

# Distance sampling online workshop

Analysis in R: Analysis of multi-species surveys

CREEM, Univ of St Andrews – October 2018

## 1 More complex analyses

This practical is based on the Montrave songbird case study in [Buckland et al. \(2015, Section 5.2.2.3\)](#), with computer code under [Montrave songbird case study](#). Both point and line transect surveys were conducted and here we use the data from the **line transect** data, although the issues (and solutions) will be similar.

These data are provided in a 'flat file' format (i.e. it contains all the necessary columns to estimate a detection function, density and abundance). While both formats are equally valid, the 'flat file' approach has a particular idiosyncrasy which we exploit here to introduce more functions and data manipulation.

Several species of birds were identified but not all species were detected on all transects. If a simple data selection is performed to select records for a particular species, then not all of the transects will be included in the resulting data (because that species may not have been seen). This doesn't matter if we are only interested in fitting detection functions, but will matter if we wish to estimate density and abundance because the effort will be too low since some of the transects are missing. To correct for this, some data frame manipulation is required. There is generally more than one way to do something in R ([R Core Team, 2018](#)) - for an alternative way see the computer code 'Montrave song bird case study' associated with [Buckland et al. \(2015\)](#), as well as Section 9 below.

## 2 Objectives of the practical

1. Data frame selection and manipulation
2. Extracting estimates from dht object
3. Customising detection function plots
4. Improve re-usability of code with functions

## 3 Importing the data

The data is in a 'flat file' format and contains the following columns:

- Region.Label - name of study
- Area - size of study region (km<sup>2</sup>)
- repeats - number of visits to transect
- Sample.Label - line transect identifier
- Effort - length of transect (km)
- distance - perpendicular distance (m)
- species - species of bird (c=chaffinch, g=great tit, r=robin and w=wren)
- visit - on which visit bird was detected.

## Solutions



European robin (*Erithacus rubecula*); one of the species in the Montrave study ([Buckland, 2006](#)).



Aerial view of Montrave study area. White diagonal lines represent transects walked for data analysed here.

Use the following command to import the data from the website associated with [Buckland et al. \(2015\)](#) and then use the head command to examine it.

```
birds <- read.csv("https://synergy.st-andrews.ac.uk/ds-manda/files/2016/11/montrave-line.csv",
  header = TRUE)
```

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

```
## Region.Label Area repeats Sample.Label
## 1 Montrave 33.2 2 1
## 2 Montrave 33.2 2 1
## Effort distance species visit
## 1 0.208 75 c 1
## 2 0.208 40 c 1
```

**Question:** Explore the data. How many transects are there?

```
length(unique(birds$Sample.Label))
```

```
## [1] 19
```

For now, save the transect labels to a new object as we will use them later on:

```
tran.lab <- unique(birds$Sample.Label)
```

The table command is a quick way to determine how many detections there are of each species:

```
table(birds$species)
```

```
##
## c g r w
## 73 32 82 156
```

As a hint of things to come, create a two-way table showing the number of detections by transect and by species. If there are zeroes in this table, it will create a challenge.

```
with(birds, table(species, Sample.Label))
```

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
c	4	7	7	5	9	7	5	1	1	3	1	0	2	2	4	3	4	7	1
g	0	2	3	5	5	1	3	2	1	1	0	0	2	0	2	0	3	1	1
r	3	8	11	5	8	7	5	7	3	1	0	2	5	6	4	0	4	3	0
w	10	11	12	11	14	12	13	11	6	3	1	4	9	12	9	2	7	6	3

Each of the line transects was visited twice which is not taken into account at present. However, it is straightforward to do so:

```
birds$Effort <- birds$Effort * birds$repeats
```

## 4 Manipulating the robin data

For the purposes of this practical, we are interested in estimating the density of robins and so we select only these records:

```
robins <- birds[birds$species == "r", ]
```

**Question:** On how many transects were robins detected?

```
length(unique(robins$Sample.Label))
```

```
## [1] 16
```

If we were to use the `robins` data as it is at present to estimate density, then density would be **incorrect** because the search effort associated with three transects is missing. Adding these missing transects to the `robins` data, requires several steps:

1. identify the missing transects,
2. select the information for the missing transects,
3. get the missing information in the correct format,
4. add the missing information to the `robins` data.

The following commands identifies the missing commands. After each command, type the name of the object which has been created to see what each command has done.

```
robin.lab <- unique(robins$Sample.Label)
miss.lab <- tran.lab[!is.element(el = tran.lab,
  set = robin.lab)]
```

To understand what this latter command has done, it can be broken down into several elements:

- elements of `tran.lab` are selected using `[]`
- the `is.element` function (without the `!` symbol) selects the elements in `tran.lab`, which are also in the `set` argument (i.e. `robin.lab`)
- the `!` is used to select the elements in `tran.lab` that are NOT in `robin.lab`.

```
## Robins were detected on the following transects:
```

```
## [1] 1 2 3 4 5 6 7 8 9 10 12 13 14
## [14] 15 17 18
```

```
##
```

```
## Therefore missing transects are:
```

```
## [1] 11 16 19
```

Now we know which transects are missing, we can select these records from the `birds` data frame:

```
miss.data <- birds[is.element(birds$Sample.Label,
                             miss.lab), ]
```

However, the information about the transects are repeated in this new data frame because we have just selected all records for these transects. A quick check of the number of rows will confirm this:

```
length(miss.data$Sample.Label)
```

To get rid of rows where Sample.Label is duplicated use the command:

```
miss.data <- miss.data[!duplicated(miss.data$Sample.Label),
                        ]
```

This command has selected the records from miss.data for which the transect label is not duplicated.

We only want to keep the information about search effort and so data in the distance, species and visit columns are set to missing:

```
miss.data$distance <- rep(NA, length(miss.lab))
miss.data$species <- rep("NA", length(miss.lab))
miss.data$visit <- rep(NA, length(miss.lab))
```

Examine miss.data.

```
miss.data

##      Region.Label Area repeats Sample.Label
## 234      Montrave 33.2        2           11
## 299      Montrave 33.2        2           16
## 339      Montrave 33.2        2           19
##      Effort distance species visit
## 234  0.078         NA      NA      NA
## 299  0.378         NA      NA      NA
## 339  0.040         NA      NA      NA
```

The final thing to do is to add the missing data (miss.data) to the robins data frame using the rbind function (this combines data frames with the same columns).

```
robins <- rbind(robins, miss.data)
```

Let's see the result of all this manipulation:

```
tail(robins, n = 4)

##      Region.Label Area repeats Sample.Label
## 334      Montrave 33.2        2           18
## 234      Montrave 33.2        2           11
## 299      Montrave 33.2        2           16
## 339      Montrave 33.2        2           19
```

```
##      Effort distance species visit
## 334  0.400      70      r      2
## 234  0.078      NA      NA     NA
## 299  0.378      NA      NA     NA
## 339  0.040      NA      NA     NA
```

If we wanted to be very tidy, then the data frame could be sorted so that the transect labels were in order:

```
robins <- robins[order(robins$Sample.Label), ]
```

## 5 Analysis

Before we fit any models, have a quick look at the histogram of distances:

```
hist(robins$distance, breaks = 20)
```

Consistent with [Buckland et al. \(2015\)](#), three detection functions are fitted using the `ds()` function in the R package *Distance* ([Miller, 2017](#)):

```
library(Distance)
robin.hn.herm <- ds(robins, truncation = 95, transect = "line",
  key = "hn", adjustment = "herm", convert.units = 0.1)
robin.uni.cos <- ds(robins, truncation = 95, transect = "line",
  key = "unif", adjustment = "cos", convert.units = 0.1)
robin.haz.simp <- ds(robins, truncation = 95,
  transect = "line", key = "hr", adjustment = "poly",
  convert.units = 0.1)
```

```
summarize_ds_models(robin.hn.herm, robin.uni.cos,
  robin.haz.simp)
```

Histogram of robins\$distance

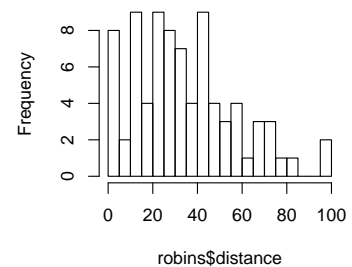


Figure 1: Perpendicular distances of robins in Montrave study.

**Question:** What is the preferred model for the robin data?

Table 2: Model selection for robin data from Montrave line transect survey.

Model	C-vM p-value	$\hat{P}_a$	$se(\hat{P}_a)$	$\Delta AIC$
robin.uni.cos	0.510	0.636	0.103	0.000
robin.haz.simp	0.732	0.679	0.053	0.565
robin.hn.herm	0.463	0.626	0.121	0.617

**Note:** All three detection function fit the data (based upon the C-vM test of exact distances). The estimated detection probability is very similar for all models, and the  $\Delta AIC$  values of all models is  $< 1$ . Hence all models will give very similar estimates of density.

## 6 Examining the dht object

The fitted model object (e.g. `robin.uni.cos`) is made up of two parts; the detection function in the `ddf` part and the estimates in the `dht` part. In this section, we look at the `dht` part.

To list the elements that are contained in `dht`, use the `names` function:

```
names(robin.uni.cos$dht)
```

```
## [1] "individuals"
```

Detections were of individual birds and so group size was not included in these data - if it had been included (in a column called size), then as well as individuals there would have been elements clusters and Expected.S.

The estimates stored in the individuals object can be listed in a similar manner:

```
names(robin.uni.cos$dht$individuals)
```

```
## [1] "bysample"      "summary"
## [3] "N"             "D"
## [5] "average.p"     "cormat"
## [7] "vc"            "Nhat.by.sample"
```

To collect together the density estimates (and estimates of precision) from all the fitted models, we can use the following command:

```
model.results <- rbind(robin.uni.cos$dht$individuals$D,
  robin.haz.simp$dht$individuals$D, robin.hn.herm$dht$individuals$D)
```

**Question:** Examine the three sets of density estimates to see if the previous suggestion (that the density estimates are similar) is confirmed.

```
model.results
```

Table 3: Density estimates for Montrave robins under three fitted detection functions.

Label	Estimate	se	cv	lcl	ucl	df
Total	0.6856845	0.1316387	0.1919814	0.4698670	1.0006305	89.83772
Total	0.6418461	0.0829831	0.1292881	0.4948932	0.8324351	41.07649
Total	0.6964545	0.1527862	0.2193772	0.4528707	1.0710540	95.26086

## 7 Goodness of fit

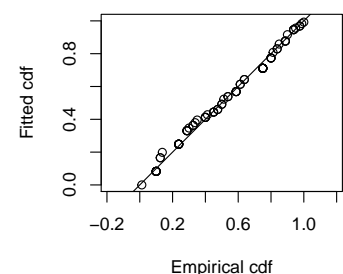
Here we look at goodness of fit test with unequal bin intervals and just consider one of the fitted models. First we specify the required bin intervals.

```
robin.brks <- c(0, 12.5, 22.5, 32.5, 42.5, 52.5,
  62.5, 77.5, 95)
```

Perform the tests using both exact distance data for the Cramer-von Mises test and specified breakpoints for  $\chi^2$  test for the uniform-cosine model that had the (slightly) smallest AIC score.

```
gof_ds(robin.uni.cos, breaks = robin.brks, chisq = TRUE,
  main = "QQ plot unif-cos for robins")
```

QQ plot unif-cos for robins



```
##
## Goodness of fit results for ddf object
##
## Chi-square tests
##           [0,12.5] (12.5,22.5] (22.5,32.5]
## Observed  11.000000  15.0000000  15.0000000
## Expected  16.553595  13.1496015  12.7439652
## Chisquare   1.863185   0.2603862   0.3993806
##           (32.5,42.5] (42.5,52.5]
## Observed   10.0000000  13.0000000
## Expected   11.7483134  10.0013304
## Chisquare   0.2601735   0.8990823
##           (52.5,62.5] (62.5,77.5]
## Observed    7.00000000  7.00000000
## Expected    7.58791419  6.35139428
## Chisquare    0.04555179  0.06623575
##           (77.5,95]      Total
## Observed    2.000000000  80.000000
## Expected    1.863886388  80.000000
## Chisquare    0.009939938  3.803935
##
## P = 0.57798 with 5 degrees of freedom
##
## Distance sampling Cramer-von Mises test (unweighted)
## Test statistic = 0.116508 p-value = 0.509846
```

**Note:** The detections fall close to the diagonal line of the qq plot, suggesting an adequate fit for the uniform cosine model. The *p*-value of the Cramer-von Mises test (at bottom of printout) confirms this. Similarly the *p*-value for the  $\chi^2$  test also suggests an adequate fit.

## 8 Customising the detection function plot

The `plot` function provides a basic plot of the fitted detection function overlaid onto the scaled distribution of distances:

```
plot(robin.uni.cos)
```

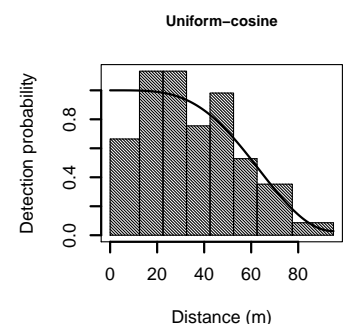
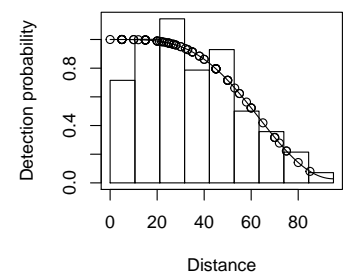
However, the plot can be customised for reporting:

```
plot(robin.uni.cos, showpoints = FALSE, black.white = TRUE,
     pl.den = 50, lwd = 2, breaks = robin.brks,
     main = "Uniform-cosine", xlab = "Distance (m)")
```

The arguments are:

- `showpoints` - logical indicating whether observed distances are shown
- `lwd` - line width (1=default)
- `pl.den` - density of shading of histogram (0=no shading)

For other options see `help(plot.ds)` (Note `plot` is a generic function which selects a relevant type of plot based the object).



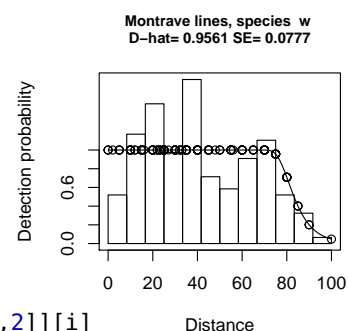
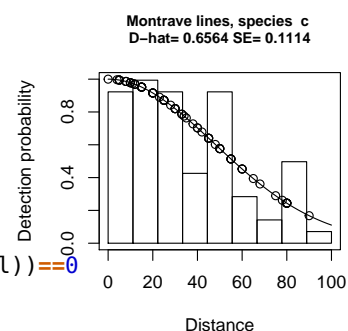
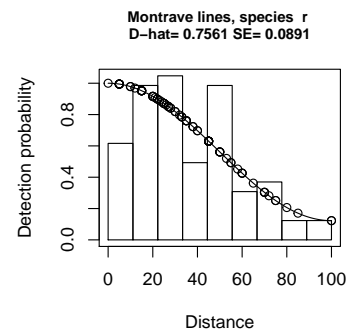
## 9 Advanced: modularising R code to work with multiple species

When analysing a multi-species survey, it is likely that the investigator will want to analyse all (or at least many) of the species encountered during the survey. This will necessitate some repetitive calculation, such as accounting for transects without detections and fitting multiple detection functions for each species.

To facilitate the repetitive nature of such analyses, it is useful to take advantage of the programmatic nature of the R language to create *functions* that can be called repeatedly with arguments to accommodate changing circumstances. The code below demonstrates such a modular approach whereby two functions `augment.empty.transects()` and `fit.hn.uni.haz()` are defined to aide in the repeated analyses.

The first function `augment.empty.transects()` performs the data manipulation described in Section 4, with two arguments: the data frame containing the full survey data and the species code on which to subset the data.

```
augment.empty.transects <- function(survey, species) {
  # Purpose: find transects on which species not detected
  #           adjust data file to correct effort
  # Input: raw data file, species on which to subset data
  # Output: data frame with correct effort for model fitting
  # Rextad August 2018
  num.transects <- length(unique(survey$Sample.Label))
  holetab <- as.matrix(table(survey$species, survey$Sample.Label))
  holes <- which(holetab, arr.ind = TRUE)
  if (length(holes[rownames(holes)==species])==0) {
    adj.survey <- survey[survey$species==species, ]
  } else {
    alltranlen <- vector(mode="numeric", length=num.transects)
    for (i in 1:num.transects) {
      alltranlen[i] <- survey$Effort[survey$Sample.Label==i][1]
    }
    empty.transects <- NULL
    for (i in 1:length(holes[rownames(holes)==species,2])) {
      empty.label <- holes[rownames(holes)==species,2][i]
      empty.length <- alltranlen[holes[rownames(holes)==species,2][i]]
      empty.record <- cbind(survey[1,1:3], empty.label, empty.length)
      empty.transects <- rbind(empty.transects, empty.record)
    }
    empty.transects[,c("a", "b", "c")] <- NA
    names(empty.transects) <- names(survey)
    adj.survey <- survey[survey$species==species, ]
    adj.survey <- rbind(adj.survey, empty.transects)
    adj.survey <- adj.survey[order(adj.survey$Sample.Label), ]
  }
  return(adj.survey)
}
```





The second function, `fit.hn.uni.haz()` fits three candidate models to a dataset provided as the first argument. The second argument is the truncation distance. The final argument determines whether the `summarize_ds_models()` table is printed.

```
fit.hn.uni.haz <- function(data, trunc, print=TRUE) {
  # Purpose: fit three key functions to transect data,
  #           perform model selection and
  #           print model selection table
  # Input: data to analyse, truncation distance, print flag
  # Output: fitted model object (class 'dsmodel')
  # Rextad August 2018
  hn.herm <- ds(data, trun=trunc, key="hn", adj="herm", con=.1)
  uni.cos <- ds(data, trun=trunc, key="unif", adj="cos", con=.1)
  haz.simp <- ds(data, trun=trunc, key="hr", adj="poly", con=.1)
  mods <- summarize_ds_models(hn.herm, uni.cos, haz.simp, output="plain")
  if(print) print(knitr::kable(mods))
  names(mods) <- c("mod", "key", "form", "fit", "pa", "sepa", "daic")
  if(mods[1,1]=="hn.herm") {
    result <- hn.herm
  } else {
    if(mods[1,1]=="uni.cos") {
      result <- uni.cos
    } else {
      result <- haz.simp
    }
  }
  return(result)
}
```

The two functions are used in tandem in the calling code below. Note the for loop that iterates through three of the four species detected in the Montrave survey (great tit not analysed because there were few detections).

```
for(species in c("r", "c", "w")) {
  best.model <- fit.hn.uni.haz(augment.empty.transects(birds, species),
                              100, print=FALSE)
  plot(best.model,
        main=paste("Montrave lines, species ", species,
                    "\nD-hat=", round(best.model$dht$individuals$D$Estimate,4),
                    "SE=", round(best.model$dht$individuals$D$se, 4)))
}
```

## References

Buckland, S. T. 2006. Point transect surveys for songbirds: robust methodologies. *The Auk*, **123**:345. URL [https://doi.org/10.1642/0004-8038\(2006\)123\[345:psfsrc\]2.0.co;2](https://doi.org/10.1642/0004-8038(2006)123[345:psfsrc]2.0.co;2).

- 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/>.