

# Vignette: Simulation studies with **ValidationExplorer** under a fixed-effort design type

Jacob Oram<sup>1</sup>   Katharine Banner<sup>1</sup>   Christian Stratton<sup>2</sup>   Kathryn M. Irvine<sup>3,\*</sup>

2024-11-19

## Abstract

Our vignette demonstrates the use of the **ValidationExplorer** package to conduct statistical simulation studies that explore the costs and inferential properties (e.g., near nominal coverage and/or minimal estimation error) of alternative validation designs. Our functions allow the user to specify a suite of candidate validation designs using either a stratified sampling procedure or a fixed-effort design type. An example of the former is provided in the manuscript entitled ‘ValidationExplorer: Streamlined simulations to aid informed management decisions using bioacoustic data in the presence of misclassification’, which was submitted to the Applications series of *Methods in Ecology and Evolution*. In this vignette, we provide an additional example of data simulation, model fitting, and visualization of simulation results when using a fixed-effort design type. Our demonstration here is intended to aid researchers and others to tailor a validation design that provides useful inference while also ensuring that the level of effort meets cost constraints.

<sup>1</sup> Department of Mathematical Sciences, Montana State University, Bozeman, MT, USA

<sup>2</sup> Department of Mathematics and Statistics, Middlebury College, Middlebury, VT, USA

<sup>3</sup> U.S. Geological Survey, Northern Rocky Mountain Science Center, Bozeman, MT, USA

\* Correspondence: Kathryn M. Irvine <kirvine@usgs.gov>

---

**Disclaimer:** This draft manuscript is distributed solely for the purposes of scientific peer review. Its content is deliberative and pre-decisional, so it must not be disclosed or released by reviewers. Because the manuscript has not yet been approved for publication by the U.S. Geological Survey (USGS), it does not represent any official USGS funding or policy. Any use of trade, firm, or product names is for descriptive purposes only and does not imply endorsement by the U.S. Government.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Conducting a simulation study with a fixed effort design type</b>	<b>3</b>
2.1	Step 0: Define measurable objectives and constraints . . . . .	3
2.2	Step 1: Installing and loading required packages . . . . .	4
2.3	Step 3: Simulate data {dataSimulation} . . . . .	5
2.4	Step 4: MCMC_tuning . . . . .	9
2.5	Step 5: Fit models to simulated data . . . . .	11
2.6	Step 6: Visualize simulations . . . . .	12
<b>3</b>	<b>Conclusion</b>	<b>13</b>
	<b>References</b>	<b>14</b>

# 1 Introduction

Automated recording units (ARUs) provide one of the main data sources for many contemporary monitoring programs that aim to provide inference about status and trends for assemblages of species (Loeb et al. 2015). As described in “**ValidationExplorer**: Streamlined simulations to aid informed management decisions using bioacoustic data in the presence of misclassification” (hereafter, “the main text”), substantial practical interest lies in identifying cost-effective validation designs that will allow a monitoring program to obtain its measurable objectives. We believe that statistical simulation studies are a valuable tool for evaluating the relative merits of candidate validation designs prior to gathering and validating ARU data, and it is our goal for **ValidationExplorer** to provide those tools.

This vignette provides detailed demonstrations of the **ValidationExplorer** package under a fixed-effort design. As described in the main text, a validation design is composed of two parts: a random mechanism for selecting observations to be manually reviewed by experts (“validated”), and a percentage or proportion that controls the sample size. We refer to the random mechanism as the *design type* and the proportion as the level of validation effort (LOVE). In our example, we assume a fixed effort (the design type) validation design, under which,  $x\%$  of recordings obtained from the first visit to a site are validated by experts. The level of validation effort is controlled by the value of  $x$ . In the following section we consider five possible LOVEs and show an example of a fixed effort design.

## 2 Conducting a simulation study with a fixed effort design type

### 2.1 Step 0: Define measurable objectives and constraints

Recall from the main text that the first step – before opening R and loading **ValidationExplorer** – is to identify and write down the set of measurable objectives that the data will be used for. Suppose that, for this example, the measurable objectives and cost constraints are the same as those Section 3 of the main text:

- The measurable objective is to estimate relative activity parameters (denoted as  $\lambda_k$ ) for each species with estimation error less than 1 call per night and with the expected width of 95% posterior intervals less than 3 calls per night.
- The monitoring program can pay their expert bat biologists to validate at most 4000 recordings. WE

Table 1: Prior knowledge of relative activity rates and occurrence probabilities for the six bat species of interest. These values will be used to simulate data.

Species	$\psi$	$\lambda$
<i>Eptesicus fuscus</i> (EPFU)	0.6331	5.9347
<i>Lasiurus cinereus</i> (LACI)	0.6122	4.1603
<i>Lasionycteris noctivagans</i> (LANO)	0.8490	14.2532
<i>Myotis californicus</i> (MYCA)	0.6972	6.1985
<i>Myotis ciliolabrum</i> (MYCI)	0.2365	11.8649
<i>Myotis evotis</i> (MYEV)	0.7036	2.4050

make the assumption that all recordings are approximately the same cost to validate (i.e., rare species autoIDs are not necessarily more expensive or time consuming than extremely common ones).

Suppose further that the species assemblage is the same as that in the main text. That is, we have six species of interest that co-occur. The existing prior knowledge (perhaps from another study) about the relative activity rates and occurrence probabilities for each species are summarized in Table 1.

## 2.2 Step 1: Installing and loading required packages

Once measurable objectives and constraints are clearly defined, the next step is to load the required packages. For first time users, it may be necessary to install a number of dependencies, as shown by Table 1 in the main text. If you need to install a dependency or update the version, run the following, with `your_package_name_here` replaced by the name of the package:

```
install.packages("your_package_name_here")
```

After installing the necessary packages, load these libraries by calling

```
library(tidyverse)
library(nimble)
library(coda)
library(rstan)
library(parallel)
library(here)
```

Finally, install and load `ValidationExplorer` by running

```
devtools::install_github(repo = "j-oram/ValidationExplorer")
library(ValidationExplorer)
```

Note: to knit this vignette, you may also need to install the `kableExtra` package.

## 2.3 Step 3: Simulate data {dataSimulation}

The first step in a simulation study is to simulate data under each of the candidate validation designs, which is accomplished with the `simulate_validatedData` function in `ValidationExplorer`. We begin by assigning values for the number of sites, visits, species, as well as the parameter values in Table 1:

```
# Set the number of sites, species and visits
nsites <- 30
nspecies <- 6
nvisits <- 4

psi <- c(0.6331, 0.6122, 0.8490, 0.6972, 0.2365, 0.7036)
lambda <- c(5.9347, 4.1603, 14.2532, 6.1985, 11.8649, 2.4050)
```

Additionally, `simulate_validatedData` requires that the user supply misclassification probabilities in the form of a matrix, subject to the constraint that rows in the matrix sum to one. An easy way to simulate a matrix of probabilities that meet these criteria is to leverage the `rdirch` function from the `nimble` package:

```
# Simulate a hypothetical confusion matrix
set.seed(10092024)
Theta <- t(apply(diag(29, nspecies) + 1, 1, function(x) {nimble::rdirch(alpha = x)}))
print(Theta)
```

```
##           [,1]      [,2]      [,3]      [,4]      [,5]      [,6]
## [1,] 0.953670138 0.015886205 0.01346950 0.0005918829 0.0153509623 0.001031315
## [2,] 0.028614016 0.914794071 0.01941594 0.0205164961 0.0065938947 0.010065584
## [3,] 0.013669368 0.006878823 0.82604797 0.0969673863 0.0005128495 0.055923603
## [4,] 0.042623395 0.028935593 0.01156918 0.8327210780 0.0599628913 0.024187867
## [5,] 0.005590593 0.034119356 0.02036092 0.0639122608 0.8635266207 0.012490252
## [6,] 0.003122368 0.003831873 0.04087851 0.0742637611 0.0632647264 0.814638762
```

Note that the above definition of `Theta` places high values on the diagonal of the matrix, corresponding to a high probability of correct classification. To lower the diagonal values, change the specification of `diag(29, nspecies)` to a smaller value. For example:

```
another_Theta <- t(apply(diag(5, nspecies) + 1, 1, function(x) {nimble::rdirch(alpha = x)}))
print(another_Theta)
```

```
##           [,1]      [,2]      [,3]      [,4]      [,5]      [,6]
## [1,] 0.55731323 0.142858110 0.053100053 0.003946041 0.04273763 0.200044944
## [2,] 0.13437231 0.499939788 0.021421570 0.014098339 0.02959848 0.300569509
## [3,] 0.10172739 0.007228338 0.515224477 0.155895124 0.13383729 0.086087378
## [4,] 0.18412846 0.210913704 0.026614757 0.479425606 0.09220850 0.006708972
## [5,] 0.27654377 0.142140440 0.004483512 0.068078975 0.32570927 0.183044035
## [6,] 0.06086223 0.043479478 0.123594690 0.083448153 0.15034635 0.538269101
```

`another_Theta` has lower values on the diagonal, and greater off-diagonal values (i.e., higher probability of misclassification). If you have specific values you would like to use for the assumed classification probabilities, these can be supplied manually:

```
manual_Theta <- matrix(c(0.9, 0.05, 0.01, 0.01, 0.02, 0.01,
                        0.01, 0.7, 0.21, 0.05, 0.02, 0.01,
                        0.01, 0.01, 0.95, 0.01, 0.01, 0.01,
                        0.05, 0.05, 0.03, 0.82, 0.04, 0.01,
                        0.01, 0.015, 0.005, 0.005, 0.95, 0.015,
                        0.003, 0.007, 0.1, 0.04, 0.06, 0.79),
                      byrow = TRUE, nrow = 6)
```

If you define the classifier manually, make sure the rows sum to 1 by running

```
all(rowSums(manual_Theta) == 1) # want this to return TRUE
```

```
## [1] TRUE
```

```
# If the above returns FALSE, see which one is not 1:
rowSums(manual_Theta)
```

```
## [1] 1 1 1 1 1 1
```

With the required inputs defined, we can simulate data using the following code:

```
sim_data <- simulate_validatedData(
  n_datasets = 50,
  nsites = nsites,
  nvisits = nvisits,
  nspecies = nspecies,
  design_type = "FixedPercent",
  scenarios = c(0.05, 0.1, 0.15, 0.3, 0.5),
  psi = psi,
```

```

lambda = lambda,
theta = Theta,
save_datasets = FALSE, # default value is FALSE
save_masked_datasets = FALSE, # default value is FALSE
directory = here::here("Vignette", "Fixed_Effort")
)

```

Note that we specified the design type through the argument `design_type = "FixedPercent"`, with the possible scenarios defined by `scenarios = c(0.05, 0.1, 0.15, 0.3, 0.5)`. These two arguments specify the set of alternative validation designs we will compare in our simulation study. Under the first validation design, 5% of recordings from the first visit to a site are validated, while in the second validation design 10% of recordings from the first visit to a site are validated, and so on.

To understand the output from `simulate_validatedData`, we can investigate `sim_data`. The output is a list, containing three objects:

- **full\_datasets**: A list of length `n_datasets` with unmasked datasets (i.e., the datasets if all recordings were validated so that every recording has an autoID and a true species label). We opted to not save these datasets by setting `save_datasets = FALSE`. If `save_datasets = TRUE`, then these will be saved individually in `directory` as `dataset_n.rds`, where `n` is the dataset number.
- **zeros**: A list of length `n_datasets` containing the true species- autoID combinations that were never observed at each site visit. For example, if, in dataset 10, there were no calls from species 1 that were classified as 3 on visit 4 to site 30, then the 10th entry of this list would contain a dataset with a row corresponding to `site = 30`, `visit = 4`, `true_spp = 1`, `id_spp = 3`, with `count = 0`. These zeros are necessary for the model to identify occurrence probabilities and relative activity rates. If `save_datasets = TRUE`, the zeros for each dataset will also be saved in `directory` individually as `zeros_in_dataset_n.rds`, where `n` is the dataset number.
- **masked\_dfs**: A nested list containing each dataset masked under each scenario. For example, `masked_dfs[[9]][[27]]` contains dataset 27, assuming validation scenario 9. If `save_masked_datasets = TRUE`, then each dataset/scenario combination is saved individually in `directory` as `dataset_n_masked_under_scenario_s.rds`, where `n` is the dataset number and `s` is the scenario number.

Examples of each are given below:

```
full_dfs <- sim_data$full_datasets
head(full_dfs[[3]]) # Dataset number 3 if all recordings were validated
```

```
## # A tibble: 6 x 10
## # Groups:   site, visit [1]
##   site visit true_spp id_spp lambda  psi  theta    z count  Y.
##   <int> <int>   <int>  <int>  <dbl> <dbl> <dbl> <int> <int> <int>
## 1     1     1     1      1     5.93 0.633 0.954     1     4    24
## 2     1     1     1      1     5.93 0.633 0.954     1     4    24
## 3     1     1     1      1     5.93 0.633 0.954     1     4    24
## 4     1     1     1      1     5.93 0.633 0.954     1     4    24
## 5     1     1     4      2     6.20 0.697 0.0289     1     1    24
## 6     1     1     3      3    14.3 0.849 0.826     1    10    24
```

```
zeros <- sim_data$zeros

# The site-visit-true_spp-autoID combinations that were never observed in
# dataset 3. Notice that count = 0 for all rows!
head(zeros[[3]])
```

```
## # A tibble: 6 x 10
## # Groups:   site, visit [1]
##   site visit true_spp id_spp lambda  psi  theta    z count  Y.
##   <int> <int>   <int>  <int>  <dbl> <dbl> <dbl> <int> <int> <int>
## 1     1     1     2      1     4.16 0.612 0.0286     0     0    24
## 2     1     1     3      1    14.3 0.849 0.0137     1     0    24
## 3     1     1     4      1     6.20 0.697 0.0426     1     0    24
## 4     1     1     5      1    11.9 0.236 0.00559     0     0    24
## 5     1     1     6      1     2.40 0.704 0.00312     1     0    24
## 6     1     1     1      2     5.93 0.633 0.0159     1     0    24
```

```
masked_dfs <- sim_data$masked_dfs

# View dataset 3 subjected to the validation design in scenario 4:
# randomly select and validate 30% of recordings from the first visit
# to each site
head(masked_dfs[[4]][[3]])
```

```
## # A tibble: 6 x 11
## # Groups:   site [1]
##   site visit true_spp id_spp lambda  psi  theta    z count  Y. scenario
##   <int> <int>   <int>  <int>  <dbl> <dbl> <dbl> <int> <int> <int>   <int>
## 1     1     1     NA      1     5.93 0.633 0.954     1     4    24         4
## 2     1     1     1      1     5.93 0.633 0.954     1     4    24         4
## 3     1     1     NA      1     5.93 0.633 0.954     1     4    24         4
## 4     1     1     1      1     5.93 0.633 0.954     1     4    24         4
## 5     1     1     NA      2     6.20 0.697 0.0289     1     1    24         4
## 6     1     1     NA      3    14.3 0.849 0.826     1    10    24         4
```

For most simulations, it will be useful to summarize the number of recordings that are validated under a given validation design and scenario. This can be accomplished using the `summarize_n_validated` function, which



outputs a vector containing the average number of recordings validated in a dataset under each scenario. Note that in our example, there are 8 scenarios (each contained in a row of the `val_scenarios` dataframe), meaning that the output is of length 8:

```
summarize_n_validated(data_list = sim_data$masked_dfs, theta_scenario = "1", scenario_names = 1:5)
```

```
## # A tibble: 5 x 3
##   theta_scenario scenario n_validated
##   <chr>           <int>      <dbl>
## 1 1               1         54.8
## 2 1               2         93.9
## 3 1               3        135.
## 4 1               4        255.
## 5 1               5        410.
```

## 2.4 Step 4: MCMC\_tuning

Running a complete simulation study can be time consuming. In an effort to help users improve the efficiency of their simulations, we provide the `tune_mcmc` function, which provides the user with information about the warmup and number of iterations required for the MCMC to reach approximate convergence. This function takes in a masked dataset and the corresponding zeros, fits a model to these data, and outputs an estimated run time for 10,000 iterations, as well as the estimated number of required warmup and total iterations.

As in the main text, we use a dataset with the lowest number of validated recordings, for which we expect the greatest number of iterations. In our example, this is scenario 1, in which an average of  $\approx 55$  recordings are validated per dataset. When we run `tune_mcmc` with dataset 39 under scenario 1, we get an error message, telling us that convergence was not reached in under 10,000 iterations. We have several options:

1. We could increase the number of iterations above 10,000 – perhaps to 20,000 and settle for a longer run time of the simulation study.
2. We could take this as a sign that the level of effort (55 calls validated per dataset of size  $\approx 3300$ ) is insufficient to identify model parameters. In this case, this scenario should not be considered.

In our experience fitting these models, the second option seems to often be the case. We adopt this approach here and run `tune_mcmc` again under scenario 2:

```
tune_list <- tune_mcmc(dataset = sim_data$masked_dfs[[2]][[39]], zeros = sim_data$zeros[[39]])
```

```
## [1] "Fitting MCMC in parallel ... this may take a few minutes"
```

We can see from the output that minimum number of iterations is 2000 with a warmup of 1000:

```
tune_list$min_iter
```

```
## [1] 1500
```

```
tune_list$min_warmup
```

```
## [1] 500
```

If we were to specify 10,000 iterations per dataset, each model fit would take approximately 2 minutes.

```
tune_list$max_iter_time
```

```
## Time difference of 2.077669 mins
```

This may seem insignificant, but over the course of an entire simulations study with 5 scenarios  $\times$  50 datasets, that corresponds to around 8 hours of run time.

To understand the MCMC in greater detail, we can examine the `convergence_matrix` output

```
tune_list$convergence_matrix
```

```
##               warmup = 500 warmup = 1000 warmup = 2000 warmup = 5000
## post-warmup iters = 1000           1           1           0           0
## post-warmup iters = 2000           1           1           1           1
## post-warmup iters = 3000           1           1           1           1
## post-warmup iters = 4000           1           1           1           1
## post-warmup iters = 5000           1           1           1           1
```

We emphasize that these results are from a single model fit; they are supplied only as guidelines, so it may be a good idea to increase the number of iterations above the minimum values output from `tune_mcmc`. This is what we do in the following section, where we specify that there should be 3000 iterations with the first 1500 discarded as warmup.

## 2.5 Step 5: Fit models to simulated data

The following code block shows our use of `run_sims` to fit models to simulated data.

```
# reduce the size of the list since validation scenario 1 is too low:
# we want scenarios 2-5 for fitting
masked_dfs_reduced <- masked_dfs[2:5]

sims_output <- run_sims(
  data_list = masked_dfs_reduced,
  zeros_list = sim_data$zeros,
  DGVs = list(lambda = lambda, psi = psi, theta = Theta),
  theta_scenario_id = "FE", # for "fixed effort"
  parallel = TRUE,
  niter = 3000,
  nburn = 1500,
  thin = 1,
  save_fits = TRUE,
  save_individual_summaries_list = FALSE,
  directory = here::here("Vignette", "Fixed_Effort")
)
```

```
## Beginning scenario 1.
```

```
## 2024-11-19 13:27:15.337427
```

```
##      |
```

```
## Beginning scenario 2.
```

```
## 2024-11-19 14:15:43.606391
```

```
##      |
```

```
## Beginning scenario 3.
```

```
## 2024-11-19 15:03:38.367523
```

```
##      |
```

```
## Beginning scenario 4.
```

```
## 2024-11-19 15:53:53.302172
```

```
##      |
```

## 2.6 Step 6: Visualize simulations

Once the simulation study is complete, you can visualize the results using several functions. `visualize_parameter_group` and `visualize_single_parameter`. These functions ensure that only converged models are included in the visualization. `visualize_parameter_group` is useful for examining an entire set of parameters, such as all relative activity parameters. The set of expected inputs for `visualize_parameter_group` are:

- `sim_summary`: A dataframe in the format of the summaries output by `run_sims`. Column names must match those of the `run_sims` output.
- `pars`: The name of the parameter “group” to be visualized (e.g, “psi”, “lambda” or “theta”).
- `theta_scenario`: The  $\Theta$  classifier ID.
- `scenarios`: Which scenarios to visualize?
- `convergence_threshold`: What value should  $\hat{R}$  be below to be considered “converged”? Default value is 1.1. This value matters because only model fits where all parameter values are below the `convergence_threshold` are used for visualization.

We can visualize the inference for the relative activity parameters in the first three scenarios in our simulation study above by running the code below. Note that we have set `convergence_threshold` to be unrealistically high for illustration purposes. Typically, the default value of 1.1 is as high of an  $\hat{R}$  value as possible for the MCMC chains to be considered “converged”.

```
visualize_parameter_group(sim_summary = sims_output,  
                          pars = "lambda",  
                          theta_scenario = 1,  
                          scenarios = 1:4, # recall that scenario 1 in the figure is actually scenario 2  
                          convergence_threshold = 1.1)
```

The features of the plot are as follows:

- Facet grids: parameters
- X-axis: Manual verification scenario
- y-axis: parameter values
- Small grey error bars: 95% posterior interval for an individual model fit where all parameters were below `convergence_threshold`.
- Colored error bars: average 95% posterior interval across all converged models under that scenario.
- Color: Coverage, or the rate at which 95% posterior intervals contain the true data-generating parameter value.

- Black dots: the true value of the parameter
- Red dots: average posterior mean

If you would like to visualize a single parameter, use `visualize_single_parameter`, which takes the same arguments as the previous visualization function:

```
visualize_single_parameter(sims_output, par = "theta[2, 1]",
                           theta_scenario = 1,
                           scenarios = 1:3,
                           convergence_threshold = 1.2)
```

Note that the scale of the y-axis is free to change from one visualization to the next. Additionally, if no datasets show evidence of convergence (i.e., no fitted models have  $\hat{R} \leq c$  for all parameters, where  $c$  is the specified convergence threshold) under a given scenario, the scenario will not appear on the x-axis.

### 3 Conclusion

We have demonstrated the use of the `ValidationExplorer` package. Note that with any simulation study results are conditional on the settings and assumptions. In the case of the count-detection model framework, the assumptions are

- The occurrence of species within a site are independent; the presence of one species carries no information about the presence or absence of another.
- For any one species, its occurrence at one location is independent of its occurrence at any other location (independence across sites).
- Visits to a site (i.e., detector nights) are independent.
- Recordings within the same site-visit are independent. This holds regardless of whether a recording is validated or remains ambiguous (only has an autoID label).
- The count of recordings generated by a species in a single detector night is a Poisson random variable.
- All species in the assemblage can co-occur and are capable of being confused (no structural zeros in the classification matrix).

The settings include the number of datasets used in the simulations, the assumed characteristics of the species assemblage (i.e., the values for each  $\psi_k$  and  $\lambda_k$ ) and classifier ( $\Theta$ ), and the number of sites and visits. Our hope is that the `ValidationExplorer` provides a useful tool for assessing the merits of various validation designs, so that effort can be thoughtfully assigned based on program goals and monitoring

objectives. Ultimately, we hope that this software streamlines the bat acoustic workflow and allows for efficient processing of information.

## References

Loeb, Susan C., Thomas J. Rodhouse, Laura E. Ellison, Cori L. Lausen, Jonathan D. Reichard, Kathryn M. Irvine, Thomas E. Ingersoll, et al. 2015. “A Plan for the North American Bat Monitoring Program (NABat).” SRS-208. USDA Forest Service.