

Distance sampling online workshop

Analysis in R: Covariates in the detection function

CREEM, Univ of St Andrews – October 2018

1 Covariates in the detection function

We illustrate fitting multiple covariate distance sampling (MCDS) models to point transect data using a bird survey in Hawaii; data on an abundant species, the Hawaii amakihi (*Hemignathus virens*) is used. This practical is based on the case study shown in Section 5.3.2 of [Buckland et al. \(2015\)](#) which duplicates the analysis presented in [Marques et al. \(2007\)](#). This set of data is included in 'Distance for Windows' as one of the Sample Projects: you can open this project (entitled amakihi.zip) in the 'Sample projects' directory in the 'My Distance projects' directory residing under 'My Documents'. We describe the analysis of these data using Distance in R.



Hawaii amakihi (*Hemignathus virens*)

2 Objectives of this practical

1. Introduce different types of plots to explore covariates
2. Add covariates to the detection function
3. Plot the detection functions.

3 Importing the data

Analysis begins by importing the data from a comma-delimited file (this file was created by copying the data from the amakihi Distance project).

```
# Import Amakihi data
amakihi <- read.csv(file = "https://synergy.st-andrews.ac.uk/ds-manda/files/2016/11/amakihi.csv")
```

Check that it has been imported correctly.

```
head(amakihi, n = 3)
```

```
## Study.Area Region.Label Sample.Label
## 1      Kana      Jul92          1
## 2      Kana      Jul92          1
## 3      Kana      Jul92          1
## Effort distance OBS MAS HAS
## 1      1      40 TJS  50  1
## 2      1      60 TJS  50  1
## 3      1      45 TJS  50  1
```

These data consist of eight columns:

- Study.Area - name of the study area
- Region.Label - survey dates which are used as 'strata'
- Sample.Label - point transect identifier
- Effort - survey effort (which is always 1 because point transects used)

- distance - perpendicular distances
- OBS - initials of the observer
- MAS - minutes after sunrise
- HAS - hour after sunrise

The latter three columns are the covariates to be considered for possible inclusion into the detection function.

There a couple of records with missing distances and so can be deleted with the following command:

```
amakihi <- amakihi[!is.na(amakihi$distance), ]
```

In this command,

- records in amakihi are selected using the square brackets []
- amakihi is a data frame and so selection can be performed on either rows or columns i.e. [rows, columns]. In this case, the selection is performed on the rows (because the selection criteria is before the comma) and all columns will be retained
- the rows selected as those where the distances (stored in amakihi\$distance) are not missing. The function is.na selects elements that are missing; the symbol ! means 'not', and so !is.na selects elements that are not missing.

4 Exploratory data analysis

It is important to gain an understanding of the data prior to fitting detection functions (Buckland et al., 2015). With this in mind, preliminary analysis of distance sampling data involves:

- assessing the shape of the collected data,
- considering the level of truncation of distances, and
- exploring patterns in potential covariates.

We begin by assessing the distribution of distances by plotting histograms with different number of bins and different truncation.

The components of the boxplot are:

- the thick black line indicates the median
- the lower limit of the box is the first quartile (25th percentile) and the upper limit is the third quartile (75th percentile)
- the height of the box is the interquartile range (75th - 25th quartiles)
- the whiskers extend to the most extreme points which are no more than 1.5 times the interquartile range.
- dots indicate 'outliers' if there are any, i.e. points beyond the range of the whiskers.

This format is probably not as useful as a histogram in a distance sampling context but boxplots can be useful to compare the distances for discrete groups in the data. Here we use boxplots to display the distribution of distances recorded by each observer and for each hour after sunrise. Note how the ~ symbol is used to define the groups.

Question: Examine the distribution of radial distances of the point transect data of the amakihi.

Basic syntax will be

```
ggplot(amakihi, aes(x=distance))
+ geom_histogram(binwidth=1)
```

Examine the full dataset, then truncate the data to 82.5m.

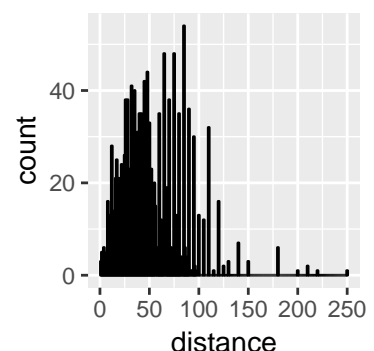
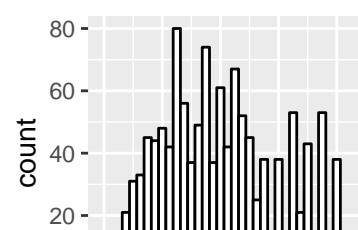


Figure 1: Three levels of detail examining distribution of detection distances.



Boxplots of distances by observer:

```
boxplot(amakihi$distance ~ amakihi$OBS, xlab = "Observer",
        ylab = "Distance (m)")
```

Boxplot of distances for each hour after sunrise:

```
boxplot(amakihi$distance ~ amakihi$HAS, xlab = "Hour",
        ylab = "Distance (m)")
```

For minutes after sunrise (a continuous variable), we create a scatterplot of MAS (on the x -axis) against distances (on the y -axis). The plotting symbol (or character) is selected with the argument `pch`:

```
plot(amakihi$MAS, amakihi$distance, xlab = "Minutes after sunrise",
     ylab = "Distance (m)", pch = 20, cex = 0.5)
```

5 Colinearity

Estimating the parameters of a detection function when covariates are involved is complex. You will recall from multiple linear regression that problems in estimation arise when two covariates in the model are highly correlated. In the exploratory data analysis, it is useful to look for colinearity in potential covariates.

```
## [1] 0.9766042
```

To alleviate the potential colinearity difficulty, hours after sunrise could be transformed to a discrete, rather than a continuous variable.

6 Adjusting the raw covariates

We would like to treat OBS and HAS as factor variables as in the original analysis; OBS is, by default, treated as a factor variable because it consists of characters rather than numbers. HAS, on the other hand, consists of numbers and so by default would be treated as a continuous variable (i.e. non-factor). That is fine if we want the effect of HAS to be monotonic (i.e. detectability either increases or decreases as a function of HAS). If we want HAS to have a non-linear effect on detectability, then we need to indicate to R to treat it as a factor as shown below.

```
amakihi$HAS <- factor(amakihi$HAS)
```

The next adjustment is to change the *reference* level of the *observer* and *hour* factor covariates - the only reason to do this is to get the estimated parameters in the detection function to match the parameters estimated by Distance for Windows. By default R uses the first factor level but by using the `relevel` function, this can be changed:

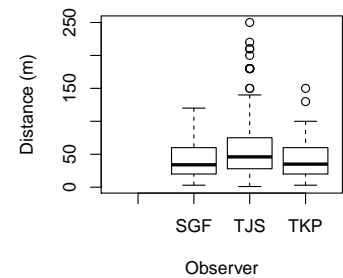


Figure 4: Detection distance distribution by Observer.

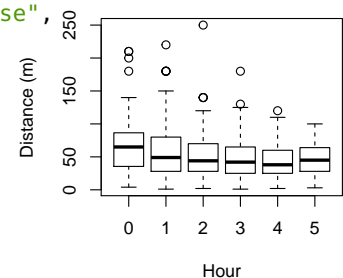


Figure 5: Detection distance distribution by hours after sunrise.

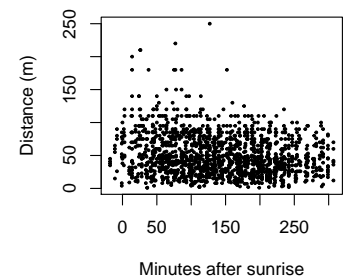


Figure 6: Detection distance distribution by minutes after sunrise.

Question: Compute the correlation of minutes after sunrise and hours after sunrise using the `cor()` function.

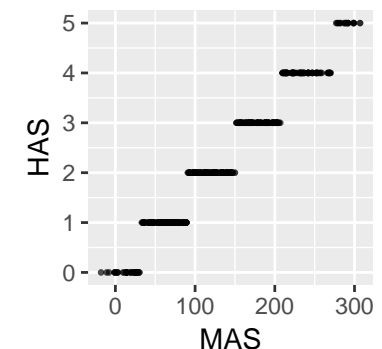


Figure 7: Diagnostics for the presence of colinearity between HAS and MAS.

```
amakihi$OBS <- relevel(amakihi$OBS, ref = "TKP")
amakihi$HAS <- relevel(amakihi$HAS, ref = "5")
```

One final adjustment, and more subtle, is a transformation of the continuous covariate, MAS. We are entertaining three possible covariates in our detection function: OBS, HAS and MAS. The first two variables, OBS and HAS, are both factor variables, and so, essentially, we can think of them as taking on values between 1 and 3 in the case of OBS, and 1 to 6 in the case of HAS. However, MAS can take on values from -18 (detections before sunrise) to >300 and the disparity in scales of measure between MAS and the other candidate covariates can lead to difficulties in the performance of the optimizer fitting the detection functions in R. The solution to the difficulty is to scale MAS such that it is on a scale (approx. 1 to 5) comparable with the other covariates.

Dividing all the MAS measurements by the standard deviation (function `sd`) of those measurements accomplishes the desired compaction in the range of the MAS covariate without changing the shape of the distribution of MAS values. The `na.rm=TRUE` argument ensures that any missing values are ignored.

```
amakihi$MAS <- amakihi$MAS/sd(amakihi$MAS, na.rm = TRUE)
```

Check what this command has done by looking at the range of the adjusted MAS:

```
range(amakihi$MAS)
```

```
## [1] -0.2364891 4.0334529
```

7 Candidate models

With three potential covariates, there are 8 possible combinations for including them in the detection function:

- No covariates
- OBS
- HAS
- MAS
- OBS + HAS
- OBS + MAS
- HAS + MAS
- OBS + HAS + MAS

Even without considering covariates there are a number of possible key function/adjustment term combinations and if all key function/covariate combinations are considered the number of potential models is large. Note that covariates are not allowed if a uniform key function is chosen and if covariate terms are included, adjustment terms are not allowed. Even with these restrictions, it is not best practice to take a scatter gun approach to detection function model fitting. [Buckland et al. \(2015\)](#) considered 13

combinations of key function/covariates. Here, we look at a subset of these as an illustration of how to incorporate covariates in the detection function.

If it is not already loaded, then first load the Distance package.

Fit a hazard rate model with no covariates or adjustment terms. By default, line transects are assumed and because our data are point transects, the argument `transect="point"` is specified:

```
hr.model0 <- ds(amakihi, transect = "point", key = "hr",
  truncation = 82.5, adjustment = NULL, order = 0)
```

The fitted model can be investigated using the `summary` function. Make a note of the AIC for this model.

```
summary(hr.model0)
```

```
##
## Summary for distance analysis
## Number of observations : 1243
## Distance range       : 0 - 82.5
##
## Model : Hazard-rate key function
## AIC   : 10807.55
##
## Detection function parameters
## Scale coefficient(s):
##           estimate      se
## (Intercept) 3.454538 0.06310866
##
## Shape coefficient(s):
##           estimate      se
## (Intercept) 0.83429 0.06533115
##
##                               Estimate
## Average p                   0.3285785
## N in covered region 3782.9624390
##                               SE          CV
## Average p                   0.02013101 0.06126697
## N in covered region 247.88663317 0.06552712
```

```
hr.obs <- ds(amakihi, transect = "point", key = "hr",
  formula = ~OBS, truncation = 82.5)
print(hr.obs$ddf$criterion)
```

```
## [1] 10778.45
```

Fit a hazard rate model with OBS and HAS in the detection function:

```
hr.obshas <- ds(amakihi, transect = "point", key = "hr",
  formula = ~OBS + HAS, truncation = 82.5)
print(hr.obshas$ddf$criterion)
```

Question: Fit a hazard rate model with OBS as a covariate in the detection function and make a note of the AIC. Has the AIC reduced by including a covariate?

Answer: Yes, AIC of the model with observer covariate is 30 AIC units smaller than the model without this covariate.

```
## [1] 10783.14
```

```
hr.has <- ds(amakihi, transect = "point", key = "hr",
  formula = ~HAS, truncation = 82.5)
hr.mas <- ds(amakihi, transect = "point", key = "hr",
  formula = ~MAS, truncation = 82.5)
hr.obsmas <- ds(amakihi, transect = "point", key = "hr",
  formula = ~OBS + MAS, truncation = 82.5)
hr.hasmas <- ds(amakihi, transect = "point", key = "hr",
  formula = ~HAS + MAS, truncation = 82.5)
hr.hasmasobs <- ds(amakihi, transect = "point",
  key = "hr", formula = ~HAS + MAS + OBS, truncation = 82.5)
```

Answer: The model with both observer and hours after sunrise had a smaller AIC than the model with observer alone, however, the reduction in AIC was only 1 AIC unit.

Question: Fit the other candidate models including covariates shown in the list above and decide which model is best in terms of AIC.

A useful function for summarising a candidate model set is `summarize_ds_models`. The arguments to the function is an enumeration of the candidate model objects.

```
summarize_ds_models(hr.model0, hr.obs, hr.has,
  hr.mas, hr.obshas, hr.obsmas, hr.has.mas,
  hr.hasmasobs)
```

Table 1: Candidate model set for Hawaii amakihi covariate analysis

	Model	Key function	Formula	C-vM p-value	\hat{P}_a	$se(\hat{P}_a)$	ΔAIC
5	hr.obsmas	Hazard-rate	~OBS + MAS	0.3890855	0.3186672	0.0201475	0.000000
2	hr.obs	Hazard-rate	~OBS	0.2707149	0.3142711	0.0204417	1.072908
7	hr.hasmasobs	Hazard-rate	~HAS + MAS + OBS	0.4559568	0.3199361	0.0200731	7.743570
4	hr.mas	Hazard-rate	~MAS	0.5579939	0.3335896	0.0201104	28.253419
1	hr.model0	Hazard-rate	~1	0.3344042	0.3285785	0.0201310	30.173452
3	hr.has	Hazard-rate	~HAS	0.5802490	0.3326696	0.0200208	30.843087
6	hr.hasmas	Hazard-rate	~HAS + MAS	0.5856382	0.3325707	0.0200345	32.841743

8 Plotting the detection functions

The detection functions can be investigated using the `plot` function as shown below. A few different plotting options are illustrated.

```
# Plot simple model
plot(hr.model0, nc = 20, main = "No covariates",
  pch = 20, pdf = TRUE)

# Plot model with OBS
plot(hr.obs, nc = 10, main = "Model with OBS covariate",
  pch = 1, cex = 0.5, pdf = TRUE)
```

What does the detection function look like for your selected model?

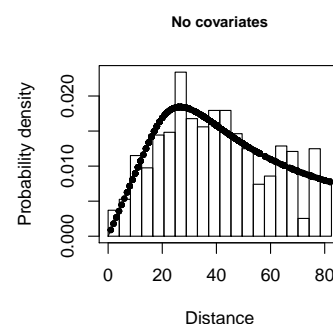
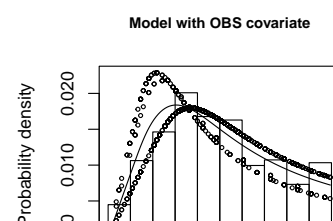


Figure 8: Detection function fit for model without covariates.



To see more sophisticated examples of plotting the detection function for the selected model, see the code accompanying (Buckland et al., 2015) [Hawaiian Amakihi case study](#).

References

- 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>.
- Marques, T. A., L. Thomas, S. G. Fancy, and S. T. Buckland. 2007. Improving estimates of bird density using multiple covariate distance sampling. *The Auk*, **124**:1229–1243. URL [https://doi.org/10.1642/0004-8038\(2007\)124\[1229:ieobdu\]2.0.co;2](https://doi.org/10.1642/0004-8038(2007)124[1229:ieobdu]2.0.co;2).