

Tradeoffs in detection efficiency and area covered in the design of optimal survey

Jake M. Ferguson, Aislyn Keyes, Michael McCartney, Katie St. Clair, Douglas Johnson, John Fieberg

September 17, 2018

Introduction

- impacts and costs of surveying invasive species: Past work has focused primarily on the early detection (e.g., Ferguson et al. 2014, Holden, Nyrop, and Ellner (2016)) and eradication (???) of invasive species or detecting rare and endangered species (???), however the complementary problem of designing efficient surveys to determine population distribution and density of animals post-detection is often a larger component of the management efforts. Determining how to optimize these efforts has received much less attention.
- Something about usefulness of Unionids as a study system for studying detection methods [Green1993, Smith2006]...
- There likely exists an empirical tradeoff between survey efficiency and survey coverage in most study designs. Surveys performed deliberately (*word choice*) will likely be time-consuming and thus cover a limited area, but have a higher probability of detecting individuals. If a survey is performed quickly it can cover more area at the expense of lower detection rates.
- Here we explored two types of surveys for zebra mussels. We used transect surveys with distance sampling, which covers a larger area but has imperfect detection, and quadrat surveys, which represent a low efficiency type of survey with high detection probability. We examined how each design performed under a range of densities. Finally, we used the empirical studies to parameterize a general model that allowed us to determine the best survey approach for a given problem based on **system size and density**.

Methods

Surveys

We conducted two different types of surveys in three Minnesota lakes in 2018. We decided which lakes to survey based on initial visits to six different lakes throughout central Minnesota (Christmas Lake, East Lake Sylvia, Lake Burgan, Lake Florida, Little Birch Lake, Sylvia Lake) that have had confirmed recent zebra mussel infestations (as determined by Minnesota Department of Natural Resources). At each lake we visited 15 different sites distributed evenly around the lake. Each site was placed in 3 to 8 m of water and determined using a bathymetry shapefile in ArcMap. We located each point in the field using a GPS unit (Garmin GPSMAP 64s). At each site our team of two divers spent 15 minutes underwater counting zebra mussels. We used these counts to determine the three lakes to conduct full surveys on, selecting lakes that displayed a range of apparent densities. **how were quadrat/distance sites differentiated**

Based on our initial 15 minute exploratory dives we conducted full surveys on Lake Florida in Kandiyohi County, Lake Burgan in Douglas County, and Little Birch Lake in Todd County. Lake Florida covers an area of 273 hectares and has a maximum depth of 12 m, Lake Burgan covers an area of 74 hectares and has a maximum depth of 13 m, Little Birch Lake covers 339 hectares and has a maximum depth of 27 m. We surveyed each of the 15 previously selected sites in each lake using two type of surveys; quadrat and distance sampling with removal.

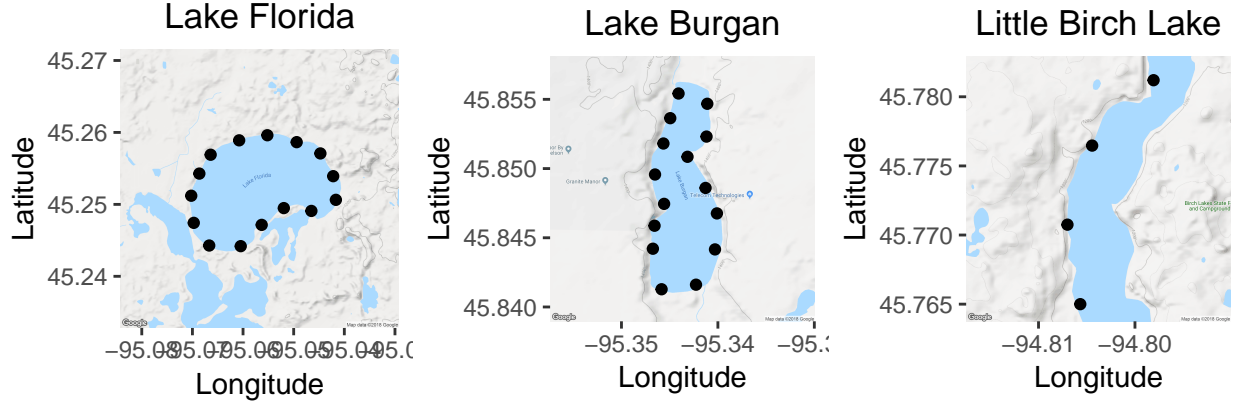


Figure 1: Map of survey sites, each point indicates the starting location of a transect. Remove place name labels. Fix x-axis on Florida.

Quadrat surveys

At each site we used the previously defined transect locations to determine the start of a transect. We ran out parallel 30 m transect lines 1 meter apart perpendicular to the shoreline, though transects were stopped earlier than 30 m if divers ran into the thermocline. Our team of two divers each took one of the transects, placing a 0.5×0.5 square meter quadrat every 2 meters along the transect. In each quadrat the diver counted all the mussels within the quadrat.

Distance sampling surveys

At each site we used the previously defined transect locations to determine the start of each survey transect. We then ran out a 30 m transect in a direction perpendicular to the shoreline, though transects were stopped earlier than 30 m if divers ran into the thermocline. Divers surveyed 1 m on either side of the transect for a transect belt that is 2 m wide.

We conducted removal surveys, which require two divers. In the removal survey, the first diver swam over the transect line marking detections for the second diver. The second diver then tried to detect animals missed by the first diver. We implemented the distance removal survey by having the primary diver swim over the transect line. Whenever the diver detected a zebra mussel or cluster of mussels, they marked the location with a survey flag then recorded the number of mussels in the cluster, the distance from the transect start to the detection (hereafter transect distance), and the perpendicular distance from the location of the detection to the transect line (hereafter detection distance). The secondary diver then looked for zebra mussels that were missed by the primary diver. Divers rotated through the primary and secondary observer roles in order to *average out* potentially innate differences between observers (Cook and Jacobson 1979).

- remove habitat survey information?

How were survey transects separated?

Statistical analysis

- Distance sampling: the counts for each transect are denoted as x_i with the total counts in the lake with k transects is denoted as $X = \sum_i^k x_i$. The length of each transect is denoted as l_i and the total length of $L = \sum_i^k l_i$. The encounter rate x_i/l_i is used to weight the survey effort. The estimated detection probability on a transect is denoted as \hat{P} and the cluster size of each detection event is denoted as s . The estimated density and variance in density (Buckland et al. 2001) are given by

$$\hat{D} = \frac{\sum_i^k x_i}{\hat{P}A} \quad (1)$$

$$\text{var}(\hat{D}) = \hat{D}^2 \left(\frac{\text{var}(X)}{X^2} + \frac{\text{var}(E(s))}{E(s)^2} + \frac{\text{var}(A\hat{P})}{A^2\hat{P}^2} \right), \quad (2)$$

with the total survey area $A = \sum_i^k a_i$, variance in counts $\text{var}(X) = (L \sum_i^k l_i (x_i/l_i - X/L)^2) / (k - 1)$, and the expected value and variance in the expected cluster size (denoted as $E(s)$ and $\text{var}(E(s))$), respectively. The variance in the detectability $\text{var}(A\hat{P})$ is **XX**. Below we describe how to determine the detection probability component using distance sampling.

In distance sampling the transect distance is used to estimate how detectability changes with distance from the transect line. We modeled detection probabilities using two model subcomponents. The distance component of the detection function, describes how distance (y) leads to changes in the probability of detecting a zebra mussel. Based on our initial inspection of the detections (Figure ??) we used the hazard-rate detectability model, which incorporates a shoulder in the detection function. A major assumption in the conventional distance sampling framework is that detection is perfect on the transect line, an assumption that has been shown to not hold for zebra mussel surveys (???). The mark-recapture component of the model is used to relax this assumption by using the removal surveys to estimate detectability on the transect line. In the mark-recapture model we assumed that each dive team collected data independently; thus, we used the point independence assumption described by Borchers et al. (2006). Point independence accounts for the effects of unmodeled covariates that can induce unexpected correlations between observers. We also included the effect of cluster size as a covariate influencing detection in the mark-recapture component of the model. We estimated detection using the mrds package (J. Laake et al. 2018), which accounts for removal survey and for the effect of distance on detectability. We obtain the transect-level detectability \hat{P} of detecting an zebra mussel cluster by integrating the detectability function over the transect strip-width (Buckland et al. 2001).

We obtained estimates of density from quadrat sampling by modifying equations 1 and 2 to assume that detection is perfect ($\hat{P} = 1$ and $\text{var}(\hat{P}) = 0$) and to account for counting individuals, not clusters ($\text{var}(E(s)) = 0$).

Finally, we looked at whether informally surveying an area can accurately represent the relative density of zebra mussels, Here, we used the initial surveys conducted at the start point of each transect for a fixed amount of time over an unknown area. We calculated the correlation between these count rates and the estimated densities on each transect.

Time budget analysis

- In order to determine the efficiency of each method, we quantified how much time each method required to complete a transect and the time required to move between transects.
- We modeled the time to conduct a transect survey by breaking up the total survey time into two components, the time to setup the transect and the time to perform the transect survey **I am not including time to do habitat survey because it is same across all types of survey, and we don't use habitat data in the analysis**. The first component, the time to setup the transect survey (hereafter referred to as the setup time), was the time required to place the transect line(s). We modeled the setup time using linear mixed-effects regression analysis with fixed-effect covariates of survey type (distance or quadrat, included as a qualitative variable), and the transect length (continuous variable). We also included lake as a random intercept term. The second component of the time budget was the time to perform the survey on each transect (hereafter referred to as the search time *word choice*). The

model structure of this component included fixed effects of transect length, an interaction between the type of survey and the detection rate (number of detections per meter of transect length) and a random slope term allowing detection time to vary with lake to account for heterogeneity that occurs due to lake conditions.

- We used the fitted models to predict the setup time and search time for a transect of length 30 m under using the estimated from the three lakes in this study.
- We modeled the time to move between consecutive transects as a function of distance between transects (hereafter, referred to as the movement time). We built the movement time model using a linear mixed-effects model to relate the time spent relocating between consecutively sampled transects to the distance between those transects for each lake in the initial density survey (for a total of six lakes). We included lake name as a random intercept term. All random effects models in this study were fit using lmer in the lme4 package (*citation*).
- We investigated the movement time for surveys with fewer transects by sequentially removing every i th transect and calculating the empirical distance between the new transect start locations. We iteratively removed every i th transect from the original 15 transects (for $i = 2$ to 8) in each of our lakes and calculated the distance between those sites. We used the travel time model described above to predict the the average time travel time for a study with n transects conducted on each lake. We modeled the travel time to conduct a survey on each lake with n transects.
- We combined the the setup time model (denoted at τ_S), search time model (τ_E), and movement time model (τ_M) to determine the expected number of transects that can be under a specific amount of time. We then determined the maximum number of transects that can be completed when a fixed total amount of time is available to complete the surveys (τ_T). The maximum number of transects that can be completed can be determined by minimizing the quantity $C(n)$ with respect to the number of transects surveyed, n , in the constrained optimization

$$C(n) = \tau_T - n \left(\tau_S + \tau_E(D, \hat{P}) + \tau_M(n) \right)$$

$$C(n) \geq 0.$$

remove equation?

- We determined the maximum number of surveys conducted, n_{\max} , obtained from solving the above optimization problem. We used n_{\max} and the empirical estimates of the standard error in lake density that were rescaled to a sample with n_{\max} transects, $\hat{\sigma} \sqrt{\frac{n}{n_{\max}}}$, where n is the number of transects from the original sample (given in Table 1). This rescaled standard error was used to compare the efficiency of the two different survey designs in each lake.

Simulation study

Look at tradeoffs in density and system size (travel time)... Options: * use empirical model to figure out time budget. * Use variance formula with Poisson/Tweedie/NB density function and ignore uncertainty in detection?. Plug into variance model.

Results

Surveys

- Quadrat surveys: The amount of area surveyed and number of zebra mussels counted in each lake are given in (Table 1). Estimates of density along with the corresponding standard errors are give in (Figure ??). We estimated that Lake Florida has the lowest density, Little Birch Lake had the highest density, and Lake Burgan had an intermediate density.

Table 1: Summary of survey results.

Design	Lake Florida			Lake Burgan			Little Birch La	
	Transects	Area surveyed	Detections	Transects	Area surveyed	Detections	Transects	Area surveyed
Quadrat	15	112.5	8	15	71.5	40	4	21.5
Distance	15	900.0	10	15	602.0	79	4	124.0

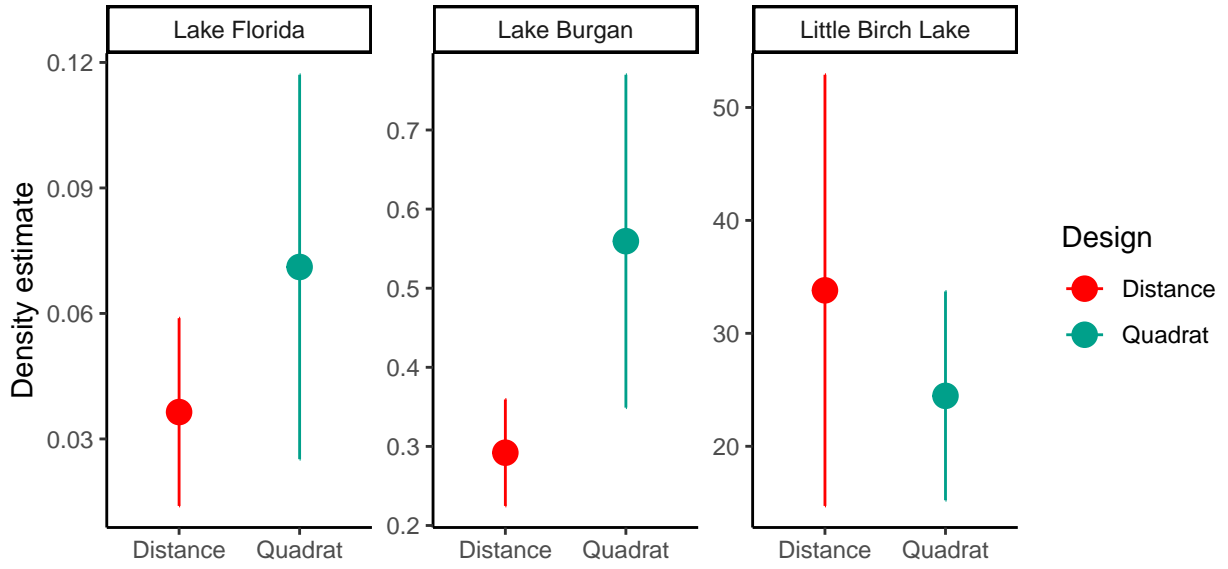


Figure 2: Density estimates and standard errors.

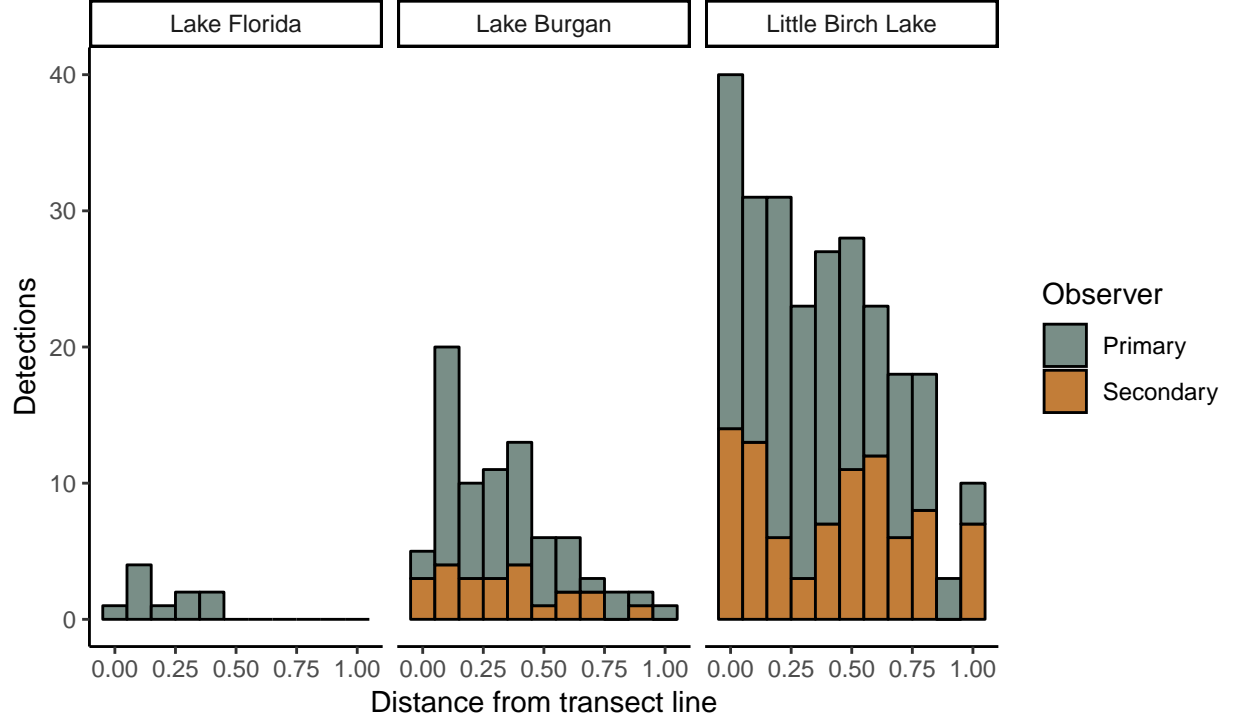


Figure 3: Stacked histogram showing the total number zebra mussel detections made by the primary and secondary diver in each lake. Distance bin widths are 0.1 m.

- Distance survey summary: The amount of area surveyed and number of zebra mussel detections made in each lake are given in (Table 1). In Lake Florida all detections were of single zebra mussels, in Lake Borgan the average cluster size was 1.18 (standard deviation = 0.13), and Little Birch Lake had the largest average cluster size with 7.83 (4.02).
- Our estimates of the detection functions indicated that the detectability of a cluster were similar between lakes with the estimated transect detectability in Lake Florida (reported as mean (standard error)) 0.31 (0.11), Lake Borgan 0.53 (0.07), and in Little Birch Lake 0.47 (0.07). Unfortunately, a lack of detections by observer 2 in Lake Florida likely leads us to be overconfident of our detection function. Despite the similarities in the overall estimates of detectability there was some differences in the shoulder width of the detection functions, potentially due to the higher vegetation levels in Little Birch Lake leading to a shorter shoulder (Figure ??).
- Estimated detection functions par estimates,
- Distance surveys: amount of area surveyed and number of detections in each lake, are given in Table 1. The average cluster size of detections in Lake Florida was , density estimate and standard error. Estimates of density along with the corresponding standard errors are give in (Figure ??). We determined that Lake Florida has the lowest density, Little Birch Lake had the highest density, and Lake Borgan had an intermediate density.

Time budget analysis

- In both the setup time and search time models we found that transect length and detection rate had postive effects on the time budget (Figure ??).
- We found that the transit time varied by transect distance, however
- We found that the transit time varied by transect distance, however

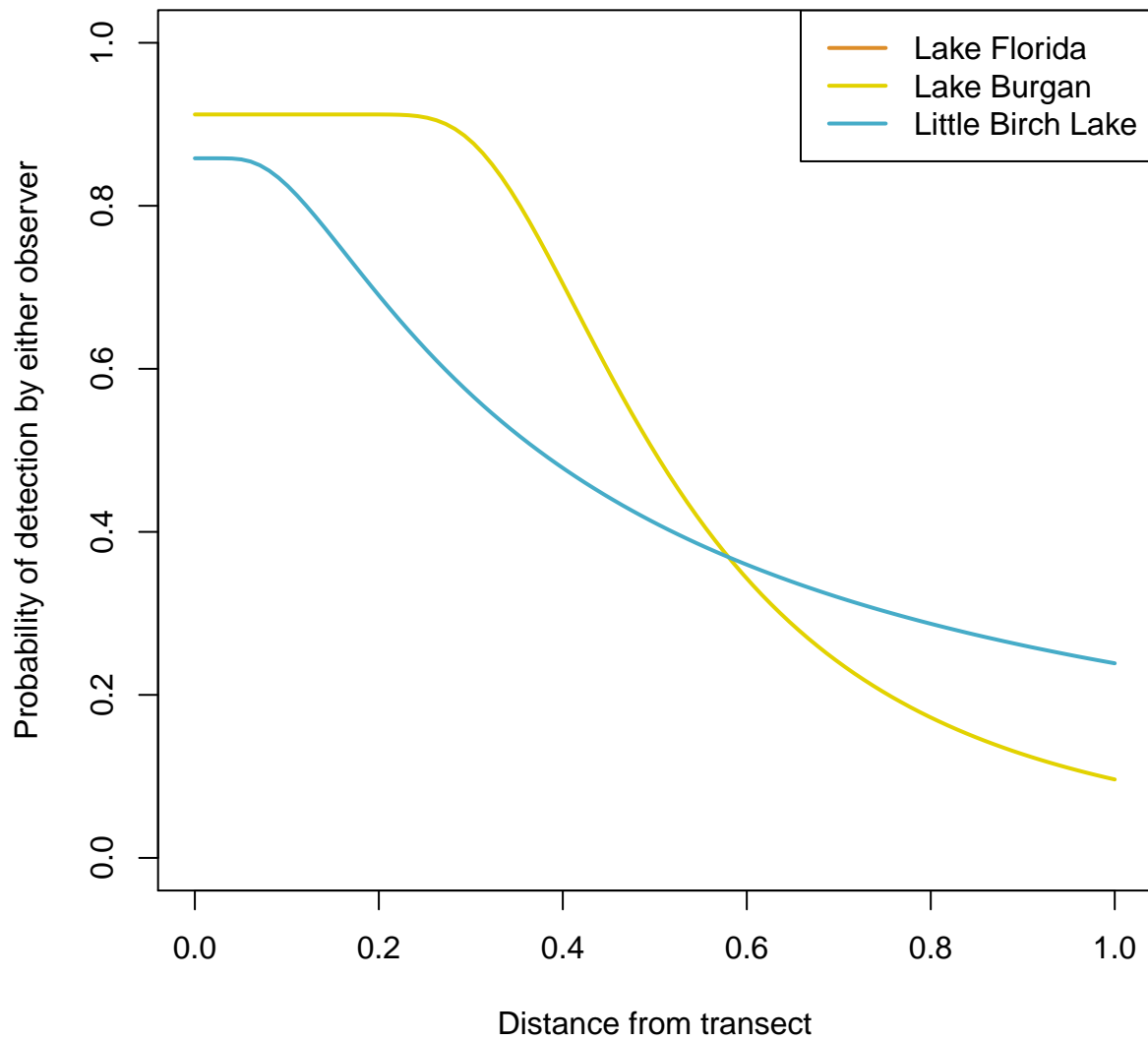


Figure 4: Average detection probability as a function of distance estimated using distance sampling for each lake. Consider removing figure...

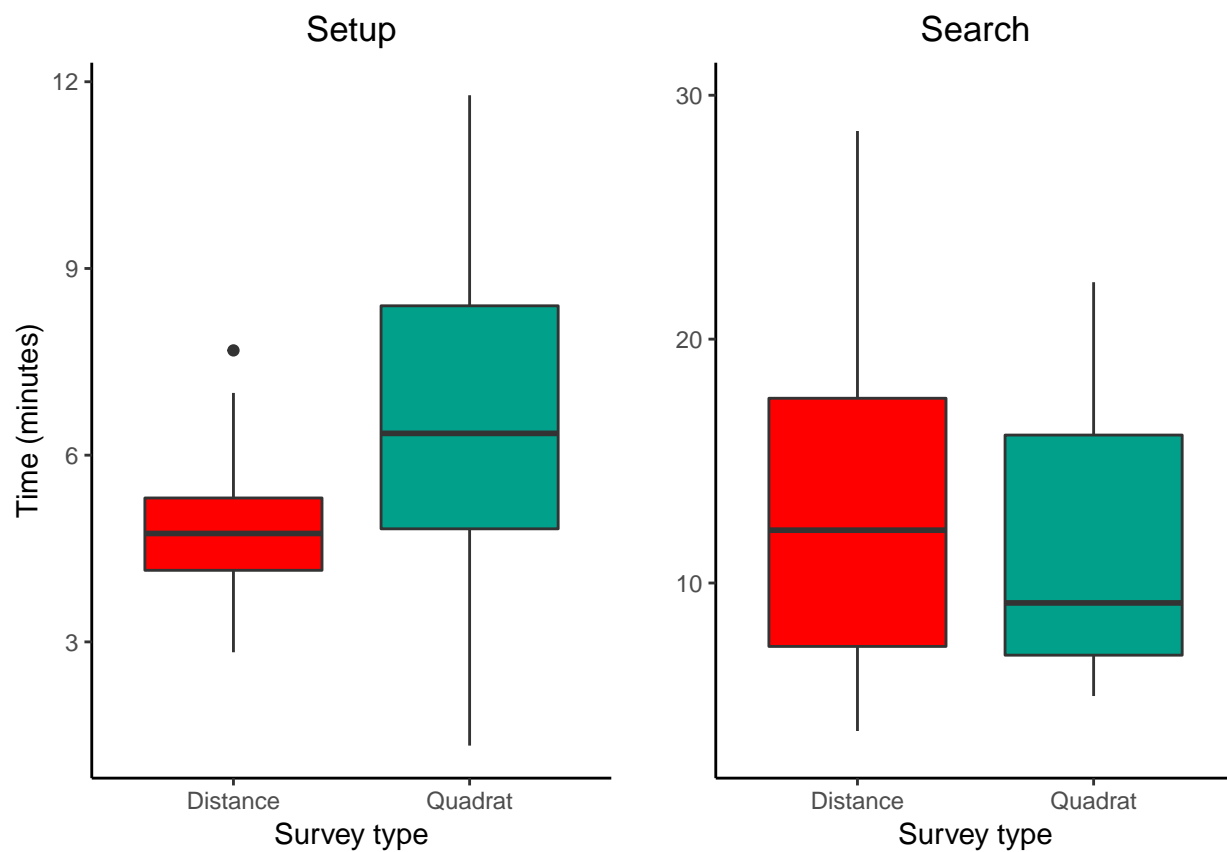


Figure 5: Empirical times for the setup and search. Consider removing figure...

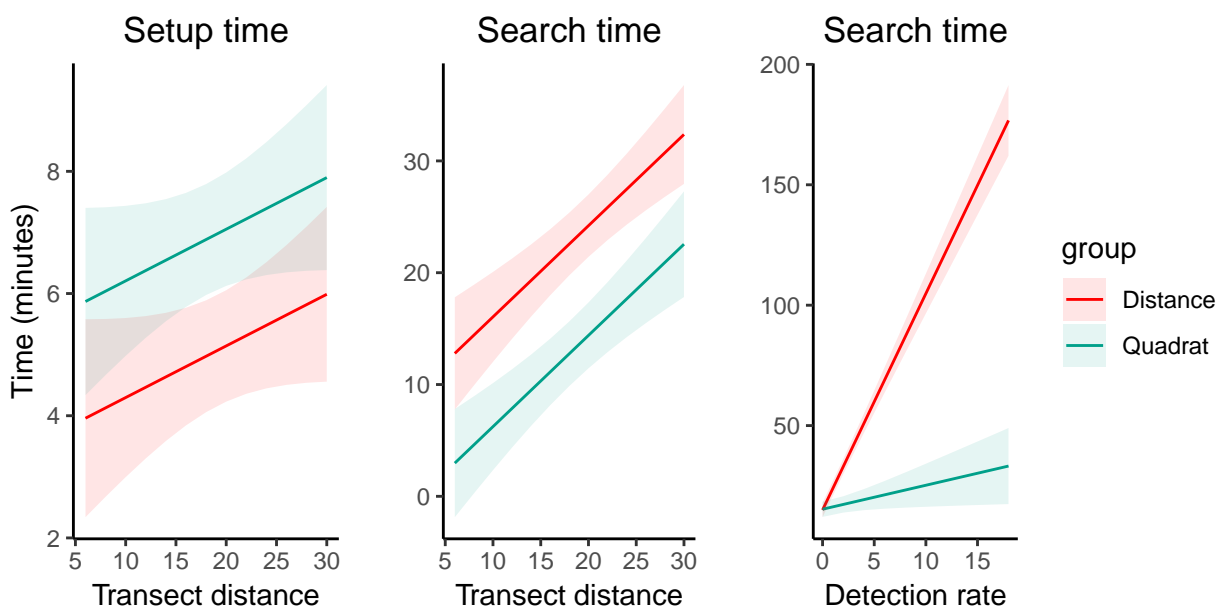


Figure 6: The impact of transect distance and detection rate on on the setup and search time. Predictions are of the fixed effects. Bands represent 95% confidence intervals of the predicted values.

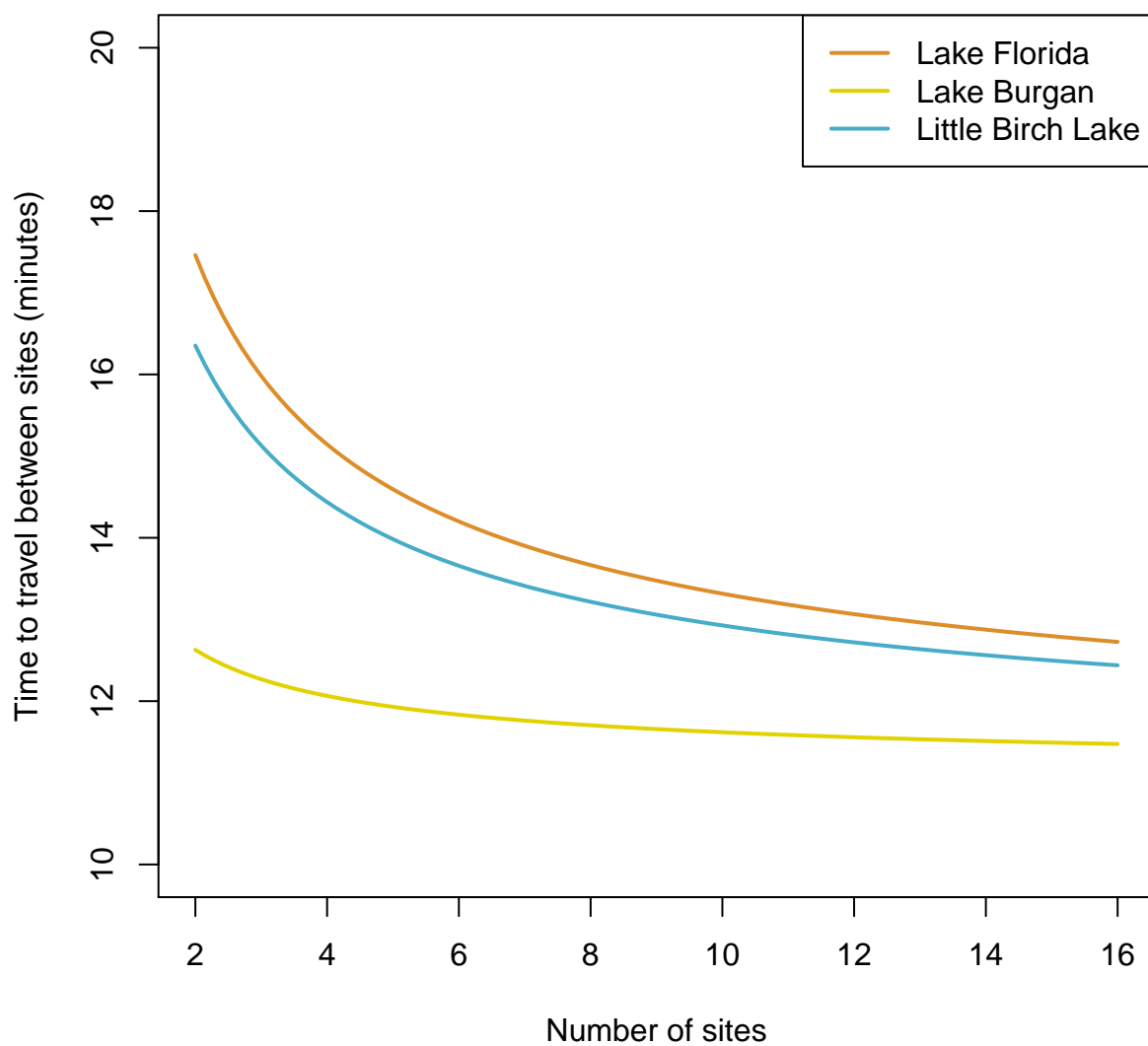


Figure 7: Predicted travel time at each lake.

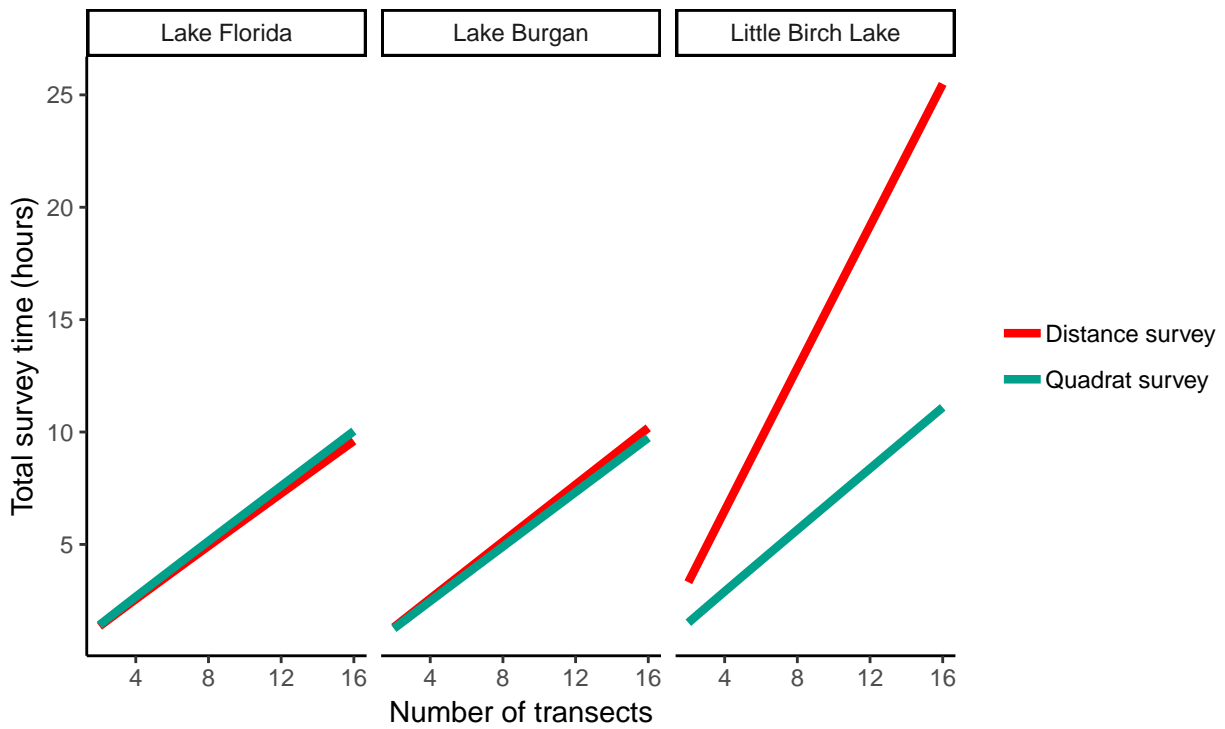


Figure 8: Maximum number of transects surveyed as a function of the total time available to conduct surveys.

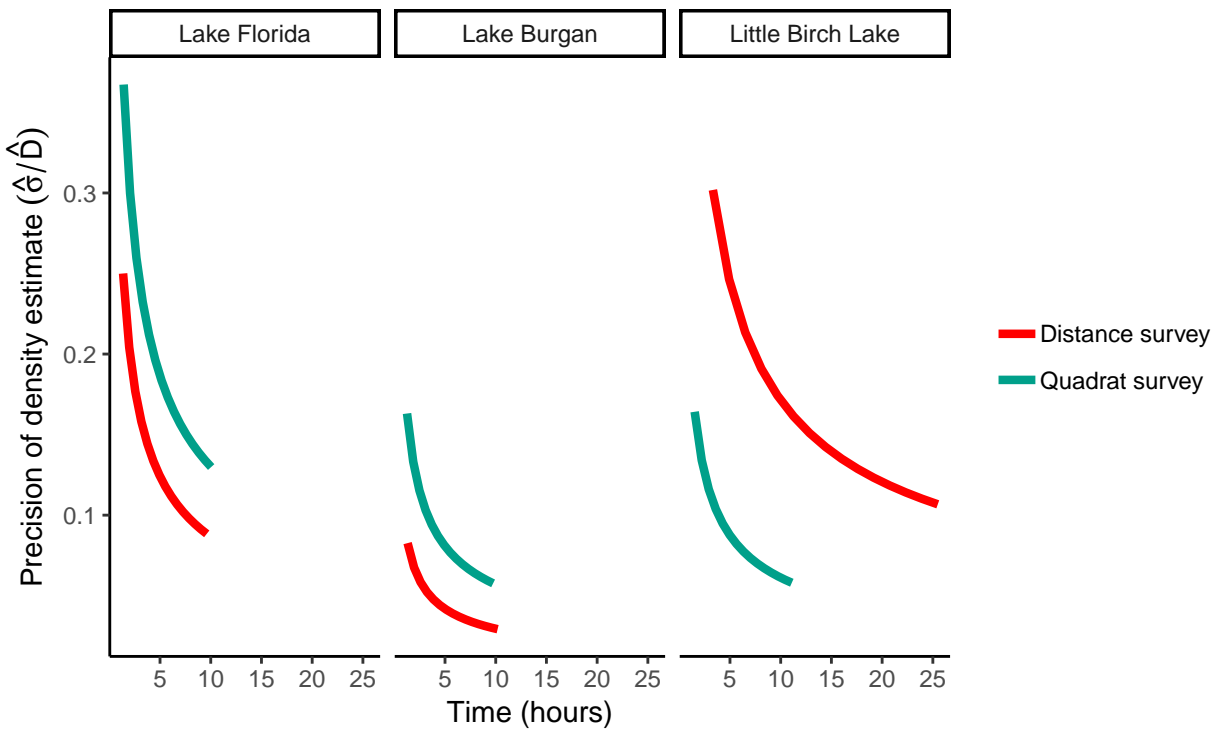


Figure 9: The predicted coefficient of variation in each lake. Fix x-axis.

Simulation study

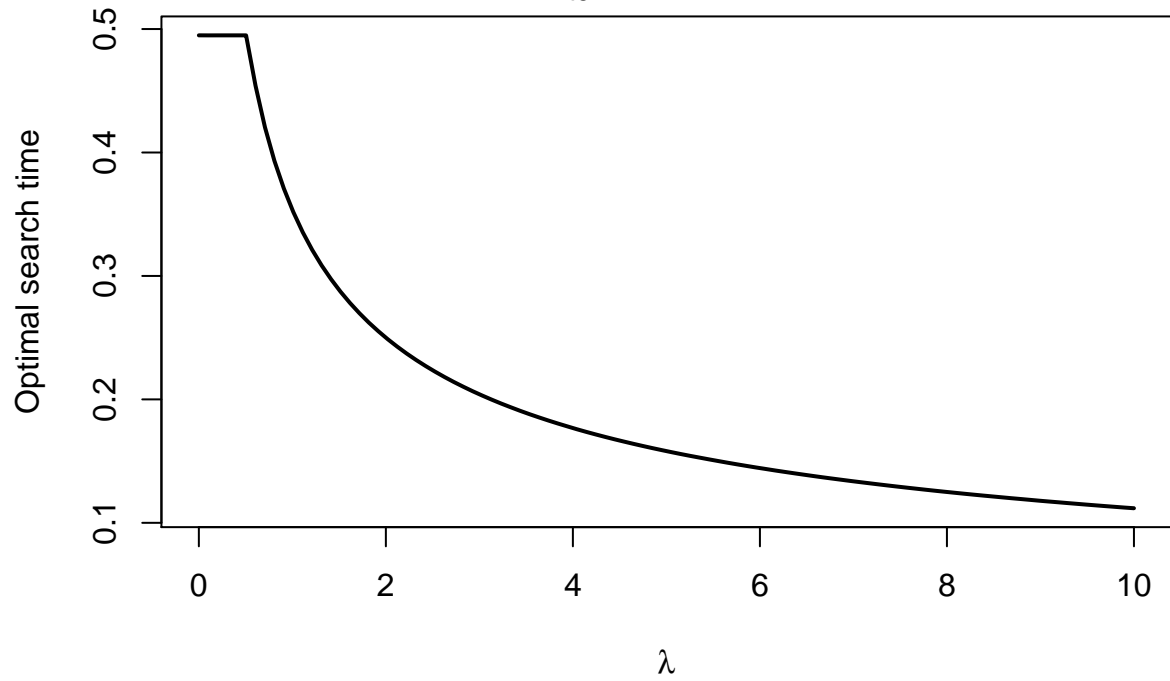
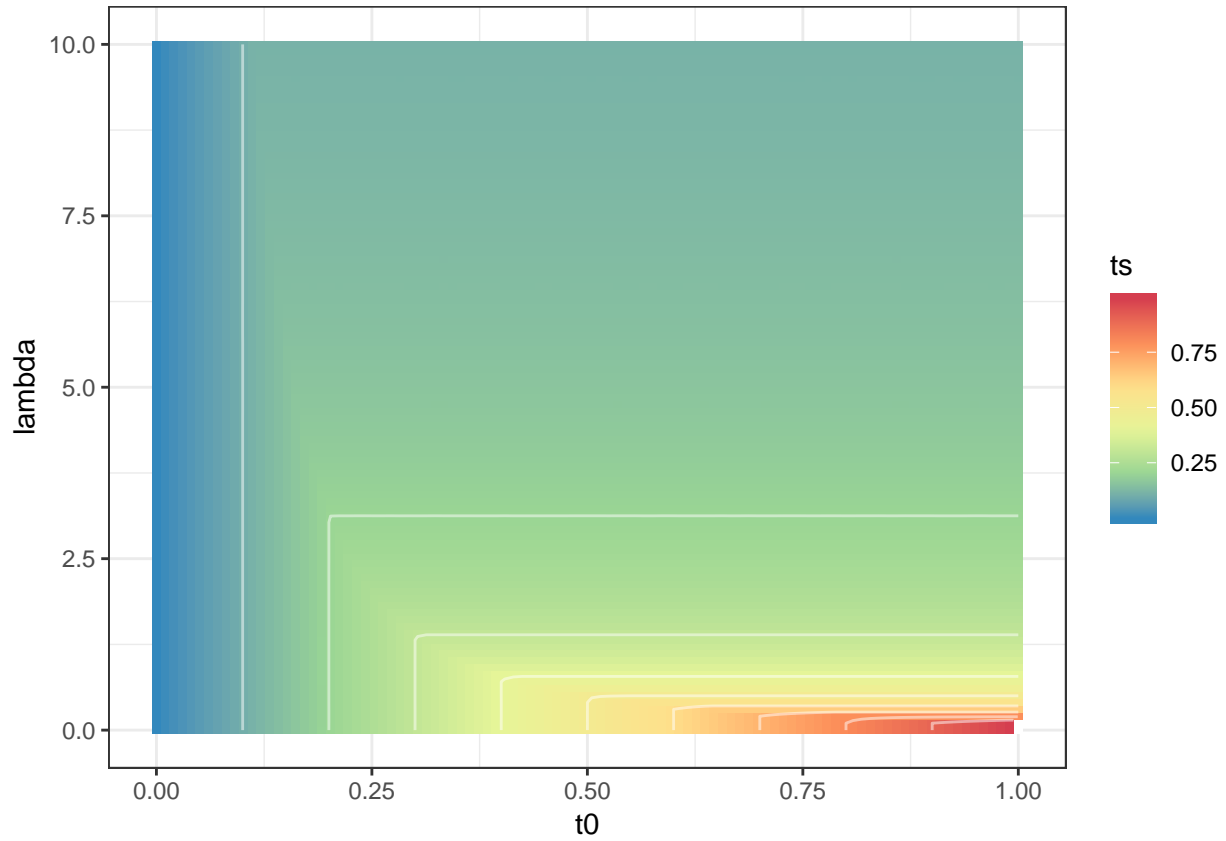
Look at tradeoffs in density and system size (travel time)... We used the results on the timebudget analysis to construct a more general model of the survey coverage/efficiency tradeoff. We start with the density variance model in equation (2), making a few simplifying assumptions. We will start by assuming that all transects are of the same length and cover and area, a . We also assume that the mean and variance of the counts follow a Tweedie distribution with mean $E(X) = \lambda n a$ and variance $\text{var}(X) = n b(a \lambda)^\phi$, where n is the number of transects completed, and ϕ is a value that typically falls between 1 and 2. We start by assuming that the detection probability is known ($\text{var}(\hat{P}) = 0$). Thus, $\text{var}(\hat{D}) = \frac{n b(a \lambda)^\phi}{(n a P)^2}$.

From the results in the previous section, we know that the total time to conduct a single survey can be modeled as $\tau = \tau_0 + \tau_S E(X)$ so the total number of transects that can be completed in a chunk of time, τ_T , is $n = \tau_T / (\tau_0 + \tau_S a \lambda)$. We can then plug in this to the variance expression above. Finally, we assume that detectability is a monotonic function of the search time with the following properties $\frac{dP}{d\tau_S} \equiv P' > 0$ and $\frac{d^2 P}{d\tau_S^2} \equiv P'' < 0$. This ensures that the function is monotonically increasing with diminishing returns on increases of effort.

We want to solve for the value of τ_S when the estimator variance is minimized $\frac{d \text{var}(\hat{D})}{d\tau_S} = 0$. Now solving:

$$\begin{aligned} \frac{d \text{var}(\hat{D})}{d\tau_S} = 0 &= \frac{d}{d\tau_S} \frac{(\tau_0 + \tau_S a \lambda)^{\phi-2}}{P^2} \\ 0 &= \frac{(\phi-2)a\lambda}{P^2} (\tau_0 + \tau_S a \lambda)^{\phi-3} - 2 \frac{P'}{P^3} (\tau_0 + \tau_S a \lambda)^{\phi-2} \\ 0 &= (\phi-2)a\lambda - 2 \frac{P'}{P} (\tau_0 + \tau_S a \lambda). \end{aligned}$$

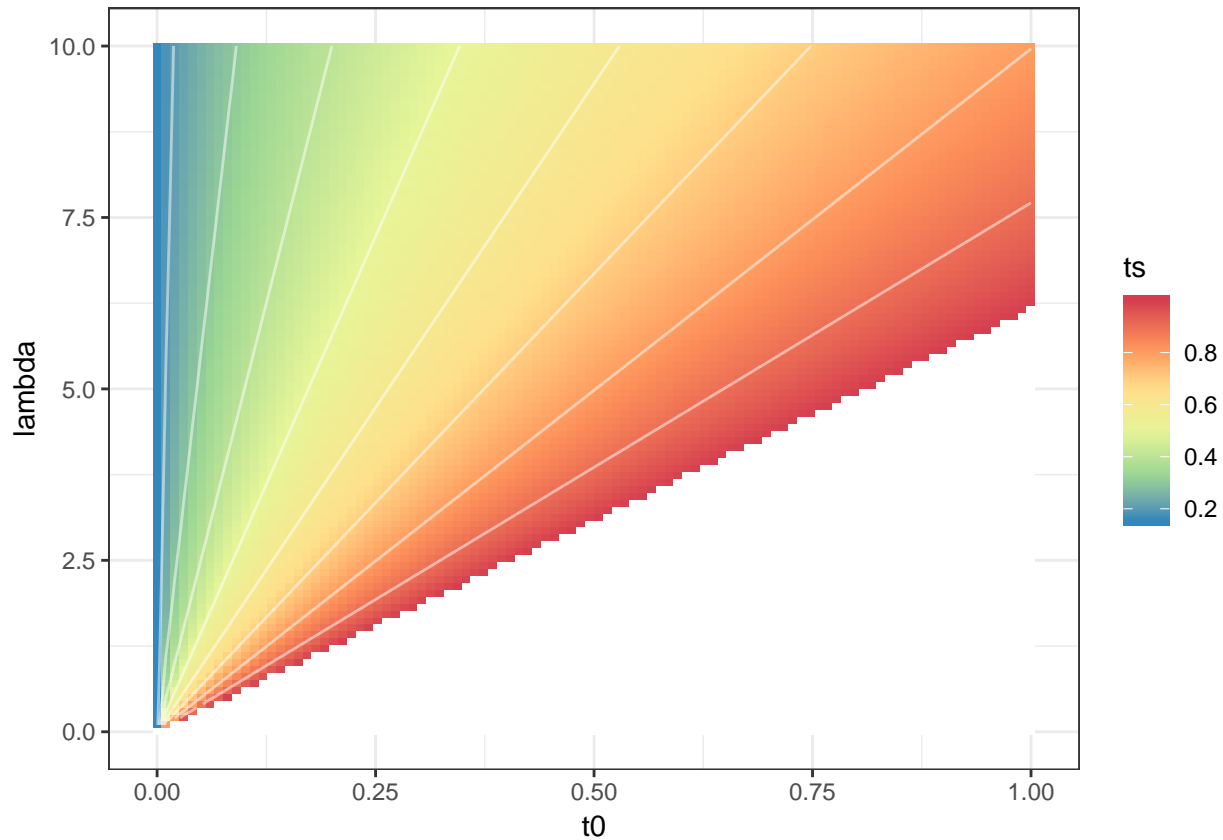
The results depends critically on the ratio P/P' , which will in general be a function of the search time τ_S . One approach is to model detection using the disc equation, $P = \tau_S / (\epsilon + \tau_S)$. In this case $P/P' = \frac{\tau_S}{\epsilon} (\epsilon + \tau_S)$. We can then get the solution for the optimal search time as,



We can do the same deal for a negative-binomial-like distribution with $\text{var}(A\hat{P}) = \alpha a\lambda + \beta a^2\lambda^2$.

Uncertainty in detectability

Now, we assume there is variation in the detectability, we assume that the estimator can be modeled with a beta-like distribution, $\text{var}(a\hat{P}) = \delta\hat{P}(1 - \hat{P})$, where δ is a scaling constant.



Discussion

- Theory tells that all else being equal, faster surveying methods are optimal as density increases. Here, we show that the practical

Borchers, D. L., J. L. Laake, C. Southwell, and C. G. M. Paxton. 2006. "Accommodating unmodeled heterogeneity in double-observer distance sampling surveys." *Biometrics* 62 (2): 372–78. doi:10.1111/j.1541-0420.2005.00493.x.

Buckland, S. T., D. R. Anderson, K. P. Burnham, J. L. Laake, D. L. Borchers, and L. Thomas. 2001. *Introduction to distance sampling: estimating abundance of biological populations*. 1st ed. Oxford: Oxford University Press.

Cook, R D, and J O Jacobson. 1979. "A design for estimating visibility bias in aerial surveys." *Biometrics* 35 (4): 735–42. doi:10.2307/2530104.

Ferguson, Jake M., Jessica B. Langebrake, Vincent L. Cannataro, Andres J. Garcia, Elizabeth A. Hamman, Maia Martcheva, and Craig W. Osenberg. 2014. "Optimal Sampling Strategies for Detecting Zoonotic Disease Epidemics." *PLoS Computational Biology* 10 (6): 1–26. doi:10.1371/journal.pcbi.1003668.

Holden, Matthew H., Jan P. Nyrop, and Stephen P. Ellner. 2016. "The economic benefit of time-varying surveillance effort for invasive species management." *Journal of Applied Ecology* 53 (3): 712–21. doi:10.1111/1365-

2664.12617.

Laake, Jeff, David Borchers, Len Thomas, David Miller, and Jon Bishop. 2018. “mrds: Mark-Recapture Distance Sampling.”