

1 Resampling-Based Methods for Biologists

2 John Fieberg¹, Kelsey Vitense¹, and Douglas H. Johnson¹

3 ¹Department of Fisheries, Wildlife and Conservation Biology, University of Minnesota,
4 2003 Upper Buford Circle, Suite 135, Saint Paul MN 55108, USA

5 Corresponding author:

6 John Fieberg¹

7 Email address: jfieberg@umn.edu

8 ABSTRACT

Ecological data often violate common assumptions of traditional parametric statistics (e.g., that residuals are Normally distributed, have constant variance, and cases are independent). Modern statistical methods are well equipped to handle these complications, but they can be challenging for non-statisticians to understand and implement. Rather than default to increasingly complex statistical methods, resampling-based methods can sometimes provide an alternative method for performing statistical inference, while also facilitating a deeper understanding of foundational concepts in frequentist statistics (e.g., sampling distributions, confidence intervals, p -values). Using simple examples and case studies, we demonstrate how resampling-based methods can help elucidate core statistical concepts and provide alternative methods for tackling challenging problems across a broad range of ecological applications.

9 INTRODUCTION

10 Falsifiable hypotheses and replication are two cornerstones of science (Johnson, 2002; Popper, 2005).
11 Replication is also critical for understanding key statistical concepts in frequentist statistics (Box 1). Yet,
12 as researchers, we typically encounter and analyze one data set at a time, making it difficult to understand
13 the concept of a *sampling distribution* (i.e., the distribution of a sample statistic across repeated samples
14 from the same population). Fully appreciating sampling distributions requires conceptualizing the process
15 of repeatedly collecting new data sets, of the same size and using the same methods of data collection,
16 and analyzing these data in the same way as in the original analysis. Similarly, *null distributions*, used to
17 test statistical hypotheses and calculate p -values, require that we consider the distribution of our statistics
18 (means, regression coefficients, etc.) across repeated hypothetical data collection and analysis efforts,
19 while adding a further constraint that the null hypothesis is true.

20 Rather than serving as a unifying concept, the importance of variability across different replicate
21 samples is often lost upon students when they take their first statistics course. The traditional formula-
22 based approach to teaching gives the impression that statistics is little more than a set of recipes, each
23 one suited to a different table setting of data. Most of these recipes rely on large-sample assumptions
24 that allow us to derive an appropriate (albeit approximate) sampling distribution. For example, the
25 sampling distribution of many statistics will be well approximated by a Normal distribution for large
26 samples. This result makes it possible to derive analytical formulas for calculating confidence intervals
27 and p -values for simple one-sample problems (e.g., involving means or proportions) and two-sample
28 problems (e.g., differences between group means or proportions). Additionally, sampling distributions
29 of coefficients in linear regression models follow t -distributions when observations are independent and
30 residuals are Normally distributed with constant variance. Introductory courses may also introduce χ^2
31 and F distributions, and mostly rely on analytical formulas for performing statistical inference.

32 Most problems facing ecologists are more complicated than the univariate problems explored in
33 introductory statistics courses; further, assumptions of linear regression are often untenable. Each
34 new challenge requires a new recipe, chosen from a more specialized and difficult-to-follow cookbook.
35 Modern statistical methods, such as generalized linear models and random-effects models, allow one
36 to relax assumptions that the residuals are Normally distributed, have constant variance, and that cases
37 are independent. These methods are powerful, but to understand them fully requires a background in
38 mathematical statistics, which most biologists lack. Consider, for example, generalized linear mixed

models (GLMMs). These methods are widely available in modern statistical software, and biologists are routinely encouraged to analyze their data using these methods. Yet, we frequently encounter biologists who have fit GLMMs but do not know how to determine if the models are appropriate for their data. They struggle to describe fitted models using equations or text. And they frequently do not know how to create effect plots to visualize how the mean response changes with changes in predictor values, which for models with a non-linear link function requires integrating over the distribution of the random effects (Fieberg et al., 2009). Simply put, few biologists have taken the requisite coursework in calculus and mathematical statistics to understand these methods, and many are unfamiliar with common statistical terminology (e.g., expected value) or probability distributions other than the Normal distribution.

Resampling methods, including permutation procedures, often offer an attractive alternative to Normal-based inferential methods; they are adaptable and require fewer assumptions. Historically, the main limitations in applying these methods were that they were computationally intensive and often required custom-written computer code (Cobb, 2007). These limitations are no longer a significant concern except for extremely large data sets, given the availability of personal computers and open-source statistical software with packages for implementing resampling-based methods (e.g., Davison and Hinkley, 1997; Pruim et al., 2017; Canty and Ripley, 2019; Simpson, 2019; Oksanen et al., 2019; R Core Team, 2019).

Our objectives of this paper are to:

1. Illustrate how resampling-based approaches can facilitate a deeper understanding of core concepts in frequentist statistics (e.g., standard errors, confidence intervals, p -values).
2. Demonstrate through simple examples and case studies how resampling-based methods can provide solutions to a range of statistical inference problems.

For resampling-based methods to work appropriately, researchers need to be able to generate data sets that preserve the structure of the original data set (e.g., any clustering or other forms of correlation). This requirement can easily be met with relatively simple computer code when addressing most univariate problems encountered in a first semester statistics course. The benefits of using resampling-based methods are often greatest, however, when analyzing complex, messy data. Determining appropriate solutions in these situations can be more challenging. Nonetheless, solutions are often available in open-source software and rely on the same set of core principles. Our case studies provide a few examples of the types of problems that can be addressed using resampling-based methods, but in truth, they barely scratch the surface of what is possible. For a more in-depth treatment of bootstrapping and permutation tests, we refer the reader to Davison and Hinkley (1997) and Manly (2006).

This paper is written primarily for the applied biologist with a rudimentary understanding of introductory statistics, but we also expect it will be of interest to instructors of introductory statistics courses. In particular, our first objective can be seen as an argument for replacing traditional approaches to teaching introductory statistics with an approach that relies heavily on computational methods, an argument that is increasingly supported by data (Tintle et al., 2011, 2012, 2015; Chance et al., 2016; VanderStoep et al., 2018). To that end, several educators have developed applets for explaining and understanding key concepts in frequentist statistics (e.g., <http://www.rossmanchance.com/ISIapplets.html> and <http://www.lock5stat.com/StatKey/index.html>). Here, we provide a tutorial review emphasizing the `mosaic` package in R, which was developed specifically to facilitate teaching resampling-based methods in introductory statistics courses (Pruim et al., 2017). We have chosen this approach because R has become a sort of lingua franca among ecologists. We begin by introducing key foundational concepts (sampling distributions, confidence intervals, null distributions and p -values) using simple examples that can be analyzed with a few lines of code. We then consider a series of increasingly complex case studies that demonstrate how these concepts apply more broadly. To facilitate learning, we have archived all data and R code needed to re-create our examples at <https://bootstrapping4biologists.netlify.com>. In addition, the data and code have been curated and included in the Data Repository of the University of Minnesota (Fieberg et al., 2020).

87 UNDERSTANDING FOUNDATIONAL CONCEPTS USING SIMULATION- 88 BASED METHODS

89 Sampling Distributions

90 The sampling distribution of a statistic tells us about the values we might expect for the statistic if we were
91 to repeatedly collect data sets of the same size from the same population and using the same sampling
92 protocols. Again, this is a difficult concept to grasp because we usually only get to collect one data set.
93 However, it is easy to use a computer to generate multiple samples from a population, and hence calculate
94 multiple statistics to approximate the sampling distribution. To illustrate, we explore here the sampling
95 distribution of a sample mean using a data set containing the count of the number of polar bears within
96 each of 164 roughly 3 x 3 km quadrats on Rowley Island in northern Foxe Basin, Nunavut (Fig. 1A,
97 Stapleton et al., 2014). First, we read in the data set which contains two variables, an id for each quadrat
98 (`Quadrat`) and a count of the number of bears in each quadrat (`Num.Bears`), making use of the `here`
99 package to make it easy to read in data from a subdirectory (Müller, 2017). We use `<-` to assign the data
100 to an object called `bdat`. We then use the `head` function to look at the first five rows of the data set (here
101 and throughout the paper, we provide R code followed by its associated output).

```
bdat <- read.csv(here("data", "bears.csv"))  
head(bdat, n=5)
```

```
102   Quadrat Num.Bears  
103 1         1         2  
104 2         2         1  
105 3         3         0  
106 4         4         0  
107 5         5         1
```

108 We use the `mean` function to calculate the mean number of bears counted per quadrat in the population
109 (i.e., the set of 164 quadrats). This function, like most in the `mosaic` package, follows a common syntax,
110 which for statistical summaries of a single variable is: `goal(~variable name, data=Mydata)`
111 (Pruim et al., 2017).

```
mean(~Num.Bears, data=bdat)
```

```
112 [1] 0.561
```

113 In this case, we can calculate the population mean, $\mu = 0.561$, exactly because we have polar bear
114 counts for all sample units. Nonetheless, we will use these data to illustrate the concept of the sampling
115 distribution. Specifically, we consider a hypothetical situation in which we have access to the counts of
116 bears in each of 164 quadrats (i.e., the population), to look at what *could* have happened if we had only
117 sampled a random subset of 75 quadrats. Our sample statistic, the sample mean (\bar{x}), provides an estimate
118 of the population mean (μ) bear count per quadrat. Below, we repeatedly take simple random samples of
119 75 quadrats using the `sample` function and calculate the sample mean for each of the resulting data sets.
120 We store the resulting sample means in an object named `samp.mean`. This feat can be accomplished
121 with essentially one line of code using the `do` function in the `mosaic` package, where its argument
122 (10,000 in this case) tells R how many times to repeatedly execute the code within the `{ }`:

```
samp.mean<-do(10000)*{  
  mean(~Num.Bears, data=sample(bdat, size = 75))  
}
```

123 The `do` function captures the 10,000 sample means and stores them in a variable named `result`
124 within the `samp.mean` object. We can look at the first five sample means using the `head` function in R:

```
head(samp.mean, n=5)
```

```
125   result
126 1  0.507
127 2  0.600
128 3  0.453
129 4  0.480
130 5  0.520
```

131 We can easily visualize the sampling distribution using a histogram or density plot (Fig. 1C), and we
 132 can calculate the standard error using the `sd` function in the `mosaic` library (the standard error is just the
 133 standard deviation of the sampling distribution; Box 1):

```
(se.mean<-sd(~result, data=samp.mean))
```

```
134 [1] 0.0772
```

135 When statistics have bell-shaped (i.e., approximately Normal) sampling distributions, we expect 95%
 136 of sample statistics to be within roughly 2 standard deviations of the mean of the sampling distribution
 137 (1.96 standard deviations to be more exact); we will refer to this result as the **2 SE rule** (Box 1). We can
 138 use the above simulation results (i.e., the 10,000 sample means contained in `samp.mean`), to verify this
 139 claim. First, we find this interval by taking the mean of the sampling distribution \pm 2 times the SE of the
 140 mean (we stored this SE in an object called `se.mean` and use the `c` operator to create a vector containing
 141 the values -2 and 2).

```
(rule.2SE<-mean(~result, data=samp.mean)+c(-2,2)*se.mean)
```

```
142 [1] 0.407 0.715
```

143 We then determine the proportion of sample means that fall in this interval (i.e., the proportion of
 144 sample means that are greater than the lower threshold and less than the upper threshold) using the `tally`
 145 function. We use `I(result >= 0.407 & result <= 0.715)` to create an indicator variable that
 146 is 1 when `result` is in the interval and zero otherwise.

```
tally(~I(result >= 0.407 & result <= 0.715),
      data=samp.mean, format="proportion")
```

```
147 I(result >= 0.407 & result <= 0.715)
148 TRUE FALSE
149 0.9548 0.0452
```

150 We see that roughly 95% of the sample means fall within 2 standard errors of the mean.

151 Central Limit Theorem

152 The sampling distribution of means (or sums) approaches that of a Normal distribution as the sample
 153 size increases. This result, given by the Central Limit Theorem, forms the basis of many formula-based
 154 methods of statistical inference and can be illustrated through simulation (Kwak and Kim, 2017). For
 155 example, the sampling distribution in our bear example is right-skewed for small sample sizes (e.g., $n = 10$;
 156 Fig. 1B); however, for samples of size 75, the sampling distribution is bell-shaped and well-approximated
 157 by a Normal distribution. Thus, rather than blindly trusting that the Central Limit Theorem applies,
 158 we can demonstrate it first hand by sampling repeatedly using different sample sizes. Additionally, the
 159 Central Limit Theorem guarantees for sufficiently large samples that the sampling distribution of \bar{x} will
 160 be approximately Normal with 95% of the sample means falling within 2 standard errors of μ in the
 161 polar bear example. This result holds even though the *population distribution* is right-skewed and highly
 162 discrete, with individual bear counts taking on only integer values. Further, many estimators (including all
 163 maximum likelihood estimators) are calculated via a sum of independent measurements, which ensures
 164 that their sampling distributions are asymptotically Normal. Understanding this powerful result helps
 165 uncover why many back-of-the-envelope calculations use *estimate* $\pm 2SE$ to form approximate 95%
 166 confidence intervals.

Bootstrap Confidence Intervals

In real applications, we do not have access to data from the full population. Instead, we have but a single sample. When making inference to the population, via traditional parametric methods or with resampling-based methods, we must assume our sample data are representative of the population, as is assumed when a large number of observations are selected by simple random sampling. Then, we can use the distribution of values in our sample to approximate the distribution of values in the population. For example, we can make many, many, copies of our sample data and use the resulting data set as an estimate of the whole population. With this estimated population in place, we could repeatedly sample from it, forming many data sets that are the same size as our original data set, and calculate statistics for each sample to approximate sampling distributions (Fig. 2). In practice, we do not actually need to make multiple copies of our sample data to estimate the population; instead, we form new data sets that are the same size as our original data set by sampling our original data with replacement, which effectively does the same thing. Sampling with replacement means that we select cases one at a time, and after each selection, we put the selected case back in the population so it can be chosen again. Thus, each observation in the original data set can occur zero, one, two, or more times in the generated data set, whereas it occurred exactly once in the original data set.

This process allows us to create a bootstrap sampling distribution to determine how much our estimates vary among repeated samples. We quantify this variability using the standard deviation of our estimates across repeated samples and refer to this standard deviation as our bootstrap standard error (BSE; Fig. 2). If the bootstrap sampling distribution is bell-shaped and centered on our sample statistic (indicating there is no sampling bias), then we can use this estimated standard error and the 2 SE rule for Normally distributed data to calculate a 95% confidence interval for our parameter of interest: $\text{estimate} \pm 2\text{BSE}$. We can also think about repeating this process many times (i.e., collect data, estimate a parameter, use the bootstrap to calculate a SE and confidence interval), in which case we expect 95% of our confidence intervals to include the true population parameter of interest. That is, the bootstrap method for calculating a 95% confidence interval is well calibrated. In practice, we only have one data set, leading to a single confidence interval that either does or does not contain the true value of the population parameter, but knowing the method is well calibrated allows us to state that we are *95% confident that our interval contains the population parameter*.

We illustrate this approach by considering data collected by the Minnesota Department of Natural Resources (MN DNR) to explore the potential impact of changing fishing regulations on the size distribution of northern pike (*Esox lucius*) in Medicine Lake, an approximately 460-acre lake in Beltrami County, MN. In 1989, the MN DNR instituted a slot limit of 22-30 inches in this lake (i.e., all caught fish within this size interval had to be released). We consider length data from 73 and 81 fish collected in trap nets in 1988 and 1993, respectively (before and after the fishing regulation was put in place). Importantly, these data come from only one lake, and many other factors may have changed between 1988 and 1993. Therefore, we must be cautious when interpreting any changes in the distribution of fish sizes (i.e., attributing the cause of length changes to the management regulation or generalizing results to other lakes). Nonetheless, we can ask, “How much did fish length, on average, change between 1988 and 1993 in Medicine Lake?” To address this question, we estimate the difference in the mean length of fish in the two samples. We also quantify our uncertainty in this estimated difference in means, recognizing that we would get a different estimate if we could go back and collect other samples of fish in those two years.

We begin by calculating the sample size, mean length in the two years, and the difference in sample means using the `tally`, `mean`, and `diffmean` functions in the `mosaic` library. Again, note that the functions in the `mosaic` library have a common syntax for data summaries involving two variables `goal(y~x, data=Mydata)`. For example, `mean(y~x, data=Mydata)` will calculate the mean of `y` for each level of a categorical variable `x`, and `diffmean(y~x, data=MyData)` will calculate the difference in sample means when `x` is a categorical variable taking on only two levels. Here, we store the difference in means in an object named `effect.size`:

```
pikedat<-read.csv(here("data", "Pikedata.csv"))
tally(~year, data=pikedat)
```

```
216 year
217 1988 1993
218    73    81
```

```
mean(length.inches~year, data=pikedat)
```

```
219 1988 1993
220 18.6 21.2
```

```
(effect_size <- diffmean(length.inches~year, data=pikedat))
```

```
221 diffmean
222      2.58
```

223 We estimate that the mean size of pike increased by roughly 2.6 inches between 1988 and 1993. To
 224 evaluate uncertainty in our estimated effect size, we explore the variability in the difference in means
 225 across bootstrapped samples. Here, we use the sample data from 1988 and 1993 as our estimate of
 226 the distribution of lengths in the population in 1988 and 1993, respectively. We repeatedly sample
 227 (10,000 times) from these estimated populations and calculate our sample statistic (the difference in
 228 means) for each of these bootstrapped data sets, making use of the `do` and `resample` functions in the
 229 `mosaic` library. Whereas the `sample` function, by default, samples without replacement, the `resample`
 230 function samples with replacement; the `group` argument is used to ensure that bootstrapping is conducted
 231 separately for each year.

```
bootdist<-do(10000)*{
  diffmean(length.inches~year, data=resample(pikedat, group=year))
}
```

232 Above, we stored the 10,000 different differences in means in a variable named `diffmean` contained
 233 in the `bootdist` object that we created. We can look at the first five differences in means using the
 234 `head` function.

```
head(bootdist, n=5)
```

```
235 diffmean
236 1      3.27
237 2      1.19
238 3      2.32
239 4      3.11
240 5      2.92
```

241 Because the bootstrap distribution is bell-shaped and centered on our sample statistic (Fig. 3), we can
 242 estimate a confidence interval for the difference in mean length in 1993 relative to 1988 using the 2 SE
 243 rule (i.e., based on an unstated Normal distributional assumption):

```
(SE<-sd(~diffmean, data=bootdist))
```

```
244 [1] 0.673
```

```
(CIa<-effect_size + c(-2, 2)*SE)
```

```
245 [1] 1.24 3.93
```

246 Alternatively, we could use quantiles of the bootstrap distribution, which does not require Normality
 247 (Fig. 2). The 0.25 quantile (or, equivalently, 25th percentile) refers to the value of x , such that 25% of
 248 the distribution falls below x . Here, we use the `qdata` function in the `mosaic` library to determine the
 249 0.025 and 0.975 quantiles. These quantiles capture the middle 95% of the bootstrap distribution for the
 250 difference in means:

```
(CIb<-qdata(~diffmean, data=bootdist, p=c(0.025, 0.975)))
```

```
251      quantile      p
252 2.5%        1.26 0.025
253 97.5%       3.90 0.975
```

254 For bell-shaped bootstrap distributions, Normal-based and percentile-based bootstrap confidence
 255 intervals work well and will usually give similar results. For small samples or skewed distributions, better
 256 methods exist (Davison and Hinkley, 1997; Hesterberg, 2015; Puth et al., 2015). Nonetheless, simple
 257 examples like this can facilitate understanding confidence intervals and other measures of uncertainty.

258 Null Distributions and p -values

259 Although null hypothesis testing has largely fallen out of favor in many disciplines, including ecology
 260 (Johnson, 1999; Hobbs and Hilborn, 2006), p -values, when fully understood, can play an important role
 261 in assessing the plausibility of statistical hypotheses (De Valpine, 2014; Murtaugh, 2014; Dushoff et al.,
 262 2019). Further, it is critically important for biologists to be able to interpret the p -values they see in
 263 papers and in output generated by statistical software. To understand p -values, we must consider the
 264 sampling distribution of our statistic across repeated samples in the case where the null hypothesis is true.
 265 Simulation-based methods are also useful for understanding this concept, provided we can generate data
 266 sets consistent with the null hypothesis. To illustrate, we will consider data from an experiment used to
 267 test various hypotheses about the mating preferences of female sage crickets (*Cyphoderris strepitans*)
 268 (Johnson et al., 1999). Males of this species will often allow a female to eat part of their hind wings when
 269 mating, which decreases the male's attractiveness to future potential mates. To explore how diet might
 270 influence mating behaviors of females, Johnson et al. (1999) randomly assigned 24 female crickets to
 271 either a low-nutrient (*starved*) or high-nutrient (*fed*) diet before putting them in a cage with a male
 272 cricket to mate. They then measured the waiting time before mating occurred. These data are contained in
 273 the *abd* package and were used by Whitlock and Schluter (2009) to introduce the Mann-Whitney U-test.

274 We can access the data using the *data* function after first loading the *abd* package:

```
library(abd)
data("SagebrushCrickets")
```

275 We then use the *str* function to explore the structure of the data set, finding that there are 2 variables:
 276 *treatment*, a factor variable with values *fed* or *starved*, and *time.to.mating*, a numerical
 277 variable containing the mating times for each cricket.

```
str(SagebrushCrickets)
```

```
278 'data.frame':   24 obs. of  2 variables:
279 $ treatment      : Factor w/ 2 levels "fed","starved": 2 2 2 2 2 2 2 2 2 2 ...
280 $ time.to.mating: num  1.9 2.1 3.8 9 9.6 13 14.7 17.9 21.7 29 ...
```

281 Our objective is to test whether the mean waiting times depend on the diet of the female: H_0 :
 282 $\mu_{fed} = \mu_{starved}$ versus H_A : $\mu_{fed} \neq \mu_{starved}$ where H_0 and H_A stand for the null and alternative hypotheses,
 283 and μ_{fed} and $\mu_{starved}$ represent the mean waiting times of crickets fed the high- and low-nutrient diets,
 284 respectively. It is natural to consider using the difference in sample means, $\bar{x}_{fed} - \bar{x}_{starved}$ to conduct
 285 our hypothesis test. Given the small sample sizes (we have 13 and 11 cases in the *fed* and *starved*
 286 treatment groups, respectively), we cannot safely assume the Central Limit Theorem ensures that the
 287 sampling distribution of $\bar{x}_{fed} - \bar{x}_{starved}$ is Normal. Further, the distribution of waiting times within each
 288 treatment group is far from Normal, and the measurements do not appear to be equally variable within
 289 each treatment group (Fig. 4). Thus, the assumptions of a parametric t-test are likely not met.

290 Although we could use a rank-based non-parametric test, such as the Mann-Whitney U-test (Whitlock
 291 and Schluter, 2009), we can gain additional insights into key statistical concepts (null distributions
 292 and p -values) if we construct our own randomization test, which also allows us to relax the Normality
 293 assumption.

294 We begin by estimating the mean waiting times for the two treatment groups (*fed* and *starved*) as
 295 well as the difference in sample means.


```

mean(time.to.mating~treatment, data=SagebrushCrickets)

296     fed starved
297     36.0     17.7

(effect_size2<-diffmean(time.to.mating~treatment, data=SagebrushCrickets))

298 diffmean
299     -18.3

```

300 We see that the mean time to waiting is 18.3 hours longer for the crickets in the `fed` treatment group.
 301 A skeptic might point out that each time you repeat this experiment you will get a different value for this
 302 statistic (i.e., the difference in sample means). Perhaps the longer mean waiting time was just a fluke.
 303 Could that be the case? To find out, we need to determine what values we might expect to see for the
 304 difference in means when the null hypothesis is true. If the null hypothesis is true, then the labels (`fed`
 305 and `starved`) should not matter. So, to simulate a sample statistic (difference in means) that we might
 306 get if the null hypothesis were true, we can: 1) randomly shuffle the treatment variable among cases to
 307 form a new data set, and 2) calculate the difference in sample means. Using the `do`, `diffmean`, and
 308 `shuffle` functions in the `mosaic` library, we can repeat this process many times:

```

randomization_dist<-do(10000)*{
  diffmean(time.to.mating~shuffle(treatment), data=SagebrushCrickets)
}

```

309 Here, the `shuffle` function permutes the `treatment` variable, ensuring that, on average, the
 310 difference in sample means between the two treatment groups is equal to zero. Again, the `do` function
 311 stores the 10,000 differences in means in a variable named `diffmean` in the `randomization_dist`
 312 object that we created. Plotting this randomization distribution (Fig. 5), we find the range of statistics
 313 (difference in sample means) that we would expect to see if the null hypothesis were true. We can then
 314 use this information to determine the probability of observing an absolute difference in sample means
 315 as large or larger than we saw in our experiment, $\bar{x}_{starved} - \bar{x}_{fed} = -18.3$, if the null hypothesis is true.
 316 This probability is our two-sided p -value. We can approximate it using our randomization distribution by
 317 counting the number of times the difference in sample means was ≤ -18.3 or ≥ 18.3 (equivalent to the
 318 number of times the absolute value of the difference was ≥ 18.3).

```

(p.value<- tally(~abs(diffmean)>=18.3, data=randomization_dist,
  format="proportion"))

319 abs(diffmean) >= 18.3
320 TRUE FALSE
321 0.126 0.874

```

322 Given a p -value of 0.126, we expect that roughly 13% of the time we would get a difference in mean
 323 waiting times as large or larger than 18.3 if the null hypothesis were true. Thus, our result is perhaps not
 324 all that surprising. Still, our sample sizes were small, leading to differences in sample means that were
 325 highly variable from sample to sample. It would behoove us to also report a confidence interval to go with
 326 our result, so the reader has an understanding of the range of effect sizes that are plausibly supported by
 327 the data. Using a bootstrap, we find that we are 95% confident that the difference in population means is
 328 between -38.9 and 3.6.

```

boot_dist<-do(10000)*{
  diffmean(time.to.mating~treatment, data=resample(SagebrushCrickets))
}
qdata(~diffmean, data=boot_dist, p=c(0.025, 0.975))

329     quantile     p
330 2.5%      -38.91 0.025
331 97.5%       3.61 0.975

```


332 Generalizations to Other Sample Statistics

333 The methods illustrated here easily generalize to other statistics and data sets, provided cases can be
334 assumed to be independent (e.g., there are no repeated measures on the same sample unit). For example,
335 we could calculate a bootstrap distribution, and hence confidence interval, for a correlation coefficient
336 using:

```
do(10000)*{cor(y~x, data=resample(Mydata))}
```

337 We could determine an appropriate null distribution for testing whether a set of regression coefficients
338 are all zero versus an alternative hypothesis that at least one is non-zero by shuffling the response variable:

```
do(10000)*{lm(shuffle(y)~x1 + x2 + x3, data=Mydata)}
```

339 Shuffling the response variable breaks any association between y and the predictor variables, thus
340 meeting the assumption of the null hypothesis. It is tempting to think that we could also determine an
341 appropriate null distribution for testing whether a regression coefficient for x_2 is zero, while adjusting for
342 two other variables (x_1 and x_3) using:

```
do(10000)*{lm(y~x1 + shuffle(x2) + x3, data=Mydata)}
```

343 Although shuffling the values of x_2 for each case will break any relationship between x_2 and y , this
344 process also breaks any association between x_2 and the other two predictor variables, x_1 and x_3 . Thus,
345 shuffling x_2 imposes additional unintended restrictions on the data beyond those required to ensure the null
346 hypothesis is true. Alternative solutions have been proposed, which involve permuting residuals, either
347 residuals of the response (Freedman and Lane, 1983; Anderson and Robinson, 2001), or the predictor at
348 stake (Collins, 1987; Dekker et al., 2007). The latter approach is perhaps simplest, has good statistical
349 properties, and is described below (see Supplementary Appendix A for a simple coded example):

- 350 1. Fit a linear regression model relating x_2 to the other predictor variables in the model:
351 `lm(x2~x1+x3, data=Mydata)`. The residuals from this model, which we label x_{2r} , capture
352 the part of x_2 that is not explained by the other variables in the model (here, x_1 and x_3). If we
353 replace x_2 in our original model with x_{2r} , fitting the following model: `lm(y~x1+x2r+x3,`
354 `data=Mydata)`, the coefficient, standard error, t-statistic, and p -value for x_{2r} will be identical
355 to the coefficient for x_2 in our original model.
- 356 2. Create an appropriate null distribution of t-statistics associated with the coefficient for x_{2r} by shuf-
357 fling x_{2r} : `do(10000)*{lm(y~x1 + shuffle(x2r) + x3, data=Mydata)}`. The
358 test is based on the t-statistic from the above randomization distribution rather than the coefficient
359 for x_{2r} itself because the t-statistic ensures the permutation distribution does not depend on
360 additional unknown parameters, whereas the regression coefficient does not have this property.

361 As the above example illustrates, solutions to more complex applications may require custom-written
362 code (and additional care), but the underlying concepts remain the same. For confidence intervals, we
363 need to generate bootstrapped data sets that mimic how we obtained the original sample data, calculate
364 our sample statistic for each of these data sets, and then use the variability in those statistics to form
365 our confidence interval. For testing null hypotheses, we need to generate data sets for which we have
366 assured compliance with the null hypothesis and then compare our observed statistic with those simulated
367 under that null hypothesis. In the words of Cobb (2007), “I like to think of [the core logic of statistical
368 inference] as three Rs: randomize, repeat, reject. Randomize data production; repeat by simulation to see
369 what’s typical and what’s not; reject any model that puts your data in its tail.”

370 For computationally intensive applications it may not be possible to generate as many as 10,000
371 random data sets as in our examples. In such cases, it is customary and important to include the
372 original sample statistic in the permutation or bootstrap distribution (Davison and Hinkley, 1997; Phipson
373 and Smyth, 2010). Doing so ensures statistical hypothesis tests are exact even for small numbers of
374 permutations (Phipson and Smyth, 2010), i.e., across repeated tests with a significance level of $\alpha = 0.05$,
375 we will reject at most 5% of the time when the null hypothesis is true.

376 Bootstrapping versus Permutations

377 Our simple examples above demonstrate how to form confidence intervals by bootstrapping and how to
378 test statistical hypotheses using permutations. In truth, one can use bootstrapping to perform hypothesis
379 tests and permutations to form confidence intervals (e.g., by “inverting” a hypothesis test so that the
380 confidence interval includes only parameter values for which we would fail to reject the null hypothesis).
381 Ideally, the chosen approach (resampling or permutation) should reflect the original study design. For
382 randomized experiments, it makes sense to use permutations that reflect alternative outcomes of the
383 randomization procedure. By contrast, bootstrap resampling may be more appropriate for observational
384 data where the “randomness” is reflected by the sampling procedure (Lock et al., 2013, p. 274).

385 Consider again, the fish length example used to demonstrate bootstrap confidence intervals. If we
386 wanted to conduct a statistical hypothesis test that mean fish length was the same in both years, we could
387 make the null hypothesis true (by adding the difference between the two sample means to the length
388 of fish in the first sample), then bootstrap. This approach is illustrated below. We begin by creating a
389 new variable, `length.null`, by adding the difference in sample means to all fish lengths collected in
390 1988. We can accomplish this task using the `mutate` function in the `dplyr` library (Wickham et al.,
391 2019). The `ifelse` function evaluates the first expression (`year=="1988"`) and executes the next
392 argument when true and the last argument when false. The difference in the mean of this new variable is
393 zero (reflecting our null hypothesis).

```
delta<-diffmean(length.inches~year, data=pikedat)
pikedat<-mutate(pikedat,
                length.null=ifelse(year=="1988", length.inches+delta,
                                   length.inches))
diffmean(length.null~year, data=pikedat)
```

```
394 diffmean
395      0
```

396 We can then use the bootstrap to quantify variability in difference in sample means *when the null*
397 *hypothesis is true*, from which we can calculate the *p*-value for our hypothesis test.

```
randfish<-do(10000)*{
  diffmean(length.null~year, data=resample(pikedat, group=year))
}
(p.value2<-tally(~abs(diffmean)>=delta, data=randfish, format="proportion"))
```

```
398 abs(diffmean) >= delta
399   TRUE  FALSE
400     0      1
```

401 None of the bootstrap differences in sample means were as large or larger than the difference we
402 observed using our original data, resulting in a *p*-value ≈ 0 .

403 This example raises another important point of consideration when choosing an appropriate method for
404 testing null hypotheses, namely differences between bootstrap- and permutation-based null distributions.
405 In the sagebrush cricket example, we permuted treatment labels to create our null distribution for the
406 difference in sample means. By permuting treatment labels, we ensure the *distribution* of the response
407 variable, *y*, is the same for both treatment groups. Thus, means, variances, and higher-order moments will
408 all, on average, be the same for the two groups. By contrast, the bootstrap-based test ensures only that
409 the sample means of the two groups are, on average, the same under the null hypothesis. Both methods
410 should reject the null hypothesis in cases where the population means differ, provided our sample size is
411 large enough. However, the permutation-based test may also reject the null hypothesis if the population
412 means are the same in the two groups, but the population variances differ. Thus, the bootstrap-based test
413 may be preferred if the goal is to detect only differences in population means, allowing for the possibility
414 that the population variances differ.

CASE STUDIES

We next consider a series of three case studies that demonstrate how the bootstrap can be used to calculate confidence intervals when faced with more challenging data analytic problems. Specifically we show how to use the bootstrap for more robust inference in the face of assumption violations (Case Study 1), to estimate complex quantities that combine results from multiple analyses (Case Study 2), and to quantify model-selection uncertainty (Case Study 3). The data and R code needed to reproduce these examples are available at <https://bootstrapping4biologists.netlify.com>. In addition, we use the site to host the output that results from running all of the coded examples.

Case Study I: Relaxing the Assumptions of Linear Regression

Zuur et al. (2009) lead off their popular book, *Mixed Effects Models and Extensions in Ecology with R*, with the statement (p. 19), “In Ecology, the data are seldom modeled adequately by linear regression models. If they are, you are lucky.” Here, we consider one of the data sets from their book, containing measurements of species richness at 45 sampling stations in the Netherlands (five stations located within each of nine beaches, Janssen and Mulder, 2004, 2005). When plotting species richness against an index of exposure of each beach to waves and surf, we see that mean species richness decreases linearly with exposure level (Fig. 6A), but the distribution of residuals is right skewed and, thus, far from Normal (Fig. 6B). In addition, observations were clustered in space so we might expect two observations from the same beach to be more alike than two observations from two different beaches.

The rest of Zuur et al.’s book is devoted to describing regression methods that allow one to relax these assumptions; e.g., using random effects to account for correlation and different statistical distributions to relax the Normality and constant variance assumptions. Without knowledge of these tools, however, is there a way to estimate uncertainty in the regression line as a summary of the relationship between species richness and exposure (over the observed range of the data)? Yes! We just need to adapt the bootstrap to mimic the way data were collected (i.e., cluster sampling, with multiple observations from each of several beaches). In this case, we need to sample *beaches* with replacement, keeping all observations from a beach when selected (we will refer to this approach as a cluster-level bootstrap). Doing so allows us to evaluate how much our estimates of regression parameters would change if we were to collect another sample of observations from a different set of nine beaches. Our estimate of uncertainty would not require assumptions of Normality or constant variance. Further, we would only need to assume that data from different beaches are independent (not that observations from the same beach are independent). Although it would be nice to have a larger sample of beaches to represent the population, alternative analysis methods (e.g., models appropriate for clustered data) would face a similar challenge of having to estimate among-beach variability from nine beaches.

Conceptually, this example is no more difficult to understand than the fisheries slot limit example. We are trying to determine how much our regression parameter estimates would vary if we were able to collect another sample of observations using the same sampling effort (nine beaches, each with five stations). Implementation becomes slightly more challenging than the previous examples because we can no longer rely on the `resample` function in the `mosaic` package (it does not accommodate cluster-level bootstrapping, but see Deen and de Rooij, 2019, for an alternative using the `ClusterBootstrap` package). The bootstrap distributions of the regression coefficients in this example are not symmetric, so it is beneficial to consider alternative bootstrap confidence interval procedures. It is important to note that not all confidence intervals perform equally well. Normal- and percentile-based intervals are simple but less accurate than intervals that attempt to account for bias and skewness of the sampling distribution (e.g., the “BCa” interval illustrated in the Supplementary Appendix); accuracy is measured in terms of an interval estimator’s *coverage rate*, defined as the frequency with which the interval estimator includes the true population parameter. We include, as supplementary material, code necessary to implement the cluster-level bootstrap, demonstrating a range of confidence interval procedures implemented using the `boot` package (Canty and Ripley, 2019, Supplementary Appendix B).

Case Study II: Combining Data Sources: Estimating Uncertainty When the Sampling Distribution is Unknown

As mentioned above, the Central Limit Theorem guarantees that the sampling distribution of many estimators, including all maximum likelihood estimators, will be Normally distributed as sample sizes become large (Casella and Berger 2002). In practice, however, it may be difficult to derive an analytical

expression for the standard error to use in conjunction with a Normal approximation. This is particularly true for estimation problems that require fitting multiple sub-models to different response variables collected from the same or different sets of sample units. As an illustrative example, we consider data from Zicus et al. (2006), used to evaluate the relative cost-effectiveness of single- and double-cylinder nesting structures for mallard (*Anas platyrhynchos*) ducks. To measure duckling productivity, Zicus et al. (2006) fit a linear regression model to duckling counts collected from 110 nesting structures over an 8-year period. To estimate costs, they fit an additional discrete-time survival model to quantify the probability that a structure would fall over and hence require a visit (with associated maintenance costs) for continued use. Both models included predictors capturing cylinder type (single or double) and size of the wetland holding the structure. Zicus et al. (2006) combined the output from the two models and estimated, for a range of wetland sizes and for both cylinder types, the expected cost of maintaining the structure divided by the expected number of ducklings produced as a measure of cost-effectiveness. Estimating the standard error of this cost-effectiveness measure was complicated for several reasons: 1) the quantity of interest was a ratio and therefore not a linear function of model parameters, 2) there was a complex dependence structure resulting from repeated measures on the same set of structures, which were also shared across the two models; and 3) the residuals from the duckling model were not Normally distributed.

Fortunately, using a bootstrap to estimate a confidence interval for cost-effectiveness was relatively straightforward. We just generated new bootstrapped data sets, mimicking the original sampling design, and repeated the analysis using these new data sets. As illustrated in Supplementary Appendix C, we: 1) resampled nest structures with replacement; 2) fit both models to the bootstrapped data sets; and 3) used the fitted models to estimate the expected number of ducklings, expected structure survival (and hence cost), and the ratio of expected cost to the expected number of ducklings for each cylinder type at different wetland sizes. This process allowed us to calculate confidence intervals for the cost effectiveness of both types of nest structures across a range of wetland sizes (Supplementary Appendix C). We have used this same general approach elsewhere to quantify uncertainty in estimated growth rates from population models when parameters are estimated from multiple data sources (e.g., Ellner and Fieberg, 2003; Fieberg et al., 2010; Lenarz et al., 2010).

One important consideration in this example is that our estimator of cost-effectiveness, a ratio of expected means, may be biased. We can also use the bootstrap to estimate $\text{Bias} = E[\hat{\theta}] - \theta$, where $E[\hat{\theta}]$ refers to the mean of the sampling distribution of $\hat{\theta}$, our estimator of the unknown parameter θ . To understand how this works, consider that we use the bootstrap to mimic sampling from our population: the bootstrap samples relate to the original sample in much the same way as the original sample relates to the population (Hall, 1988, Fig. 2). Thus, to estimate bias, we can compare the mean of our bootstrap statistics, θ_b^* , to the same statistic calculated using our original sample as though it were the true population, $\hat{\theta}$. In Supplementary Appendix C, we estimate the bias of our cost-effectiveness measure for both deployment types and across a range of wetland sizes, showing that it is small relative to the standard error.

Case Study III: Evaluating Model-Selection Uncertainty

The bootstrap can be used to quantify model uncertainty in addition to sampling variation (Buckland et al., 1997; Fieberg and Johnson, 2015; Harrell, 2016). Consider, for example, the common situation in which a researcher desires to develop a predictive model by selecting a limited number of explanatory variables from some larger set (e.g., using AIC statistics associated with various sub-models). Several models may provide nearly equally good fits to the data, in which case we might expect a different model to come out as “best” each time we collect a new data set (where “best” is determined via backwards stepwise selection using AIC). We can explore whether this is the case by fitting the same set of candidate models to each of several bootstrapped data sets.

In this case study, we explore model-selection uncertainty arising from backwards stepwise selection using AIC (Supplementary Appendix D). We begin by fitting a linear regression model to abundance data of longnose dace (*Rhinichthys cataractae*) collected from the Maryland Biological Stream Survey (downloaded from <http://www.biostathandbook.com/multipleregression.html>). Starting with six predictors (measurements of various stream attributes), we use the `stepAIC` function in the MASS package in R (Venables and Ripley, 2013) to sequentially eliminate predictors until we can no longer reduce AIC further by dropping any of the remaining predictors. Implementing this process resulted in a reduced model containing three variables: acreage drained by the stream, maximum depth, and NO3. We can then use functions in the `rms` (regression modeling strategies) package in R (Harrell

2018) to explore variability in the best model as determined by backwards stepwise selection. First, we refit the full model using `ols` (this function is equivalent to `lm`, but it stores additional information with the model fit). We then use the `validate` function to choose a best model by applying backwards selection to each of 100 bootstrapped data sets (i.e., `validate` fits a series of models, using backwards selection to choose the best model for each bootstrapped data set). Applying this process resulted in 31 different best models. The most frequently chosen model, selected in 24 of the 100 cases, included only acreage and NO3. The model originally chosen using stepAIC rose to the top in only nine of the 100 bootstrapped data sets. One bootstrap sample resulted in a model with all six predictors, and five bootstrap samples led to a model with a single predictor (maximum depth in four instances and NO3 in one instance). In addition to demonstrating the high degree of uncertainty associated with choosing a best model, these bootstrap results could be used to calculate model-averaged predictions (using the frequency with which a model was selected to determine model weights; Buckland et al., 1997).

DISCUSSION

Biologists need to understand, quantify, and communicate measures of effect sizes and their uncertainty. With frequentist statistics, we conceptualize and measure uncertainty by considering how statistics (e.g., means, proportions, regression coefficients) vary from sample to sample; yet, in reality, we typically have only one sample at our disposal. Resampling-based methods provide a natural way to understand foundational concepts related to uncertainty, including sampling distributions, standard errors, confidence intervals, and *p*-values. This understanding, in turn, makes it possible to develop custom analyses for a range of messy data scenarios using the same set of core concepts.

Consistent with the general call to provide estimates of effect sizes and their uncertainty (e.g., Johnson, 1999), our case studies emphasized applications of the bootstrap in this vein rather than resampling-based methods for conducting hypothesis tests. The bootstrap can be viewed as a Swiss-army knife of statistical tools, providing estimates of uncertainty for a wide range of problems; we just need to: 1) use our sample data to estimate the distribution of variable values in the population; 2) generate new data sets by repeatedly sampling from this approximated population; and 3) analyze each set of data in a consistent manner. To some readers, it may feel like bootstrapping is cheating because it seems to make up data. In the fisheries example, we had a total of 154 observations of fish length, but with a couple of lines of R code and a computer, we produced 10,000 such data sets. Something definitely seems fishy. In fact, the name bootstrapping derives from the phrase “to pull oneself up by one’s bootstraps,” a physical impossibility. Are we really getting something for nothing? No, we are not generating any new data. What bootstrapping does is explore the internal variability of a single data set. Formula-based estimates of uncertainty operate in a very similar manner, but are only available in a few select cases where we can derive the sampling distribution from our underlying model assumptions or via the Central Limit Theorem. For example, the formula for the standard error of a sample mean, $SE = s/\sqrt{n}$, also uses a measure of internal variability (the sample standard deviation, *s*) to quantify variability among repeated samples. A mathematical justification for the bootstrap procedure was provided by its developer, Bradley Efron (Efron, 1979), extending work on the jackknife, a linear approximation to the bootstrap (Quenouille, 1949).

Important to note is that bootstrapping can only be as good as the original data. This means that the sample data must be representative of the population from which they were drawn, preferably randomly, and that the sample size should be adequate to capture the variation in the population. When generating new data sets, it is also important to mimic the way the original data were collected and to identify independent sample units for resampling. In our first case study, we resampled beaches rather than individual observations because we expected observations from the same beach to be correlated. Similarly, one can create “blocks” of observations close in time or space and treat these as independent sample units when data are spatially or temporally correlated (Chernick and LaBudde, 2011). Although we focused on methods that rely on case resampling, which approximates the distribution of values in the population by making many copies of the sample data, parametric bootstrapping is also possible. In the latter case, we must assume the population values follow a particular parametric distribution (e.g., gamma, beta, etc.), use the sample data to estimate parameters of this distribution, and then use this parameterized distribution to simulate new data sets.

In addition to the examples presented in conjunction with our case studies, bootstrapping can be used to provide accurate estimates of model fit (e.g., R^2) and calibrate models (Harrell, 2016, Supplementary

Appendix D), to reduce prediction variance when using “greedy analysis methods” with many unknown parameters (e.g., random forests as a solution to the instability of regression trees; Breiman, 2001), and to quantify uncertainty across a diverse set of applications (Davison and Hinkley, 1997; Manly, 2006). We find the bootstrap particularly attractive for decision analysis because it can provide a distribution of possible outcomes across a set of potential management actions while accounting for model and parameter uncertainty (e.g., Ellner and Fieberg, 2003). Ecologists have a long history of using resampling-based methods to analyze multivariate response data using ordination methods (ter Braak, 1990; Vendrig et al., 2017; Van den Brink and ter Braak, 1999; ter Braak and Smilauer, 2018). The `vegan` and `permute` packages provide methods for restricted permutations (e.g., permuting observations within blocks for randomized complete block designs) that can be used with multivariate data (e.g., Anderson and ter Braak, 2003; Oksanen et al., 2019; Simpson, 2019); the `mvabund` package also provides methods for bootstrapping generalized linear models fit to multivariate response data (Wang et al., 2019, 2012). Another interesting example in biology is the use of the bootstrap to quantify the level of confidence associated with different clades on a phylogenetic tree (Felsenstein Joseph, 1985; Efron et al., 1996).

We acknowledge that fully model-based alternatives exist for addressing many of the challenges encountered in our case studies. For example, random effects could be used to account for potential correlation among observations from the same beach, and an appropriate count distribution (e.g., Poisson or negative binomial) could be used to model the counts rather than a Normal distribution. Hierarchical models with random effects “feel good” because they match the underlying structure of many ecological data sets. Further, mixed effect models offer significant advantages for modeling hierarchical data as these approaches allow one to consider both within- and between-cluster variability. For example, we could have allowed each beach to have its own intercept drawn from a Normal distribution. Rather than relax assumptions, however, this hierarchical approach would add even more assumptions (Murtaugh, 2007). Before trusting inference from the model, we would need to evaluate whether the Normal distribution was appropriate for describing variability among beaches. Adding random effects also implies that observations from the sample cluster (in this case beach) are *positively* correlated (Fieberg et al., 2009). Although this assumption is reasonable for these data, it may not be appropriate in other situations. For example, the lead author recently collaborated on a survival analysis of moose calves that included data from several twins (Severud et al., 2019). The fate of twins may be positively correlated due to genetics, mother’s health, and environmental effects, but it is also possible that losing a calf might increase maternal investment for, and the survival probability of, its twin. Using a cluster-level bootstrap allowed the fates of these individuals to be *positively or negatively correlated*, whereas using random effects would not (Smith and Murray, 1984).

A simple cluster-level bootstrap provided a reasonable solution to non-independence and non-Normal data in our first case study because we had balanced data (i.e., equal numbers of observations for each cluster). Similar to Generalized Estimating Equations that use a working-independence assumption, this approach may be sub-optimal when applied to unbalanced data and potentially problematic depending on the mechanisms causing variability in sample sizes among clusters (e.g., if the size of the cluster is in some way related to the response of interest; Williamson, 2014). Also of note, we were interested in the effect of a predictor variable, *exposure*, that did not vary within a cluster (i.e., *exposure* was constant across all measurements at a beach). Models that use random effects to model correlation offer substantial advantages when interest lies in predictors that vary both within and among clusters (Muff et al., 2016). A parametric bootstrap (i.e., simulating from a fitted model) is always possible with random-effects models (e.g., using the `bootMer` function in the `lme4` package; Bates et al., 2015). This approach does not allow one to relax model assumptions, but could prove useful for estimating confidence intervals for functions of model parameters in situations where one is uncomfortable assuming the sampling distribution is Normal or where it is impossible to derive an appropriate standard error. Alternatively, Warton et al. (2017) proposed bootstrapping probability integral transform (PIT-) residuals as a general method appropriate for non-Normal, and possibly clustered or multivariate, data. Model-based and resampling-based solutions to regression problems, particularly those involving dependent data (e.g. repeated measures, time series, spatial data), tend to be fairly complex. Therefore, it is important to consult with someone who has expertise in these areas and to recognize that statisticians may not agree on a best solution.

We occasionally hear from colleagues that it is hopeless to try to teach frequentist statistics to biologists – e.g., it is impossible for non-statisticians to truly understand and interpret confidence intervals – and therefore, we should just teach Bayesian methods. Although replication plays a much lesser role in

631 Bayesian statistics, we note that Bayesian p -values are frequently used to perform goodness-of-fit tests
632 (Kery, 2010), and understanding the properties of estimators across repeated samples remains critically
633 important for evaluating Bayesian estimators (Rubin, 1984; Little, 2006). In particular, demonstrating
634 that Bayesian procedures are well calibrated (e.g., 95% credible intervals contain parameters used to
635 simulate data roughly 95% of the time) can help overcome the criticism that Bayesian methods rely on
636 a subjective form of probability, with results dependent on one's chosen priors. This point cannot be
637 overstated, particularly in today's highly polarized world in which individuals have strongly held prior
638 beliefs that differ among groups.

639 In our own work, we take a pragmatic approach to data analysis and use a variety of tools, including
640 frequentist and Bayesian model-based inference as well as resampling-based methods. Yet, we find that
641 resampling-based methods often provide an easier entry point to appropriate statistical analyses when
642 consulting with less experienced and mathematically savvy users. We have also found resampling-based
643 methods helpful for teaching foundational concepts in frequentist statistics and that undergraduate biology
644 majors are able to adapt resampling-based methods to new problems. Thus, we see great opportunity in
645 initiatives to use resampling-based methods to improve statistical thinking in the biological sciences (e.g.,
646 www.causeweb.org/STUB). In summary, we believe resampling-based methods should be used more
647 frequently, both in practice and in the classroom.

648 **ACKNOWLEDGMENTS**

649 We thank S. Stapleton and the Government of Nunavut for sharing the polar bear data. We thank A. Gray,
650 C. ter Braak, and N. Tintle for many helpful suggestions that greatly improved the manuscript.

Box 1: Key Concepts in Frequentist Statistics

- A *sampling distribution* is the distribution of sample statistics computed using different samples of the same size from the same population.
- A *bootstrap distribution* is a distribution of statistics computed using different samples of the same size from the same *estimated* population formed by merging many copies of the original sample data. Alternatively, the sample data may be used to estimate parameters of a statistical distribution, and then this distribution can be used to generate new samples. This alternative is termed the *parametric bootstrap*.
- A *null* or *randomization distribution* is a collection of statistics from samples simulated assuming the null hypothesis is true.
- The *standard error* of a statistic is the standard deviation of the sampling distribution. When forming confidence intervals, we estimate the standard error using the standard deviation of a bootstrap distribution. When calculating p-values, we estimate the standard error using the standard deviation of the randomization distribution.
- *2 SE rule*: when statistics have bell-shaped (i.e., approximately Normal) sampling distributions, we expect roughly 95% of sample statistics to be within 2 standard deviations of the mean of the sampling distribution
- A *confidence interval* for a parameter is an interval computed from data using a method that will capture the parameter for a specified proportion of all samples (e.g., 95% of the time for a 95% confidence interval)
- The *p-value* is the chance of obtaining a sample statistic as extreme (or more extreme) than the observed sample statistic, if the null hypothesis is true.

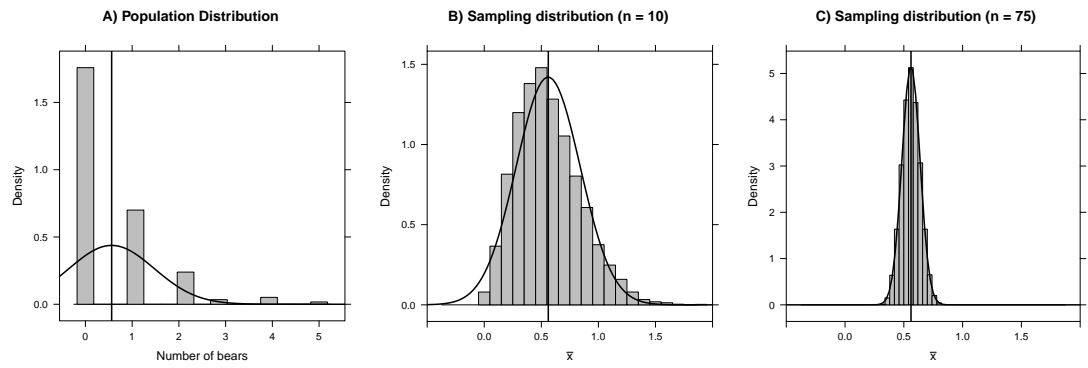


Figure 1. Histograms depicting A) the number of bears in each of 164 roughly 3 x 3 km quadrats on Rowley Island in northern Foxe Basin, Nunavut (Stapleton et al. 2014), B) the sampling distribution of the mean number of bears in 10,000 simple random samples of size 10 plots, and C) the sampling distribution of the mean number of bears in 10,000 simple random samples of size 75 plots. The vertical line gives the population mean and the smooth curves depict Normal approximations to the distributions in each panel.

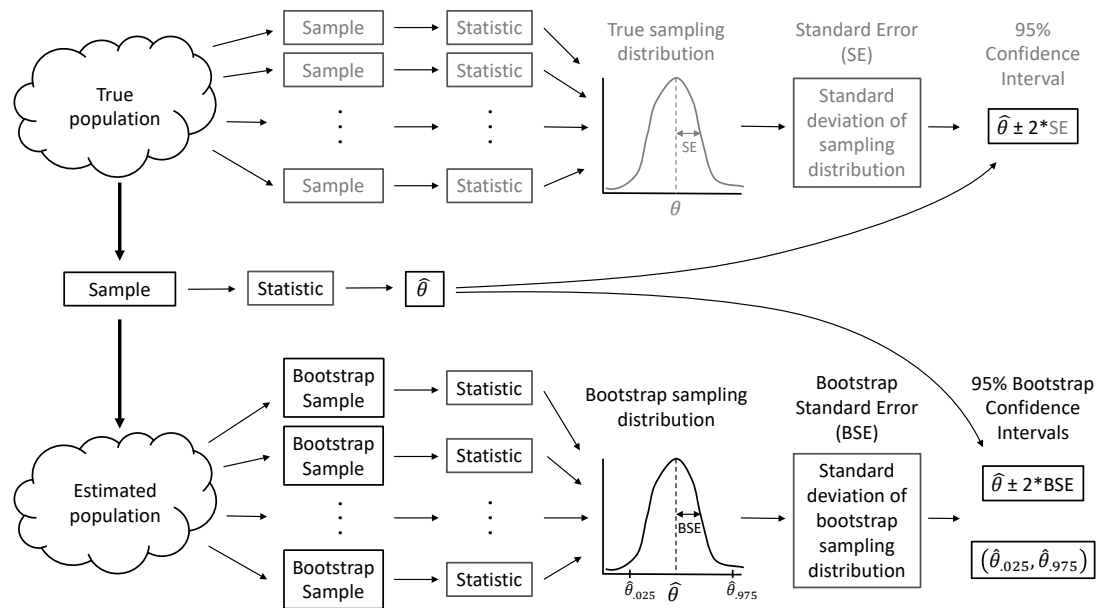


Figure 2. A sampling distribution is the distribution of sample statistics computed using different samples of the same size from the same population. We can estimate characteristics of the sampling distribution (e.g., its standard deviation) using bootstrapping, in which we repeatedly sample from an estimated population. Each bootstrap sample should be the same size as the original sample, and bootstrap samples should be formed in such a way that they preserve the structure of the original data set (e.g., any clustering or other forms of correlation).

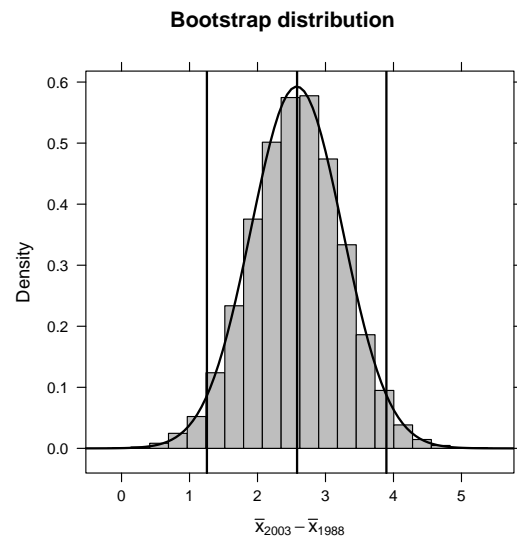


Figure 3. Bootstrap distribution of the difference in mean length of fish (year 2003 relative to year 1988). Vertical lines indicate the 95% confidence interval using the percentile-based method, and the smooth curve illustrates the best-fit Normal distribution.

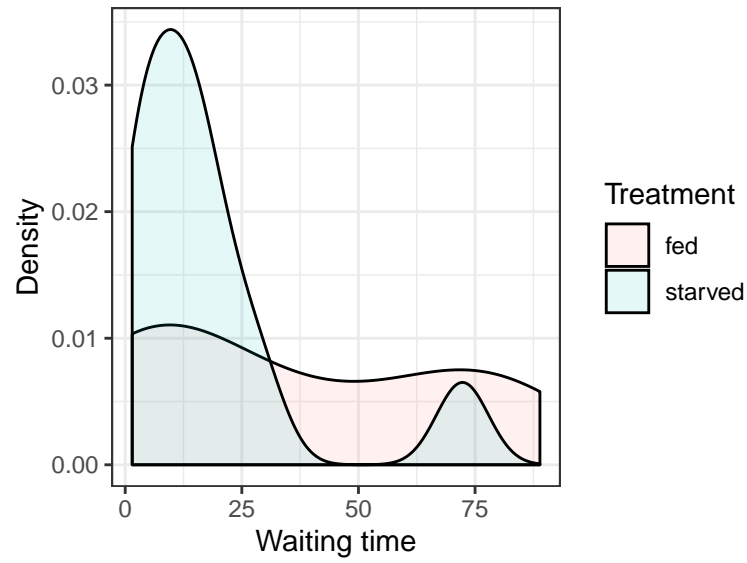


Figure 4. Distribution of waiting times to mating for female sage brush (*Cyphoderris strepitans*) crickets randomly assigned to either a low-nutrient (*starved*) or high-nutrient (*fed*) diet.

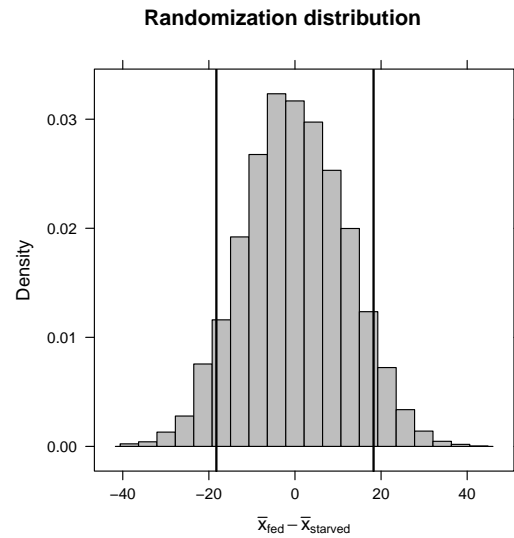


Figure 5. Randomization distribution for the difference in mean waiting time until mating formed by shuffling treatment assignments (fed versus starved) among cases. This distribution informs us of the types of statistics (difference in means) that we would expect to see if the null hypothesis were true. Our observed sample statistic, $(\bar{x}_{starved} - \bar{x}_{fed} = -18.3$, is indicated by the leftmost black vertical line. When the null hypothesis is true, we would expect to get a difference in sample means ≤ -18.3 or ≥ 18.3 13 % of the time (i.e., the p -value is 0.13).

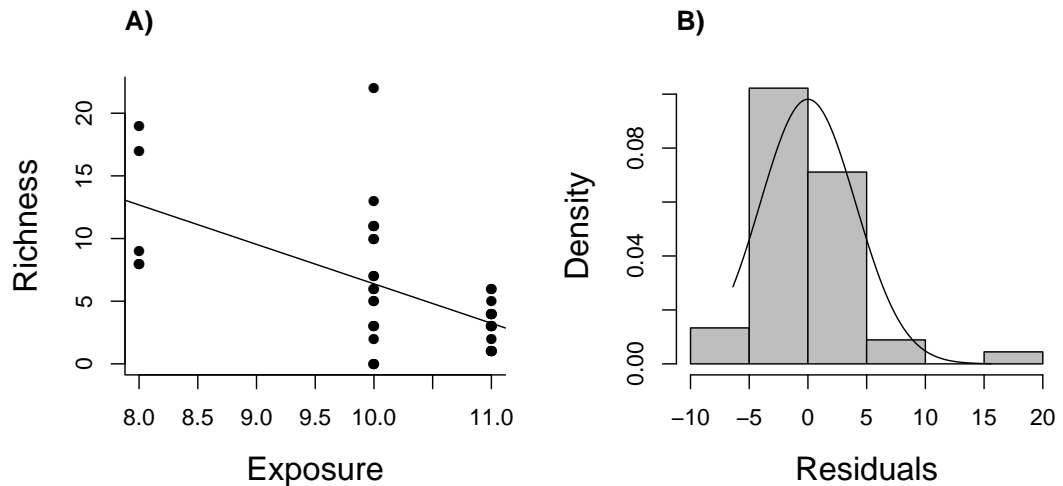


Figure 6. Regression line and diagnostics for linear model relating species richness to the level of exposure of the beach measured at 45 stations in the Netherlands.

4

REFERENCES

- Anderson, M. and ter Braak, C. (2003). Permutation tests for multi-factorial analysis of variance. *Journal of Statistical Computation and Simulation*, 73(2):85–113.
- Anderson, M. J. and Robinson, J. (2001). Permutation tests for linear models. *Australian & New Zealand Journal of Statistics*, 43(1):75–88.
- Bates, D., Mächler, M., Bolker, B., and Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1):1–48.
- Breiman, L. (2001). Random forests. *Machine Learning*, 45(1):5–32.
- Buckland, S. T., Burnham, K. P., and Augustin, N. H. (1997). Model selection: An integral part of inference. *Biometrics*, 53:603–618.
- Canty, A. and Ripley, B. D. (2019). *boot: Bootstrap R (S-Plus) Functions*. R package version 1.3-22.
- Chance, B., Wong, J., and Tintle, N. (2016). Student performance in curricula centered on simulation-based inference: A preliminary report. *Journal of Statistics Education*, 24(3):114–126.
- Chernick, M. R. and LaBudde, R. A. (2011). *An introduction to Bootstrap Methods with Applications to R*. Wiley.
- Cobb, G. W. (2007). The introductory statistics course: A ptolemaic curriculum? *Technology Innovations in Statistics Education*, 1(1).
- Collins, M. F. (1987). A permutation test for planar regression. *Australian Journal of Statistics*, 29(3):303–308.
- Davison, A. C. and Hinkley, D. V. (1997). *Bootstrap Methods and their Applications*. Cambridge University Press.
- De Valpine, P. (2014). The common sense of p values. *Ecology*, 95(3):617–621.
- Deen, M. and de Rooij, M. (2019). Clusterbootstrap: An r package for the analysis of hierarchical data using generalized linear models with the cluster bootstrap. *Behavior Research Methods*, pages 1–19.

- Dekker, D., Krackhardt, D., and Snijders, T. A. (2007). Sensitivity of MRQAP tests to collinearity and autocorrelation conditions. *Psychometrika*, 72(4):563–581.
- Dushoff, J., Kain, M. P., and Bolker, B. M. (2019). I can see clearly now: Reinterpreting statistical significance. *Methods in Ecology and Evolution*, 10(6):756–759.
- Efron, B. (1979). Bootstrap methods: Another look at the jackknife. *The Annals of Statistics*, 7(1):1–26.
- Efron, B., Halloran, E., and Holmes, S. (1996). Bootstrap confidence levels for phylogenetic trees. *Proceedings of the National Academy of Sciences*, 93(23):13429–13429.
- Ellner, S. P. and Fieberg, J. (2003). Using pva for management despite uncertainty: effects of habitat, hatcheries, and harvest on salmon. *Ecology*, 84(6):1359–1369.
- Felsenstein Joseph (1985). Confidence limits on phylogenies: An approach using the bootstrap. *Evolution*, 39(4):783–791.
- Fieberg, J. and Johnson, D. H. (2015). MMI: Multimodel inference or models with management implications? *Journal of Wildlife Management*, 79(5):708–718.
- Fieberg, J., Rieger, R. H., Zicus, M. C., and Schildcrout, J. S. (2009). Regression modelling of correlated data in ecology: Subject-specific and population averaged response patterns. *Journal of Applied Ecology*, 46(5):1018–1025.
- Fieberg, J., Vitense, K., and Johnson, D. H. (2020). Data and r code supporting: Resampling-based methods for biologists. *University of Minnesota Digital Conservancy*, <https://doi.org/10.13020/wn56-9y75>.
- Fieberg, J. R., Shertzer, K. W., Conn, P. B., Noyce, K. V., and Garshelis, D. L. (2010). Integrated population modeling of black bears in Minnesota: Implications for monitoring and management. *PLoS ONE*, 5(8).
- Freedman, D. and Lane, D. (1983). A nonstochastic interpretation of reported significance levels. *Journal of Business & Economic Statistics*, 1(4):292–298.
- Hall, P. (1988). Theoretical comparison of bootstrap confidence intervals. *The Annals of Statistics*, 16(3):927–953.
- Harrell, F. (2016). *Regression Modeling Strategies: With Applications to Linear Models, Logistic and Ordinal Regression, and Survival Analysis*. Springer.
- Hesterberg, T. C. (2015). What teachers should know about the bootstrap: Resampling in the undergraduate statistics curriculum. *American Statistician*, 69(4):371–386.
- Hobbs, N. T. and Hilborn, R. (2006). Alternatives to statistical hypothesis testing in ecology: A guide to self teaching. *Ecological Applications*, 16(1):5–19.
- Janssen, G. and Mulder, S. (2004). De ecologie van de zandige kust van nederland: Inventarisatie van het marcobenthos van zand en brandingszone. *Rapportnr.: 2004.033*.
- Janssen, G. and Mulder, S. (2005). Zonation of macrofauna across sandy beaches and surf zones along the dutch coast. *Oceanologia*, 47(2).
- Johnson, D. H. (1999). The insignificance of statistical significance testing. *The Journal of Wildlife Management*, 63(3):763–772.
- Johnson, D. H. (2002). The importance of replication in wildlife research. *Journal of Wildlife Management*, 66(4):919–932.
- Johnson, J. C., Ivy, T. M., and Sakaluk, S. K. (1999). Female remating propensity contingent on sexual cannibalism in sagebrush crickets, *cyphoderris strepitans*: A mechanism of cryptic female choice. *Behavioral Ecology*, 10(3):227–233.

- 726 Kery, M. (2010). *Introduction to WinBUGS for Ecologists: A Bayesian Approach to Regression, ANOVA,*
727 *Mixed Models, and Related Analyses.* Academic Press.
- 728 Kwak, S. G. and Kim, J. H. (2017). Central limit theorem: The cornerstone of modern statistics. *Korean*
729 *Journal of Anesthesiology*, 70(2):144.
- 730 Lenarz, M. S., Fieberg, J., Schrage, M. W., and Edwards, A. J. (2010). Living on the edge: Viability of
731 moose in northeastern minnesota. *Journal of Wildlife Management*, 74(5):1013–1023.
- 732 Little, R. J. (2006). Calibrated Bayes: A Bayes/frequentist roadmap. *American Statistician*, 60(3):213–
733 223.
- 734 Lock, R. H., Lock, P. F., Morgan, K. L., Lock, E. F., and Lock, D. F. (2013). *Statistics: Unlocking the*
735 *Power of Data.* John Wiley & Sons.
- 736 Manly, B. F. (2006). *Randomization, Bootstrap and Monte Carlo Methods in Biology.* CRC Press.
- 737 Muff, S., Held, L., and Keller, L. F. (2016). Marginal or conditional regression models for correlated
738 non-normal data? *Methods in Ecology and Evolution*, 7(12):1514–1524.
- 739 Murtaugh, P. A. (2007). Simplicity and complexity in ecological data analysis. *Ecology*, 88(1):56–62.
- 740 Murtaugh, P. A. (2014). In defense of P values. *Ecology*, 95(3):611–617.
- 741 Müller, K. (2017). *here: A Simpler Way to Find Your Files.* R package version 0.1.
- 742 Oksanen, J., Blanchet, F. G., Friendly, M., Kindt, R., Legendre, P., McGlinn, D., Minchin, P. R., O’Hara,
743 R. B., Simpson, G. L., Solymos, P., Stevens, M. H. H., Szoecs, E., and Wagner, H. (2019). *vegan:*
744 *Community Ecology Package.* R package version 2.5-6.
- 745 Phipson, B. and Smyth, G. K. (2010). Permutation p-values should never be zero: calculating exact
746 p-values when permutations are randomly drawn. *Statistical Applications in Genetics and Molecular*
747 *Biology*, 9(1).
- 748 Popper, K. (2005). *The Logic of Scientific Discovery.* Routledge.
- 749 Pruim, R., Kaplan, D. T., and Horton, N. J. (2017). The mosaic package: Helping students to “think with
750 data” using R. *The R Journal*, 9(1):77–102.
- 751 Puth, M. T., Neuhäuser, M., and Ruxton, G. D. (2015). On the variety of methods for calculating
752 confidence intervals by bootstrapping. *Journal of Animal Ecology*, 84(4):892–897.
- 753 Quenouille, M. H. (1949). Problems in plane sampling. *The Annals of Mathematical Statistics*, 20(3):355–
754 375.
- 755 R Core Team (2019). *R: A Language and Environment for Statistical Computing.* R Foundation for
756 Statistical Computing, Vienna, Austria.
- 757 Rubin, D. B. (1984). Bayesianly justifiable and relevant frequency calculations for the applied statistician.
758 *The Annals of Statistics*, 12(4):1151–1172.
- 759 Severud, W. J., Obermoller, T. R., Delgiudice, G. D., and Fieberg, J. R. (2019). Survival and cause-specific
760 mortality of moose calves in northeastern minnesota. *Journal of Wildlife Management*, 83(5):1131–
761 1142.
- 762 Simpson, G. L. (2019). *permute: Functions for Generating Restricted Permutations of Data.* R package
763 version 0.9-5.
- 764 Smith, D. W. and Murray, L. W. (1984). An alternative to Eisenhart’s Model II and Mixed Model in the
765 case of negative variance estimates. *Journal of the American Statistical Association*, 79(385):145–151.
- 766 Stapleton, S., LaRue, M., Lecomte, N., Atkinson, S., Garshelis, D., Porter, C., and Atwood, T. (2014).
767 Polar bears from space: Assessing satellite imagery as a tool to track arctic wildlife. *PLoS ONE*, 9(7).

- 768 ter Braak, C. (1990). Update notes: canoco, version 3.10. agricultural mathematics group, wageningen.
769 *The Netherlands*, 35.
- 770 ter Braak, C. J. and Smilauer, P. (2018). Canoco reference manual and user's guide: software for ordination
771 (version 5.10). Technical report, Biometris, Wageningen University & Research.
- 772 Tintle, N., Chance, B., Cobb, G., Roy, S., Swanson, T., and VanderStoep, J. (2015). Combating anti-
773 statistical thinking using simulation-based methods throughout the undergraduate curriculum. *American*
774 *Statistician*, 69(4):362–370.
- 775 Tintle, N., VanderStoep, J., Holmes, V. L., Quisenberry, B., and Swanson, T. (2011). Development
776 and assessment of a preliminary randomization-based introductory statistics curriculum. *Journal of*
777 *Statistics Education*, 19(1).
- 778 Tintle, N. L., Topliff, K., VanderStoep, J., Holmes, V.-L., and Swanson, T. (2012). Retention of statistical
779 concepts in a preliminary randomization-based introductory statistics curriculum. *Statistics Education*
780 *Research Journal*, 11(1):21.
- 781 Van den Brink, P. J. and ter Braak, C. (1999). Principal response curves: Analysis of time-dependent
782 multivariate responses of biological community to stress. *Environmental Toxicology and Chemistry: An International Journal*, 18(2):138–148.
- 784 VanderStoep, J. L., Couch, O., and Lenderink, C. (2018). Assessing the association between quantitative
785 maturity and student performance in an introductory statistics class: Simulation-based vs non simulation-
786 based. In *Proceedings of the International Conference on Teaching Statistics*, volume 10.
- 787 Venables, W. and Ripley, B. (2013). *Modern Applied Statistics with S-PLUS*. Springer.
- 788 Vendrig, N. J., Hemerik, L., and ter Braak, C. J. (2017). Response variable selection in principal response
789 curves using permutation testing. *Aquatic Ecology*, 51(1):131–143.
- 790 Wang, Y., Naumann, U., Eddelbuettel, D., Wilshire, J., and Warton, D. (2019). *mvabund: Statistical*
791 *Methods for Analysing Multivariate Abundance Data*. R package version 4.0.1.
- 792 Wang, Y., Naumann, U., Wright, S. T., and Warton, D. I. (2012). mvabund—an r package for model-based
793 analysis of multivariate abundance data. *Methods in Ecology and Evolution*, 3(3):471–474.
- 794 Warton, D. I., Thibaut, L., and Wang, Y. A. (2017). The pit-trap—a “model-free” bootstrap procedure for
795 inference about regression models with discrete, multivariate responses. *PloS one*, 12(7):e0181790.
- 796 Whitlock, M. C. and Schluter, D. (2009). *The Analysis of Biological Data*. Roberts and Company.
- 797 Wickham, H., Francois, R., Henry, L., and Müller, K. (2019). *dplyr: A Grammar of Data Manipulation*.
798 R package version 0.8.3.
- 799 Williamson, J. M. (2014). Informative cluster size. *Wiley StatsRef: Statistics Reference Online*, pages
800 1–2.
- 801 Zicus, M. C., Rave, D. P., and Fieberg, J. R. (2006). Cost-effectiveness of single- versus double-cylinder
802 over-water nest structures. *Wildlife Society Bulletin*, 34(3):647–655.
- 803 Zuur, A., Ieno, E., Walker, N., Saveliev, A., and Smith, G. (2009). *Mixed Effects Models and Extensions*
804 *in Ecology with R*. Springer.