

CREEM, Univ of St Andrews: Distance sampling on-line workshop

Analysis in R: Analysis of classic duck nest data

March 2022

1 Using survey data

We use field data to fit different detection function models and estimate density and abundance. The data were collected during line transect surveys of duck nests in Monte Vista National Wildlife Refuge, Colorado, USA in 1967 and 1968. Twenty transects, each 25.75km in length were walked 5 times over the two years. Total transect length of 128.75km (25.75×5) and a distance out to 2.4m was searched. Consult [Anderson and Pospahala \(1970\)](#) for a description of the survey. Distances of detected nests have been provided in a 'csv' text file in a basic format required by 'Distance'. The columns in the the file are:

- Study.Area - name of the study region (Monte Vista NWR)
- Region.Label - identifier of regions or strata (in this case there is only one region and it is set to 'Default')
- Area - size of the stratum
- Sample.Label - line transect identifier
- Effort - length of each transect
- distance - perpendicular distances (m).

The distances allow different key functions/adjustments to be fitted in the detection function model and, by including the transect lengths and area of the region, density and abundance can be estimated.

2 Objectives of the practical

1. Import a text file
2. Understand the structure of a data frame
3. Fit different key functions/adjustments in the detection function model using `ds`
4. Examine the results of an analysis, i.e. `ddf` and `dht` components of a `dsmodel` object

3 Importing the data

The file containing the duck nest survey data is located on the online workshop website. Either download the file following [this link](#) or the content of the data file can be read directly from the web into an R ([R Core Team, 2018](#)) data frame named `nests` via the following command

```
nests <- read.csv(file="https://workshops.distancesampling.org/online-course/exercisepdfs/Ch7/dataset", header=TRUE)
```

(Solutions)



The Monte Vista National Wildlife Refuge.

The URL of the file location is quite long, we repeat it here so it is more legible.

```
{https://workshops.distancesampling.org/online-course/exercisepdfs/Ch7/datasets/ducks-area-effort.csv}
```

This command is made up of several components:

- `read.csv` is a function to read data files of type 'csv' (comma-separated values),
- the function has two arguments specified; `file` specifies the name of the data file and `header=TRUE` specifies that the first row of the data file contains the names of the data columns.
- the `<-` symbol has assigned the data set to an object called `nests`. Note that there is now an object called `nests` listed on the 'Environment' tab.

To check that the data file has been read into R correctly, use the `head` and `tail` 'functions' to look at the top and bottom rows of the data, respectively. To look at the first few rows of `nests` type the following command.

```
head(nests, n=2)
```

```
##   Region.Label Area Sample.Label Effort distance
## 1      Default 40.47           1 128.75    0.06
## 2      Default 40.47           1 128.75    0.07
```

The `head` function as used above displays the first 6 records of the named object. The argument `n` controls the number of rows to display. To display the *last* 2 records in the data, type the command:

```
tail(nests, n=2)
```

```
##   Region.Label Area Sample.Label Effort distance
## 533      Default 40.47           20 128.75    2.38
## 534      Default 40.47           20 128.75    2.13
```

The object `nests` is a dataframe object made up of rows and columns. Use the function `dim` to find out the dimensions of the data set (i.e. the total number of rows and columns):

```
dim(nests)
```

```
## [1] 534  5
```

Another way to look at a data frame is to move to the 'Environment tab' in R-Studio and click on the rectangle (with the grid); this opens a new tab showing the data.

4 Summarising the perpendicular distances

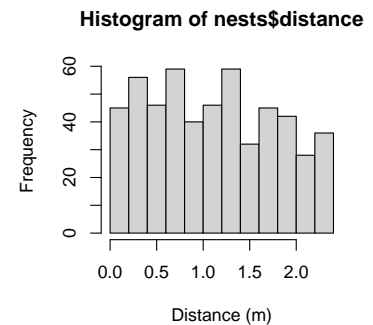
To access an individual column within a data frame use the `$` symbol, for example to summarise the distances:

```
summary(nests$distance)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  0.010  0.540   1.080   1.117   1.670   2.400
```

Similarly to plot the histogram of distances, the command is:

```
hist(nests$distance, xlab="Distance (m)")
```



5 Fitting different models

To use the `ds` function, load the `Distance` package (Miller, 2017).

The function `ds` requires a data frame to have a column called `distance`, we specify the name of the data frame as follows:

```
library(Distance)
```

```
## Loading required package: mrds
```

```
## This is mrds 2.2.6
```

```
## Built: R 4.1.3; ; 2022-03-17 18:30:31 UTC; windows
```

```
##
```

```
## Attaching package: 'Distance'
```

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

```
##
```

```
##      create.bins
```

```
conversion <- convert_units("meter", "kilometer", "square kilometer")
```

```
nest.model1 <- ds(nests, key="hn", adjustment=NULL, convert_units = conversion)
```

```
## Fitting half-normal key function
```

```
## Key only model: not constraining for monotonicity.
```

```
## AIC= 928.134
```

The `convert_units` argument ensures that the correct units are specified - in this example, distances are in metres, lengths in km and the area in km^2 . Think of this argument as a divider used to transform units of transect effort into units of perpendicular distance (e.g., $1\text{km} / 0.001 = 1000\text{m}$). The `convert_units` function performs these calculation for you if you provide the units in which perpendicular distances, transects and study area size are recorded in the data.

This call to `ds` fits a half-normal key function with no adjustment terms. Summarise the fitted model:

```
summary(nest.model1$ddf)
```

```
##
## Summary for ds object
## Number of observations : 534
## Distance range       : 0 - 2.4
## AIC                  : 928.1338
##
## Detection function:
## Half-normal key function
##
## Detection function parameters
## Scale coefficient(s):
##           estimate      se
## (Intercept) 0.9328967 0.1703933
##
##           Estimate      SE      CV
## Average p      0.8693482 0.03902053 0.04488481
## N in covered region 614.2533225 29.19683065 0.04753223
```

Plot the detection function with the histogram having 12 bins:

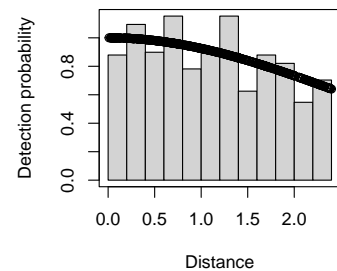
```
plot(nest.model1, nc=12)
```

To fit different detection functions, change the key and adjustment arguments. For example to fit a half-normal key function with cosine adjustment terms, use the command:

```
nest.model2 <- ds(nests, key="hn", adjustment="cos", convert_units = conversion)
```

By default, AIC selection will be used to fit adjustment terms of up to order 5.

```
##
## Summary for ds object
## Number of observations : 534
## Distance range       : 0 - 2.4
## AIC                  : 928.1338
##
## Detection function:
## Half-normal key function
##
## Detection function parameters
## Scale coefficient(s):
##           estimate      se
## (Intercept) 0.9328967 0.1703933
##
##           Estimate      SE      CV
## Average p      0.8693482 0.03902053 0.04488481
## N in covered region 614.2533225 29.19683065 0.04753223
```



Question: Have any adjustment terms been selected?

Answer: No adjustment terms have been included in the preferred model, because there is only a single row in the list of Detection function parameters.

To fit a hazard rate key function with Hermite polynomial adjustment terms, then use the command:

```
nest.model3 <- ds(nests, key="hr", adjustment="herm", convert_units = conversion)
```

```
## Error :
## gosolnp-->Could not find a feasible starting point...exiting
```

```
summary(nest.model3$ddf)
```

```
##
## Summary for ds object
## Number of observations : 534
## Distance range       : 0 - 2.4
## AIC                  : 929.7934
##
## Detection function:
## Hazard-rate key function
##
## Detection function parameters
## Scale coefficient(s):
##           estimate      se
## (Intercept) 0.9190194 0.2081124
##
## Shape coefficient(s):
##           estimate      se
## (Intercept) 0.2899024 0.6393472
##
##           Estimate      SE      CV
## Average p      0.8890651 0.04957233 0.05575782
## N in covered region 600.6309174 34.59069491 0.05759060
```

Use the `help` command to find out what other key functions and adjustment terms are available.

6 The `ds` object

The objects created with `ds` (e.g. `nest.model1`) are made up of two parts. We can list them using the `names` function as below:

```
names(nest.model1)
```

```
## [1] "ddf" "dht" "call"
```

The detection function information is in the `ddf` part and the density and abundance estimates are stored in the `dht` part. To access each part, then the `$` can be used (as with columns in a data frame). For example to see what information is stored in the `ddf` part, we can use the `names` function again:

Question: Is there much difference in the probability of detecting a nest given it is within the maximum detection distance between the three models fitted to the duck nest data?

Answer: No, there is barely any difference between the estimated probability of detection. For the first two models, $\hat{P}_a = 0.869$ and for the hazard rate model, $\hat{P}_a = 0.889$.

```
names(nest.model1$ddf)
```

```
## [1] "call"          "data"          "model"         "meta.data"    "control"
## [6] "method"       "ds"           "par"          "lnl"         "hessian"
## [11] "dsmodel"     "criterion"    "fitted"       "Nhat"        "name.message"
```

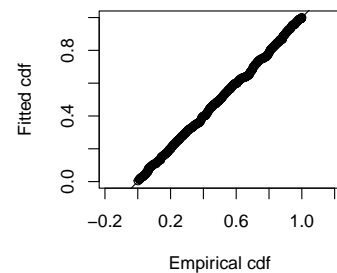
7 Goodness of fit

Before making inference from the detection function we have fitted, we should evaluate the model. First assessment is goodness of fit, accomplished using the function `gof_ds`:

```
gof_ds(nest.model1)
```

```
##
## Goodness of fit results for ddf object
##
## Distance sampling Cramer-von Mises test (unweighted)
## Test statistic = 0.0353634 p-value = 0.955416
```

Calling the function `gof_ds` with the default arguments and exact distance data, a *q-q* plot is produced along with the unweighted Cramer-von Mises goodness of fit test.



Question: Interpret the *q-q* plot and CvM test results for the duck nest data.

Answer: The half-normal detection function fits the duck nest data very well. All points of the *q-q* plot fall on the diagonal line and the *p*-value associated with the CvM test statistic is very large ($\gg 0.05$) indicating a good fit of the model to the data.

8 Estimating density and abundance

So far, we have concentrated on the detection function but, with more information such as transect lengths and the area of the region, we can estimate density and abundance. The second component of a `ds` object, contains this additional information. This information can be viewed with:

```
str(nest.model1$dht$individuals, max=1)
```

```
## List of 8
## $ bysample      :'data.frame':  20 obs. of  9 variables:
## $ summary      :'data.frame':  1 obs. of  9 variables:
## $ N            :'data.frame':  1 obs. of  7 variables:
## $ D            :'data.frame':  1 obs. of  7 variables:
## $ average.p    : num 0.869
## $ cormat       : num [1, 1] 1
## $ vc          :List of 3
## $ Nhat.by.sample:'data.frame':  20 obs. of  9 variables:
```

This `dht` object contains considerable information. However, focus upon three tables generated by the `summary()` function: `summary`, `abundance` and `density`. Dig more deeply into the content of these tables.

8.1 Summary information

This provides information about the survey:

- size of study area,
- area covered by sampling effort
- length of all transects
- number of detections
- number of transects
- encounter rate (ER) number of detections per unit transect length and its associated variability

```
nest.model1$dht$individuals$summary
```

Region	Area	CoveredArea	Effort	n	k	ER	se.ER	cv.ER
Default	40.47	12.36	2575	534	20	0.207	0.008	0.038

8.2 Abundance estimates

Estimated density multiplied by the size of the study area.

```
nest.model1$dht$individuals$N
```

Label	Estimate	se	cv	lcl	ucl	df
Total	2011.23	118.85	0.06	1788.91	2261.19	99.56

8.3 Density estimates

Density estimated using the formula

$$\hat{D} = \frac{n}{a\hat{P}_a}$$

where n (number of nests)=534, a (covered area)=12.36 and \hat{P}_a (probability of detection)=0.8693

```
nest.model1$dht$individuals$D
```

Label	Estimate	se	cv	lcl	ucl	df
Total	49.7	2.94	0.06	44.2	55.87	99.56

References

Anderson, D. R. and R. S. Pospahala. 1970. Correction of bias in belt transect studies of immotile objects. The Journal of Wildlife Management, **34**:141–146. URL <http://www.jstor.org/stable/3799501>.

Question: Compute (by hand) the density estimate resulting from the estimated probability of detection arising from the hazard rate detection function: $\hat{P}_a = 0.8891$.

Answer: $\hat{D} = \frac{534}{12.36 \cdot 0.8891} = 48.59 \text{ nests} \cdot \text{km}^{-2}$ compared to 49.70 nests·km⁻² (a difference of 2.3%). In other words, the difference in estimated density is quite small; particularly when uncertainty is taken into account.

- Buckland, S. T., E. A. Rexstad, T. A. Marques, and C. S. Oedekoven. 2015. Distance Sampling: Methods and Applications. Springer. URL <https://www.springer.com/gb/book/9783319192185>.
- Miller, D. L. 2017. Distance: Distance Sampling Detection Function and Abundance Estimation. URL <https://CRAN.R-project.org/package=Distance>. R package version 0.9.7.
- R Core Team. 2018. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.