

MONTANA STATE UNIVERSITY
DEPARTMENT OF MATHEMATICAL SCIENCES
WRITING PROJECT

TITLE

Author:
MEAGHAN WINDER

Supervisor:
DR. ANDREW HOEGH

Spring 2020



A writing project submitted in partial fulfillment
of the requirements for the degree

Master's of Science in Statistics

APPROVAL

of a writing project submitted by

Meaghan Winder

This writing project has been read by the writing project advisor and has been found to be satisfactory regarding content, English usage, format, citations, bibliographic style, and consistency, and is ready for submission to the Statistics Faculty.

Date

Andrew Hoegh
Writing Project Advisor

Date

Mark C. Greenwood
Writing Projects Coordinator

Contents

1	Introduction	2
2	Data	7
3	Methods	13
3.1	Bayesian Modeling Background	13
3.1.1	Gibbs Sampler	16
3.1.2	Metropolis-Hastings Algorithm	18
3.1.3	Bayesian Hierarchical Modeling	19
3.2	Occupancy Models	20
3.2.1	Multi-scale Occupancy Models	22
3.3	Implementation	24
3.3.1	msocc package	25
3.3.2	*implementation for traditional occupancy model* . . .	25
3.3.3	Alternative Options	25
4	Analysis	26
5	Results	34
5.1	Ecological Implications	34
5.2	Sampling Recommendations	35
6	Future Work	35
7	Acknowledgements	35
8	References	36
9	Appendix - R Code	40

Abstract

abstract text here

1 Introduction

In early 2020, the City of Austin, Texas approved the spending of four million dollars over the next five years in an attempt to remove zebra mussels from the city’s source of drinking water with a liquid copper sulfate pentahydrate released into the water intake pipes (Bontke, 2020). This is one of many pursuits to remove dreissenid mussels¹ from water bodies across the United States, and four million dollars is merely a small fraction of what is spent annually on control and mitigation efforts.

Zebra mussels are native to the Caspian and Black Seas, but have become widespread in both Europe and the United States; they were discovered in the Great Lakes in the late 1980s and have since spread rapidly across the continental U.S. The United States National Park Service stated that “[o]nce a population of zebra mussels has become established in a water body, there is very little to be done to remove them. Prevention, therefore, is the best way to keep a water body clean of zebra mussels” (U.S. National Park Service, 2017); hence, early detection of invasive species, such as dreissenid mussels, has become a priority in order for organizations to plan, budget, and install necessary technologies before colonization has occurred (Holser, 2017).

¹Zebra mussels (*Dreissena polymorpha*) and quagga mussels (*Dreissena rostriformis bugensis*) collectively.

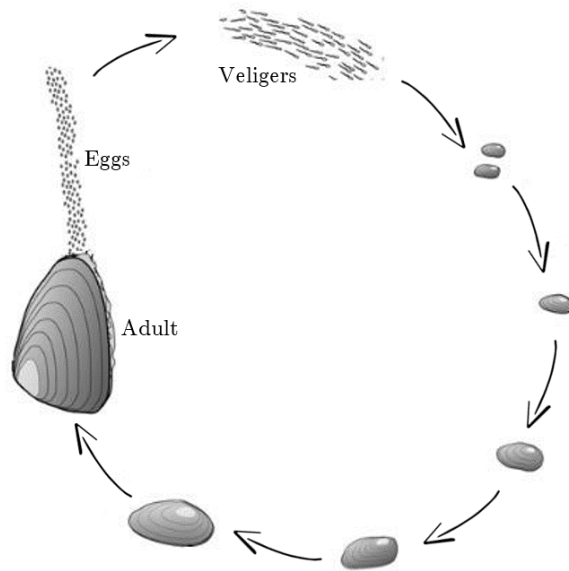


Figure 1: A diagram depicting the life cycle of dreissenid mussels.

Zebra mussels live between two and five years; they start as microscopic veligers but mature to thumbnail sized adults; they begin reproduction at two years of age, after which, females can release up to one million eggs per year (U.S. National Park Service, 2017). A diagram depicting the life cycle of dreissenid mussels is shown in Figure 1. Dreissenid mussel veligers free-swim in the water; often, they travel to uninfested waters on boats or through other aquatic recreational activities, however, sometimes they are moved by nature and travel downstream to uninfested waters. Adult dreissenids attach and colonize hard surfaces in the water, this process of accumulation of adult zebra mussels on rocks, native mussels, docks, boats, or other hard surfaces is referred to as “biofouling,” and objects that are in the water for long periods

of time become difficult and costly to clean.

Once a water body is colonized by adult dreissenid mussels, water supply and delivery facilities, water recreation sites, and other water dependent economies in that body of water become much more expensive to maintain and operate (Bureau of Reclamation, 2019). For example, after dreissenid mussels were detected in Montana’s Tiber Reservoir in 2016, with suspected detections in Canyon Ferry Reservoir as well, the state of Montana highlighted the economic impacts that a full dreissenid mussel colonization could have on the state; in January 2019, the Montana Department of Natural Resources and Conservation published a report titled *Enumeration of Potential Economic Costs of Dreissenid Mussels Infestation in Montana*. The report states that if dreissenid mussels were to infest every lake and river in Montana, it could cost the state more than 230 million dollars annually in mitigation costs, as well as lost revenue (Nelson, 2019).

Dreissenid infestations result not only in economic impacts, but in environmental ones as well. Dreissenid mussels are filter feeders and siphon plankton from the water, which can lead to changes the water body ecosystem by increasing water clarity; a single adult dreissenid can filter about a liter of water per day, which reduces the availability of algae for native mussels and bottom feeding fish (Bureau of Reclamation, 2019). Additionally, “biofouling” can prevent native mussels from moving, feeding, reproducing, or regulating the water system. Several actions, such as the 2017 initiative, *Safeguarding the West from Invasive Species*, by the Department of the In-

terior, have been taken to protect water bodies in the western United States from the economic and ecological threats posed by the invasive dreissenid mussels. Early detection of dreissenid mussels can reduce the economic and ecological repercussions of dreissenid infestations, however there are issues with the available early detection methods.

The established standard for early detection of dreissenids in the western United States is plankton tow sampling for mussel veligers. Using a fine mesh net, water and debris are collected at multiple sampling sites within each water body; the debris from each net collected at the same sampling site on the same day is aggregated and examined, using cross-polarized light microscopy, for the free-swimming veligers. Following the microscopic examination, positive species identification is confirmed using polymerase chain reaction (PCR) chemistry. This early detection method requires a breeding population, so is limited to the weeks immediately following a spawning event (Nichols, 1996); spawning begins at water temperatures above 10°C for quagga mussels and above 12°C for zebra mussels (McMahon, 1996; Mills et al., 1996). This suggests that veliger availability in northern latitude water bodies is typically limited to warmer months (Sepulveda et al., 2019).

An alternative method, growing in popularity, for detection of rare, endangered, or invasive species, is environmental DNA (eDNA) surveys (Schmeltz and Kinziger, 2016). Environmental DNA methods can detect DNA diffused by the target species in water sampled from a water body. Multiple water samples are collected from each sampling site within a lake, the sam-

ples are then analyzed using one of several different types of PCR chemistry. Sepulveda, Amberg, and Hanson (2019) suggest the use of eDNA surveys may widen the seasonal sampling window of plankton tow methods, since eDNA does not rely on a breeding population. Environmental DNA surveys are more time and cost effective than traditional sampling methods for species of low abundance (Rees et al., 2014). However, a positive eDNA result does not necessarily indicate that the target species is present or alive at the site; positive eDNA results can be obtained from “a failed introduction, from external sources, or from field contamination, rather than fresh DNA from mussel colonization” (Sepulveda et al., 2019). One criticism of the detection of dreissenid mussels using eDNA is there is a possibility of obtaining false-positive results; since control efforts for invasive species are costly, there is some hesitation in using eDNA surveys as the sole decision-making tool for the management of dreissenid mussels. When the two early detection methods result in conflicting answers, decision making can be even more complicated, since a positive eDNA result suggests only that the DNA of the target species is present, regardless of whether the species is alive or even present at all, but when veligers are detected, positive eDNA results indicate a potential colonization, which is useful to managers (Holser, 2017).

Occupancy models allow the occurrence of a species to be accurately estimated, even when the species is imperfectly detected. For both the plankton tow and eDNA surveys, there is a non-zero probability of a false negative result. Since plankton tow survey methods are restricted to capturing only

mussel veligers, dreissenid mussels can be present in the lake and not be captured in the plankton tow nets at one or more of the sampling sites, either because the veligers are missed with the nets or because there are no free-swimming veligers available in the water; even if the veligers are captured in the net, there is a possibility that they are not detected using cross-polarized light microscopy. Similarly for eDNA methods, dreissenid mussels could be present in the lake, but their DNA could be missed in one or more of the samples from each of the sampling sites; even if dreissenid mussel DNA is present in the sample, it could be missed in the PCR replicate. Replication in the survey design can help researchers learn about the detection probabilities for each of the early detection methods.

Some research question about the detection probability. Furthermore, sampling recommendations such as how many samples to take at low concentration sites and how many samples are necessary to ensure a high probability of at least one detection at the sites with low concentrations are of interest.

A few sentences here about the structure of the rest of the paper.

2 Data

The eDNA data were collected during the spring, summer, and fall of 2019 as part of a study on the detection of dreissenid mussels and Asian clams in northeastern lakes using environmental DNA. There were six lakes surveyed across Maine, New Hampshire, New York, and Vermont; each lake was sam-

pled for dreissenid mussels, Asian clams, or both. The data available are from three of the lakes sampled for dreissenid mussels. There are multiple sites within each lake, from each, five one-liter water samples were collected. Most sites were visited on more than one day throughout the course of the study. Each of the samples were analyzed using a DNA amplification technique, known as digital droplet PCR (ddPCR), on a Biorad Digital Droplet PCR system using dreissenid mussel assays. This PCR technique divides the sample aliquot into micro droplets to be analyzed separately. A field blank of clean lab water accompanied the set of five samples from each site on a given day to ensure no cross contamination in the lab. Droplet replications for dreissenid mussel DNA were negative for all field blank samples; these samples were removed from the data set.

The lakes are coded “BOM,” “LG,” and “MG,” and these names will be used throughout. There are three sites in lake BOM (BOM1, BOM2, BOM3), each of which was sampled five times on three different days (May 28, July 8, and October 21), for a total of 15 samples per site and 45 total samples from lake BOM. There are five sites in lake LG; the first three sites (LG1, LG2, LG3) were sampled five times on both May 28 and July 8, and the final two sites (LG4, LG5) were each sampled five times on September 25, for a total of 40 samples from lake LG. There were originally five sites in lake MG, but one (MG2) was inaccessible on the study day, so there are data for four sites in lake MG (MG1, MG3, MG4, MG5), each of which was sampled five times on August 23, for a total of 20 samples from lake MG. In all, there are twelve

Table 1: 10 sample rows of the eDNA data.

Lake	Site	Sample	Date	Water Temperature	Concentration	Positive Droplets
BOM	BOM1	BOM1w0528195	5/28/2019	11.0	8.51	60
BOM	BOM1	BOM1w0708194	7/8/2019	25.5	0.81	7
BOM	BOM2	BOM2w1021191	10/21/2019	NA	1.14	8
BOM	BOM3	BOM3w0528193	5/28/2019	6.0	178.72	1154
BOM	BOM3	BOM3w0708193	7/8/2019	24.5	175.99	1388
BOM	BOM3	BOM3w1021191	10/21/2019	NA	2.04	13
LG	LG1	LG1w0528194	5/28/2019	12.0	0.00	0
LG	LG2	LG2w0708194	7/8/2019	23.0	0.00	0
LG	LG2	LG2w0708195	7/8/2019	23.0	0.00	0
MG	MG5	MG5w082319D	8/23/2019	18.3	0.00	0

BOM3 all have numbers positive droplets near 0. It is noted that samples with fewer than three positive droplets should be treated with caution, as they have potential to be false positives. This is immediately noticeable in lakes LG and MG, as all the values of positive droplets per sample are low for these two lakes. Of the 40 samples from lake LG, 38 of them resulted in no positive droplets, and the remaining two samples resulted in two or three positive droplets per sample. Similarly for lake MG, 16 of the 20 samples resulted in no positive droplets, and the remaining samples resulted in one or two positive droplets per sample. While larger numbers of positive droplets were observed in samples from lake BOM, there are still three samples which resulted in some, but fewer than three, positive droplets. To limit the possibility of false positive samples for these analyses, to be considered a positive sample, the recommended threshold of three or more positive droplets will be used. The top plot in Figure 4 displays the proportion of positive samples by site if all samples with positive droplets are treated as positive samples, and the bottom plot displays the proportion of positive samples by site using the

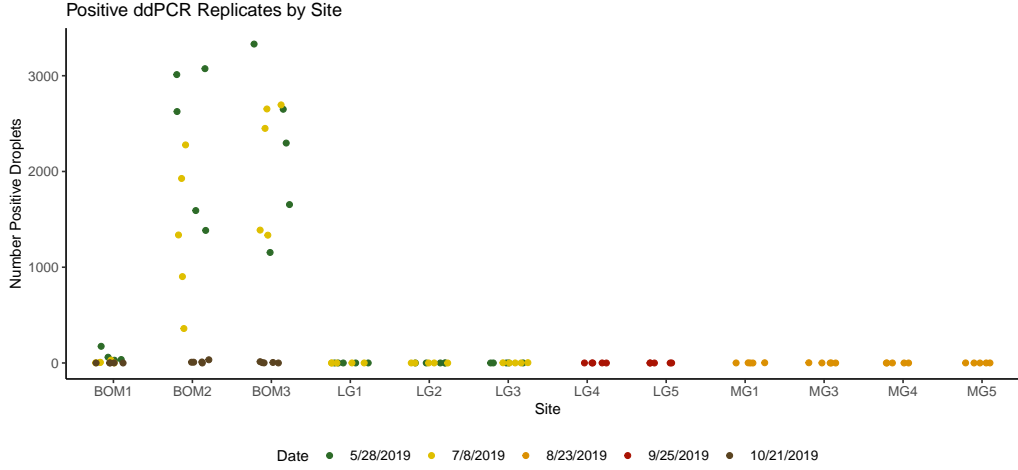


Figure 3: A plot of the positive droplets from each sample by site, colored by the date on which the sample was taken.

above threshold to limit the positive sites to those with three or more positive droplets. Figure 4 can be used to visualize the changes in the proportion of positive samples from each site when the threshold is used. The use of the threshold changes the proportion of positive samples in six of the twelve sites; the proportion of positive samples are: 0.5333 for site BOM1, 0.9333 for site BOM2, 0.8666 for lake BOM3, 0.0000 for sites LG1, LG2, LG4, and LG5, 0.1000 for site LG2, and finally 0.0000 for all sites in lake MG.

The concentrations in the study range between 0.00 and 733.80 copies/ μ L. Based on the top plot in Figure 5, it is clear that the sites with the highest concentrations occur in samples from sites BOM2 and BOM3, with average concentrations of 206 and 235 copies/ μ L in each of the respective sites. Sites LG1, LG2, LG4, LG5, and MG4 all have average concentrations of 0 copies/ μ L. The samples from the remaining sites tend to have low, but non-zero

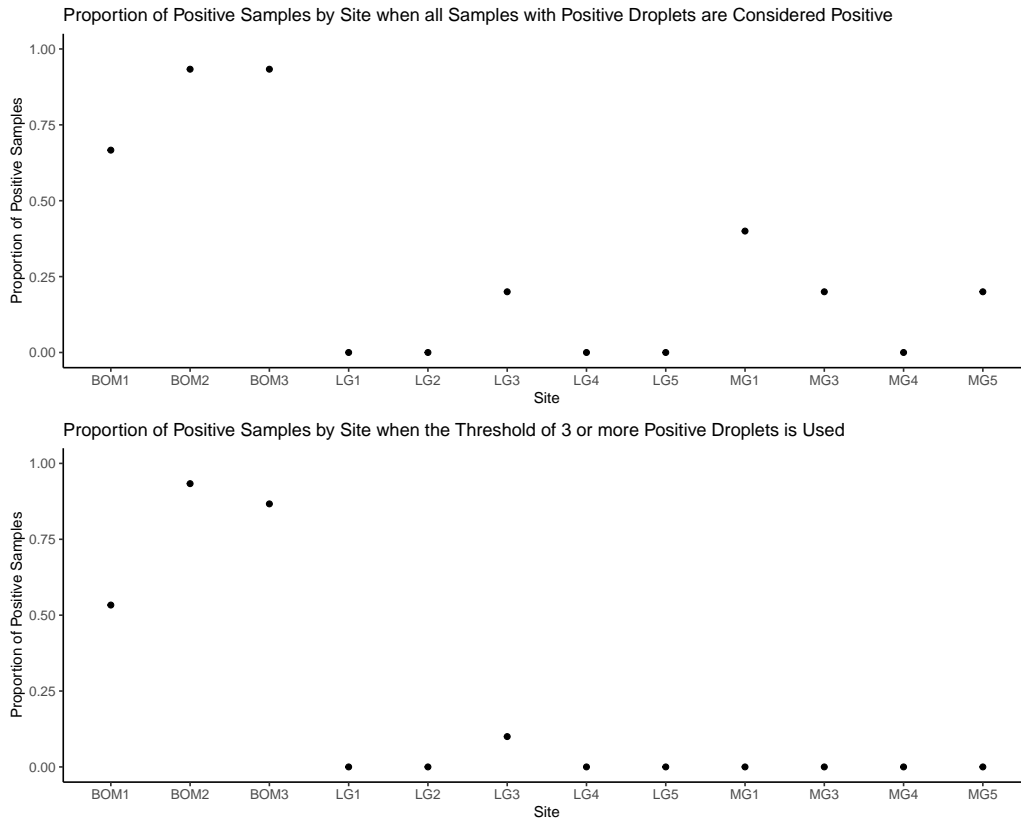


Figure 4: A plot of the proportion of positive droplets by site if all samples with any positive droplets are treated as positive detections (top), and a plot of the proportion by site if samples with fewer than three positive droplets are treated as negative (bottom).

concentrations; the bottom plot in Figure 5 displays the samples by site for non-zero concentrations below 10 copies/ μ L.

Water temperature could play a role in the availability of dreissenid mussel eDNA in the water, or the probability that the DNA is able to be detected using ddPCR. The water temperatures observed in lake BOM range from 6°C to 25.5°C, but temperatures were not recorded at any of the sites in lake BOM on October 21. The temperatures from a USGS gauge at a nearby lake, were positively adjusted by 2°C to obtain an estimated temperature of 14°C for all sites at lake BOM on October 21, 2019. The observed water temperatures in lake LG range from 12°C to 23°C, and from 17.2°C to 18.9°C in lake MG. Based on Figure 6, it seems that the water temperatures on the same day are similar across sites within the same lake. For sites that were visited more than once, it is clear that water temperatures were cooler in the spring but warmed up in the summer. *A sentence or two to transition to the methods section.*

3 Methods

3.1 Bayesian Modeling Background

In Bayesian statistics, inferences are based on the posterior distributions of the unknown parameters. The posterior distribution is a function of the likelihood of the observed data given a sampling model and some prior belief about the unknown parameters. Let $p(\mathbf{y}|\boldsymbol{\theta})$ denote the likelihood for some

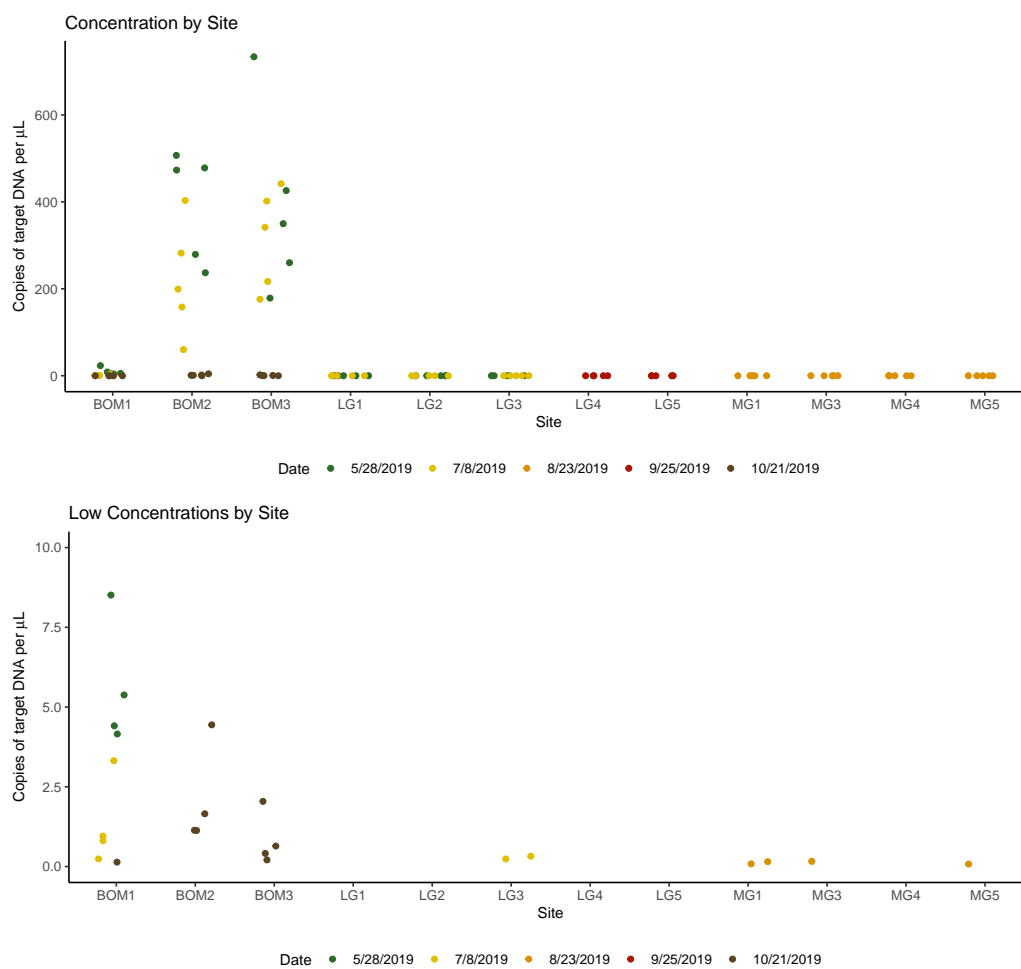


Figure 5: A plot displaying the number of copies of target DNA per μL per sample by site (top), and a plot displaying the number of copies of target DNA per μL for samples with low, but non-zero, concentrations.

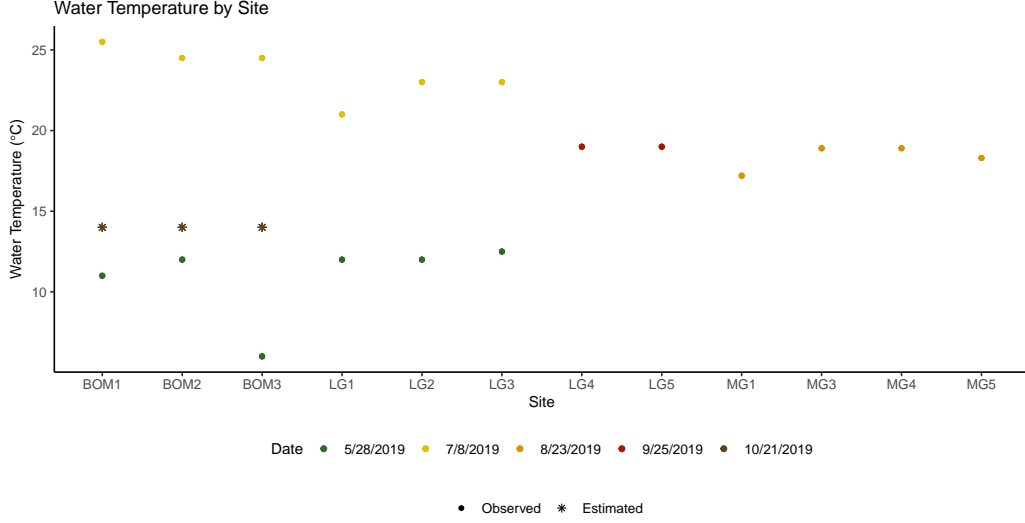


Figure 6: A plot of the water temperatures by site, colored by the date on which the sample was taken. The estimated water temperatures are denoted with stars.

sampling model, and $p(\boldsymbol{\theta})$ denote the prior distribution on the unknown parameters $\boldsymbol{\theta}$. Then, using Bayes' rule, the posterior distribution, $p(\boldsymbol{\theta}|\mathbf{y})$, is:

$$p(\boldsymbol{\theta}|\mathbf{y}) = \frac{p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{\int_{\Theta} p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})}$$

If the resulting posterior distribution is of the same distributional family as the prior distribution for the unknown parameter, then the prior is said to be conjugate for the sampling model. For example, if the observed data are n independent binomial trials, with some unknown probability of success, θ , and a beta prior distribution is used for θ , it can be shown that the resulting posterior distribution is also a beta distribution. The complete model definition is:

$$p(y_1, \dots, y_n | \theta) \sim \text{binomial}(n, \theta)$$

$$p(\theta) \sim \text{beta}(\alpha, \beta)$$

Using this, it can be shown that,

$$p(\theta | y_1, \dots, y_n) \sim \text{beta}\left(\sum_{i=1}^n y_i + \alpha, n - \sum_{i=1}^n y_i + \beta\right)$$

Therefore, since the prior distribution and the posterior distribution are of the same distributional family, then the beta prior is conjugate for the binomial sampling model.

However, the resulting posterior distribution is not always a named statistical distribution. In these situations, it is impossible to sample directly from the posterior distribution, but iterative sampling mechanisms can be used to approximate the unknown posterior distribution in order to make inferences.

3.1.1 Gibbs Sampler

The most efficient iterative sampling technique available is the Gibbs sampler. In order to draw samples from an approximation of the joint posterior distribution using a Gibbs sampler, the full conditional posterior distribution of each of the unknown parameters must have a closed-form solution. The full conditional posterior distribution of a parameter is the distribution of that parameter, conditional on all other unknown parameters, the data,

and the prior distributions. Once a full conditional posterior distribution is calculated for each unknown parameter, the iterative sampling can proceed in the following manner.

Gibbs Sampler

Let $\boldsymbol{\theta} = \{\theta_1, \dots, \theta_p\}$ denote the vector of unknown parameters and $p(\theta_i|.)$ denote the full conditional distribution of θ_i . The Gibbs Sampler generates the s^{th} iteration as follows:

1. Sample $\theta_1^{(s)} \sim p(\theta_1|\theta_2^{(s-1)}, \theta_3^{(s-1)}, \dots, \theta_p^{(s-1)}, y_1, \dots, y_n)$.
2. Sample $\theta_2^{(s)} \sim p(\theta_2|\theta_1^{(s)}, \theta_3^{(s-1)}, \dots, \theta_p^{(s-1)}, y_1, \dots, y_n)$.
- \vdots
- p . Sample $\theta_p^{(s)} \sim p(\theta_p|\theta_1^{(s)}, \theta_2^{(s)}, \dots, \theta_{p-1}^{(s)}, y_1, \dots, y_n)$.

This process generates a dependent sequence of $\boldsymbol{\theta}$ vectors for each iteration, which together, converge to the joint posterior distribution, $p(\boldsymbol{\theta}|\mathbf{y})$. The Gibbs sampler is a basic Markov Chain Monte Carlo (MCMC) algorithm, where the current state only depends on the previous state, and additionally, the results should not depend on the starting values of $\boldsymbol{\theta}^{(0)}$.

One common use of a Gibbs sampler is for a normal sampling with unknown mean and variance. With this sampling model, a normal prior on the mean term and an inverse-gamma prior on the variance term enable the use of a Gibbs sampler to approximate the joint posterior distribution. However, in some cases, such as generalized linear models, semi-conjugate priors or closed-form solutions for the full conditional distributions are not available,

and in those cases a Gibbs sampler cannot be used to sample from the joint posterior distribution of the unknown parameters.

3.1.2 Metropolis-Hastings Algorithm

In cases that do not permit the use of a Gibbs sampler, a Metropolis-Hastings algorithm is often used to sample from the target distribution. Unlike a Gibbs sampler, a Metropolis-Hastings algorithm requires tuning, but proceeds in the following manner.

Metropolis-Hastings Algorithm

Let $\boldsymbol{\theta}^{(s)}$ be the current set of parameter estimates, $\boldsymbol{\theta}^*$ denote a new proposed set of parameters, and $J(\boldsymbol{\theta}^*|\boldsymbol{\theta}^{(s)})$ denote the proposal distribution, which is usually a random walk distribution. For example, $J(\boldsymbol{\theta}^*|\boldsymbol{\theta}^{(s)}) \sim MVN(\boldsymbol{\theta}^{(s)}, \gamma^2 \mathbf{I}_p)$ where p is the number of unknown parameters, and γ is thought of as the step size, or the average distance the proposed set of parameters falls from the current set. To complete an iteration:

1. Sample $\boldsymbol{\theta}^*|\boldsymbol{\theta}^{(s)} \sim J(\boldsymbol{\theta}^*|\boldsymbol{\theta}^{(s)})$.
2. Calculate the acceptance ratio $r = \frac{p(y|\boldsymbol{\theta}^*)p(\boldsymbol{\theta}^*)}{p(y|\boldsymbol{\theta}^{(s)})p(\boldsymbol{\theta}^{(s)})}$.
3. Set $\boldsymbol{\theta}^{(s+1)}$:
 - If $r \geq 1$ then the proposed set, $\boldsymbol{\theta}^*$, is more attractive than the current set, $\boldsymbol{\theta}^{(s)}$, so $\boldsymbol{\theta}^{(s+1)} = \boldsymbol{\theta}^*$.

- If $r < 1$ then the proposed set, $\boldsymbol{\theta}^*$, is less attractive than the current set, $\boldsymbol{\theta}^{(s)}$, however the relative frequency of samples of $\boldsymbol{\theta}^*$ to $\boldsymbol{\theta}^{(s)}$ should be r , so with probability r , $\boldsymbol{\theta}^{(s+1)} = \boldsymbol{\theta}^*$.

The Metropolis algorithm is a generalization of the Metropolis-Hastings algorithm in which the proposal distribution is symmetric, meaning that $J(\boldsymbol{\theta}^*|\boldsymbol{\theta}^{(s)}) = J(\boldsymbol{\theta}^{(s)}|\boldsymbol{\theta}^*)$. The step size portion of the algorithm requires tuning and trace plots or other convergence tools should be used to ensure that the algorithm has efficiently explored the entire parameter space and converged to the true joint posterior distribution.

3.1.3 Bayesian Hierarchical Modeling

Oftentimes data is collected in a hierarchical structure, that is data with a multilevel structure, such as students within classes, or sampling sites within a larger area of interest. For example, with two levels, the responses are independent observations from group j and follow some distribution with unknown parameter θ_j , meaning that $y_{1,j}, \dots, y_{n_j,j} | \theta_j \sim p(y | \theta_j)$; but the unknown parameters in each group are related to each other, such that they are independent samples from some distribution with parameter ϕ , such that $\theta_1, \dots, \theta_m | \phi \sim p(\theta | \phi)$. With this representation, $p(y | \theta)$ is the variability among measurements within a group, and $p(\theta | \phi)$ is the sampling variability across groups. In order to completely specify the model, a prior distribution, $\phi \sim p(\psi)$, is needed for ϕ . These models can be adapted to include covariates at each level, or to describe a more complicated structure with more levels.

3.2 Occupancy Models

In ecological studies, there are several different state variables which may be of interest: abundance, vital rates, and occupancy are a few examples. Both mark-recapture and occupancy studies can be used to learn about the previously mentioned state variables, but there are advantages and limitations with each study type. In general, mark-recapture surveys are studies where individuals are captured or observed, given a unique mark, their identities recorded, and finally released; on subsequent occasions, both marked and unmarked individuals are captured, their identities recorded, unmarked individuals are marked, and they are all released. From this, capture histories, or encounter histories (typically a series of 0's and 1's) are recorded for each individual that is captured over the course of the study. There are several variations of mark-recapture studies, in which apparent survival, abundance, or both can be estimated depending on the study design. However, mark-recapture methods cannot be used when individuals of a species cannot be marked or uniquely identified. In this scenario, occupancy of a particular species can be recorded at each occasion. Occupancy methods are useful when studying a species over a large spatial scale for many years, particularly when the sampling sites are large relative to the species of interest. Though there is a loss in information about individuals of a species when occupancy studies are used, sometimes researchers are interested in the range of a species or the spread of invasion and therefore occupancy is the state variable of interest, and occupancy studies are most appropriate. Occupancy

surveys are often used to track the mitigation efforts or spread of invasive species, but are plagued with imperfect detection probabilities.

Occupancy models are a way to move past logistic regression for occupancy data in order to account for imperfect detection of the target species. However, multiple samples per site must be taken in order to learn about detection probabilities. For example, suppose five samples are taken from a site and two return a positive result, indicating that the species is present at the site, but the species was missed in three of the five samples, which provides information about the detectability of the species. There are several different types of occupancy models, but the traditional model is defined as follows. At the site level, Z_i is defined as the presence ($Z_i = 1$) or absence ($Z_i = 0$) of the target species at the i^{th} site ($i = 1, \dots, M$). Then, $Z_i \sim \text{Bernoulli}(\psi_i)$, where ψ_i denotes the probability that the target species occupies the i^{th} site. At the sample level, Y_{ij} is defined as whether the target species is detected or not in the j^{th} replicate at the i^{th} site ($j = 1, \dots, J_i$). It follows that, $Y_{ij}|z_{ij} \sim \text{Bernoulli}(z_i p_{ij})$, where z_i is the latent occupancy state for the i^{th} site and p_{ij} denotes the probability that the target species is detected on the j^{th} survey, given the species occupies site i . Figure 7 illustrates the structure of a traditional occupancy study, where each sample results in a binary response.

One assumption of occupancy models is that there is no unmodeled heterogeneity in the occupancy probabilities, or the detection probabilities, therefore covariates can be included at the site and sample levels. At the site

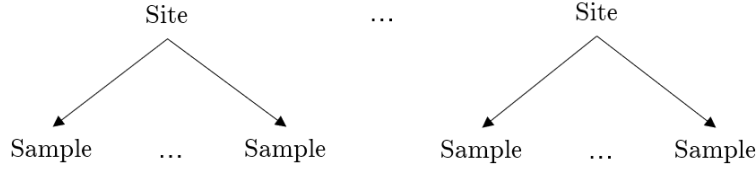


Figure 7: A diagram displaying the structure of a traditional occupancy study.

level, \mathbf{x}_i are the site specific covariates and $\boldsymbol{\beta}$ is the corresponding vector of regression parameters for ψ_i . At the sample level, \mathbf{w}_{ij} are the sample specific covariates and $\boldsymbol{\alpha}$ is the corresponding vector of regression parameters for p_{ij} . It is also assumed that the sites are closed to changes in occupancy, so a site that is occupied at the beginning of the study should remain occupied throughout; if an occupancy study takes place over many years, there are models to account for the potential changes in occupancy status at the sites. Another assumption of these models is that the species is not misidentified, however, there are models that can account for false negative and false positive detections. The final assumption of these models is that the detection of the species is independent across sites. Violations of these assumptions can result in biased estimates of the occupancy and detection probabilities.

3.2.1 Multi-scale Occupancy Models

For more elusive species, the samples from each site are split into subsamples, or replicates, and analyzed on a smaller scale. This introduces a hierarchical structure in the data creating dependencies which need to be

accounted for (MacKenzie et al., 2018). One modeling approach to these types of data is the hierarchical latent state multi-scale occupancy model defined by Dorazio and Erickson which accounts for false negatives at the site, sample, and replicate levels (Dorazio and Erickson, 2018). The model is defined similarly to the traditional occupancy model described above, with an additional replicate level, the notation is as follows. At the site level, Z_i is defined as the presence ($Z_i = 1$) or absence ($Z_i = 0$) of the target species at the i^{th} site ($i = 1, \dots, M$). Then, $Z_i \sim \text{Bernoulli}(\psi_i)$, where ψ_i denotes the probability that the target species occupies the i^{th} site. At the sample level, A_{ij} is defined as the presence or absence of the target species in the j^{th} sample from the i^{th} site ($j = 1, \dots, J_i$). Then, $A_{ij}|z_i \sim \text{Bernoulli}(z_i\theta_{ij})$, where z_i is the latent occupancy state for the i^{th} site, and θ_{ij} is the conditional probability that the target species is present in the j^{th} sample from the i^{th} site, given the target species is present at the site. At the replicate level, Y_{ijk} is defined as whether the target species is detected or not in the k^{th} replicate of the j^{th} sample collected at the i^{th} site ($k = 1, \dots, K_{ij}$). It follows that, $Y_{ijk}|a_{ij} \sim \text{Bernoulli}(a_{ij}p_{ijk})$, where a_{ij} is the latent occupancy state for sample j from site i , and p_{ijk} is the conditional probability that the target species is detected in the k^{th} replicate of the j^{th} sample collected at the i^{th} site, given the target species is present in that sample. If p_{ijk} does not differ among the replicates and the replicates are statistically independent, the replicate level can be simplified in the following manner: $Y_{ij}|a_{ij} \sim \text{binomial}(K_{ij}, a_{ij}p_{ij})$, where $Y_{ij} = \sum_{k=1}^{K_{ij}} Y_{ijk}$, and p_{ij} is the condi-

tional probability of detection in each replicate of the j^{th} sample collected at the i^{th} location, given that the target species is present in that sample. Similarly to above, an assumption of these models is that there is no unmodeled heterogeneity in the detection and occupancy probabilities, so covariates can be included in each level of the model to account for differences between individual sites, samples, or replicates; the notation for the covariates and corresponding regression parameters for the multi-scale occupancy model are as follows. At the site level, \mathbf{x}_i are the site specific covariates and $\boldsymbol{\beta}$ is the corresponding vector of regression parameters for ψ_i . At the sample level, \mathbf{w}_{ij} are the sample specific covariates and $\boldsymbol{\alpha}$ is the corresponding vector of regression parameters for θ_{ij} . At the replication level, where there are no replicate level covariates, \mathbf{v}_{ij} are the sample specific covariates and $\boldsymbol{\delta}$ is the corresponding vector of regression parameters for p_{ij} . Site specific covariates can be included for θ_{ij} and p_{ij} . The remaining assumptions from traditional occupancy models, described above, hold for multi-scale occupancy models as well.

3.3 Implementation

There are several options for fitting occupancy models; the **R** package used for this analysis as well as some alternative methods for fitting occupancy models are described below.

3.3.1 msocc package

get rid of this subsection?

The multi-scale occupancy model described above (Dorazio and Erickson, 2018) can be easily fit in a Bayesian framework with the `msocc` package (Stratton, 2020). This package implements the Pólya Gamma data augmentation strategy described by Polson et. al (2013) to enable the use of a Gibbs Sampler in order to speed up the computation and avoid the tuning necessary for a Metropolis Algorithm. The model is specified in a familiar way to users of R at each level of the hierarchy in the model fitting function `msocc_mod()`. To make the package more computationally efficient, there is an option to use the Beta-Binomial sampler when possible. Built into the package are posterior summary tools (`posterior_summary()`) and a function for generating credible interval plots (`cred_plot()`) at all levels of the hierarchy.

3.3.2 *implementation for traditional occupancy model*

3.3.3 Alternative Options

There are several alternative options for fitting occupancy models. For example, the `eNDAoccupancy` (Dorazio and Erickson, 2018) package fits the Bayesian multi-scale occupancy model defined above using a Metropolis-Hastings algorithm, which is less efficient and requires tuning. These models can also be written in STAN or JAGS and run in R using the packages `rstan` (Stan Development Team, 2019) and `rjags` (Plummer, 2019) respectively,

but again these methods are less computationally efficient than the methods described above. There are various frequentist methods for fitting occupancy models as well. A commonly used R package for fitting frequentist occupancy models is `unmarked` (Fiske and Chandler, 2011). Ecologists often utilize Program MARK, or access its functionalities through the R package `RMark` (Laake, 2013) to fit occupancy models, and other various models popular in ecology.

4 Analysis

A paragraph about the motivation to analyze the data this way, traditional occupancy model framework with the threshold

```
# format data for use in PGocc4 function
## number of samples per sites

J <- rep(0, M)

for(j in 1:length(J)){
  for(i in 1:nrow(eDNA)){
    if(eDNA$Site[i] == Site.names[j]){
      J[j] <- J[j] + 1
    }
  }
}
```

```

## number the from each site 1:J[i]

samp <- matrix(NA, nrow = M, ncol = max(J))

for(i in 1:M){
  if(J[i] == ncol(samp)){
    samp[i, ] <- seq(1:J[i])
  } else {
    samp[i, 1:J[i]] <- seq(1:J[i])
  }
}

samp <- na.omit(as.vector(t(samp)))

## detection/non-detection matrix

y <- matrix(NA, nrow = M, ncol = max(J))

for(i in 1:M){
  for(k in 1:nrow(eDNA)){
    if(eDNA$Site[k] == Site.names[i] & eDNA$Detect3[k] == 1){
      y[i, samp[k]] <- 1
    }
  }
}

```

```

    } else if(eDNA$Site[k] == Site.names[i] & eDNA$Detect3[k] == 0){
      y[i, samp[k]] <- 0
    }
  }
}

y <- matrix(as.integer(y), nrow = M, ncol = max(J))

## site level covariates
avg.water.temp <- eDNA %>%
  group_by(Site) %>%
  summarise(mean(Water.Temp))

X <- data.frame(Site.names)
names(X) <- "site"
X[, 2] <- X$site %>%
  as.character() %>%
  removeNumbers() %>%
  factor()
X <- cbind(X, avg.water.temp[, 2])

```

```

X <- as.matrix(X)

colnames(X)[2:3] <- c("lake", "avg.water.temp")

## sample level covariates
### water temperature
W1 <- matrix(NA, nrow = M, ncol = max(J))

for(i in 1:M){
  if(J[i] == ncol(W1)){
    W1[i, ] <- eDNA$Water.Temp[eDNA$Site == Site.names[i]]
  } else {
    W1[i, 1:J[i]] <- eDNA$Water.Temp[eDNA$Site == Site.names[i]]
  }
}

### date
W2 <- data.frame(matrix(NA, nrow = M, ncol = max(J)))

for(i in 1:M){
  if(J[i] == ncol(W2)){
    W2[i, ] <- eDNA$Date.Collect[eDNA$Site == Site.names[i]]
  }
}

```

```

    } else {
      W2[i, 1:J[i]] <- eDNA$Date.Collected[eDNA$Site == Site.names[i]]
    }
  }

W2 <- as.matrix(W2)

### site
W3 <- data.frame(matrix(NA, nrow = M, ncol = max(J)))

for(i in 1:M){
  if(J[i] == ncol(W3)){
    W3[i, ] <- eDNA$Site[eDNA$Site == Site.names[i]]
  } else {
    W3[i, 1:J[i]] <- eDNA$Site[eDNA$Site == Site.names[i]]
  }
}

W3 <- as.matrix(W3)

W <- list(W1 = W1, W2 = W2, W3 = W3)

```



```
# merge for use in PGocc4
```

```
data <- vb_Designs(W = W, X = X, y = y)
```

```
post_summary <- function(parameter.matrix, param, plot.options = NULL, plot = T){  
  mean <- apply(parameter.matrix, 1, mean)  
  median <- apply(parameter.matrix, 1, median)  
  lwr <- apply(parameter.matrix, 1, quantile, probs = 0.025)  
  upr <- apply(parameter.matrix, 1, quantile, probs = 0.975)  
  out <- cbind("mean" = mean, "median" = median, "2.5%" = lwr , "97.5%" = upr)  
  if(missing(param)){  
    out <- out  
  } else if(param == "psi"){  
    out <- cbind.data.frame(out, rownames = Site.names)  
  } else if(param == "p"){  
    out <- cbind.data.frame(Site = eDNA$Site, Sample = samp, out)  
  } else if(param == "alpha"){  
    rownames(out) <- paste("alpha", 0:(nrow(out) - 1))  
  } else if(param == "beta"){  
    rownames(out) <- paste("beta", 0:(nrow(out) - 1))  
  }  
  if(plot == T){  
    for(i in 1:nrow(parameter.matrix)){
```

```

df <- data.frame(parameter.matrix[i, ])
tr.plot <- ggplot(df, aes(x = 1:ncol(parameter.matrix), y = parameter.matrix[, i]))
  geom_path() +
  xlab("Post Burn-in Iterations") +
  ylab(rownames(out)[i]) +
  theme_bw() +
  theme(title = element_text(size = 10),
        panel.border = element_blank(),
        panel.grid.major = element_blank(),
        panel.grid.minor = element_blank(),
        axis.line = element_line(colour = "black"))
if(is.null(plot.options)){
  tr.plot <- tr.plot
} else {
  tr.plot <- tr.plot + plot.options}
print(tr.plot)
}
}
return(out)
}

```

```

cred_plot <- function(post_summary_out, parameter){
  df <- data.frame(post_summary_out)
  if(nrow(df) == 12){
    rownames(df) <- Site.names
  }
  if(parameter == "psi"){
    cred.plot <- ggplot(df, aes(x = rownames(df), y = mean)) +
      geom_point() +
      xlab("") +
      geom_errorbar(aes(ymax = X97.5., ymin = X2.5.), width = 0) +
      theme_bw() +
      theme(title = element_text(size = 10),
            panel.border = element_blank(),
            panel.grid.major = element_blank(),
            panel.grid.minor = element_blank(),
            axis.line = element_line(colour = "black"),
            axis.text.x = element_text(angle = 45, hjust = 1)) +
      ylab(expression(paste("Estimated ", psi)))
  } else if(parameter == "p"){
    cred.plot <- ggplot(df, aes(x = eDNA$Date.Collect, y = mean)) +
      geom_point() +
      xlab("") +

```

```

geom_errorbar(aes(ymax = X97.5., ymin = X2.5.), width = 0) +
theme_bw() +
theme(title = element_text(size = 10),
      panel.border = element_blank(),
      panel.grid.major = element_blank(),
      panel.grid.minor = element_blank(),
      axis.line = element_line(colour = "black"),
      axis.text.x = element_text(angle = 45, hjust = 1)) +
ylab("Estimated p") +
facet_wrap(~Site)
}
print(cred.plot)
}

```

5 Results

5.1 Ecological Implications

when analyzed at the droplet level, less confident about negative result

5.2 Sampling Recommendations

6 Future Work

arbitrary amplitude cutoff for choosing the droplets that are positive and negative

threshold for the number of positive droplets to be considered a positive sample

what is a false positive and how does it bias the results

discussion of droplet level and how multi-scale occupancy model can

7 Acknowledgements

8 References

- Bontke, J. (2020). City spends \$4 million on liquid compound to stop spread of zebra mussels. <https://cbsaustin.com/news/local/city-spends-4-million-on-liquid-compound-to-stop-spread-of-zebra-mussels>. Date accessed: January 29, 2020.
- Bureau of Reclamation (2019). Invasive mussels. <https://www.usbr.gov/mussels/index.html>. Date accessed: February 11, 2020.
- Clark, A. (2018). *Rcppocc: Fits various Bayesian occupancy models assuming all regression effects uses a logit link function*. R package version 1.0.
- Dorazio, R. M. and Erickson, R. A. (2018). **eDNAoccupancy**: An R package for multiscale occupancy modelling of environmental DNA data. *Molecular Ecology Resources*, 18(2):368 – 380.
- Fiske, I. and Chandler, R. (2011). *unmarked: An R Package for Fitting Hierarchical Models of Wildlife Occurrence and Abundance*.
- Holser, D. M. (2011). Early detection of dreissenid species: Zebra/Quagga mussels in water systems. *Aquatic Invasions*, 8(2):217 – 222.
- Holser, D. M. (2017). Where is the body? Dreissenid mussels, raw water

- testing, and the real value of environmental DNA. *Management of Biological Invasions*, 8(3):335 – 341.
- Laake, J. (2013). *RMark: An R Interface for Analysis of Capture-Recapture Data with MARK*.
- MacKenzie, D. I., Nichols, J. D., Royle, J. A., and Pollock, K. H. (2018). *Occupancy estimation and modeling: inferring patterns and dynamics of species occurrence*. Academic Press.
- McMahon, R. F. (1996). The physiological ecology of the zebra mussel, *Dreissena polymorpha*, in North America and Europe. *American Zoologist*, 36(3):339 – 363.
- Mills, E. L., Rosenberg, G., Spidle, A. P., Ludyanskiy, M., Pligin, Y., and May, B. (1996). A review of the biology and ecology of the quagga mussel (*Dreissena bugensis*), a second species of freshwater dreissenid introduced to North America. *American Zoologist*, 36(3):271 – 286.
- Nelson, N. M. (2019). Conservation and resource development.
[urlhttp://dnrc.mt.gov/divisions/cardd](http://dnrc.mt.gov/divisions/cardd). Date Accessed: February 15, 2020.
- Nichols, S. J. (1996). Variations in the reproductive cycle of *Dreissena polymorpha* in Europe, Russia, and North America. *American Zoologist*, 36(3):311 – 325.

- Plummer, M. (2019). *rjags: Bayesian Graphical Models using MCMC*. R package version 4-10.
- Polson, N. G., Scott, J. G., and Windle, J. (2013). Bayesian inference for logistic models using pólya–gamma latent variables. *Journal of the American Statistical Association*, 108(504):1339–1349.
- R Core Team (2019). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Rees, H. C., Maddison, B. C., Middleditch, D. J., Patmore, J. R., and Gough, K. C. (2014). The detection of aquatic animal species using environmental dna – a review of edna as a survey tool in ecology. *Journal of Applied Ecology*, 51:1450 – 1459.
- Schmelzle, M. C. and Kinziger, A. P. (2016). Using occupancy modelling to compare environmental DNA to traditional field methods for regional-scale monitoring of an endangered aquatic species. 16:1 – 14.
- Sepulveda, A. J., Amberg, J. J., and Hanson, E. (2019). Using environmental DNA to extend the window of early detection for dreissenid mussels. *Managment of Biological Invasions*, 10(2):342 – 358.
- Stan Development Team (2019). *RStan: the R interface to Stan*. R package version 2.19.2.
- Stratton, C. (2020). *msocc: Package for Fitting and Analyzing*

Computationally Efficient Multi-scale Occupancy Models. R package
version 1.1.0.

U.S. National Park Service (2017). Invasive Zebra Mussels.

<https://www.nps.gov/articles/zebra-mussels.htm>. Date

Accessed: February 11, 2020.

9 Appendix - R Code

```
# packages used
library(car)
library(dplyr)
library(tidyr)
library(dataRetrieval)
library(kableExtra)
library(ggplot2)
library(gridExtra)
library(tm)
library(readxl)
## library(devtools)
## devtools::install_github("StrattonCh/msocc")
library(msocc)
library(Rcppocc)

# load eDNA data
eDNA <- read.csv(
  "C:/Users/mwind/OneDrive/Writing Project_EXTRA/eDNA.csv")

# rename lakes
levels(eDNA$Lake) <- c("BOM", "LG", "MG")

# remove field blank samples
eDNA <- eDNA %>%
  filter(Site != "tb")
eDNA$Site <- droplevels(eDNA$Site)

# reorder dates in chronological order
eDNA$Date.Collecte d <- factor(eDNA$Date.Collecte d,
```

```

                                levels = c("5/28/2019",
                                              "7/8/2019",
                                              "8/23/2019",
                                              "9/25/2019",
                                              "10/21/2019"))

Site.names <- unique(eDNA$Site)

M <- length(Site.names)

lake <- eDNA$Site %>%
  as.character() %>%
  removeNumbers() %>%
  factor()

# generate table of 10 sample rows of eDNA data
set.seed(03142020)
knitr::kable(some(eDNA), 'latex', booktabs = T, linesep = "",
              caption = "\\label{tab:eDNA_data}
              10 sample rows of the eDNA data.",
              align = 'c', row.names = F,
              col.names = c("Lake", "Site", "Sample", "Date",
                            "Water Temperature", "Concentration",
                            "Positive Droplets")) %>%
kable_styling(latex_options =
              c("scale_down", "hold_position"))

```