

# Sampling regimes for rare species with imperfect detection

*Jens Åström*

*2018-10-26*

## Intro

I here explore some effects of rare occurrence and imperfect detection in species observation and explore alternative survey regimes. This is mostly relevant for rare species, and not for the general survey of insects. Early detection of alien species is another use case.

The objective is to identify suitable sampling strategies for a couple of different scenarios, where we maximise the chances of detecting a species at least once while considering the economic costs. The basic principle here is that a species can be present in a sampling location with a given probability of occurrence. The observers then have a given chance (probability) of detecting a species on a visit, if it is present. The total probability of detecting a species will be dependent on both the probability of occurrence and the probability of detection, how many locations are visited, and how many visits we make at each location.

If  $\psi$  represents the occurrence probability and  $\theta$  is the detection probability, the probability of observing the species at least once in  $J$  locations visited  $K$  times is:

$$probObserv = 1 - (1 - \psi * (1 - (1 - \theta)^K))^J.$$

The costs can be assumed to scale roughly linearly to the total amount of visits, i.e.  $totalCost = J * K * c$ , where  $c$  is the cost of one visit/survey.

## First look

We can explore how the overall probability of observing the species, and the associated costs depend on  $\psi, \theta, J, K, c$  graphically. I have made a convenience function `obsProb` to calculate the values. A lot of the results seem pretty obvious in retrospect, but it is still worth plotting them to get a better feel for the possibilities. At the moment, we only consider one and the same cost for each survey visit. It may be worth while to code up using a higher initial cost of the first visit, and lower costs for subsequent visits.

```
require(InsectSurvPower)

## Loading required package: InsectSurvPower
##custom library, available through devtools::install_github("NINAnor/InsectSurvPower")
require(dplyr)

## Loading required package: dplyr
##
## Attaching package: 'dplyr'
## The following objects are masked from 'package:stats':
##
##   filter, lag
```

```
## The following objects are masked from 'package:base':
##
## intersect, setdiff, setequal, union
```

```
require(gridExtra)
```

```
## Loading required package: gridExtra
```

```
##
```

```
## Attaching package: 'gridExtra'
```

```
## The following object is masked from 'package:dplyr':
```

```
##
```

```
## combine
```

```
require(ggplot2)
```

```
## Loading required package: ggplot2
```

```
require(xtable)
```

```
## Loading required package: xtable
```

We specify one or a range of values for occurrence probability in each location, detection probabilities, the number of locations visited, the number of visits per location, and the cost of each visit. To illustrate, we here use occurrence probabilities ranging from 0.01 to 1, two different detection probabilities as 0.5, and 0.8. We specify 30 locations, each visited 4 times, at an individual visit cost of 5000.

```
obsDf <- obsProb(occProb = seq(0.01, 1, by = 0.05),
                 detectProb = c(0.5, 0.8),
                 locations = 30,
                 visits = 4,
                 visitCost = 5000)
```

```
obsDf
```

```
## # A tibble: 40 x 7
```

```
##   occProb detectProb locations visits visitCost obsProb totCost
##   <dbl>    <dbl>    <dbl>  <dbl>    <dbl>    <dbl>    <dbl>
## 1  0.0100      0.500      30.0    4.00      5000    0.244  600000
## 2  0.0600      0.500      30.0    4.00      5000    0.791  600000
## 3  0.110      0.500      30.0    4.00      5000    0.909  600000
## 4  0.160      0.500      30.0    4.00      5000    0.932  600000
## 5  0.210      0.500      30.0    4.00      5000    0.937  600000
## 6  0.260      0.500      30.0    4.00      5000    0.937  600000
## 7  0.310      0.500      30.0    4.00      5000    0.937  600000
## 8  0.360      0.500      30.0    4.00      5000    0.937  600000
## 9  0.410      0.500      30.0    4.00      5000    0.937  600000
## 10 0.460      0.500      30.0    4.00      5000    0.937  600000
```

```
## # ... with 30 more rows
```

The probability of observing the species at least once is found in the `obsProb` column, and the total cost in the column `totCost`. The function returns an object of a specific class with a custom plotting function. This makes repeated plottings easier. We can specify a grouping variable to split up the lines.

```
plot(obsDf, group = "detectProb",
     xVar = "occProb",
     yVar = "obsProb",
     titleVar = c("locations", "visits"),
```

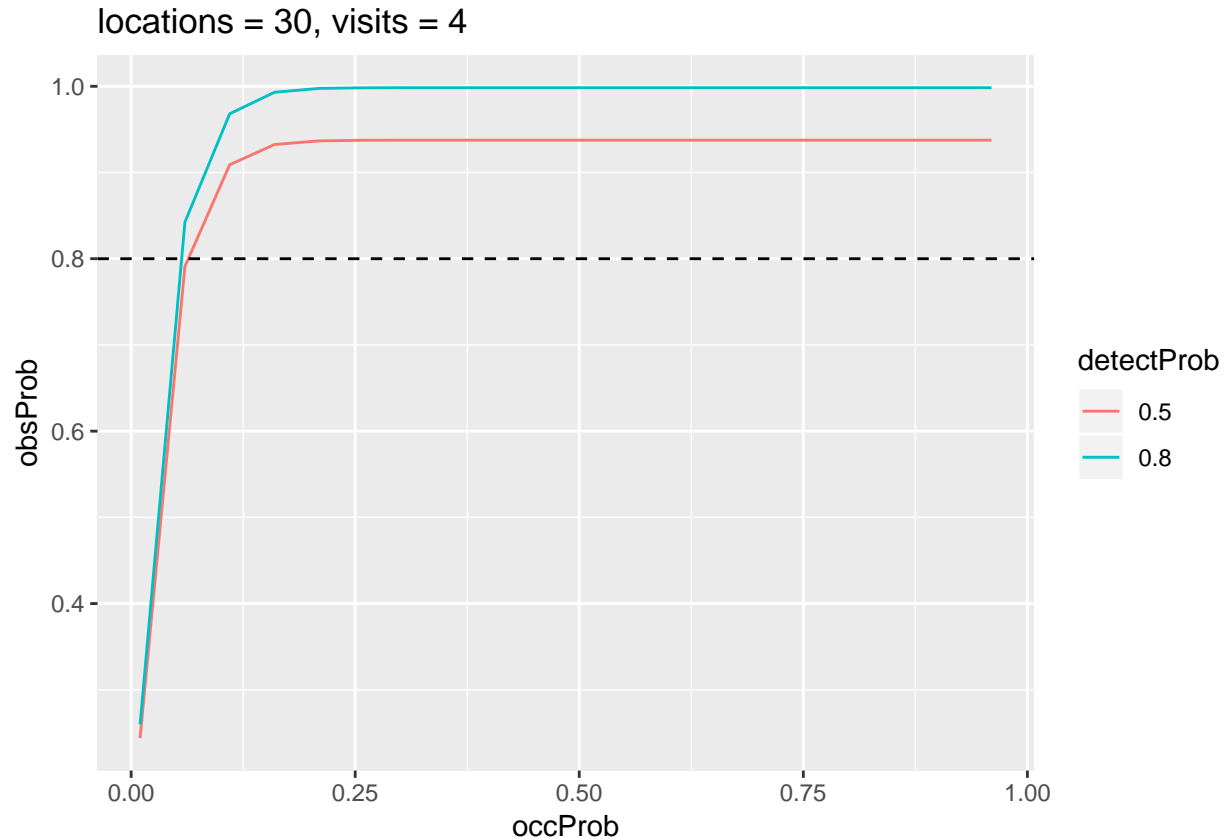


Figure 1: Observation probability for surveying a species with a detection probability of 0.5 in 30 locations, each visited 4 times.

```
hline = 0.8)
```

Figure 1 shows how the overall observation probability is dependent on both occurrence and detection probability. A threshold of a total observation probability of 0.8 is added for comparison. The threshold is reached in this case when the probability of occurrence in each location is 6%. With 30 locations, the overall observation probability rises quite sharply as a result of increased occurrence probability. With only 4 visits per location, the detection probability limits the overall achievable observation probability.

However, these occurrence probabilities are probably unreasonably high for rare species in the real world. If we assume that we search for a truly rare species with an occurrence probability of only 0.001, we find that reaching overall detectabilities above 80% is challenging (Figure 2).

```
rareLocDf <- obsProb(occProb = 0.001,
  detectProb = seq(0.4, 0.8, by = 0.2),
  locations = seq(30, 1000, by = 20),
  visits = seq(4, 16, by = 4),
  visitCost = 5000)

rareLocDf

## # A tibble: 588 x 7
##   occProb detectProb locations visits visitCost obsProb totCost
##   <dbl>    <dbl>    <dbl>  <dbl>    <dbl>  <dbl>  <dbl>
## 1 0.00100      0.400      30.0    4.00     5000  0.0257 600000
```

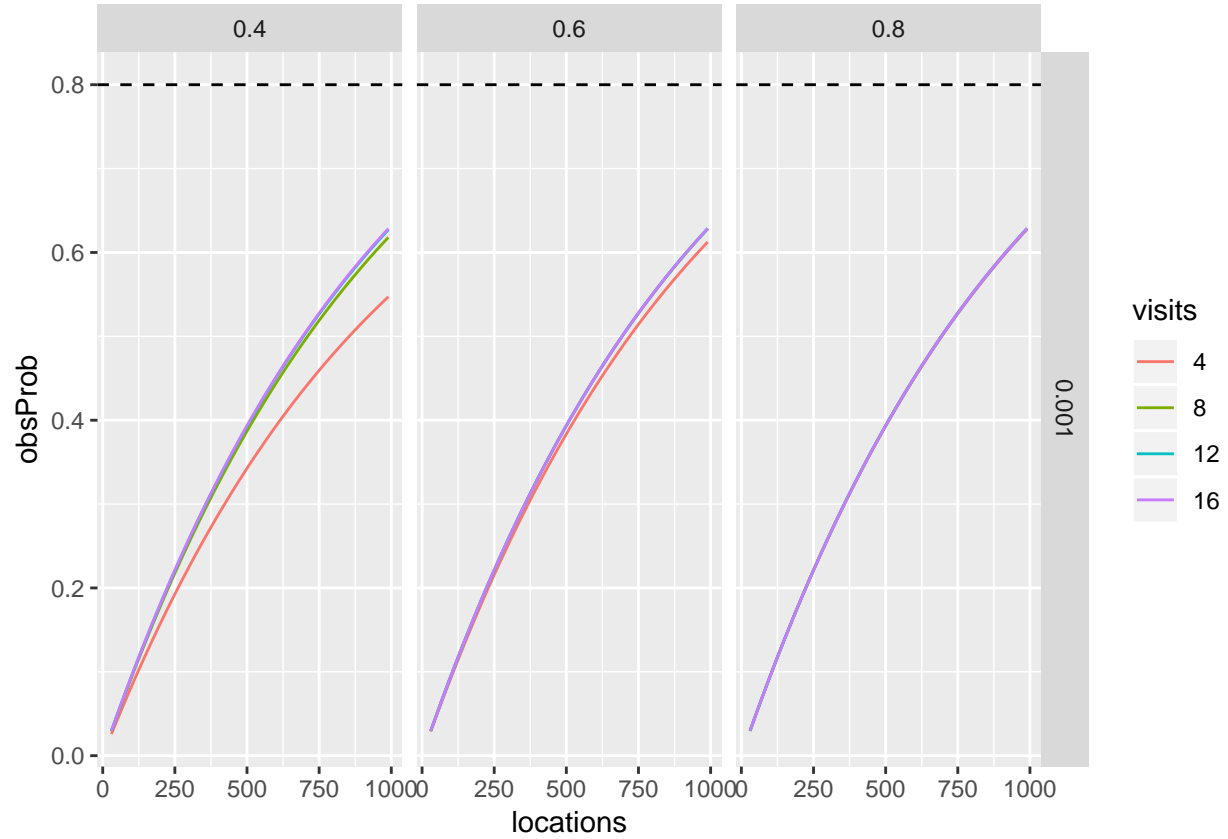


Figure 2: Probability of detecting a rare species (occurrence probability = 0.001) as a function of the number of visited locations and visits per location for different detection probabilities (detectProb = 0.4, 0.6, and 0.8).

```
## 2 0.00100      0.600      30.0    4.00      5000  0.0288  600000
## 3 0.00100      0.800      30.0    4.00      5000  0.0295  600000
## 4 0.00100      0.400      50.0    4.00      5000  0.0425 1000000
## 5 0.00100      0.600      50.0    4.00      5000  0.0475 1000000
## 6 0.00100      0.800      50.0    4.00      5000  0.0487 1000000
## 7 0.00100      0.400      70.0    4.00      5000  0.0589 1400000
## 8 0.00100      0.600      70.0    4.00      5000  0.0659 1400000
## 9 0.00100      0.800      70.0    4.00      5000  0.0675 1400000
## 10 0.00100     0.400      90.0    4.00      5000  0.0750 1800000
## # ... with 578 more rows
```

```
plot(rareLocDf, group = "visits",
     xVar = "locations",
     yVar = "obsProb",
     hline = 0.8) +
  facet_grid(occProb ~ detectProb)
```

In figure 2, we see that observing a very rare species with a high certainty is difficult even with a very large number of visited locations. In these cases, it doesn't really help to visit each location many times, as the overall probability is limited by the number of locations we visit. Figure 3 shows the results of maximising the number of visits in fixed, but large number of locations, for a very rare species.

```
rareVisitDf <- obsProb(occProb = 0.001,
  detectProb = seq(0.2, 0.8, by = .2),
  locations = 250,
  visits = seq(1, 11, by = 5),
  visitCost = 5000)

rareVisitDf

## # A tibble: 12 x 7
##   occProb detectProb locations visits visitCost obsProb totCost
##   <dbl>     <dbl>     <dbl>  <dbl>    <dbl>   <dbl>   <dbl>
## 1 0.00100     0.200       250    1.00     5000  0.0443 1250000
## 2 0.00100     0.400       250    1.00     5000  0.0885 1250000
## 3 0.00100     0.600       250    1.00     5000  0.133   1250000
## 4 0.00100     0.800       250    1.00     5000  0.177   1250000
## 5 0.00100     0.200       250    6.00     5000  0.163   7500000
## 6 0.00100     0.400       250    6.00     5000  0.211   7500000
## 7 0.00100     0.600       250    6.00     5000  0.220   7500000
## 8 0.00100     0.800       250    6.00     5000  0.221   7500000
## 9 0.00100     0.200       250   11.0     5000  0.202  13750000
## 10 0.00100     0.400       250   11.0     5000  0.220  13750000
## 11 0.00100     0.600       250   11.0     5000  0.221  13750000
## 12 0.00100     0.800       250   11.0     5000  0.221  13750000

plot(rareVisitDf, group = "detectProb",
  xVar = "visits",
  yVar = "obsProb",
  titleVar = c("occProb", "locations"),
  hline = 0.8)
```

Although the overall cost increases linearly with the total number of samples (figure 4), in cases with very rare species, this doesn't mean that the overall detection probability continues to increase indefinitely (figure 3)

```
rareCostDf <- obsProb(occProb = 0.001,
  detectProb = 0.4,
  locations = seq(50, 250, by = 50),
  visits = seq(1, 21, by = 5),
  visitCost = 5000)

rareCostDf

## # A tibble: 25 x 7
##   occProb detectProb locations visits visitCost obsProb totCost
##   <dbl>     <dbl>     <dbl>  <dbl>    <dbl>   <dbl>   <dbl>
## 1 0.00100     0.400      50.0    1.00     5000  0.0195  250000
## 2 0.00100     0.400     100     1.00     5000  0.0381  500000
## 3 0.00100     0.400     150     1.00     5000  0.0557  750000
## 4 0.00100     0.400     200     1.00     5000  0.0725 1000000
## 5 0.00100     0.400     250     1.00     5000  0.0885 1250000
## 6 0.00100     0.400     50.0    6.00     5000  0.0465 1500000
## 7 0.00100     0.400     100     6.00     5000  0.0908 3000000
## 8 0.00100     0.400     150     6.00     5000  0.133   4500000
## 9 0.00100     0.400     200     6.00     5000  0.173   6000000
## 10 0.00100     0.400     250     6.00     5000  0.211   7500000
## # ... with 15 more rows
```

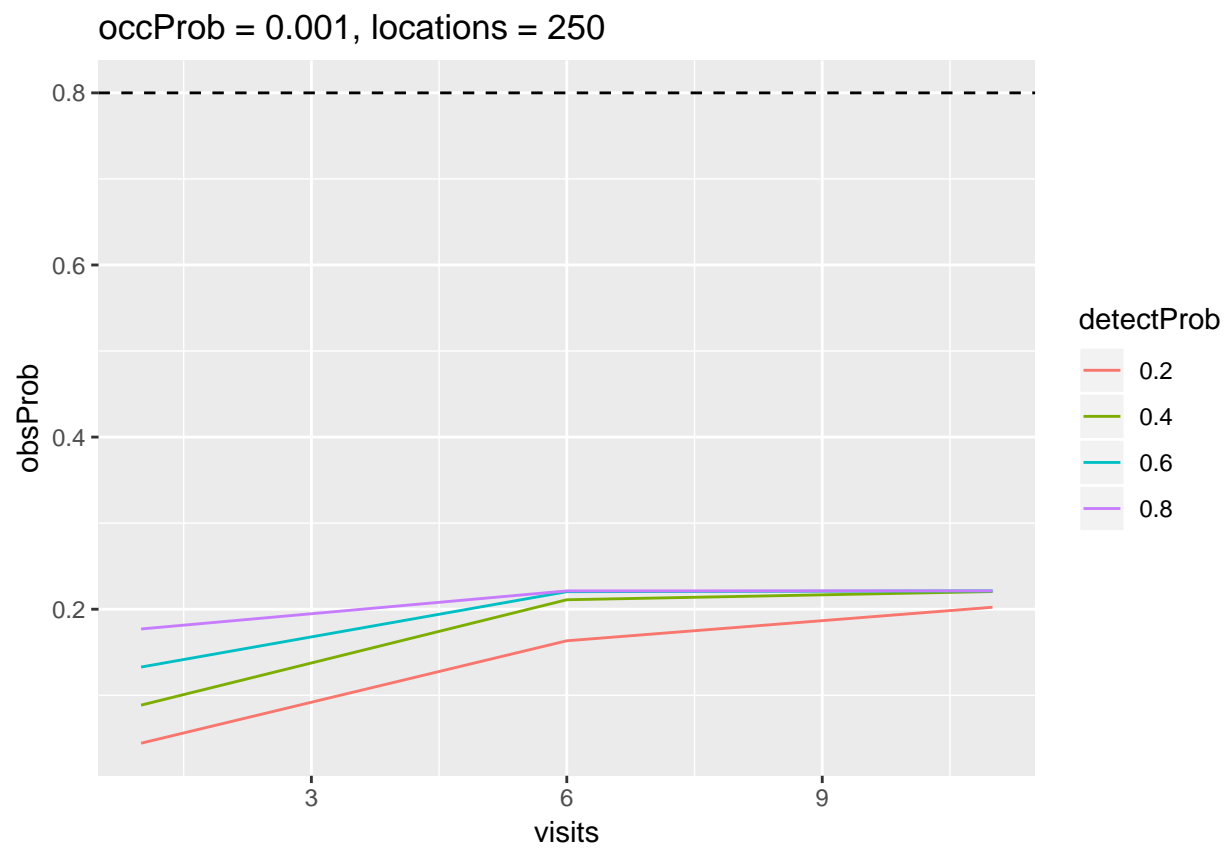


Figure 3: Observation probability for surveying a very rare species as function of sampled locations.

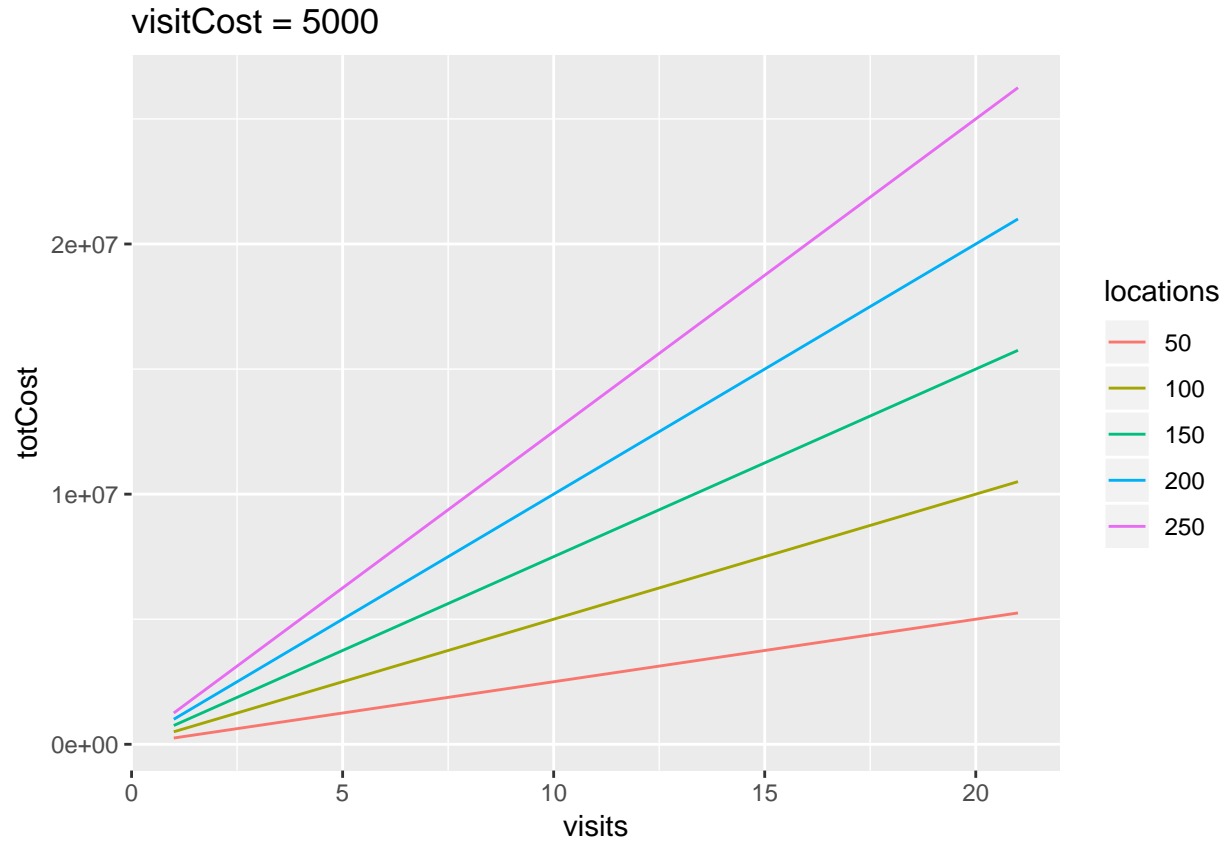


Figure 4: Total survey cost as a function of the total number of samples

```
plot(rareCostDf, group = "locations",
     xVar = "visits",
     yVar = "totCost",
     titleVar = "visitCost")
```

## Some strategies

As seen in figure 5, when we deal with a very rare species, it is little use increasing the number of visits to each location (or to spend money maximizing the detection probability), if we can't at the same time span a very large number of locations. We must in these cases concentrate on increasing the number of visited locations. Still, for a very rare species such as displayed in figure 5, with an occurrence probability of 0.01%, reaching an overall observation probability of 80% requires more at least 1650 locations, which would be unfeasible for most survey programs.

Alternatively, if detectability is low but presence is relatively high, we should focus on increasing the number of revisits per location, instead of trying to cover many locations (figure 6).

```
lowOccurrDf <- obsProb(occProb = 0.001,
                      detectProb = 0.8,
                      locations = seq(50, 2000, by = 50),
                      visits = c(1, seq(5, 20, by = 5)),
                      visitCost = 5000)
```

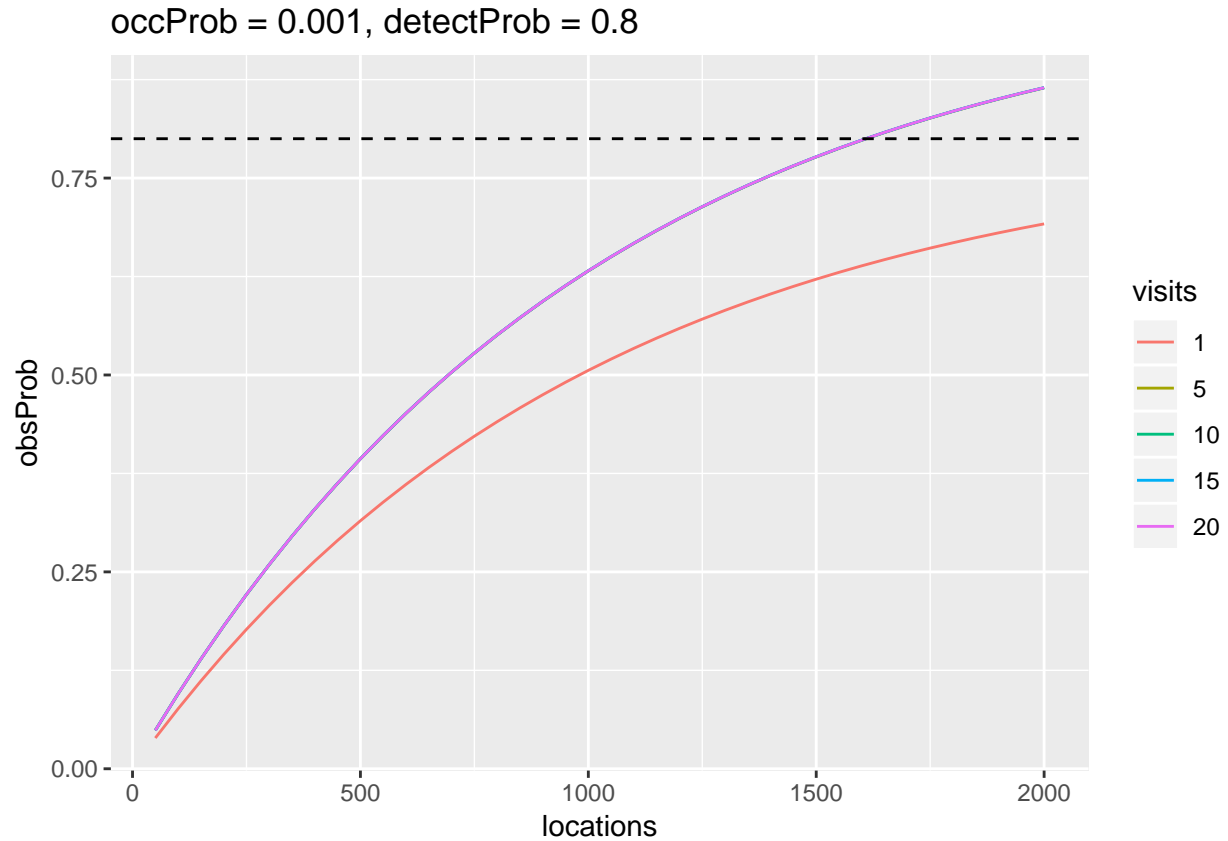


Figure 5: Observation probability for surveying a rare species as function of sampled locations.

```
lowOccurrDf

## # A tibble: 200 x 7
##   occProb detectProb locations visits visitCost obsProb totCost
##   <dbl>     <dbl>     <dbl>  <dbl>    <dbl>   <dbl>   <dbl>
## 1 0.00100     0.800      50.0    1.00     5000  0.0390 250000
## 2 0.00100     0.800     100    1.00     5000  0.0762 500000
## 3 0.00100     0.800     150    1.00     5000  0.111  750000
## 4 0.00100     0.800     200    1.00     5000  0.145 1000000
## 5 0.00100     0.800     250    1.00     5000  0.177 1250000
## 6 0.00100     0.800     300    1.00     5000  0.207 1500000
## 7 0.00100     0.800     350    1.00     5000  0.236 1750000
## 8 0.00100     0.800     400    1.00     5000  0.264 2000000
## 9 0.00100     0.800     450    1.00     5000  0.290 2250000
## 10 0.00100     0.800     500    1.00     5000  0.315 2500000
## # ... with 190 more rows

plot(lowOccurrDf, group = "visits",
     xVar = "locations",
     yVar = "obsProb",
     titleVar = c("occProb", "detectProb"),
     hline = 0.8)
```



```
lowDetectDf <- obsProb(occProb = 0.05,
                      detectProb = 0.2,
                      locations = seq(50, 250, by = 50),
                      visits = c(1, seq(5, 20, by = 5)),
                      visitCost = 5000)

lowDetectDf

## # A tibble: 25 x 7
##   occProb detectProb locations visits visitCost obsProb totCost
##   <dbl>    <dbl>    <dbl>  <dbl>    <dbl>    <dbl>    <dbl>
## 1 0.0500    0.200    50.0   1.00    5000    0.185  250000
## 2 0.0500    0.200   100    1.00    5000    0.199  500000
## 3 0.0500    0.200   150    1.00    5000    0.200  750000
## 4 0.0500    0.200   200    1.00    5000    0.200 1000000
## 5 0.0500    0.200   250    1.00    5000    0.200 1250000
## 6 0.0500    0.200   50.0   5.00    5000    0.621 1250000
## 7 0.0500    0.200   100    5.00    5000    0.668 2500000
## 8 0.0500    0.200   150    5.00    5000    0.672 3750000
## 9 0.0500    0.200   200    5.00    5000    0.672 5000000
## 10 0.0500    0.200   250    5.00    5000    0.672 6250000
## # ... with 15 more rows

plot(lowDetectDf, group = "visits",
     xVar = "locations",
     yVar = "obsProb",
     titleVar = c("occProb", "detectProb"),
     hline = 0.8)
```

## Plausible values

It is difficult to guess plausible values for occurrence and detectability for real world species, but it is reasonable to assume that we only have to consider rather low occurrence probabilities, since we are working on early detections. We can explore our possibilities of observing a species that occur in between 0.1 to 1% of all studied locations. We can set the number of locations to 100 and with two visits, as a reasonable possibility.

```
guessDf <- obsProb(occProb = c(0.001, 0.005, 0.01, 0.02, 0.04, 0.05, 0.1),
                  detectProb = c(0.1, seq(0.2, 0.8, by = 0.2)),
                  locations = c(50, seq(100, 300, by = 100)),
                  visits = seq(2, 6, by = 2),
                  visitCost = 5000)

guessDf
```

```
## # A tibble: 420 x 7
##   occProb detectProb locations visits visitCost obsProb totCost
##   <dbl>    <dbl>    <dbl>  <dbl>    <dbl>    <dbl>    <dbl>
## 1 0.00100    0.100    50.0   2.00    5000 0.00927  500000
## 2 0.00500    0.100    50.0   2.00    5000 0.0421  500000
## 3 0.0100    0.100    50.0   2.00    5000 0.0750  500000
## 4 0.0200    0.100    50.0   2.00    5000 0.121  500000
## 5 0.0400    0.100    50.0   2.00    5000 0.165  500000
```

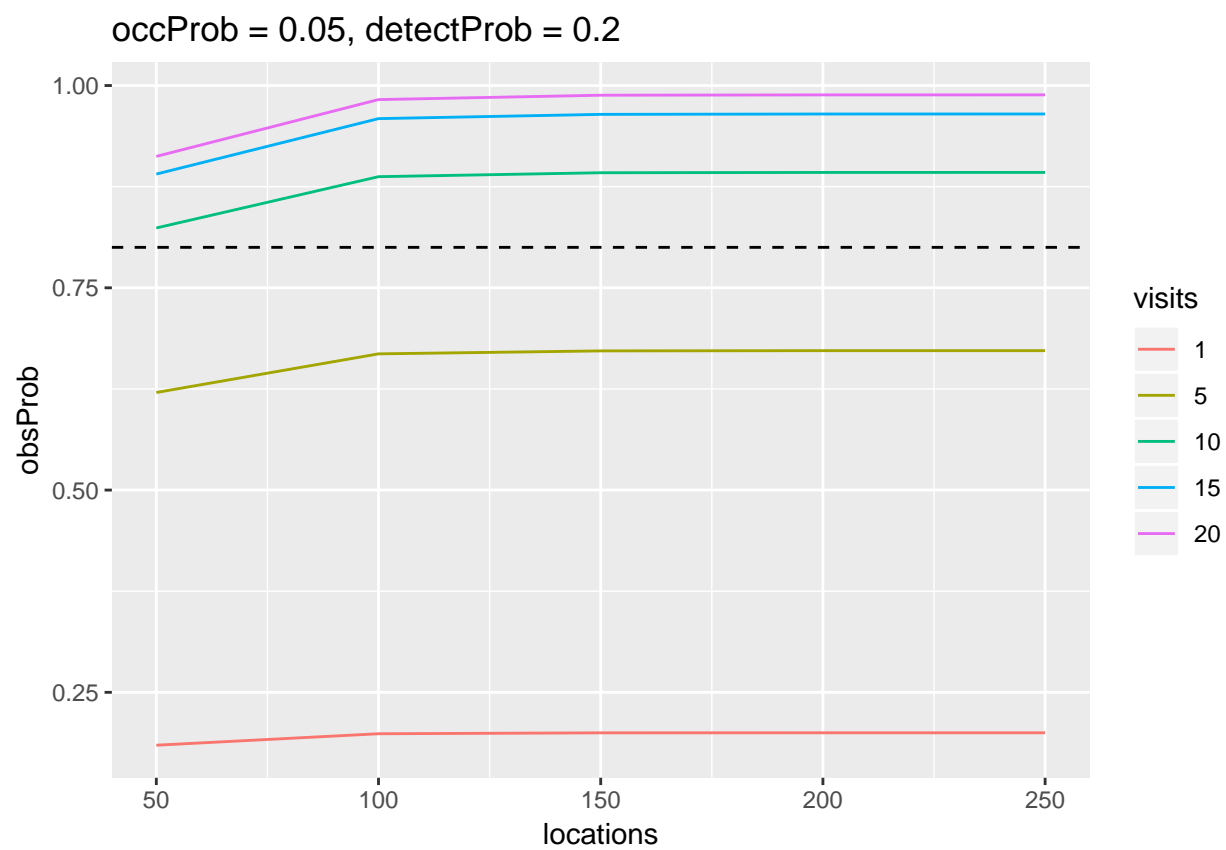


Figure 6: Observation probability for surveying a cryptic species as function of sampled locations.

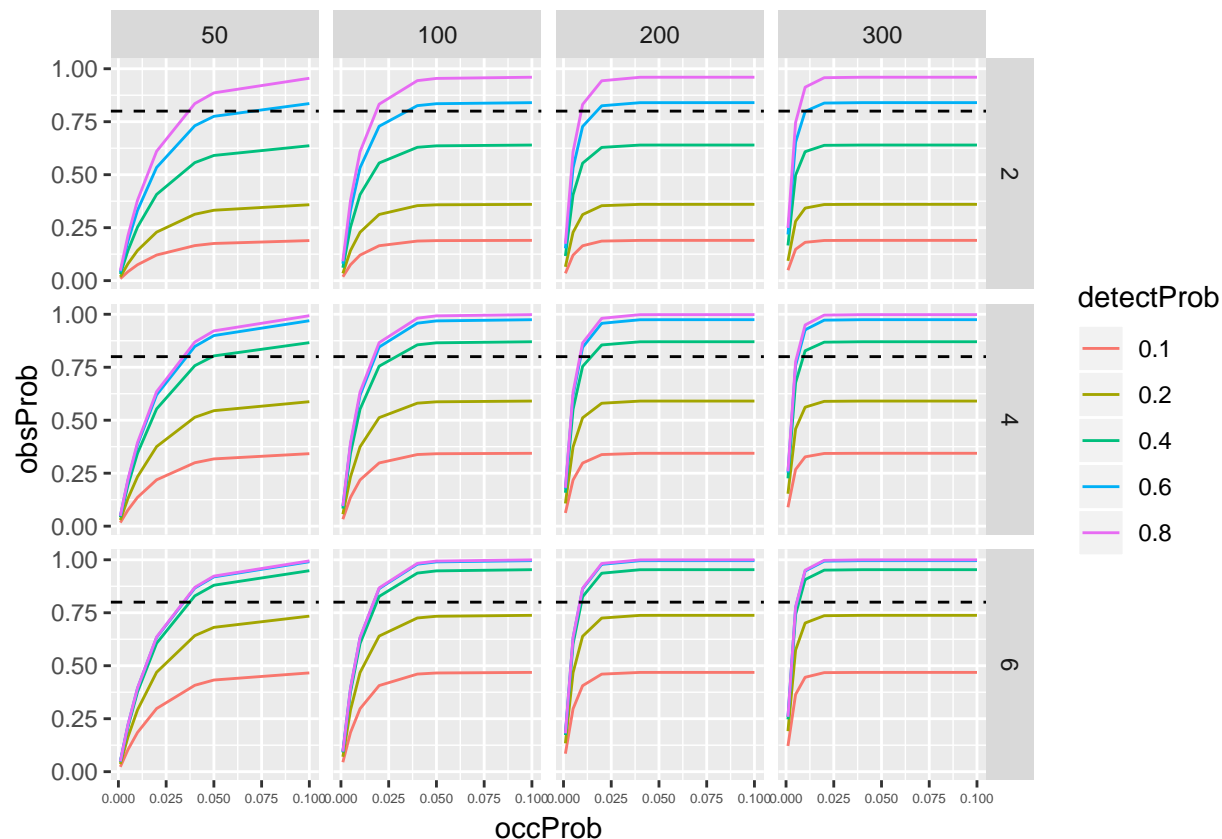


Figure 7: Observation probability for a plausible range of values. The plots are divided by number of sample locations in columns, and by number of visits in rows.

```
## 6 0.0500      0.100      50.0    2.00      5000 0.175      500000
## 7 0.100      0.100      50.0    2.00      5000 0.189      500000
## 8 0.00100    0.200      50.0    2.00      5000 0.0176     500000
## 9 0.00500    0.200      50.0    2.00      5000 0.0798     500000
## 10 0.0100    0.200      50.0    2.00      5000 0.142      500000
## # ... with 410 more rows
```

```
plot(guessDf, group = "detectProb",
     xVar = "occProb",
     yVar = "obsProb",
     hline = 0.8) +
  facet_grid(visits ~ locations) +
  theme(axis.text.x = element_text(size = 5))
```

```
plot(guessDf, group = "detectProb",
     xVar = "occProb",

     yVar = "obsProb",
     hline = 0.8) +
  facet_grid(visits ~ locations) +
  scale_y_log10() +
  scale_x_log10() +
  theme(axis.text.x = element_text(size = 7))
```

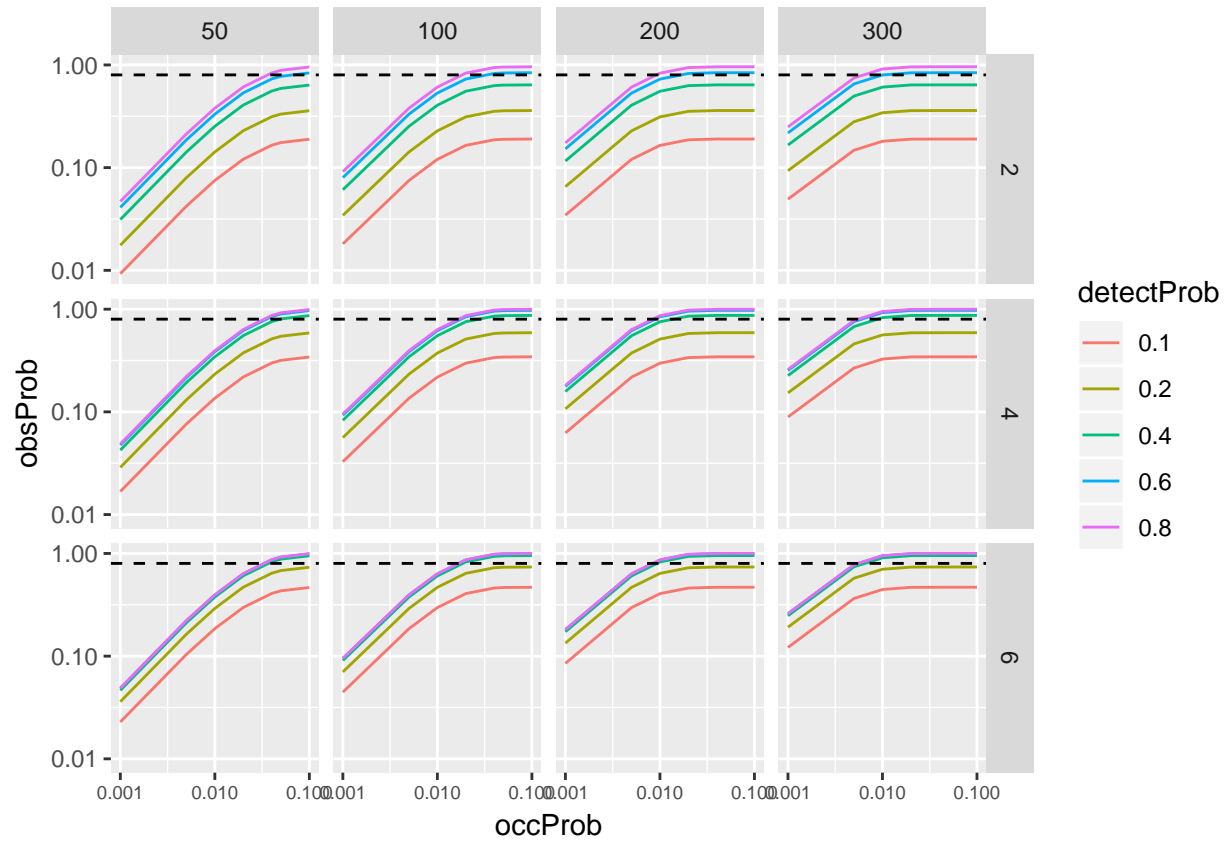


Figure 8: Observation probability for a plausible range of values. The plots are divided by number of sample locations in columns, and by number of visits in rows. Note the logarithmic scale.

## Unequally distributed probabilities

\*\* There is no definition for early detection. The earlier we want to detect (no occurrences is low) the more costly to detect. We can use the alien/native map as a relative risk map and use that to “model” a weighted occurrence map. This we visit in a weighted fashion as well, a number of times. Need to work out the math. We can then calculate the costs for different number of alien occurrences.\*\*

In the equation  $probObserv = (1 - (1 - \psi)^J) * (1 - (1 - \theta)^K)$ ,  $\psi$  designates the probability that a site we visit contains a certain species. So far, we have only considered situations where the occurrence probabilities ( $\psi$ ) are the same for all sites. In reality, however, the probability that a specific site that you visit contains a certain species will vary between sites. It will depend both on the probability of a site containing a species, and the probability that you visit that site. We can view the probability of a specific site containing a species as a weighted sampling without replacement. For example, if we know there are a 100 sites containing species x, we can calculate the probability that each site contains species x if we know the probability weights. This probability can be written, following Eframidis & Spirakis 2006, as  $p_i(k) = \frac{w_i}{\sum_{S_j \in V-S} w_j}$ . For

brevity, however, we will simply designate this probability as  $Pr[w_i, n]$ , where the weights  $w_i$  sums to 1, and  $n$  is the number of sites with the species. Similarly, we choose which sites to visit as a sampling without replacement, so that the probability of visiting a site  $i$  is  $Pr[u_i, v]$ , where  $u_i$  is the visitation weights, which sums to 1, and  $v$  is the number of sites you visit. The probability of visiting a site with an alien species then becomes  $\psi_i = Pr[w_i, n] * Pr[u_i, v]$ , and the probability of visiting any location inhabited by an alien species is  $PoccurrVisit = (1 - \prod_i (1 - (Pr[w_i, n] * Pr[u_i, v])))$ .

For simplicity, we can assume that the probability of detection is equal for each site and visit. The probability of detecting an alien species is then  $Pdetect = (1 - \prod_i (1 - (Pr[w_i, n] * Pr[u_i, v]))) * (1 - (1 - d)^K)$ , where  $d$  is the detection probability, and  $K$  is again the number of visits per site.

In the best of worlds,  $w_i$  and  $u_i$  would match up well, so that we visit the most likely sites to contain a species of interest. The overall probability of detection will decrease as the difference between these variables increase, or in simpler terms we visit the wrong places. Also, the probability will go down the more spread out the probability weights are. In the extreme case, with no spread in these probabilities, there is 100 % certainty that a species will be present in location 1, and 100 % certainty that we will visit just that. In this case, the probabilities are 1 for both occurrence and visitation, so that we are certain we visit a site with a presence. In the other end, there might be no information in the weights, so that the occurrences are randomly spread out, and we visit random sites. In that case,  $w_i$  are all the same and  $u_i$  are all the same, and we end up with the first equation.

In reality, we don't know the true weights that the sites will contain a specific species, and we might choose to visit the wrong sites. For this example calculation, we will assume we know the occurrence weights, and visit them accordingly. In other words, that  $w_i = u_i$ . If we stipulate the desired detection probability (at e.g 0.8), we can calculate how many sites  $v$  we need to visit to be able to detect a species that occur in  $n$  sites with weighted probabilities  $w_i$ , with a specific certainty.

So far, the best estimate for the occurrence weights are the occurrence modelling of alien vascular plants from Olsen et al. 2017. Using this as input, and some simple assumptions, we can calculate the needed number of sites.

To get test it out, we can use the 10km scale, which isn't so resource intensive. The 1km scale isn't really interactive since it takes to long time to calculate.

For now, we get rid of the geometry column to increase speed.

```
predWeights <- predMap %>%
  select(sites = ssbid,
         weights = pred) %>%
  sf::st_set_geometry(NULL)
```

But we start with a situation where the occupied patches are distributed randomly, and we visit the sites

randomly. In other words, where the occurrence and visitation weights all are equal.

```
predWeightsZero<- predWeights
predWeightsZero$weights <- 1
```

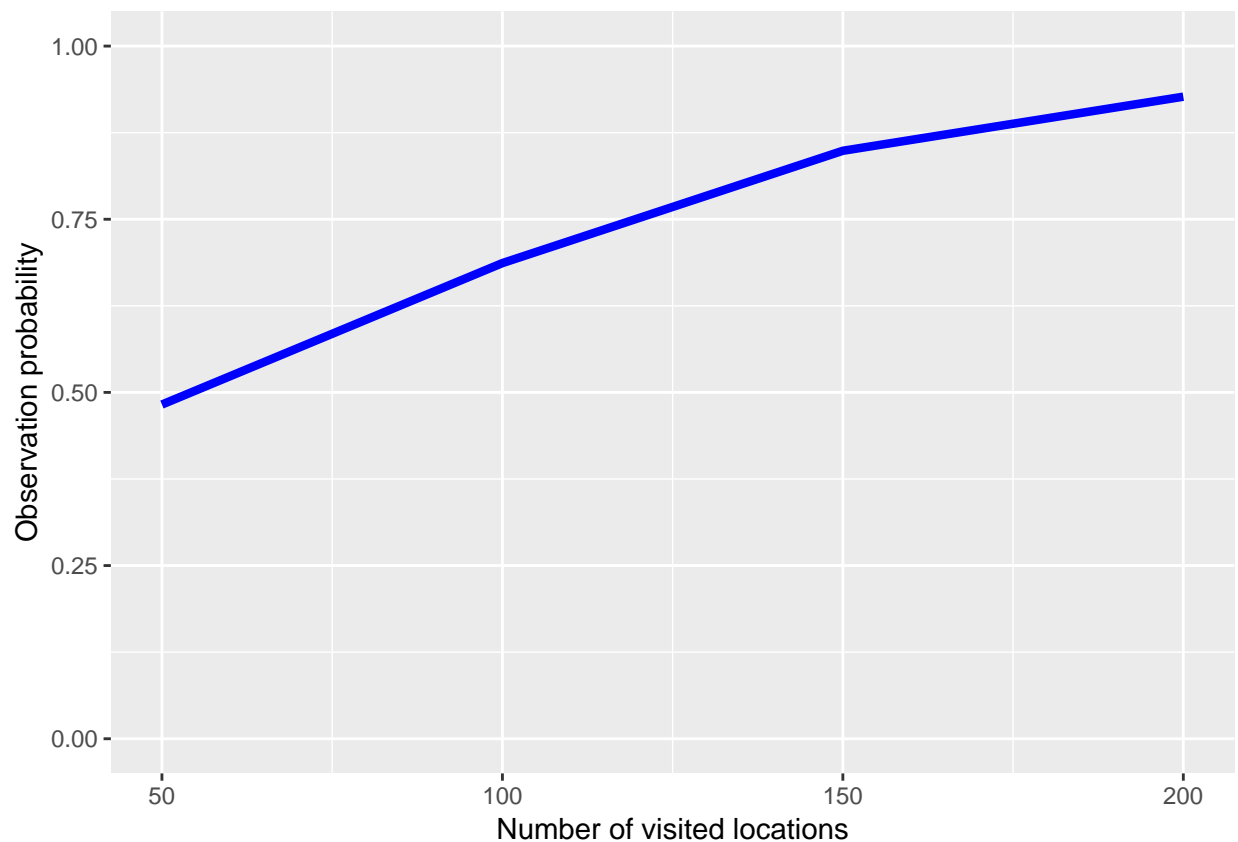
```
system.time(predProb50Zero <- weightedDetection(occWeights = predWeightsZero,
visWeights = predWeights,
noOccur = 50,
noLocations = seq(50, 200, by = 50),
noVisits = 1,
detectProb = 1))
```

```
##      user  system elapsed
##    3.660   0.080   3.743
```

```
predProb50Zero
```

```
## # A tibble: 4 x 2
##   noLocations probObs
##       <dbl>   <dbl>
## 1      50.0   0.482
## 2     100    0.687
## 3     150    0.849
## 4     200    0.927
```

```
plot(predProb50Zero)
```



We can see that the probability of visiting a randomly occupied cell in this case starts from about 0.45 and

approaches 1 as we increase our number of visited locations from 50 to 200. We can quality check this with a simpler function.

```
test <- function(noOccur = 50,
                 noLocations = 50,
                 nIter = 999){

  prop <- function(noLocations. = noLocations,
                   noOccur. = noOccur){
    visited <- sample(1:4057, noLocations., replace = F) #number of 10km cells
    occupied <- sample(1:4057, noOccur, replace = F)

    any(visited %in% occupied)
  }

  sum(replicate(nIter, prop())/nIter)
}

test()
```

```
## [1] 0.4834835
```

This can seem as a high number, but it appears to check out. We get the same results if the occurrence probabilities is distributed according to informative weights, and only the visitations are random (not shown).

But what happens when we have information about the occurrence of the species? In effect, we limit the number of potential sites we visit to a smaller value, which have higher probability of housing the species. We use the prediction map to set the occurrences, and visitation probabilities. We continue with the detection probability set to 1, with just 1 visit per site.

```
system.time(predProb <- weightedDetection(occWeights = predWeights,
                                           visWeights = predWeights,
                                           noOccur = 50,
                                           noLocations = seq(50, 200, by = 50),
                                           noVisits = 1,
                                           detectProb = 1,
                                           nIter = 999))
```

```
##      user  system elapsed
##  4.184    0.020    4.206
```

```
predProb
```

```
## # A tibble: 4 x 2
##   noLocations probObs
##       <dbl>   <dbl>
## 1       50.0   0.999
## 2       100    1.00
## 3       150    1.00
## 4       200    1.00
```

```
plot(predProb)
```

The result is virtually certainty that we will observe the species. This of course depend on the quality of the predictions, i.e. our weights. This prediction map is actually quite informative. From the histogram of

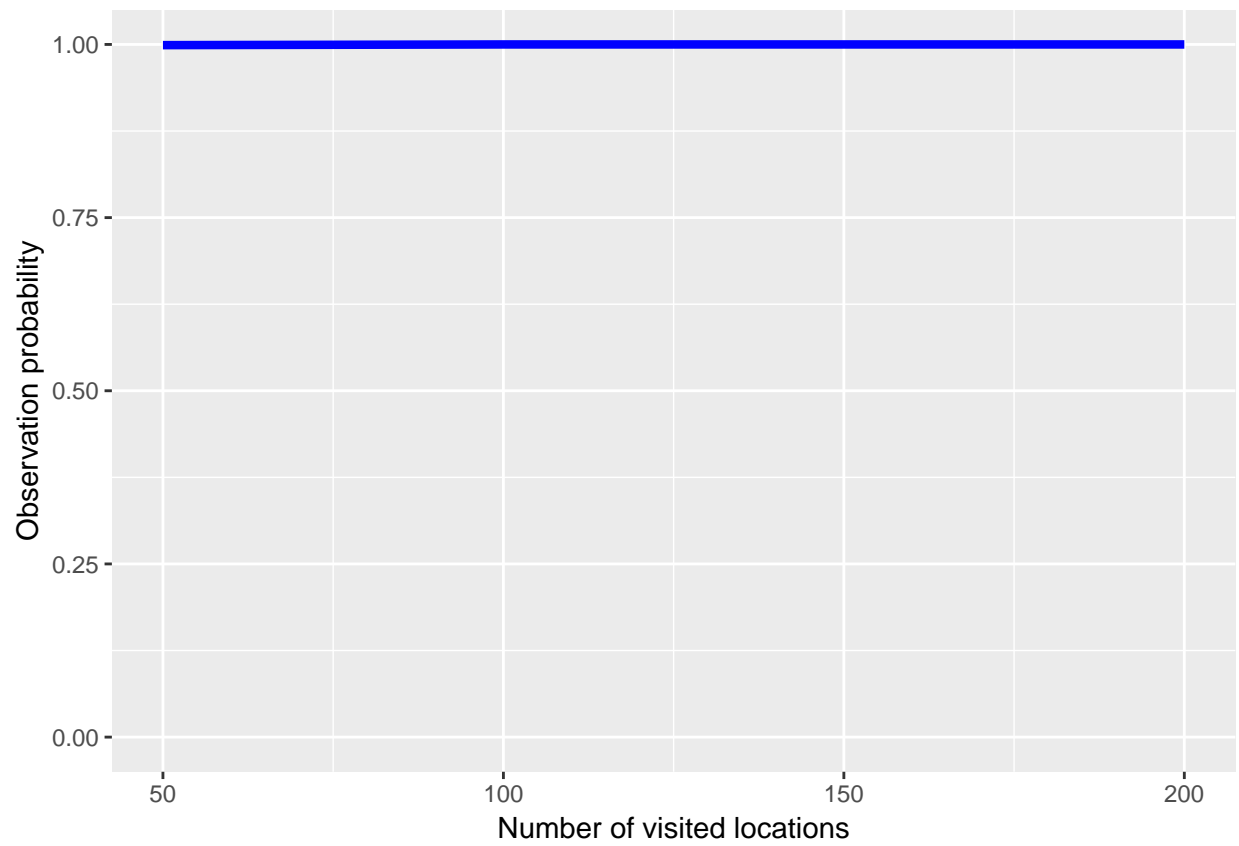
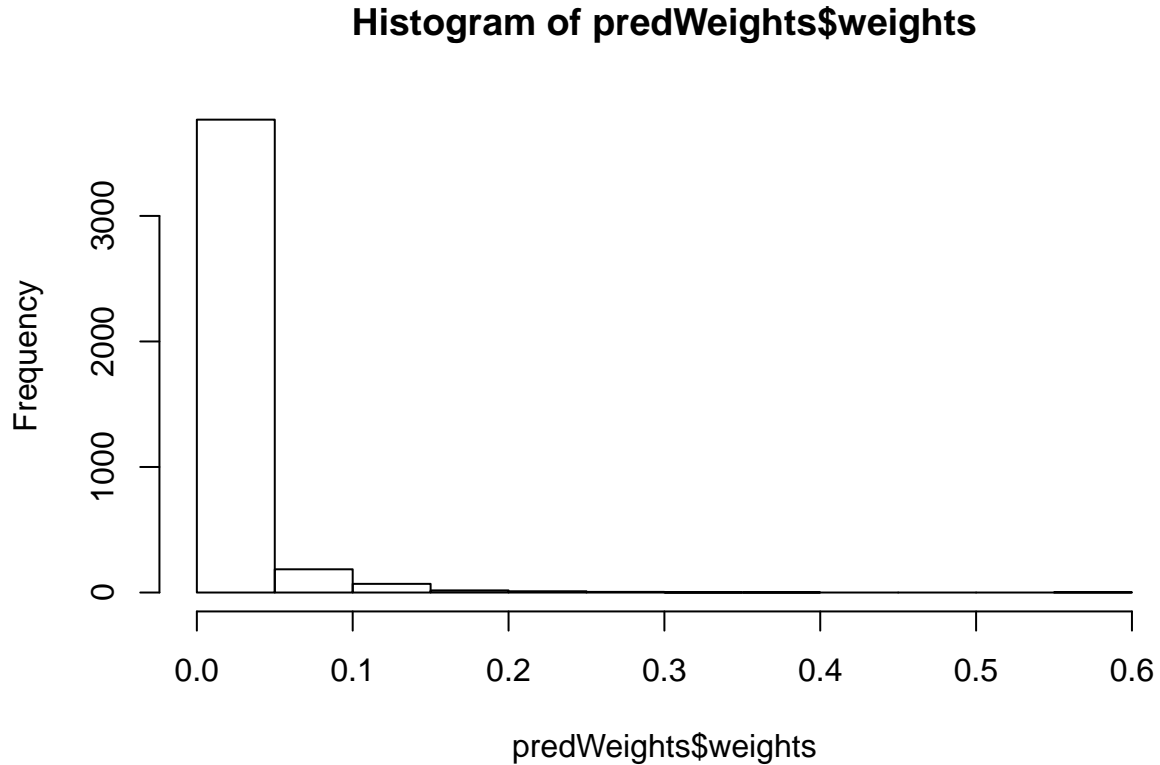


Figure 9: 50 out of 4057 grid cells occupied according to weights, 1 visit with 1 detection probability



weights, we can see that a small number of sites have a proportionally high weight. This limits the realised occurrences quite a bit.

```
hist(predWeights$weights)
```



Test with 1km map

```
require(RPostgreSQL)
require(DBI)
require(rpostgis)
require(SurveyPower)
require(tidyverse)
data(map1km)
```

```
con <- dbConnect(RPostgreSQL::PostgreSQL(), host = "gisdata-db.nina.no", dbname = "gisdata", user = "postgres", password = "postgres")
#pred <- pgGetRast(con, name = c("hotspot_ias", "bigPred1km"), rast = "rast", bands = 1, boundary = NULL)
```

```
predQ <- "SELECT ssbid, ST_Value(pred.rast, ST_Centroid(ssb.geom)) pred
FROM ssb_data_utm33n.ssb_1km ssb,
hotspot_ias.\"evenintbigpred1km\" pred
WHERE ST_Intersects(ssb.geom, pred.rast)
"
```

```
pred <- dbGetQuery(con, predQ)
```

```
pred <- pred %>%
```

```

mutate(ssbid = as.character(ssbid))

predMap1km <- map1km %>% left_join(pred, by = c("ssbid" = "ssbid"))

predMap1km$pred[is.na(predMap1km$pred)] <- 0

#plot(predMap1km["pred"]) #Slow

predWeights1km <- predMap1km %>%
  select(sites = ssbid,
         weights = pred) %>%
  sf::st_set_geometry(NULL)

devtools::use_data(predWeights1km)

```

Since we now have about 100 times as many potential sites, we would need to increase the occurrences accordingly. But while 50 out of 4057 10x10km squares sounds reasonable for an “early” detection, multiplying this with 100 yields 5000 locations in a 1x1km grid. This sounds like a lot for an early detection. But it puts the 50 occurrences above into perspective. Surveying 10x10km cells with good detection probability is a tall order.

For the 1x1km analysis, the calculations takes to much time to do on the fly, so we pre-calculate them and load the results.

```
data("predWeights1km")
```

It is reasonable to assume that we won’t reach a higher detection probability than 0.8 for a single visit, and even that is probably high for a truly novel species. But we can explore the range of occupied sites that are reasonable to handle with such a good detection probability.

```

system.time(predOcc500Det0.8 <- weightedDetection(occWeights = predWeights1km,
                                                  visWeights = predWeights1km,
                                                  noOccur = 500,
                                                  noLocations = seq(50, 300, by = 50),
                                                  noVisits = 1,
                                                  detectProb = 0.8,
                                                  nIter = 999))

save(predOcc500Det0.8, file = "predOcc500Det0.8.Rdata")

load(file = "predOcc500Det0.8.Rdata")
predOcc500Det0.8

```

```

## # A tibble: 6 x 2
##   noLocations probObs
##   <dbl>      <dbl>
## 1      50.0    0.343
## 2     100     0.538
## 3     150     0.691
## 4     200     0.798
## 5     250     0.866
## 6     300     0.913

plot(predOcc500Det0.8,
     threshold = 0.8 )

```

So for the case of 500 occupied cells, we reach the target observation probability after visiting about 200 sites.

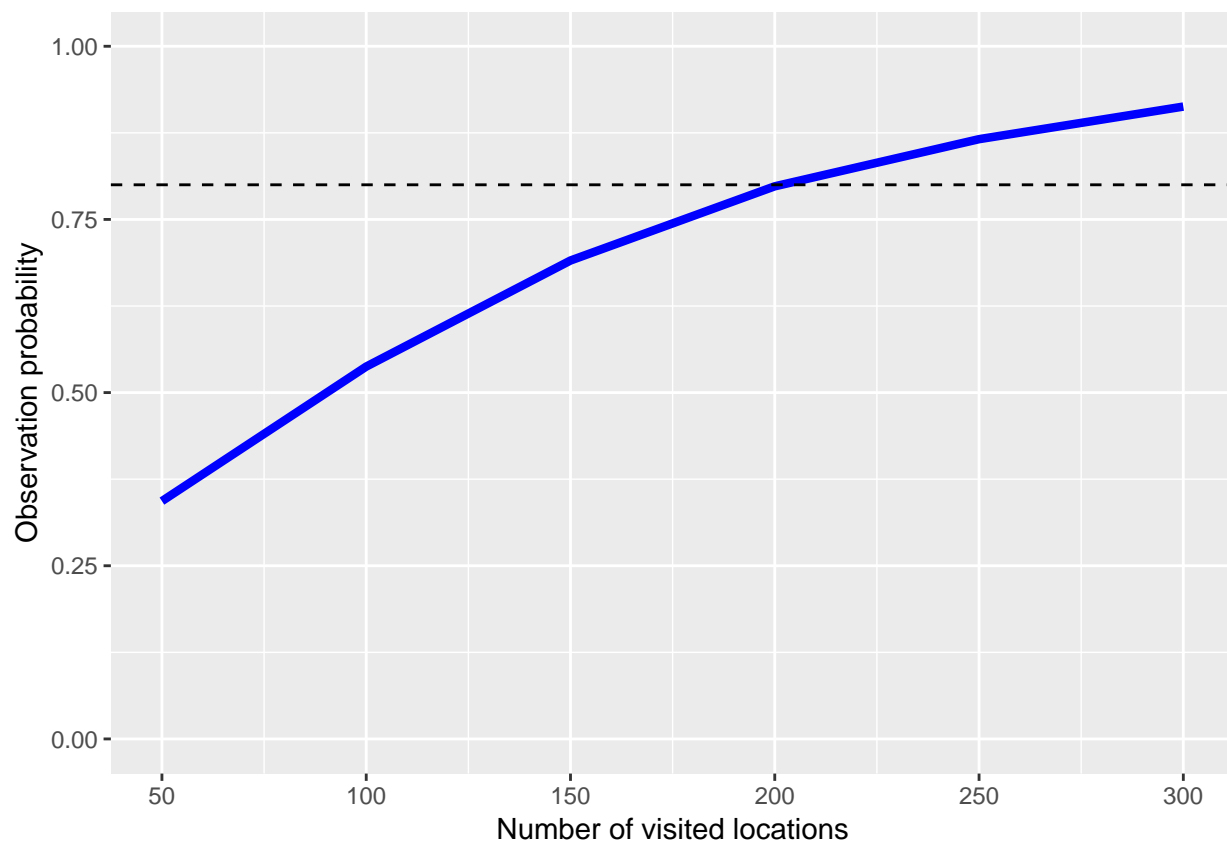


Figure 10: Estimated observation probability of an alien vascular plant species occurring in 500 1x1km grid cells as a function of the number of visited locations. Occurrences and location selection is based on the same weights, modelled from actual alien vascular plant species occurrences. Each visit has a 0.8 probability of detecting the species if present. The threshold of the desired observation probability of 0.8 is shown as a dashed line.

In similar fashion, we can explore the case with 200 and 100 occupied cells, respectively.

```
system.time(predOcc200Det0.8 <- weightedDetection(occWeights = predWeights1km,
  visWeights = predWeights1km,
  noOccur = 200,
  noLocations = seq(50, 600, by = 50),
  noVisits = 1,
  detectProb = 0.8,
  nIter = 999))

save(predOcc200Det0.8, file = "predOcc200Det0.8.Rdata")

load(file = "predOcc200Det0.8.Rdata")
xtable(predOcc200Det0.8)
```

noLocations	probObs
50.00	0.13
100.00	0.28
150.00	0.40
200.00	0.48
250.00	0.52
300.00	0.63
350.00	0.68
400.00	0.73
450.00	0.76
500.00	0.81
550.00	0.83
600.00	0.88

```
plot(predOcc200Det0.8,
  threshold = 0.8)

system.time(predOcc100Det0.8 <- weightedDetection(occWeights = predWeights1km,
  visWeights = predWeights1km,
  noOccur = 100,
  noLocations = seq(50, 600, by = 50),
  noVisits = 1,
  detectProb = 0.8,
  nIter = 999))

save(predOcc100Det0.8, file = "predOcc100Det0.8.Rdata")

load(file = "predOcc100Det0.8.Rdata")

xtable(predOcc100Det0.8)

plot(predOcc100Det0.8,
  threshold = 0.8)
```

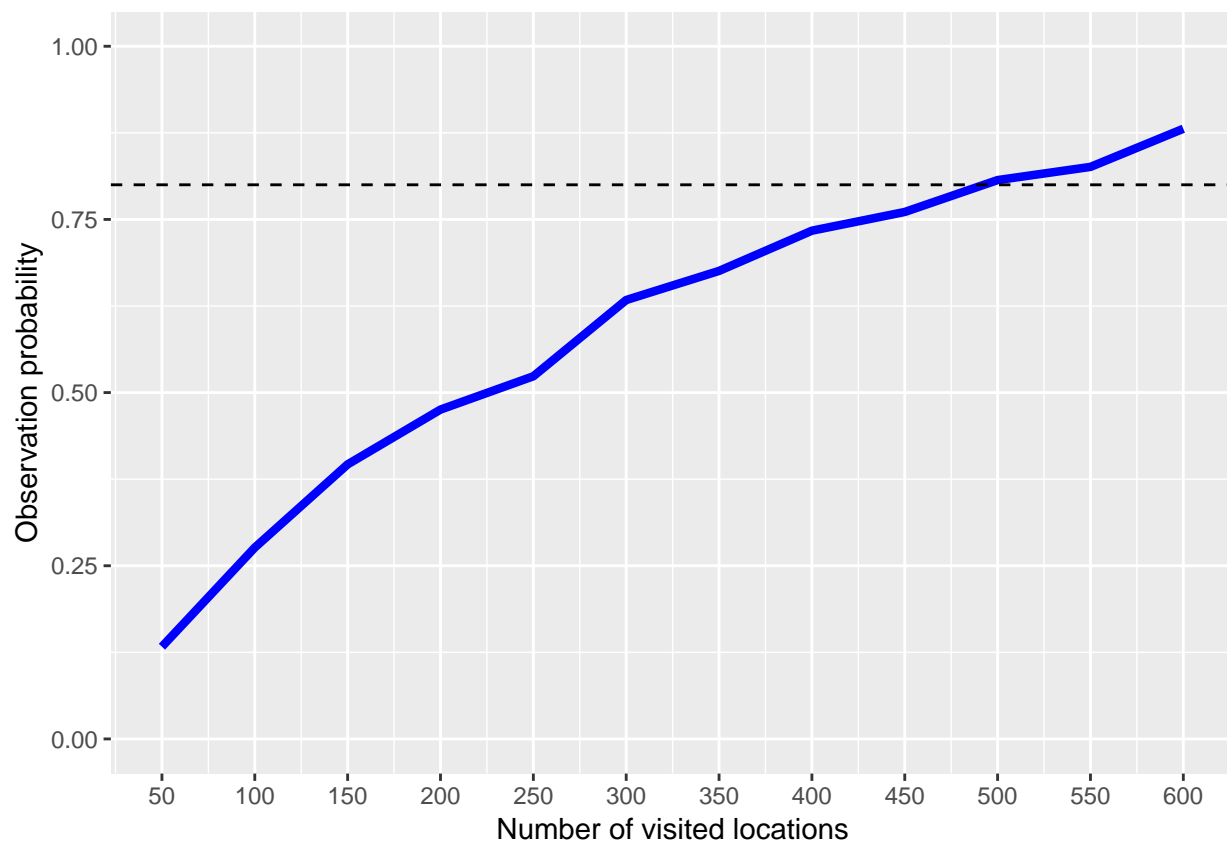
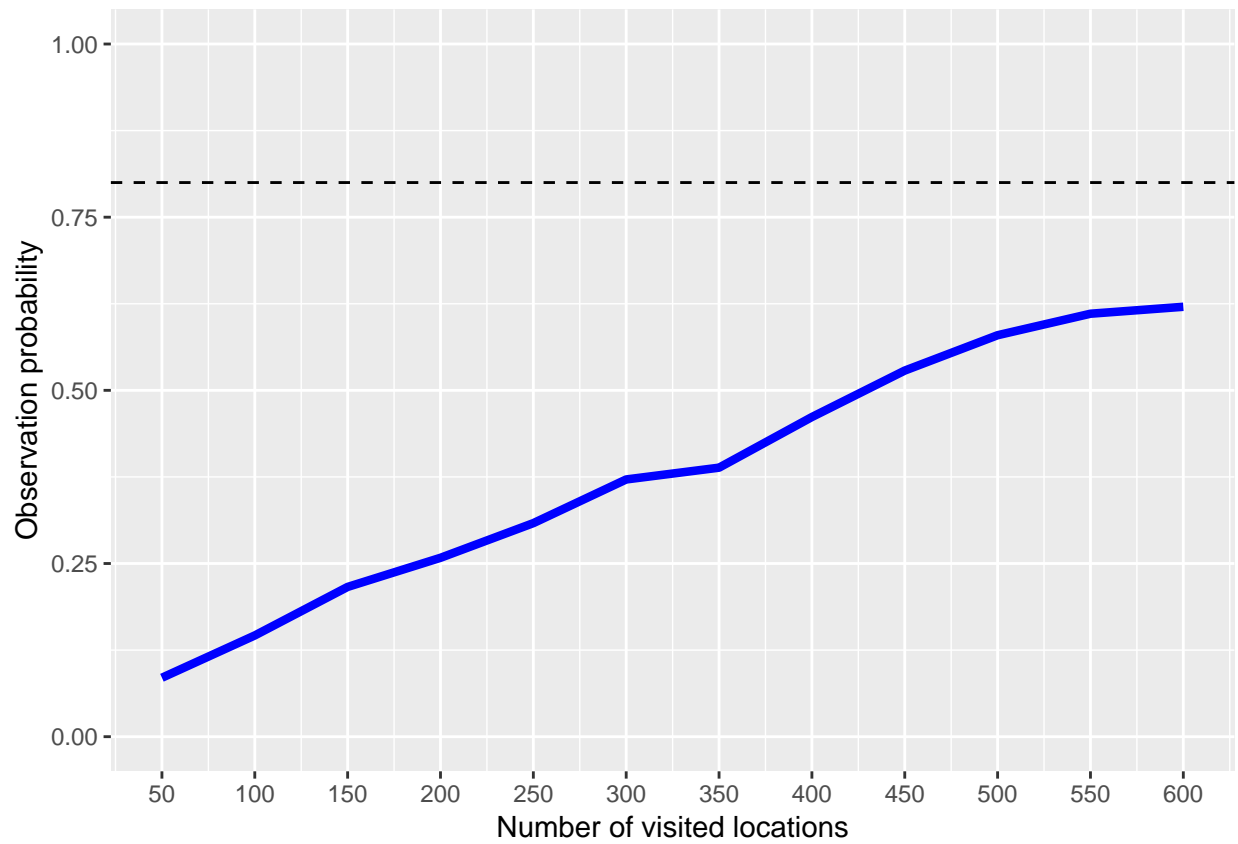


Figure 11: Estimated observation probability of an alien vascular plant species occurring in 200 1x1km grid cells as a function of the number of visited locations. Occurrences and location selection is based on the same weights, modelled from actual alien vascular plant species occurrences. Each visit has a 0.8 probability of detecting the species if present. The threshold of the desired observation probability of 0.8 is shown as a dashed line.

noLocations	probObs
50.00	0.09
100.00	0.15
150.00	0.22
200.00	0.26
250.00	0.31
300.00	0.37
350.00	0.39
400.00	0.46
450.00	0.53
500.00	0.58
550.00	0.61
600.00	0.62



In the case of 200 occupied cells, we reach an overall observation probability of 0.8 after about 500 visited sites with an observation probability of 0.8. With only 100 occupied cells, we don't reach the target of 0.8 even after 600 visited locations.

We can see the effect of a lower detection probability.

## INSERT 0.3 PROB

### Surveying 250x250m squares

In practice, it can be challenging to survey even a 1x1km grid with any respectable observation probability. We might therefore subdivide the squares in smaller units, with the result that our detection probability decreases from simply not covering the place where the species occupies. If a species just occurs in one out of 16 250x250 squares within a 1x1km square and we visit only one such smaller square, our detection probability drops to 1/16 of the former level. How much the detection probability drops of course depend on the aggregation pattern of the species, i.e. how much of the 1x1km cell it is present in. On a tangent, the number of occurrences we consider to be acceptable within an “early detection” framework depends on the scale of the grid cells considered, and the way that the species are aggregated. 500 out of ca 50 000 1x1km cells constitutes occurrences in 1% of the cells. If we accept a 1% occurrence in the about 800 000 250x250 cells, that amounts to 8000 occurrences. As long as these are not extremely aggregated, this could hardly be seen as an early establishment phase.

```
load("predOcc500Det0.05.Rdata")
load("predOcc500Det0.05Vis4.Rdata")
load("predOcc500Det0.05Vis10.Rdata")
load("predOcc500Det0.025Vis4.Rdata")
load("predOcc500Det0.025Vis10.Rdata")
load("predOcc100Det0.05Vis10.Rdata")
load("predOcc100Det0.05Vis4.Rdata")
load("predOcc100Det0.05Vis1.Rdata")
```

### Case of 500 1x1km cells occupied, but we survey only 250x250m subsquares

We start with the case of 500 occurrences spread out in the 50 000 1x1km cells, but when we visit just a 16th of these cells once.

```
plot(predOcc500Det0.05,
      threshold = 0.8) +
  ggtitle("500 occurrences, 0.05 observation probability, 1 visit per site")
```

It is clear that these conditions does not let us reach the desired detection probability of 0.8 even with a great number of sampled locations.

```
plot(predOcc500Det0.05Vis4,
      threshold = 0.8) +
  ggtitle("500 occurrences, 0.05 observation probability, 4 visit per site")
```

```
plot(predOcc500Det0.05Vis10,
      threshold = 0.8) +
  ggtitle("500 occurrences, 0.05 observation probability, 10 visit per site")
```

We can continue to explore the possibilities with a lower total occurrence, for example when a species is present in 100 out of the 50 000 1x1km grid cells.

```
plot(predOcc100Det0.05Vis1,
      threshold = 0.8) +
  ggtitle("100 occurrences, 0.05 observation probability, 1 visit per site")
```

This situation leaves us with very slim chances of detecting the species. But what happens when we increase the number of visits per cell?

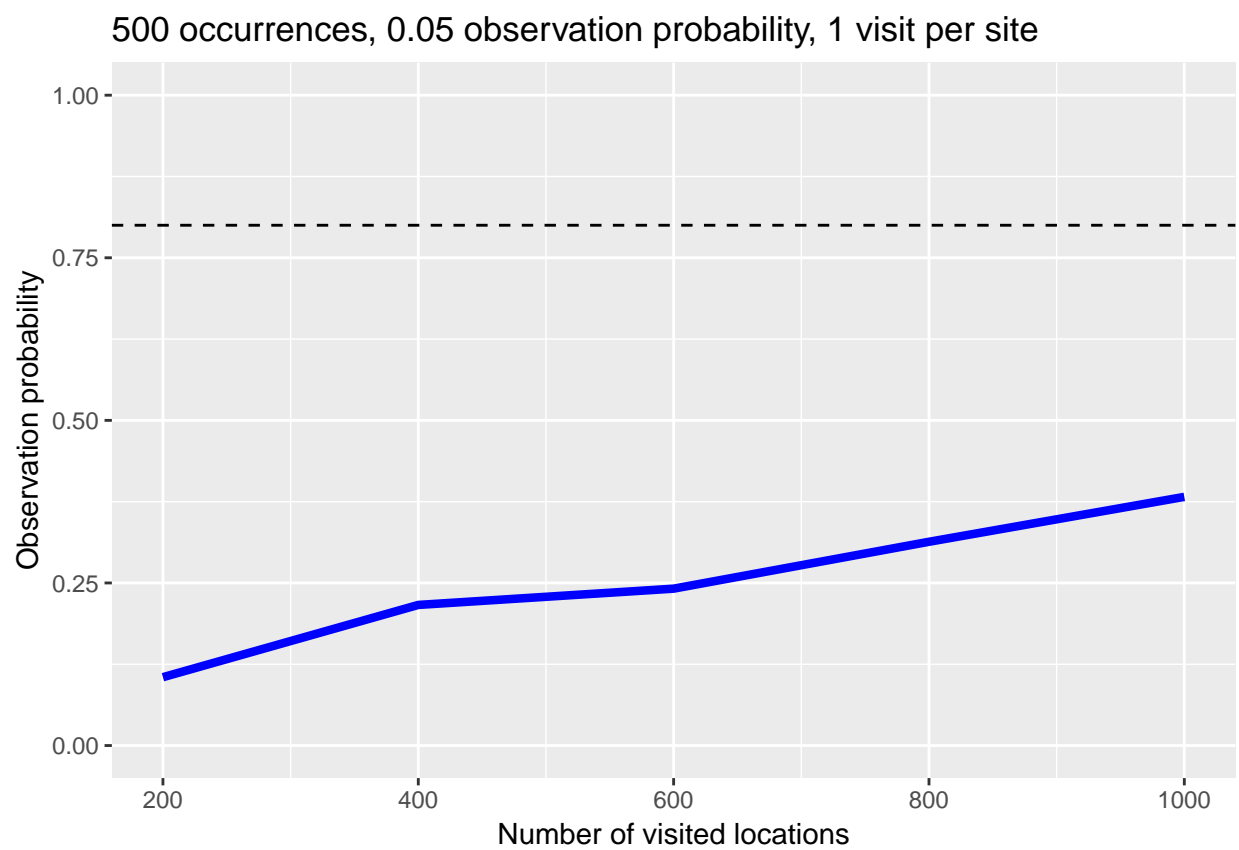


Figure 12: Probability of detecting a species at least once that is present in 500 of the 1x1km cells, but we survey only a 1/16 of the cell.



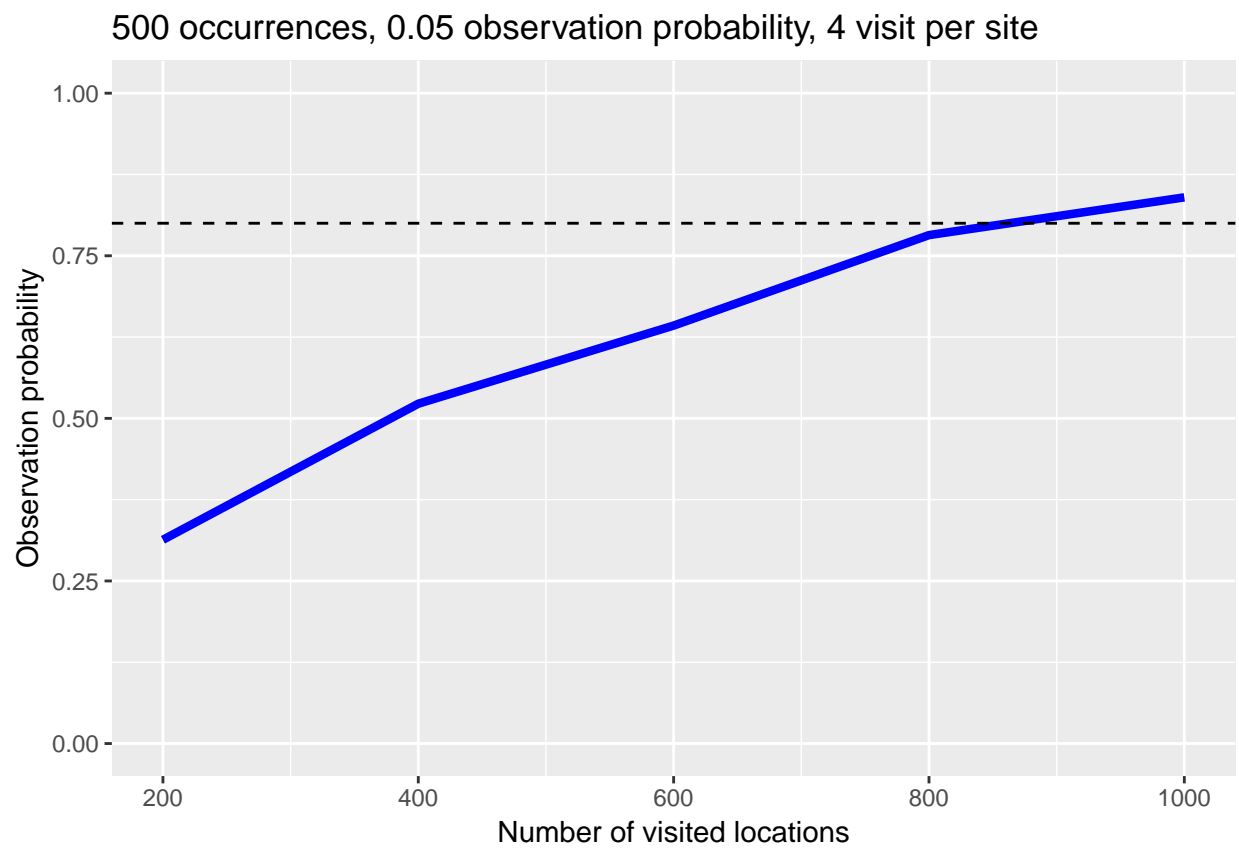


Figure 13: Probability of detecting a species at least once that is present in 500 of the 1x1km cells, but we survey only a 1/16 of the cell. Here we visit the 1x1km cell 4 times.

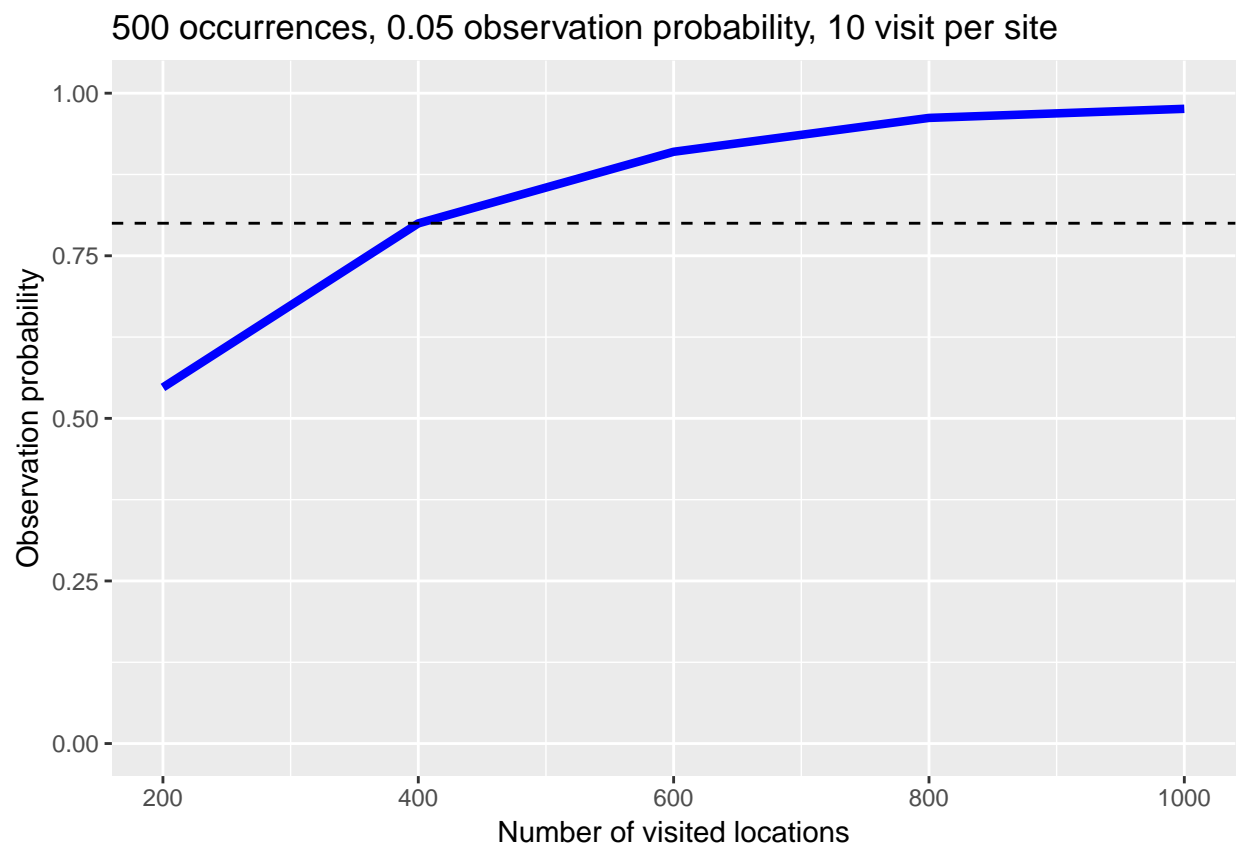


Figure 14: Probability of detecting a species at least once that is present in 500 of the 1x1km cells, but we survey only a 1/16 of the cell. Here we visit the 1x1km cell 10 times.

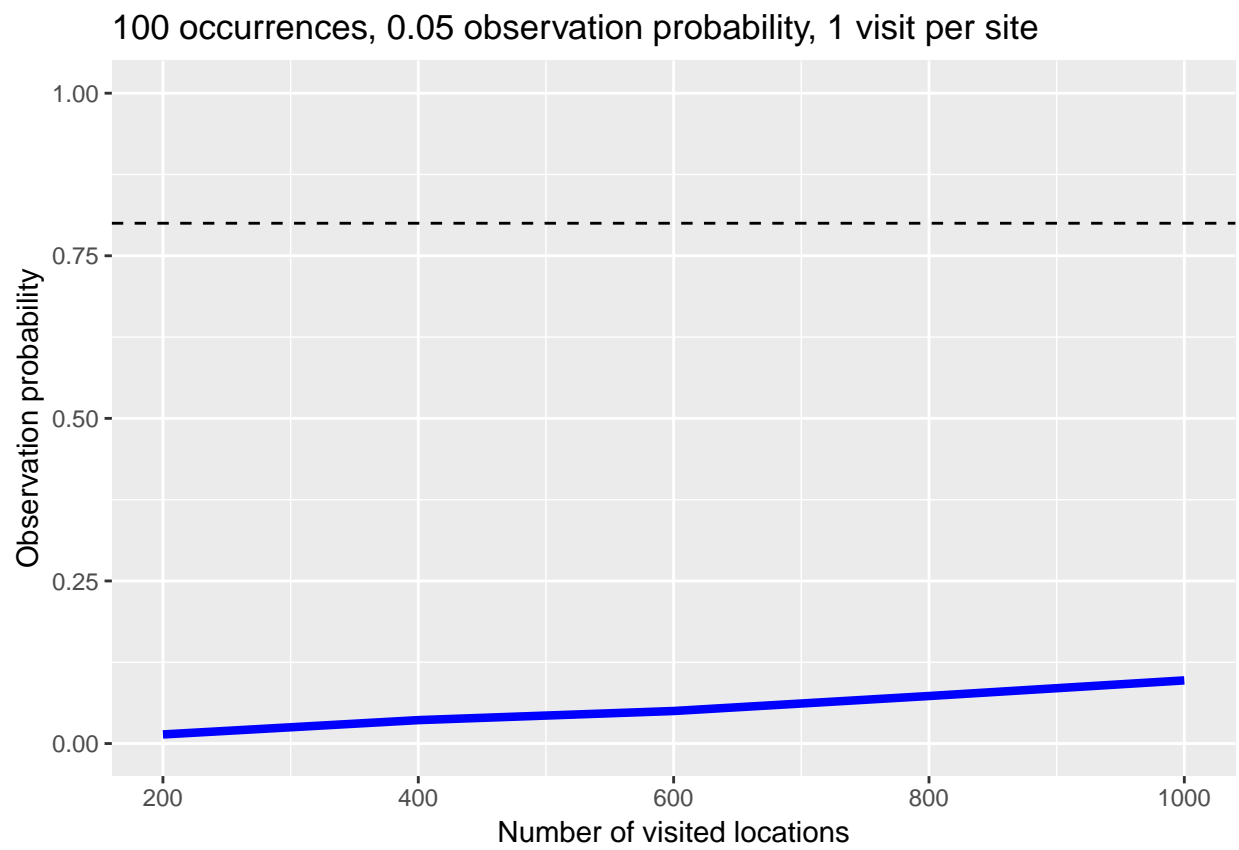


Figure 15: Probability of detecting a species at least once that is present in 500 of the 1x1km cells, but1 we survey only a 1/16 of the cell.

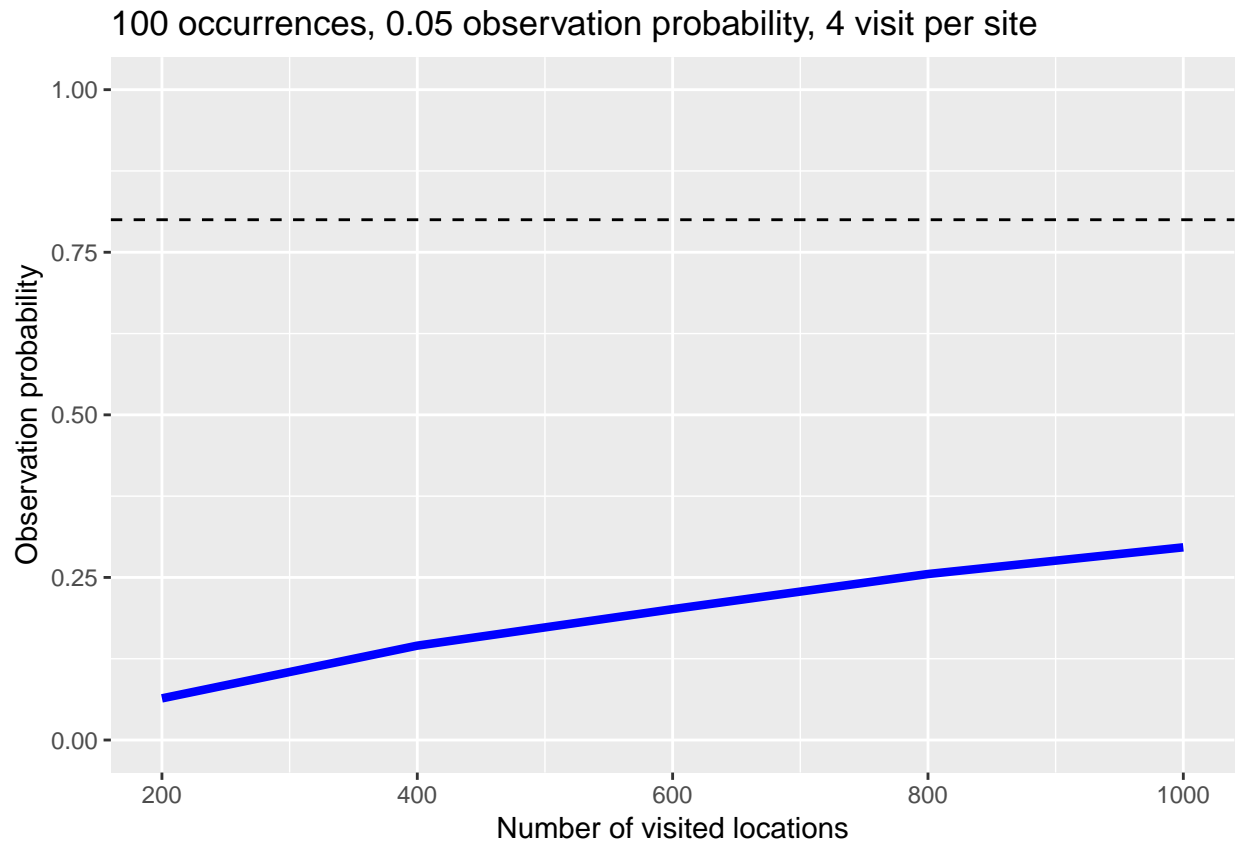
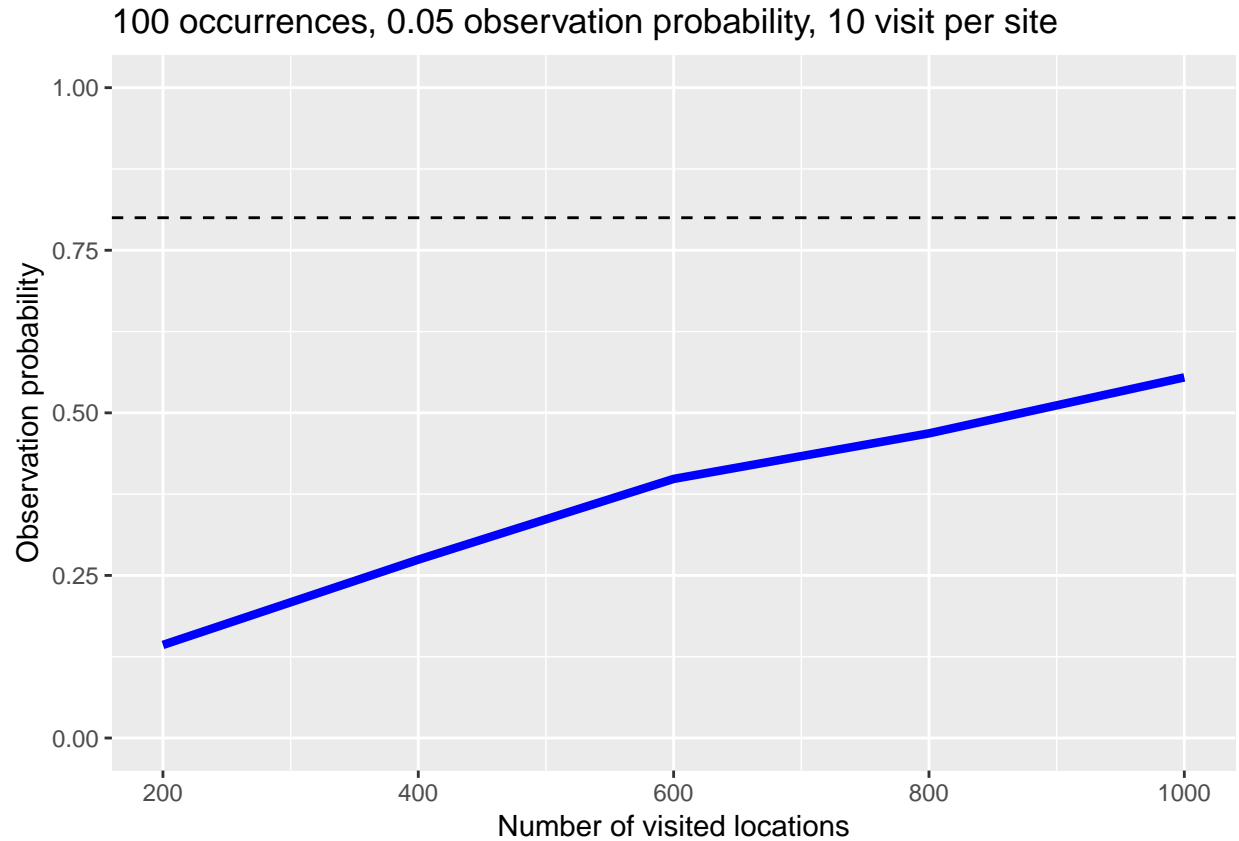


Figure 16: Probability of detecting a species at least once that is present in 100 of the 1x1km cells, but we survey only a 1/16 of the cell. Here we visit the 1x1km cell 4 times.

```
plot(predOcc100Det0.05Vis4,  
      threshold = 0.8) +  
  ggtitle("100 occurrences, 0.05 observation probability, 4 visit per site")  
  
plot(predOcc100Det0.05Vis10,  
      threshold = 0.8) +  
  ggtitle("100 occurrences, 0.05 observation probability, 10 visit per site")
```



With 4 or even 10 visits, we still don't reach the desired detection probability with so few number of occurrences, when we only manage to sample 1/16 of the squares at a time.

## References

Pavlos S. Efraimidis, Paul G. Spirakis, Weighted random sampling with a reservoir, Information Processing Letters, Volume 97, Issue 5, 16 March 2006, Pages 181-185, ISSN 0020-0190, 10.1016/j.ipl.2005.11.003.