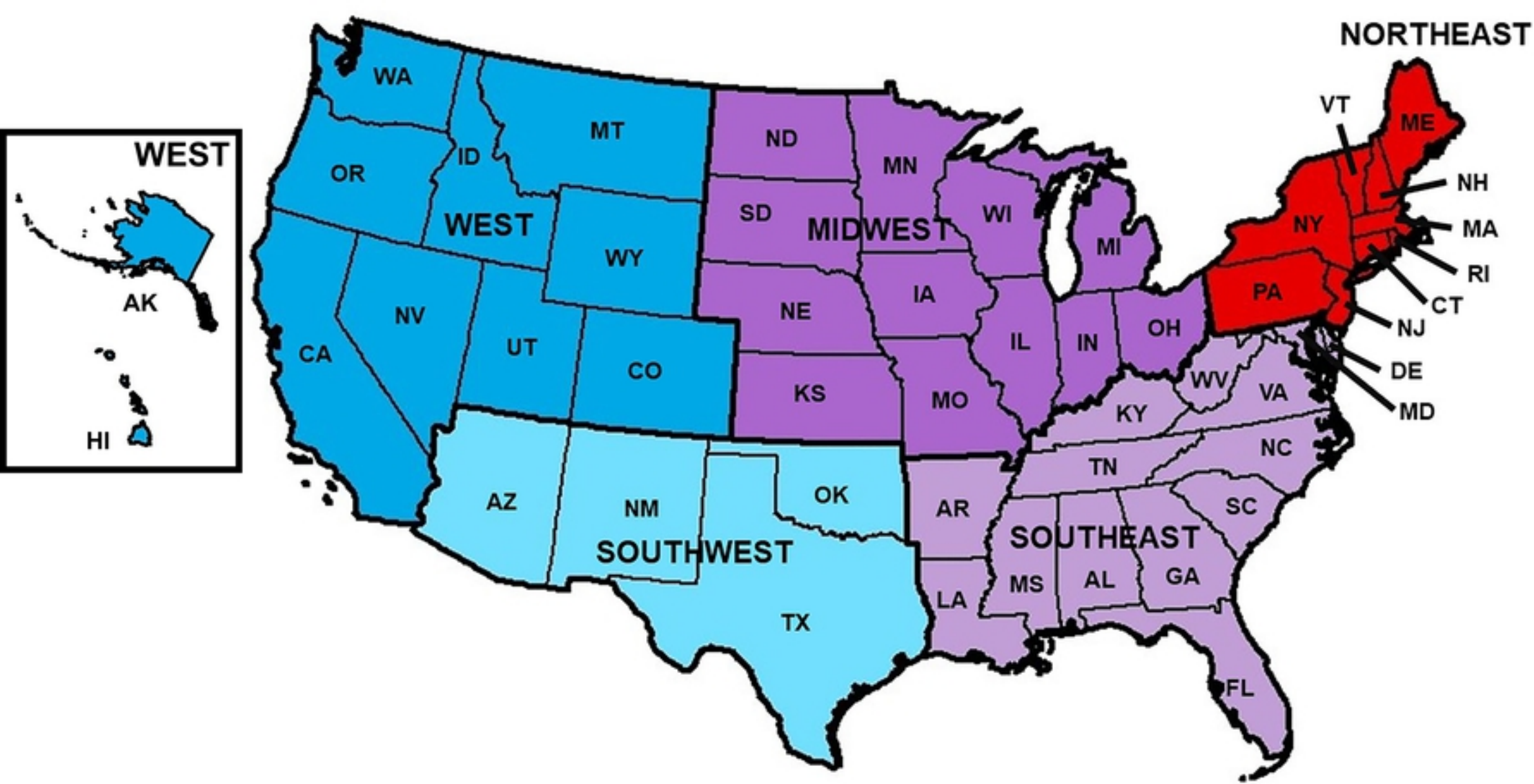


Figure 1

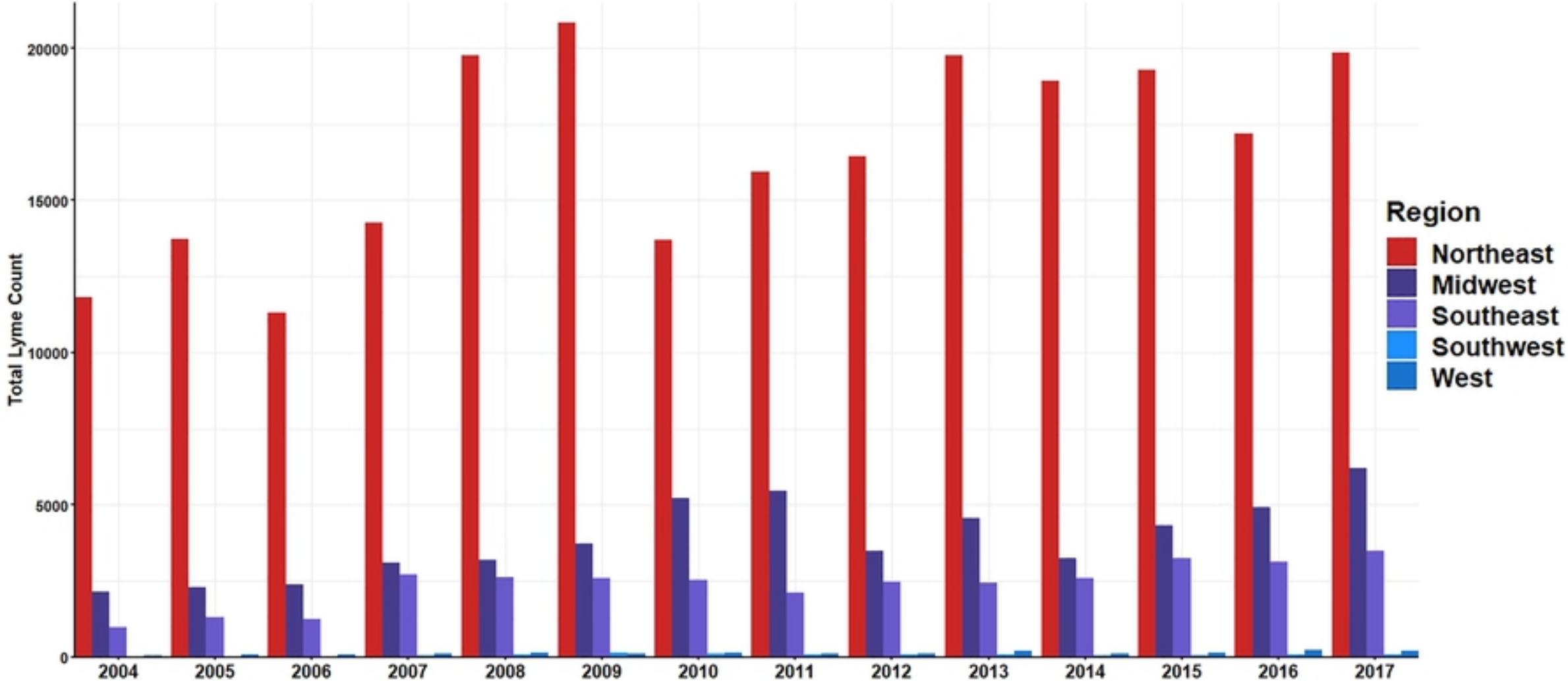


Figure 2



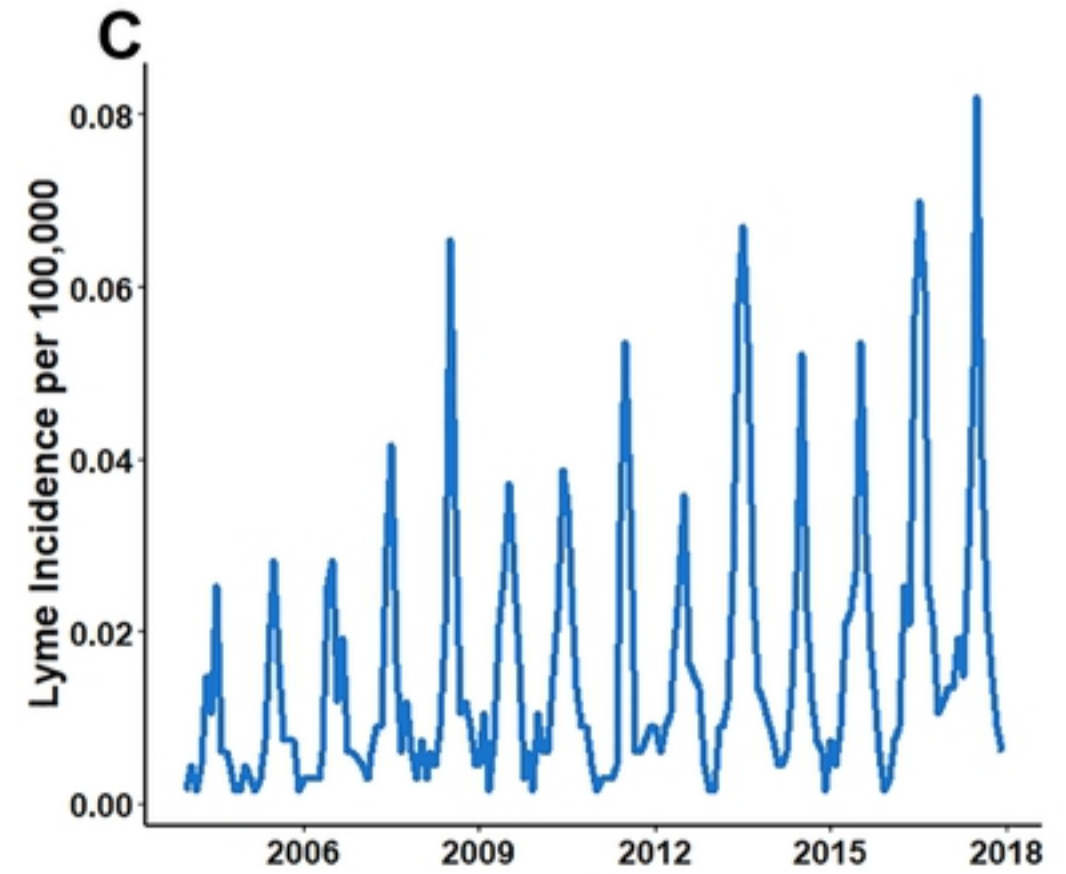
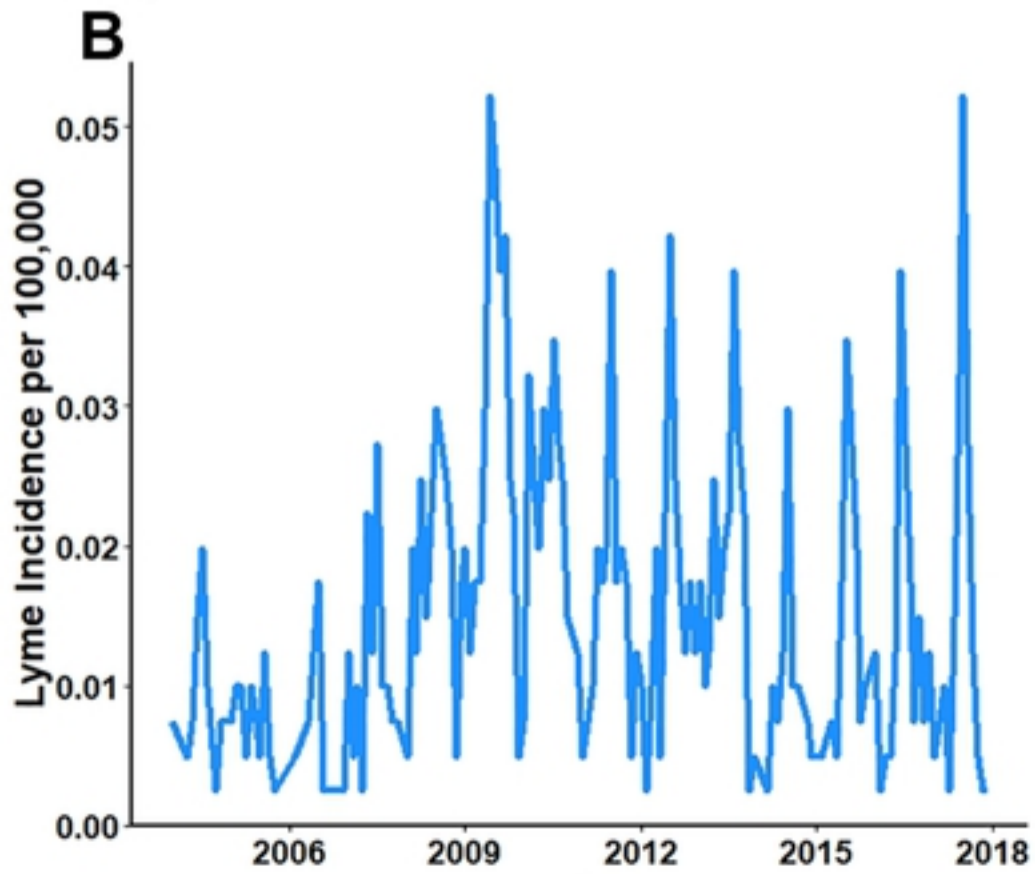
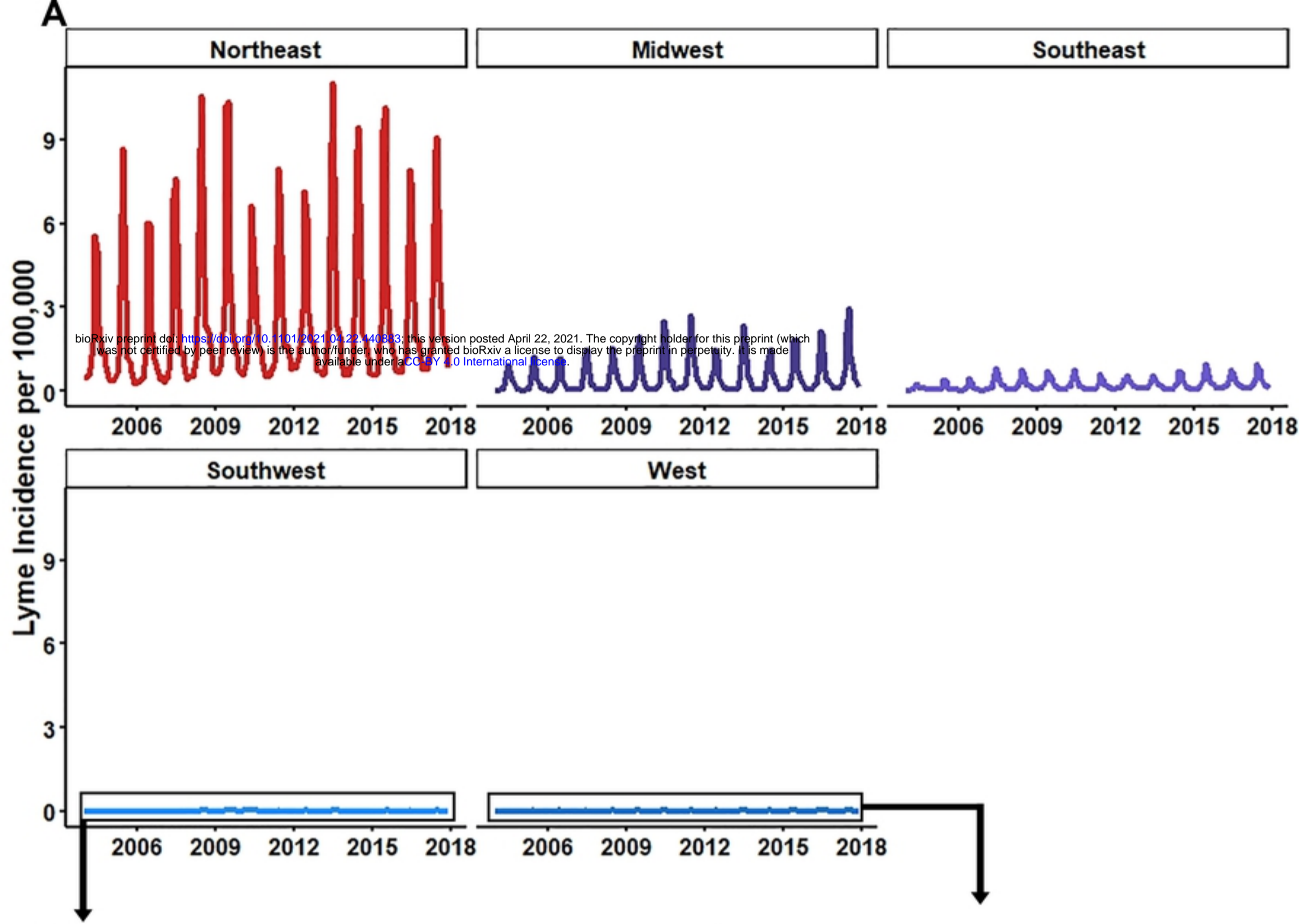


Figure 3

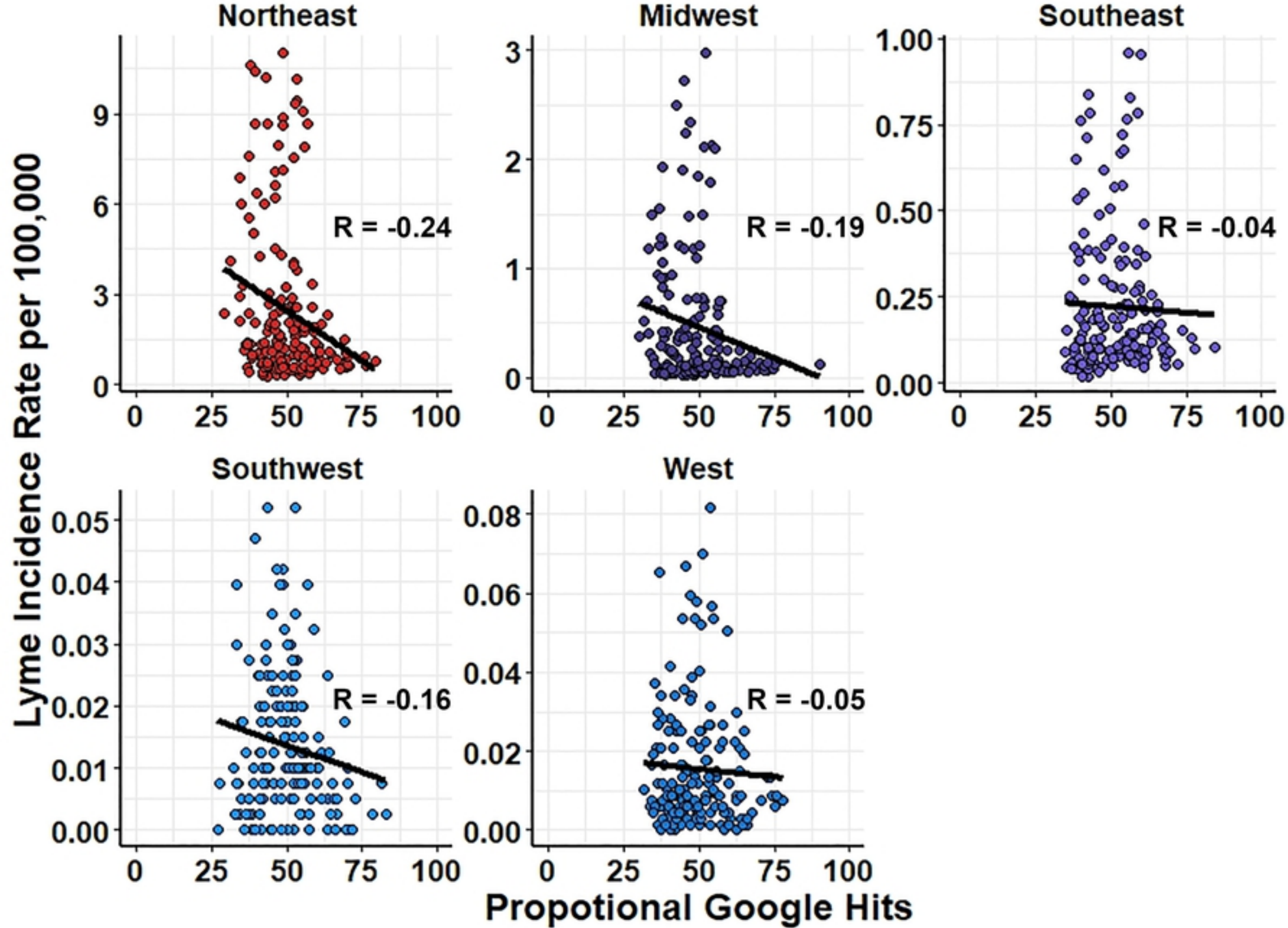


Figure 4

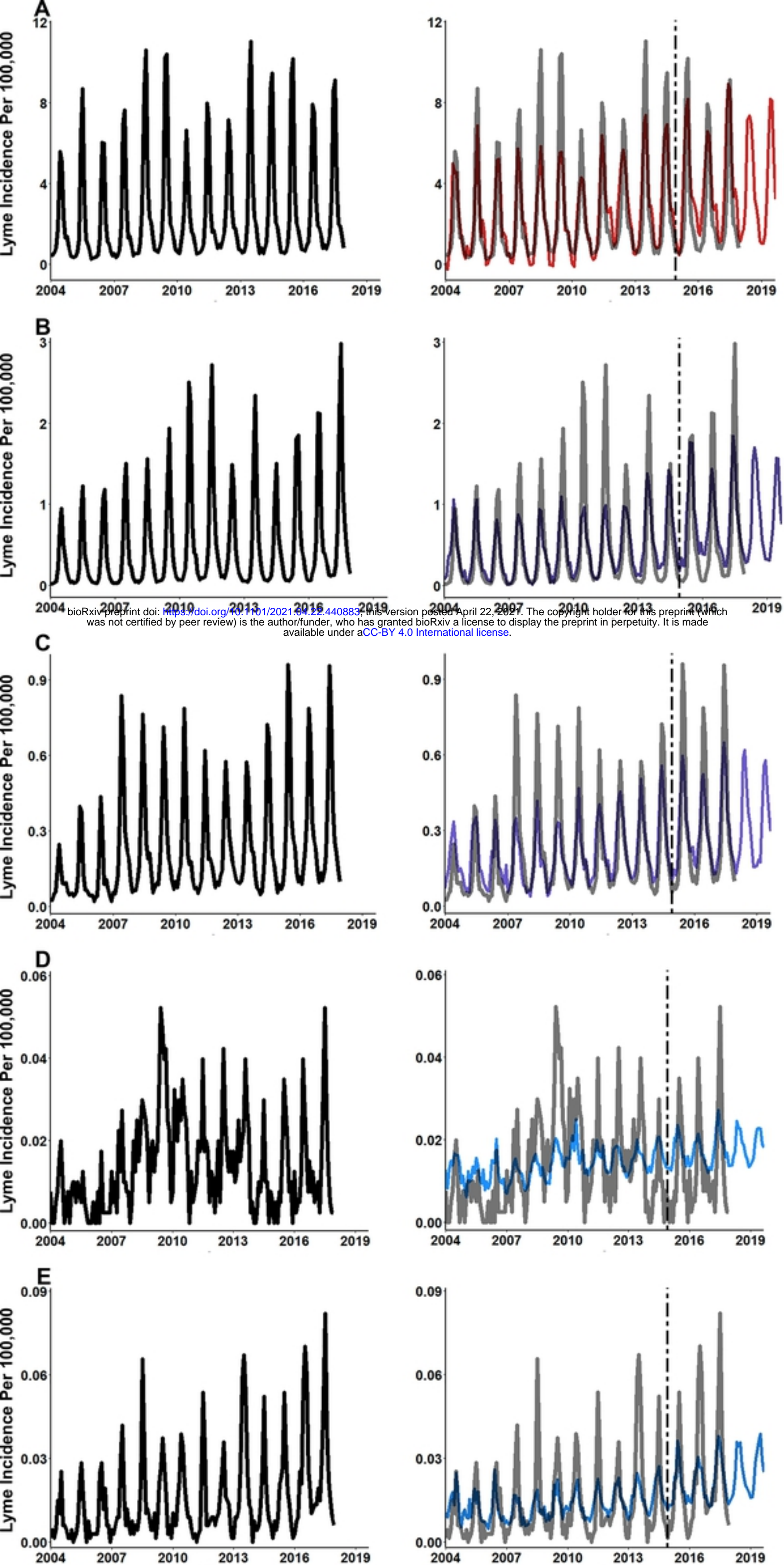


Figure 5

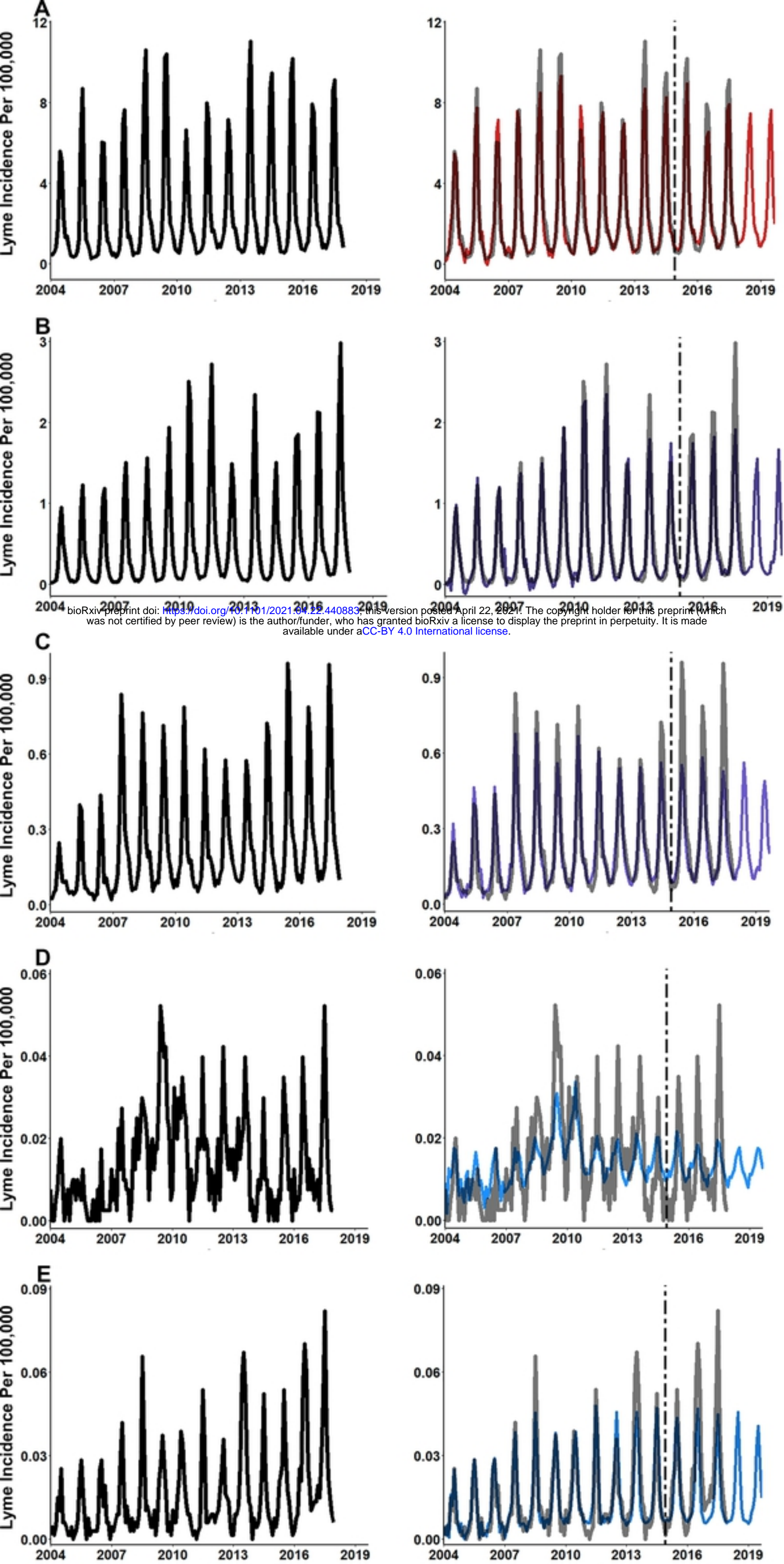


Figure 6

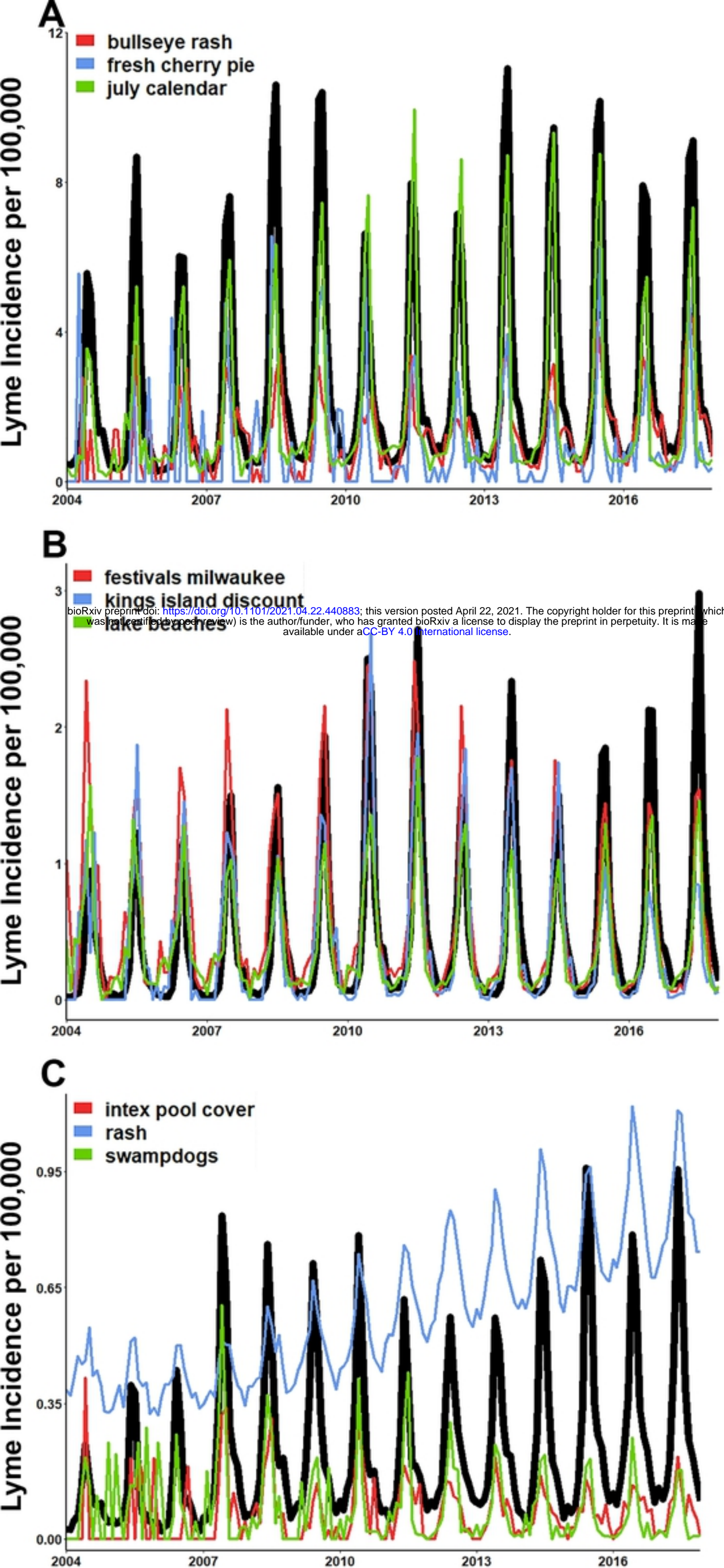


Figure 8

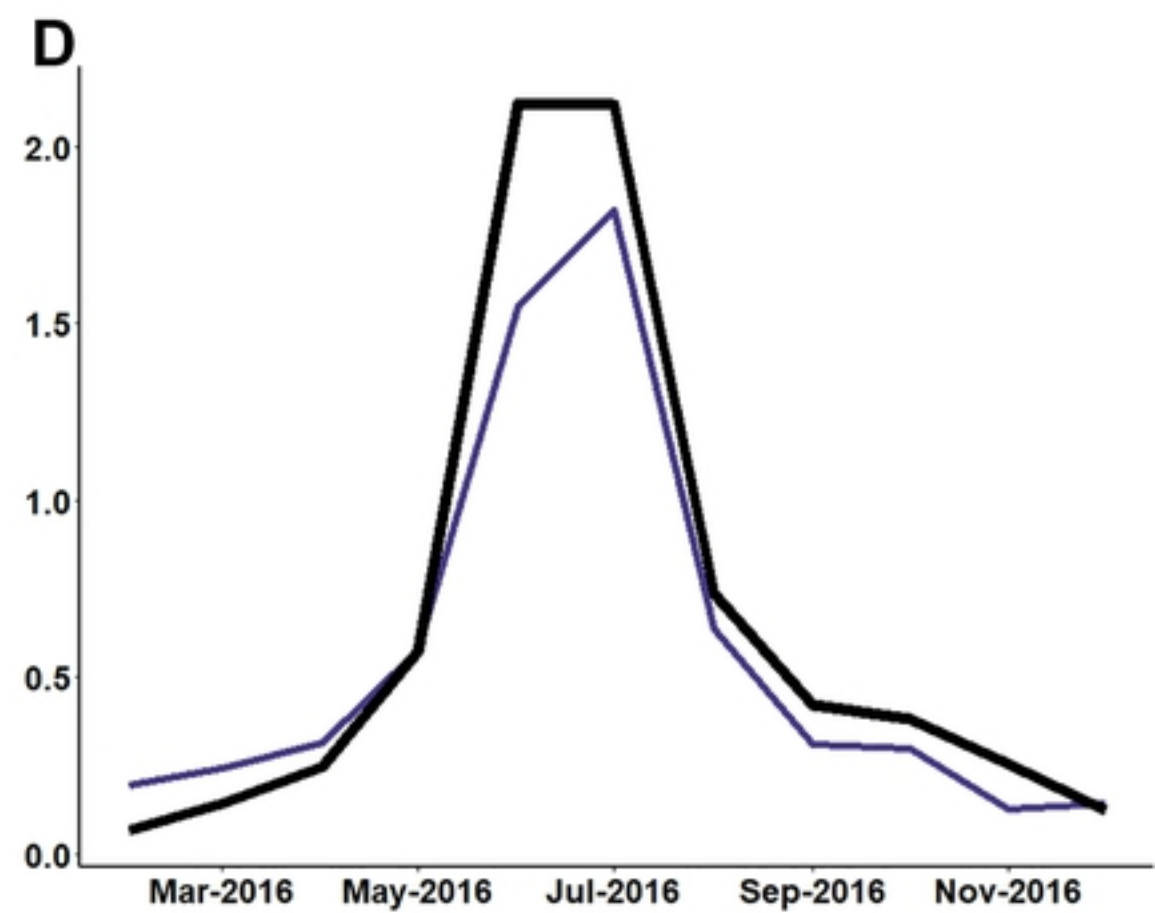
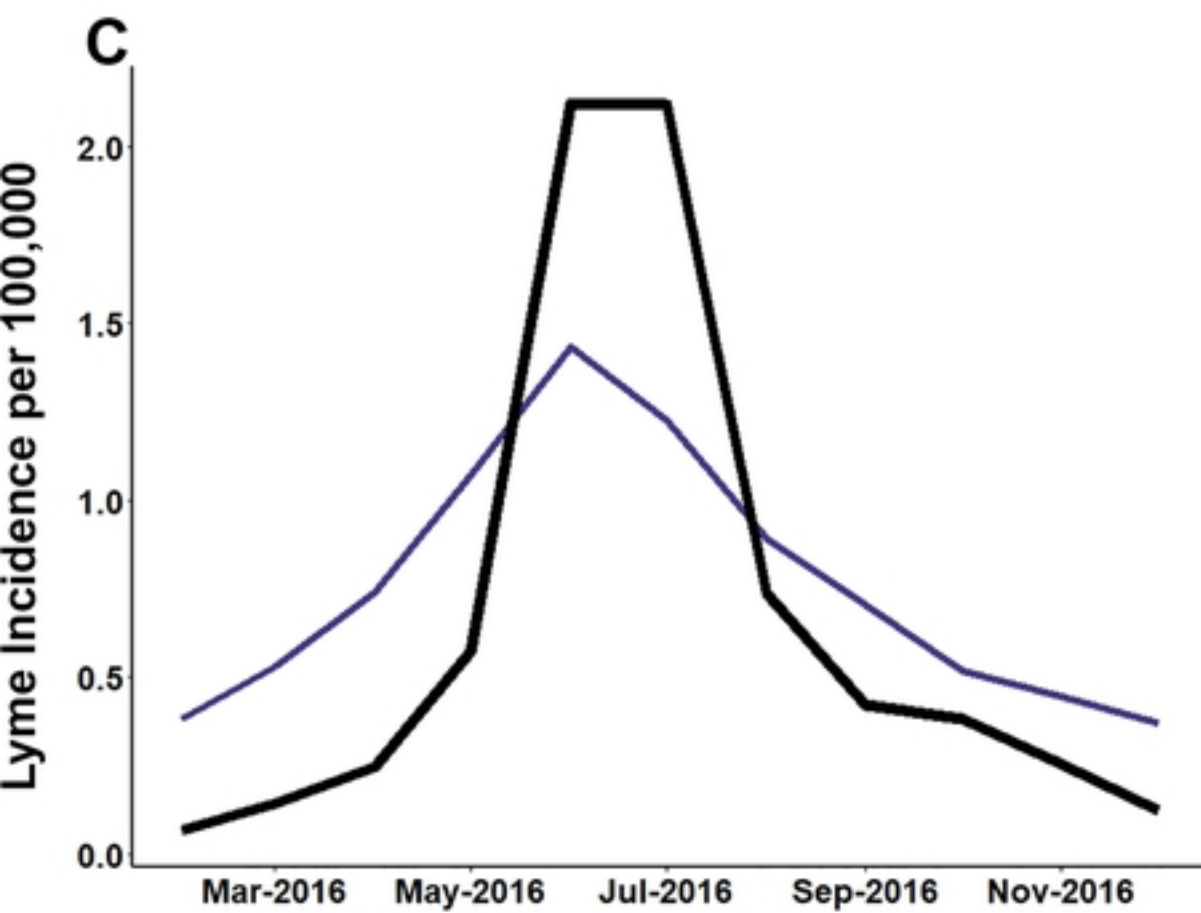
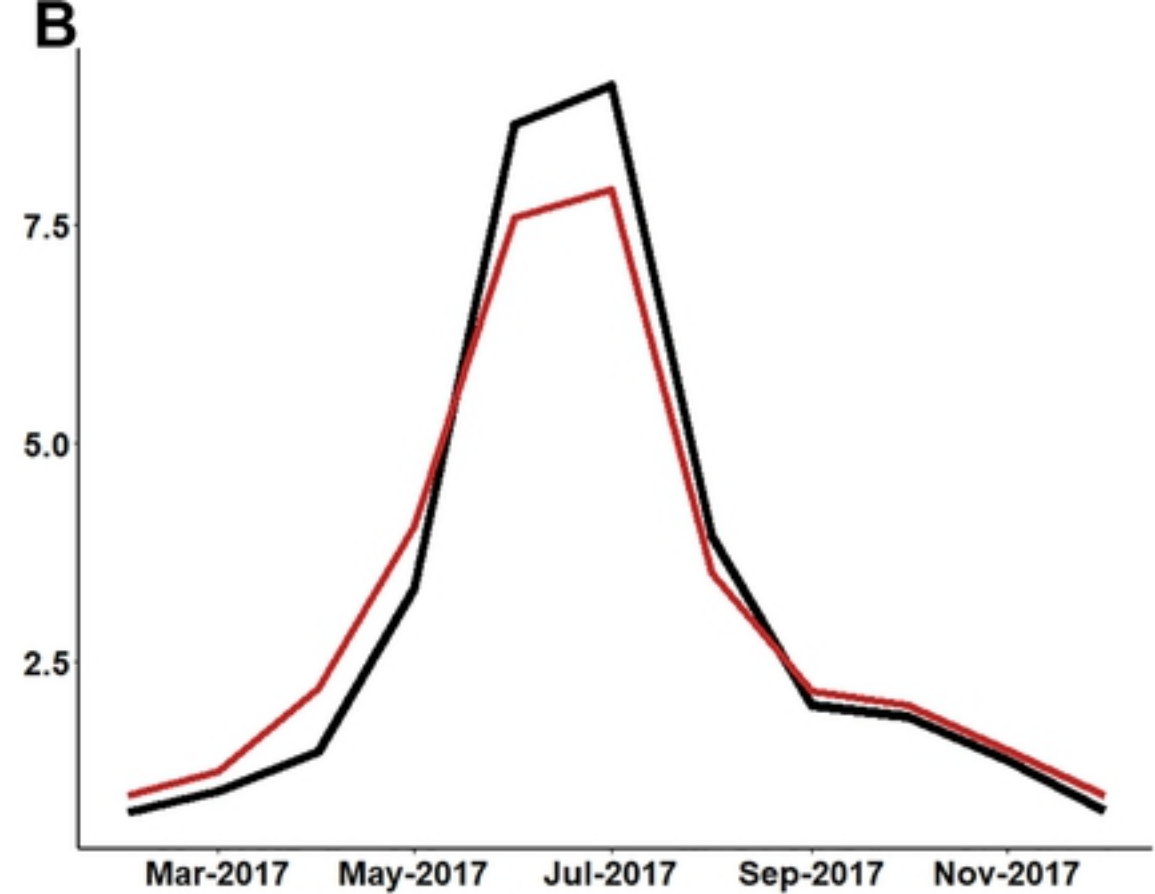
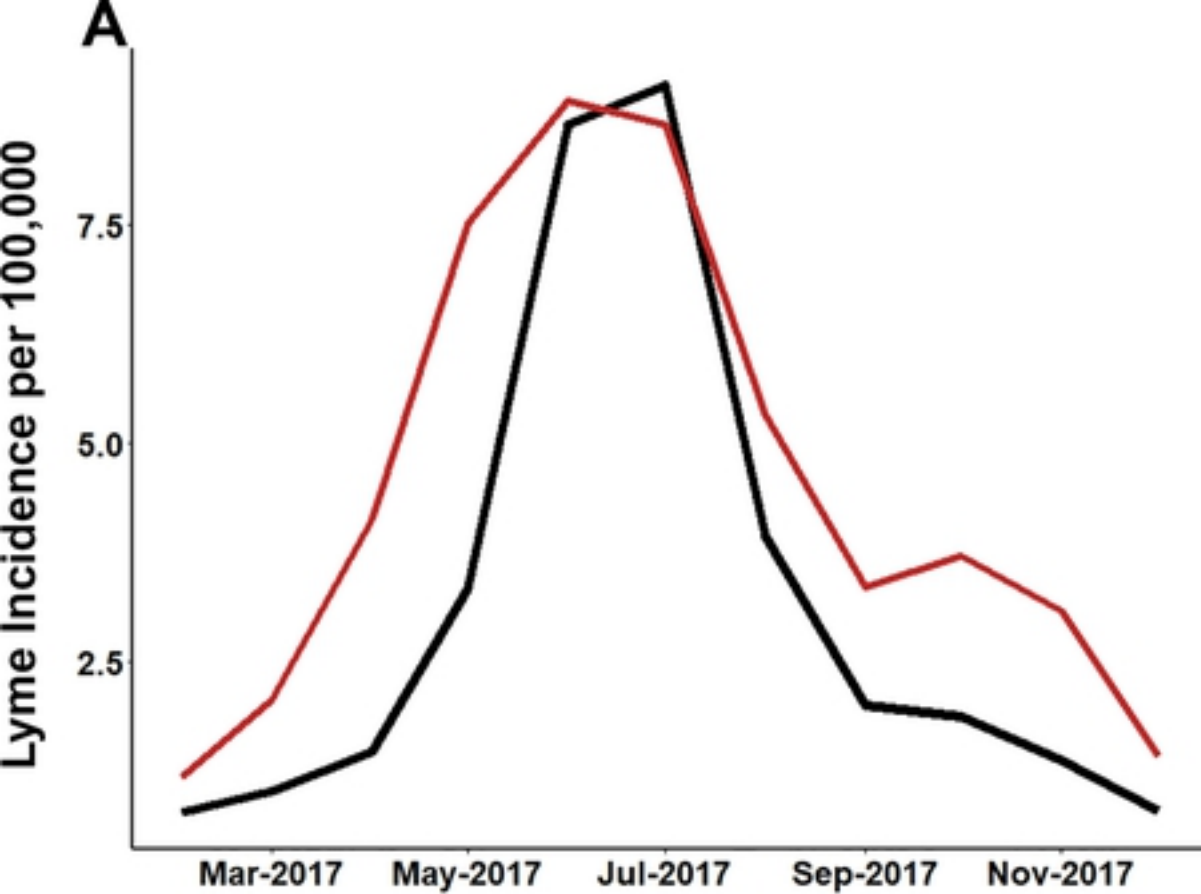


Figure 7