# hormLong: An R package for longitudinal data analysis in wildlife endocrinology studies

Benjamin Fanson, Kerry V Fanson

The growing number of wildlife endocrinology studies have greatly enhanced our understanding of comparative endocrinology, and have also generated extensive longitudinal data for a vast number of species. However, the extensive graphical analysis required for these longitudinal datasets can be time consuming because there is often a need to create tens, if not hundreds, of graphs. Furthermore, routine methods for summarising hormone profiles, such as the iterative baseline approach and area under the curve (AUC), can be tedious and non-reproducible, especially for large number of individuals. We developed an R package, hormLong, which provides the basic functions to perform graphical and numerical analyses routinely used by wildlife endocrinologists. To encourage its use, hormLong has been developed such that no familiarity with R is necessary. Here, we provide a brief overview of the functions currently available and demonstrate their utility with previously published Asian elephant data. We hope that this package will promote reproducibility and encourage standardization of wildlife hormone data analysis.

1    *For PeerJ*

2

3    ***hormLong*: An *R* package for longitudinal data analysis in wildlife endocrinology**

4    **studies**

5

6    Benjamin G. Fanson and Kerry V. Fanson

7    Centre for Integrative Ecology, School of Life and Environmental Sciences, Deakin

8    University, Waurn Ponds, VIC, Australia

9

10    Corresponding author:

11        Benjamin Fanson

12        Centre for Integrative Ecology, School of Life and Environmental Sciences, Deakin

13        University, 75 Pigdons Road, Waurn Ponds, VIC 3216, Australia

14        E-mail: bfanson@gmail.com

15

16

## ABSTRACT

The growing number of wildlife endocrinology studies have greatly enhanced our understanding of comparative endocrinology, and have also generated extensive longitudinal data for a vast number of species. However, the extensive graphical analysis required for these longitudinal datasets can be time consuming because there is often a need to create tens, if not hundreds, of graphs. Furthermore, routine methods for summarising hormone profiles, such as the iterative baseline approach and area under the curve (AUC), can be tedious and non-reproducible, especially for large number of individuals. We developed an *R* package, *hormLong*, which provides the basic functions to perform graphical and numerical analyses routinely used by wildlife endocrinologists. To encourage its use, *hormLong* has been developed such that no familiarity with *R* is necessary. Here, we provide a brief overview of the functions currently available and demonstrate their utility with previously published Asian elephant data. We hope that this package will promote reproducibility and encourage standardization of wildlife hormone data analysis.

**Keywords**

## INTRODUCTION

Longitudinal hormone monitoring is routinely used in wildlife endocrinology studies and provides a unique insight into endocrine physiology that cannot be obtained from single samples. The amount of longitudinal endocrine data is rapidly increasing due to the development of new techniques and advances in technology (e.g., non-invasive hormone monitoring, catheterization techniques, cheaper assays). Consequently, researchers routinely handle large endocrine datasets with an extensive number of samples. One of the greatest challenges with these large datasets is efficient and reproducible data analysis. Analysing longitudinal hormone data generally includes (1) graphical visualization of the data, (2) identification of peaks, and (3) quantifying the magnitude of the response.

Similar to other time series data (Cowpertwait & Metcalfe 2009; Montgomery et al. 2015), graphical analysis plays an important role in identifying patterns in hormone profiles. Researchers often monitor dozens of individuals, but create profiles for each individual one-at-a-time. Furthermore, temporal events (e.g. pregnancy, mating, stressors) are often added to graphs by hand. This process of creating dozens of graphs, marking events, and updating each graph separately becomes quite time-consuming. In addition, when multiple hormones are being monitored, it is useful to overlay hormone profiles in order to explore temporal correlations. However, this involves restructuring each individual dataset, which takes yet more time and can introduce error when done by hand.

Another challenge with analysing longitudinal hormone data is being able to distinguish the signal from the noise. There is a certain amount of inherent variability in any hormone profile due to both biological (e.g., pulsatile release, variability in steroid metabolism) and methodological factors (e.g., sampling design, pipetting error, assay variability). One common approach for identifying meaningful increases (peaks) in longitudinal datasets is the iterative baseline approach (Brown et al. 1996; Clifton & Steiner 1983). In this approach, hormone values exceeding the mean + ($n$ * SD) are excluded, where $n$ is the criterion for the number of standard deviations (SD) used in the calculation. The mean and SD are recalculated, and this culling processes is repeated until no points exceed the cut-off. Remaining values are considered "baseline" values and excluded points are considered "peaks". The appropriate value of $n$ needs to be adjusted depending on the characteristics of the dataset (number of samples and amount of variation). Although this approach is really useful for identifying peaks, it can be tedious to run these iterative calculations for each study

67 subject, and this becomes even more cumbersome when calculating and comparing different

68 values of $n$.

69      In addition to detecting presence/absence of peaks (above), it is often desirable to

70 quantify the magnitude of the response. One approach is to calculate the magnitude of the

71 peak using either absolute difference (peak minus baseline) or relative increase (ratio of peak

72 to baseline). A more complicated method is to calculate the area under the curve (AUC;

73 (Cockrem & Silverin 2002; Sheriff et al. 2010). An advantage of this technique is that it

74 incorporates both the magnitude of the peak as well as the duration, which are both

75 biologically meaningful. Without specialized software, the AUC can be a tedious calculation

76 and hinders reproducibility.

77      To facilitate efficient and reproducible data analysis, we developed a user-friendly $R$

78 package that provides wildlife endocrinologists with a toolkit for analysing longitudinal

79 hormone data and requires no prior programming experience. The package includes

80 functions allowing for exploratory graphical analysis (including mass production of

81 longitudinal profiles, box plots, and overlaying multiple hormones), iterative baseline

82 calculation, and AUC calculation. To demonstrate the utility of this package, we analysed a

83 previously published hormone dataset (Fanson et al. 2014). This study looked at changes in

84 circulating cortisol across the estrous cycle (i.e., relative to progesterone) in Asian elephants.

85 We included these data as an example dataset called *hormElephant* in the package.

86

## 87 DESCRIPTION

88 **(a) Philosophy**

89 The goal of this package is to provide a toolkit that facilitates efficient and reproducible

90 analysis of longitudinal hormone data commonly used by wildlife endocrinologists. With

91 that in mind, we created functions that perform routine characterization methods (e.g.

92 iterative baseline and AUC calculations), as well as a suite of data visualization functions to

93 facilitate graphical analysis.

94 To encourage researchers who are less familiar with $R$ to use these functions, we developed

95 an $R$-minimal workflow which allows users with no prior $R$ experience to be able to run the

96 functions. To this end, we created a detailed manual that includes instructions on how to

97 install $R$, load the *hormLong* package, and prepare data, in addition to detailed explanations

98 and examples of each function. We also developed an $R$ script template that can be easily

99 modified for analysis of a researcher's own data, eliminating most $R$ coding (manual is

100 located at http://hormlong.weebly.com and the package is available on GitHub at

101 https://github.com/bfanson/hormLong). Output files are in *csv* and *pdf* format. *csv* files can

102 be used in any spreadsheet or statistical software (e.g. Excel, SPSS, JMP) and *pdf* files can

103 be opened in vector-graphics programs (e.g. Illustrator, Inkscape) and modified easily for

104 manuscripts.

105

106 **(b) Typical workflow**

107 Figure 1 illustrates a standard workflow for *hormLong*. In short, data are imported and

108 date/time formatted. Then the baseline analysis (*hormBaseline*) is run, which creates a

109 *hormLong* object. This object can then be used for other functions that create graphs or

110 calculate summary data. The list of current functions in *hormLong* is in Table 1.

111

112 **(c) Data preparation and import**

113 The data needs to be organized in Excel (or similar program) prior to importing to *R*. The

114 data should be in 'long form' (i.e. one hormone concentration per row) to take advantage of

115 grouping capabilities of *hormLong*. For example, the elephant dataset has five columns: (1)

116 elephant name (e.g. 'Ele1', 'Ele2'), (2) date sample collected (e.g. '29-Apr-07', '01-May-

117 07'), (3) hormone type (e.g. 'Progesterone', 'Cortisol'), (4) hormone concentration (e.g. 0.34,

118 0.28), (5) name of an event (e.g. 'mated', 'ovulated'). At a minimum, the dataset must have

119 animal identifier, date collected (or numeric days), and hormone concentration. Please see

120 manual for detailed examples. The data must be saved as a *csv* file.

121 Once data are suitably prepared, the *csv* file can be imported into *R* using the function

122 *hormRead( )*. If dates and/or times are part of the dataset, the function *hormDate( )* handles

123 formatting of these variables so they are compatible with all *hormLong* functions.

124 Example code for import and date formatting:

```
125     hormElephant = hormRead()
126     hormElephant = hormDate(data      = hormElephant,
127                             date_var = 'Date_collected',
128                             name      = 'Date')
```

129

130 **(d) Baseline Analysis**

131 The iterative baseline calculation is a common method used for detecting peaks in

132 longitudinal datasets (Brown et al. 1996; Clifton & Steiner 1983). In this method, the mean

133 and standard deviation (SD) are calculated for the dataset. Any values that are greater than

134   the cutoff value (determined as the *mean + (n \* SD)*) are removed, and this process is

135   repeated until no values exceed the cutoff.  Values remaining at the end of this process are

136   considered "baseline", whereas those that have been excluded are classified as "peaks".

137          The *hormBaseline()* function allows users to easily run these iterative calculations using

138   a single line of code.  This function can run separate baseline calculations for multiple groups

139   (e.g., individuals, species, and/or hormones) at the same time because it allows the user to

140   define the grouping of the hormone data using the *by_var* argument.  For instance,

141   *by_var*='species, id' would perform separate calculations for each individual for each

142   species.  The function returns a *hormLong* object that is used as the basis of the other

143   functions described below. The ease of performing these calculations makes it much faster to

144   adjust criteria and identify an appropriate cutoff criteria for your dataset.  If the criteria is too

145   conservative (i.e., high value of *n*), then it is less likely to identify any peaks.  Conversely, if

146   the criteria is too low then it may result in the majority of the values being classified as

147   "peaks".

148          For the elephant dataset, we ran *hormBaseline()* in order to identify peaks in the cortisol

149   and progesterone data.  We wanted to calculate a separate baseline for each individual

150   elephant and each hormone, so we included *by_var*='Ele, Hormone', where 'Ele' is the

151   column name containing the elephant's identifier.  We tested 3 different baseline cutoff

152   criteria in order to identify an appropriate criteria for our dataset: (1) mean + 1.5 SD, (2)

153   mean + 2 SD, and (3) mean + 3 SD (Figure 2).  For this dataset, the first criteria is too liberal

154   and consequently nearly all the values are identified as peaks, which is not useful (Figure

155   2A).  On the other hand, the third criteria is too strict and no points were identified as peaks

156   (Figure 2C).  For this dataset, we decided to use a criteria of 2 SD (Figure 2B).  The

157   *hormBaseline()* function produces an object (called "*result*" in the example code below) that

158   can then be graphed to visualize the calculated baseline cutoff for each elephant.

159          Example code for mean + 1.5 SD:

```
160          result =  hormBaseline(data        = hormElephant,
161                                 by_var      = 'Ele, Hormone',
162                                 conc_var    = 'Cong_ng_ml',
163                                 time_var    = 'Date',
164                                 event_var   = 'Event',
165                                 criteria    = 1.5)
166
```

167   **(e) Data Visualization**

168    Data visualization is an essential component of identifying patterns in longitudinal hormone

169    profiles. To facilitate this process, we have developed several plotting functions. The

170    *hormPlot()* function is the basic plotting function that creates longitudinal profiles, broken up

171    according to the *by_var* statement and plotted with the baseline cutoff. Specific events (e.g.

172    mating, parturition, stressor) can be plotted onto profile graphs by adding an event column

173    into the user's dataset prior to import. If large temporal gaps exist in the data,

174    *hormPlotBreaks()* can be used remove those gaps. When considering multiple hormones,

175    *hormPlotOverlap()* overlays multiple hormone profiles, and *hormPlotRatio()* plots the ratio

176    of two specified hormones. In order to visualize differences in the distribution of multiple

177    groups, *hormBoxPlot()* creates vertical boxplots for all groups specified. All plots are

178    exported as *pdf* files and have several formatting options (e.g. plot size, number of plots per

179    page, date format, setting all x-axes/y-axes to the same range).

180    For the elephant dataset, we ran *hormPlot()* to visualize the longitudinal plots with three

181    different baseline cutoff criteria (see above; Figure 2). This produced longitudinal plots for

182    each elephant with a reference line showing the baseline cutoff and arrows indicating all

183    events. Next, we wanted to overlay cortisol and progesterone plots (Figure 3A). This allowed

184    us to identify when cortisol peaks occurred relative to progesterone peaks. Using this

185    function, it was clear that peaks in cortisol predominantly occurred during the follicular

186    phase, just before progesterone began to increase.

187    Example code for longitudinal plots with baseline cutoff:

188    ````
hormPlot(result)
````

189

190    Example code for overlaying cortisol and progesterone plots:

191    ````
hormPlotOverlap(result,
````

192    ````
              hormone_var='Hormone',
````

193    ````
              colors='green, purple' )
````

194

195    **(f) Summary Statistics**

196    After identifying peaks using baseline criteria, it is often necessary to extract summary

197    statistics from longitudinal profiles for subsequent analyses (e.g. ANOVA in the user's

198    preferred statistical software). The function *hormSumTable()* exports summary statistics into

199    a *csv* file for this purpose. For the elephant data, the exported summary statistics are shown

200    in Table 2.

201    Alternatively, the user may want run a statistical analysis (e.g. linear mixed model) on

202    the original dataset, but need each sample identified as 'baseline' or 'peak', as determined

203    from the iterative baseline method. This can be achieved by including *save_date*=TRUE in

204    *hormBaseline*() and a *csv* file will be created.

205         Example code for obtaining summary statistics:

206         ```
         hormSumTable(result)
         ```

207

208    **(g) Area Under the Curve Analysis**

209    Area under the curve (AUC) is often used to calculate the magnitude of a response.  The

210    *hormArea()* function performs this calculation using the following algorithm: 1) for

211    subsequent time points, determine whether the line crosses the lower bound cutoff threshold

212    (see below for options); 2) if it does cross, calculate the time at which the line crosses the

213    cutoff threshold; 3) using these new end time points, calculate the AUC (see below for

214    calculation methods).  As with baseline calculations, AUC can be calculated for multiple

215    groups in a single step using the *by_var* statement.

216         Three different lower bounds can be used for AUC calculations: 1) area from the x-axis

217    ('origin'); 2) area from the baseline mean ('baseline'); or 3) area from peak cutoff value

218    determined from *hormBaseline()* ('peak').  For each scenario, *hormArea()* calculates the area

219    above the reference line and counts the number of discrete peaks.  Therefore, in the origin

220    scenario, the entire profile constitutes a single peak.  Users can also choose between two

221    commonly used calculation methods: 1) trapezoid method $[\sum \frac{1}{2} * (t_i - t_{i+1}) * [(c_i + c_{i+1}) -$

222    $cutoff]]$; or 2) spline [integrating over *spline(method='natural)* from *stats* package in *R*]

223    (Adams et al. 2011; Cockrem & Silverin 2002; Littin & Cockrem 2001).  After calculating

224    AUC for each peak, the function produces a summary table that includes each peak identity

225    with its corresponding AUC value.  Longitudinal plots of the peak AUCs are also produced

226    (Figure 3B), allowing the user to match up peak identity in table with specific points on the

227    plot and, especially for the spline method, to assess the appropriateness of the fit.

228         For the elephant dataset, we ran *hormArea()* to quantify the area of each cortisol peak in

229    each longitudinal profile (Figure 3B).  This allows for comparisons of the magnitude of

230    cortisol peaks across cycles or among individuals.

231         Example code for obtaining summary statistics:

232         ```
         hormArea(result, lower_bound = 'peak')
         ```

233

## CONCLUSIONS

234

235 *hormLong* is an *R* package tailored to the analysis of longitudinal hormone data in wildlife

236 endocrinology studies. This package provides an efficient and easy method for implementing

237 the iterative baseline approach and calculating AUC for a large number of individuals.

238 Furthermore, the graphical capabilities of this package greatly reduce the time-consuming

239 process of graph creation, producing searchable *pdf* files with separate profiles for each

240 individual in seconds. We have simplified the *R* code so that minimal *R* experience is

241 required by the user, with all results exported from the *R* environment to allow the user to use

242 other software when preferred. We hope that wide-spread adoption of *hormLong* will result

243 in more reproducible hormone analysis and comparable results. The manual can be

244 downloaded from http://hormlong.weebly.com and the package is available on GitHub at

245 https://github.com/bfanson/hormLong.

246

## ACKNOWLEDGEMENTS

247

252

# REFERENCES

Adams NJ, Farnworth MJ, Rickett J, Parker KA, and Cockrem JF. 2011. Behavioural and corticosterone responses to capture and confinement of wild blackbirds (Turdus merula). *Applied Animal Behaviour Science* 134:246-255. http://dx.doi.org/10.1016/j.applanim.2011.07.001

Brown JL, Wildt DE, Wielebnowski N, Goodrowe KL, Graham LH, Wells S, and Howard JG. 1996. Reproductive activity in captive female cheetahs (*Acinoyx jubatus*) assessed by faecal steroids. *Journal of Reproduction and Fertility* 106:337-346.

Clifton DK, and Steiner RA. 1983. Cycle Detection: A Technique for Estimating the Frequency and Amplitude of Episodic Fluctuations inBlood Hormone and Substrate Concentrations. *Endocrinology* 112:1057-1064. doi:10.1210/endo-112-3-1057

Cockrem JF, and Silverin B. 2002. Variation within and between Birds in Corticosterone Responses of Great Tits (Parus major). *General and Comparative Endocrinology* 125:197-206. http://dx.doi.org/10.1006/gcen.2001.7750

Cowpertwait PS, and Metcalfe AV. 2009. *Introductory time series with R*: Springer Science & Business Media.

Fanson KV, Keeley T, and Fanson BG. 2014. Cyclic changes in cortisol across the estrous cycle in parous and nulliparous Asian elephants. *Endocrine connections* 3:57-66.

Littin K, and Cockrem J. 2001. Individual variation in corticosterone secretion in laying hens. *British Poultry Science* 42:536-546.

Montgomery DC, Jennings CL, and Kulahci M. 2015. *Introduction to time series analysis and forecasting*: John Wiley & Sons.

Sheriff MJ, Krebs CJ, and Boonstra R. 2010. Assessing stress in animal populations: Do fecal and plasma glucocorticoids tell the same story? *General and Comparative Endocrinology* 166:614-619. 10.1016/j.ygcen.2009.12.017

278 **Table 1: List of functions in *hormLong*.**

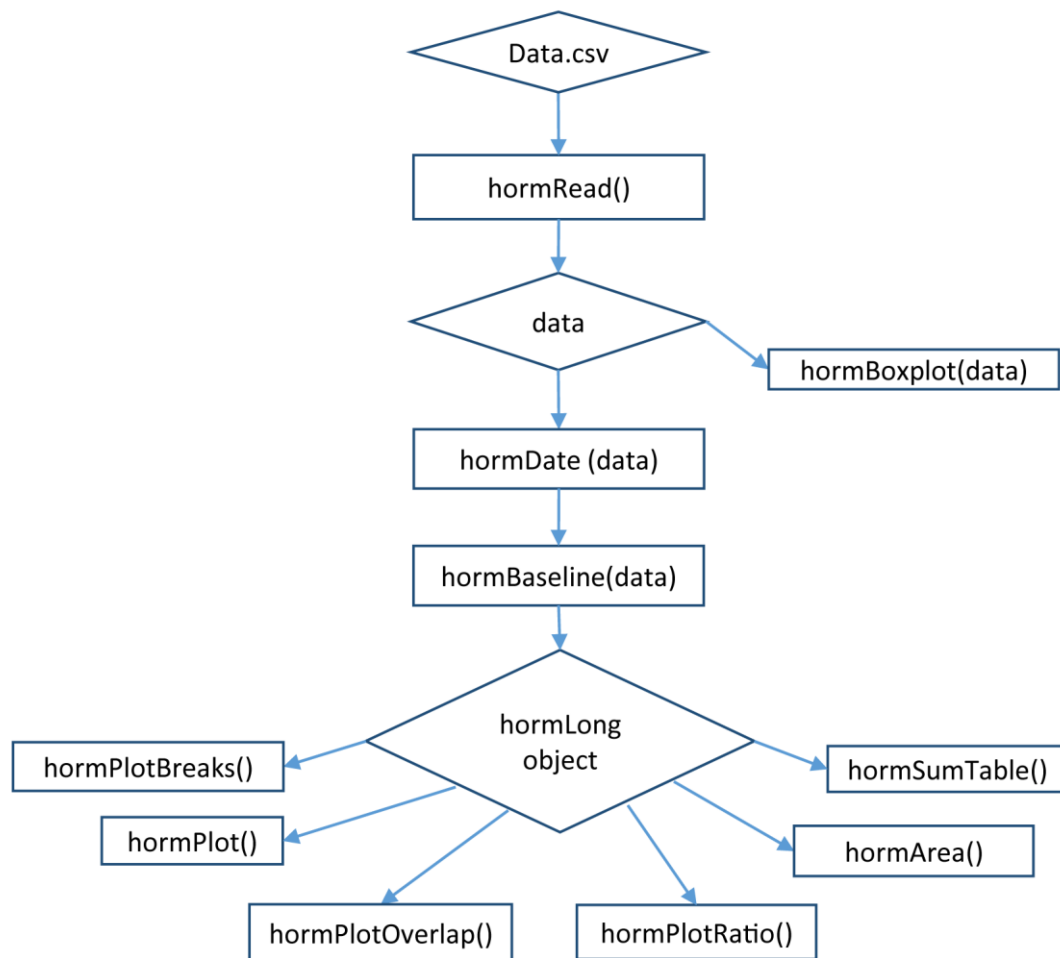| Type | Name | Description |
|---|---|---|
| **Import and data handling** | *hormRead()* | Provides a pop-up window to import file |
| | *hormDate()* | Converts character date (e.g. "2014-01-01",'01-January-2014') to numeric date field. If a time column ('18:10:01') is also supplied, then a date-time field is created. |
| **Analysis** | *hormBaseline()* | Main function that calculates peak cutoff value using iterative algorithm. Produces a *hormLong* object that is used for most other functions |
| | *hormSumTable()* | Calculates basic statistics for hormone data, such as mean, min, max, baseline mean, %CV |
| | *hormArea()* | Calculates area under the curve (AUC) for all peaks |
| **Visualization** | *hormPlot()* | Produces longitudinal plots of hormone profiles for each group specified in *by_var*. Includes baseline cutoff and individual specific events |
| | *hormPlotBreaks()* | Similar to *hormPlot()*, except that temporal gaps in endocrine profiles are removed. |
| | *hormPlotOverlap()* | Produces longitudinal plots in which multiple hormone are overlaid. |
| | *hormArea()* | Produces longitudinal plots in which AUC for peaks are delineated and numbered. This plot complements hormAUC analysis table so that numbered peaks can be assessed visually. |
| | *hormBoxplot()* | Produces simple boxplots comparing hormone concentrations using grouping function *by_var*. |

279

280

281

282 **Table 2: Example output for *hormSumTable()*.** Base_mean is the mean of baseline values from iterative process. Peak_mean is mean of all

283 peak values. Cutoff is the cutoff threshold (mean + ($n$ * SD)) determined from *hormBaseline()*. Other statistics are based on all hormone values.

| Ele | Hormone | mean | median | sd | percent_cv | min | max | cutoff | base_mean | peak_mean | peak_base |
|-----|---------|------|--------|------|------------|------|------|--------|-----------|-----------|-----------|
| Ele1 | Cortisol | 0.83 | 0.7 | 0.51 | 61.62 | 0 | 2.49 | 1.09 | 0.61 | 1.62 | 2.67 |
| Ele1 | Progesterone | 0.41 | 0.36 | 0.37 | 89.72 | 0 | 1.31 | 0.94 | 0.34 | 1.13 | 3.36 |
| Ele2 | Cortisol | 0.62 | 0.46 | 0.52 | 84.72 | 0.19 | 2.84 | 0.66 | 0.42 | 1.42 | 3.41 |
| Ele2 | Progesterone | 0.85 | 0.82 | 0.52 | 61.85 | 0.05 | 2.77 | 1.66 | 0.78 | 2.15 | 2.75 |

| | |
|---|---|
| mean | average (of all points for that set of grouping variables) |
| median | median (of all points for that set of grouping variables) |
| sd | standard deviation (of all points for that set of grouping variables) |
| percent_cv | percent coefficient of variation (SD/mean*100) |
| min, max | minimum and maximum values (of all points for that set of grouping variables) |
| cutoff | threshold value for peaks, calculated as mean+($n$*SD) for final iteration of baseline calculation (i.e., when no more points are removed). Points below this are baseline and above are peaks. |
| base_mean | average of all points classified as baseline |
| peak_mean | average of all points classified as peaks |
| peak_base | ratio of peak-to-baseline (calculated as peak_mean/base_mean) |

284

285
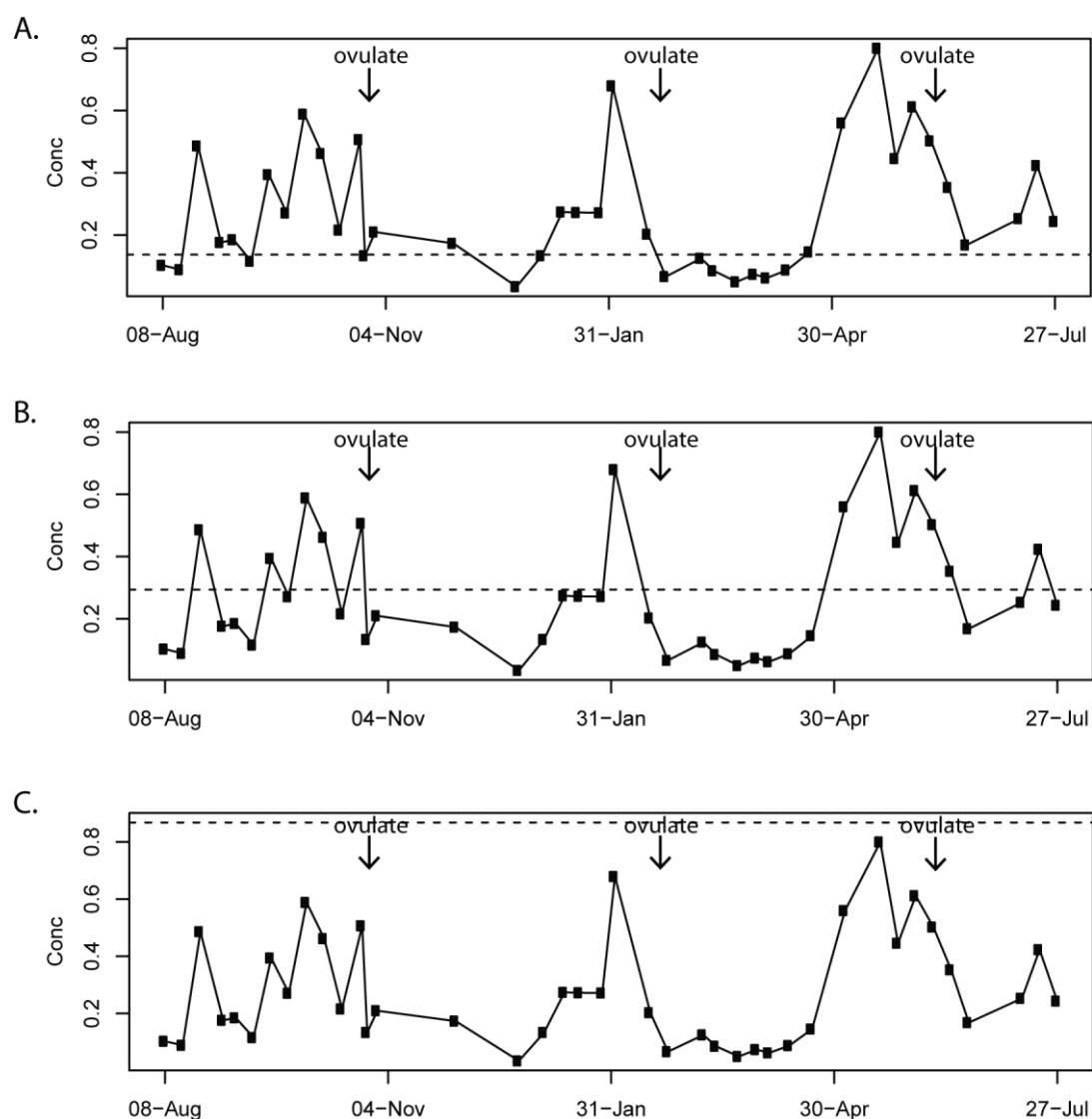


286

287 **Figure 1: Flowchart of a typical *hormLong* analysis.** Diamonds show *R* objects and boxes are

288 functions.

289

Figure 2: **Example of *hormPlot()* with varying criteria for a single individual.** The dashed line represents the cutoff criteria: A) mean + 1.5, B) mean + 2, and C) mean + 3.0. Arrows and text show the occurrence of an event.
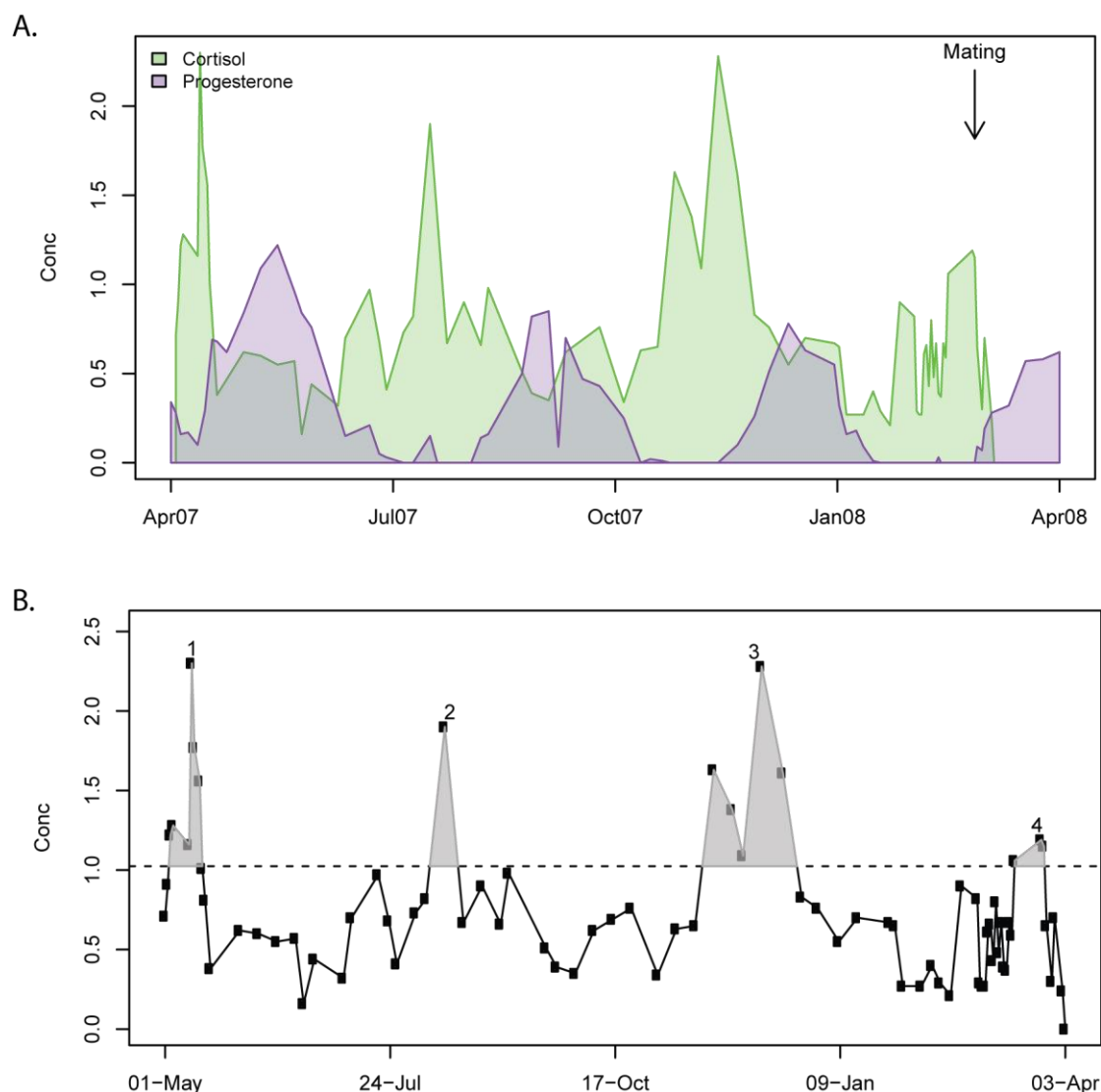
296

297

A.



B.

298

**Figure 3:** **Example of (A)** *hormPlotOverlap()* **and (B)** *hormArea()* **plot.** For (A), the

different colours represent cortisol (green) and progesterone (purple). For (B), numbers

indicate discrete peak number (matches up with outputted table) and shaded area shows the

AUC calculated in the output data table. Dashed lines is the baseline cutoff value (note –

other cutoff criteria can be used for *hormArea()*, see manual).