

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

Analysis in R: Covariates in the detection function

March 2022

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 in [Buckland et al. \(2015, Section 5.3.2\)](#) 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 ([R Core Team, 2018](#)).

(Solutions)



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.

```
amakihi <- read.csv(file="https://raw.githubusercontent.com/distanceworkshops/online-course/gh-pages/  
{file="https://raw.githubusercontent.com/distanceworkshops/  
online-course/gh-pages/exercisepdfs/Ch7/datasets/amakihi.csv"}")  
Check that it has been imported correctly.
```

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

```
##   Study.Area Region.Label Sample.Label Effort distance OBS MAS HAS  
## 1      Kana      Jul-92           1      1      40 TJS  50  1  
## 2      Kana      Jul-92           1      1      60 TJS  50  1  
## 3      Kana      Jul-92           1      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 (1 for all points because they are visited a single time)
- distance - radial distance of detection from observer
- OBS - initials of the observer
- MAS - minutes after sunrise
- HAS - hour after sunrise

Note there is no `Area` field in the dataset. Detection functions can be fitted to the data, but bird density and abundance cannot be estimated. 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:

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

In this command,

- records in `amakihhi` are selected using the square brackets `[]`
- `amakihhi` 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 `amakihhi$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. 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.

Boxplots of distances by observer:

```
ggplot(droplevels(amakihhi), aes(x=OBS, y=distance)) +
  geom_boxplot() + labs(x="Observer initials", y="Radial distance (m)")
```

Boxplot of distances for each hour after sunrise:

```
ggplot(amakihhi, aes(x=factor(HAS), y=distance)) +
  geom_boxplot() + labs(x="Hours after sunrise", y="Radial 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).

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

Basic syntax will be
`ggplot(amakihhi, aes(x=distance)) + geom_histogram(binwidth=1)`
 Examine the full dataset, then truncate the data to 82.5m.

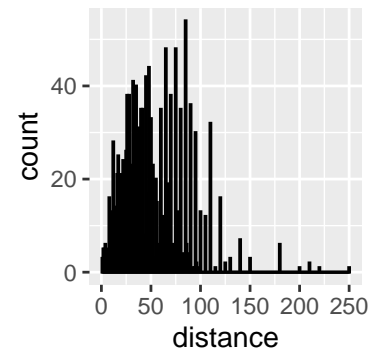


Figure 1: Two levels of detail examining distribution of detection distances.

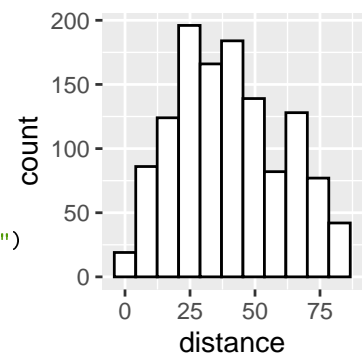
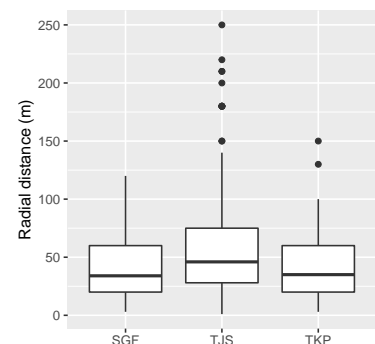


Figure 2: Two levels of detail examining distribution of detection distances.



```
ggplot(amakihi, aes(x=MAS, y=distance)) +
  geom_point(alpha=0.25, size=1.6) +
  labs(x="Hours after sunrise", y="Radial distance (m)")
```

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.

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, indicate it is a factor:

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

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.

Scaling the MAS measurements accomplishes the desired compaction in the range of the MAS covariate without changing the shape of the distribution of MAS values.

```
amakihi$MAS <- scale(amakihi$MAS)
```

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

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

```
## [1] -2.087244 2.182698
```

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

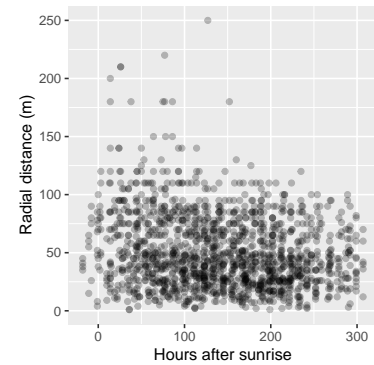


Figure 5: Detection distance distribution by minutes after sunrise.

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

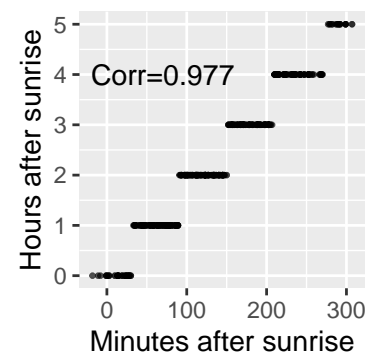


Figure 6: Diagnostics for the presence of colinearity between HAS and MAS.

- 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, load the `Distance` package ([Miller et al., 2019](#)).

We are not using the `convert_units` argument in calls to `ds` because we are not reporting density or abundance estimates. If we were reporting those estimates, we would be including a unit conversion in our calls to `ds`.

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:

```
library(Distance)
conv <- convert_units("meter", NULL, "square kilometer")
hr.model0 <- ds(amakihi, transect="point", key="hr", truncation=82.5,
               adjustment=NULL, convert_units = conv)
```

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 0.06311561
##
## Shape coefficient(s):
##           estimate      se
## (Intercept) 0.83429 0.06533064
##
##           Estimate      SE      CV
## Average p      0.3285785 0.02013402 0.06127615
## N in covered region 3782.9625812 247.91912254 0.06553571
##
## Summary statistics:
##   Region   Area CoveredArea Effort   n   k      ER   se.ER   cv.ER
## 1 Apr-93 0.8766811 0.8766811  41 231  41 5.634146 0.3289972 0.05839344
## 2 Apr-94 0.5131792 0.5131792  24 141  24 5.875000 0.2712859 0.04617632
## 3 Apr-95 0.8552986 0.8552986  40 212  40 5.300000 0.4938078 0.09317129
## 4 Dec-92 0.8125337 0.8125337  38 140  38 3.684211 0.2996092 0.08132249
```

```
## 5 Jan-94 0.8766811 0.8766811 41 172 41 4.195122 0.3088521 0.07362173
## 6 Jul-92 0.8766811 0.8766811 41 146 41 3.560976 0.1945877 0.05464449
## 7 Jul-93 0.8552986 0.8552986 40 201 40 5.025000 0.2984737 0.05939775
## 8 Total 5.6663532 5.6663532 265 1243 265 4.690566 0.1348732 0.02875414
##
```

```
## Density:
## Label Estimate se cv lcl ucl df
## 1 Apr-93 801.9204 67.87753 0.08464373 678.7270 947.4742 169.95434
## 2 Apr-94 836.2016 64.15917 0.07672691 718.8152 972.7579 165.79572
## 3 Apr-95 754.3606 84.12273 0.11151527 604.6289 941.1723 79.56603
## 4 Dec-92 524.3818 53.39462 0.10182394 428.5694 641.6143 90.07507
## 5 Jan-94 597.1009 57.19381 0.09578583 494.1049 721.5664 112.86904
## 6 Jul-92 506.8415 41.61285 0.08210229 431.1872 595.7698 193.95798
## 7 Jul-93 715.2193 61.03658 0.08533967 604.4755 846.2520 160.47187
## 8 Total 667.6186 44.65372 0.06688508 585.6130 761.1076 1419.87901
```

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?

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

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

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

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

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

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

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

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

Answer: The model with both observer and hours after sunrise had an AIC score 5 AIC units larger than the model with observer alone; suggesting the inclusion of HAS along with OBS is not a better model than the model with OBS alone.

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.hasmas, hr.hasmasobs)
```

Table 1: Candidate model set for Hawaii amakihi covariate analysis.

Model	Key func-tion	Formula	C-vM p-value	\hat{P}_a	$se(\hat{P}_a)$	ΔAIC
hr.obsmas	Hazard-rate	~OBS + MAS	0.389	0.319	0.02	0.000

Model	Key function	Formula	C-vM p-value	\hat{P}_a	se(\hat{P}_a)	Δ AIC
hr.obs	Hazard-rate	~OBS	0.271	0.314	0.02	1.073
hr.obshas	Hazard-rate	~OBS + HAS	0.450	0.320	0.02	5.760
hr.hasmasobs	Hazard-rate	~HAS + MAS + OBS	0.456	0.320	0.02	7.744
hr.mas	Hazard-rate	~MAS	0.558	0.334	0.02	28.253
hr.model0	Hazard-rate	~1	0.334	0.329	0.02	30.173
hr.has	Hazard-rate	~HAS	0.580	0.333	0.02	30.843
hr.hasmas	Hazard-rate	~HAS + MAS	0.586	0.333	0.02	32.842

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",
     cex=0.5, pdf=TRUE, showpoints=FALSE)
add_df_covar_line(hr.obs, data.frame(OBS=c("SGF", "TJS", "TKP")),
                  col=c("blue", "green", "red"), lty=2, pdf=TRUE)
legend("topright", legend=c("SGF", "TJS", "TKP"),
      col=c("blue", "green", "red"), lwd=2)
```

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

```
# Fit best model
hr.best <- ds(amakihi, transect="point", key="hr",
             truncation=82.5, quiet=TRUE, formula=~OBS+MAS)
# Plot model with OBS and MAS
plot(hr.best, nc=10, main="Model with OBS and MAS covariates",
     pch=".", pdf=TRUE)
```

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>.

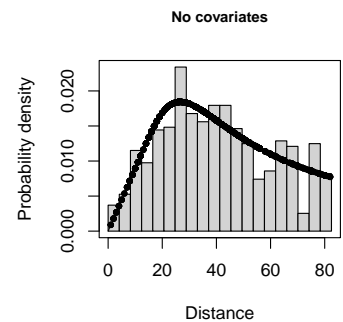


Figure 7: Detection function fit for model without covariates.

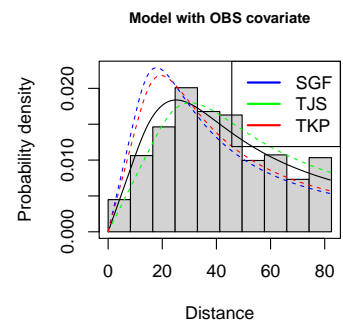
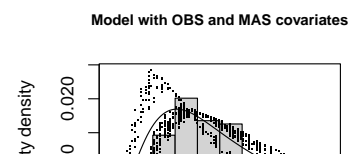


Figure 8: Detection function fit for model with observer covariate.



- 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).
- Miller, D. L., E. Rexstad, L. Thomas, L. Marshall, and J. L. Laake. 2019. Distance sampling in R. *Journal of Statistical Software*, **89**. URL <https://doi.org/10.18637/jss.v089.i01>.
- 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/>.