

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

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Data</b>	<b>6</b>
<b>3</b>	<b>Modeling Background</b>	<b>6</b>
3.1	Occupancy Models . . . . .	6
3.2	Bayesian Modeling Background . . . . .	10
<b>4</b>	<b>Methods</b>	<b>10</b>
<b>5</b>	<b>Analysis</b>	<b>10</b>
<b>6</b>	<b>Conclusion</b>	<b>11</b>
6.1	Further Investigations . . . . .	11
<b>7</b>	<b>References</b>	<b>12</b>
<b>8</b>	<b>Appendix - R Code</b>	<b>14</b>

## 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 mussels<sup>1</sup> from water bodies across the United States, and four million dollars is only 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 United States. 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, so that organizations can plan, budget, and install necessary technologies before colonization has occurred (Holser, 2017).

---

<sup>1</sup>Zebra mussels (*Dreissena polymorpha*) and quagga mussels (*Dreissena rostriformis bugensis*) collectively.

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). 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 infested with 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). 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 Depart-

ment of the Interior, 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 mussel species 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). 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 for detection of rare, endangered, or invasive species, one growing in popularity, is environmental DNA (eDNA) surveys (Schmelzle and Kinziger, 2016). Environmental DNA methods can detect DNA diffused from the target species from water sampled from a water body. Multiple water samples are collected from each sampling site within a lake,

the samples are then analyzed using one of several types of PCR chemistry. Sepulveda, Amberg, and Hanson (2019) suggest the use of eDNA surveys may widen the seasonal sampling window over plankton tow methods, since eDNA does not rely on a breeding population. This method is 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 mean 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 invasive species.

*Occupancy models allow the occurrence of a species to be accurately estimated, even when the species is imperfectly detected. For both the plankton tow surveys and the eDNA surveys, there is a non-zero probability of a false positive result. repeated surveys at the same site to learn about detection probability. generally, even if the site is target species is present at the site, not all samples from the site will contain the target species, and even if the sample contains the target species, not all subsamples from it will test positive for the target species.*

When these two methods result in conflicting answers, decision making

can be even more complicated, since a positive eDNA result only suggests 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).

Research question(s): How does the detection probability of zebra mussels using microscopy compare to the detection probability of zebra mussels using eDNA? What is the false negative rate for both methods? Given that, how many samples should we collect at a given water body?

## 2 Data

lots of data visualization here

## 3 Modeling Background

### 3.1 Occupancy Models

**I did not finish this section like I was supposed to :( but I did learn a lot more about multi-scale occupancy models and I will try to have this section cleaned up by the end of the week**

Occupancy is the presence of a particular species on a given site, this may not be the first choice of state variables to ecologists but occupancy studies are useful when there is a large spatial scale or the study is conducted over



many years, when abundance or vital rates are hard to measure. Occupancy studies are also useful over capture-recapture methods when individuals cannot be marked or uniquely identified. However, sometimes patterns of species occurrence are of interest, this happens when researchers are interested in the range of a species or the spread of invasion. The sampling units for occupancy studies are called 'sites'. We can learn about detection probabilities when multiple site visits are used. Also, when using occupancy models we need to account for imperfect detection because it is possible that the researchers could miss the species even if it is present at the site.

in the hierarchical multi-scale occupancy model (with three nested levels) defined by (Dorazio and Erickson, 2017):

- learn about the probability of species occurrence at a site
- $Z_i$  denotes the presence or absence of the target species at the  $i^{th}$  site ( $i = 1, \dots, M$ )
- $Z_i \sim \text{Bernoulli}(\psi_i)$
- $\psi_i$  is the probability that the target species is present at the  $i^{th}$  site
- $\beta$  are the site level regression parameters for  $\psi_i$
- $\mathbf{x}_i$  are the site level covariates for  $\psi_i$
- learn about the conditional probability of species occurrence in a sample of a site given the species is present at the site

- $A_{ij}$  denotes the presence or absence of the target species in the  $j^{th}$  sample from the  $i^{th}$  site ( $j = 1, \dots, J_i$ )
- $A_{ij}|z_i \sim \text{Bernoulli}(z_i\theta_{ij})$
- $z_i$  is a realized value of  $Z_i$
- $\theta_{ij}$  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 location
- $\alpha$  are the sample level regression parameters for  $\theta_{ij}$
- $w_{ij}$  are the sample level covariates for  $\theta_{ij}$
- learn about the conditional probability of detection of the speies in a subsample of a sample given that the species is present in the sample
- $Y_{ijk}$  denotes the detected or not detected in the  $k^{th}$  replicate of the  $j^{th}$  sample collected at the  $i^{th}$  site ( $k = 1, \dots, K_{ij}$ )
- $Y_{ijk}|a_{ij} \sim \text{Bernoulli}(a_{ij}p_{ijk})$
- $a_{ij}$  is a realized value of  $A_{ij}$
- $p_{ijk}$  denotes the 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 independent, then  $Y_{ij} = \sum_{k=1}^{K_{ij}} Y_{ijk}$
- $Y_{ij}|a_{ij} \sim \text{Binomial}(K_{ij}, a_{ij}p_{ij})$
- $p_{ij}$  is the conditional probability of detection of the target species in each replicate of the  $j^{\text{th}}$  sample collected at the  $i^{\text{th}}$  location, given that the target species is present in that sample
- $\delta$  are the sample level regression parameters for  $p_{ij}$
- $\mathbf{v}_{ij}$  are the sample level covariates for  $p_{ij}$

The assumptions are:

- The occupancy state of sites is constant during a single season.
- The occupancy probability is constant across sites, or is modeled appropriately using site-level covariates.
- The probability of detection given occupancy status is constant across sites, or modeled appropriately using site-level covariates.
- The species is not misidentified, no false positives.

In WILD 502 when talking about multi-season occupancy models, we talked about extirpation and colonization rates, but I think that these could be modeled with a latent variable(s)? I don't think they are of particular interest here.

## 3.2 Bayesian Modeling Background

# 4 Methods

Package options for Multi-season (should we even be using that or should we be modeling the time component in another way) single-species occupancy models:

- nimble.dynamic.occ
- STAN
- JAGS
- wiquid package??
- Frequentist Options:
  - unmarked
  - Program MARK
- write my own package?

Package for eDNA data:

msocc package

# 5 Analysis

Could do EDA here

## 6 Conclusion

### 6.1 Further Investigations

## 7 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.
- Dorazio, R. M. and Erickson, R. A. (2017). `ednaoccupancy`: An R package for multiscale occupancy modelling of environmental DNA data. *Molecular Ecology Resources*, 18(2):368 – 380.
- 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.
- 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.

- Nichols, S. J. (1996). Variations in the reproductive cycle of *Dreissena polymorpha* in Europe, Russia, and North America. *American Zoologist*, 36(3):311 – 325.
- 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.
- U.S. National Park Service (2017). Invasive Zebra Mussels. <https://www.nps.gov/articles/zebra-mussels.htm>. February 11, 2020.

## 8 Appendix - R Code

Things to think about for the simulation or questions I have:

- think about how changing the detection probabilities and occupancy probabilities impact the results:
  - High occupancy, high detection
  - High occupancy, low detection
  - Low occupancy, high detection
  - Low occupancy, low detection
- How do we/can we include sample level covariates to account for sampling effort (number of tows, if available)?
- If we are getting more data:
  - How do we account for multiple sampling seasons (years)?
  - How do the assumptions of occupancy models change when we have several sampling years versus only 1

Additional questions, not directly related to the simulation:

- How do we account for this multilevel (for lack of a better word) testing process? For example, sometimes (?) when they find them in the microscope, then they test them using polymerase chain reaction (PCR) to confirm positive identification, and sometimes gene sequencing (? –



looks like it was only used on one observation in the sample data), or scanning electron microscopy (SEM).

- One assumption of occupancy models is that there is a non-zero probability of misidentification. I read an article (Denise Holser) that suggests that there might be an issue with that assumption in this situation because there are similar looking organisms that may be present in the waters; she then suggest that the Bureau of Reclamation has attempted to mediate this issue with improved microscopic methods and improved PCR methods. I will look into this some more.

```
# SIMULATED DATA 1
set.seed(1202020)
p <- 0.6 # constant detection probability... could change to depend on covariate
psi <- 0.6 # constant occupancy probability... could change to depend on covariate
M <- 10 # number of sites
J <- 8 # constant number of samples per site... does not need to be constant

z <- rbinom(M, 1, psi) # site occupancy

y1 <- matrix(NA, nrow = M*J, ncol = 4)
y1[, 1] <- rep(1:M, each = J) # column indicating site
y1[, 2] <- rep(1:J, M) # column indicating sample within site
y1[, 3] <- rep(z, each = J) # column indicating true site occupancy

# creates a column of 1's and 0's indicating whether the species was detected
for(i in 1:(M*J)){
  y1[i, 4] <- rbinom(1, 1, p*y1[i, 3]) # if the site is not occupied the species
}

colnames(y1) <- c("Site", "Sample", "True Occupancy", "Detected")
```

```

library("car")
some(y1) # view 10 sample rows of the simulated data

# or

y2 <- rbinom(M, J, p*z) # total number of detections for the J samples within ea
# SIMULATED DATA 2
set.seed(1222020)
p <- 0.6 # constant detection probability... could change to depend on covariate
psi <- 0.6 # constant occupancy probability... could change to depend on covaria
M <- 10 # number of sites
J <- sample(1:10, M, replace = T) # number of times each of the sites were sample

z <- rbinom(M, 1, psi) # site occupancy

y <- rep(NA, M)

for(i in 1:M){
  y[i] <- rbinom(1, J[i], p*z[i]) # total number of detections for the J samples
}
#SIMULATED DATA 3
set.seed(1222020)
M <- 10 # number of sites
x <- runif(10, 0, 10)
beta1.true <- 2
beta2.true <- 0.5
beta3.true <- 0.5
beta4.true <- 0.2
p <- exp(beta1.true - beta2.true*x)/(1 + exp(beta1.true - beta2.true*x)) # detect
psi <- exp(beta3.true + beta4.true*x)/(1 + exp(beta3.true + beta4.true*x)) # occup
J <- sample(1:10, M, replace = T) # number of times each of the sites were sample

z <- rep(NA, M)
y <- rep(NA, M)

for(i in 1:M){
  z[i] <- rbinom(1, 1, psi[i]) #site occupancy

```

```
y[i] <- rbinom(1, J[i], p[i]*z[i]) # total number of detections for the J samp  
}
```