

# Butterfly Monitoring and Analyses

Reto Schmucki

2024-08-09



# Table of contents

<b>Preface</b>	<b>1</b>
NOTE - this book project is a work in progress that only started. . . . .	1
<b>I Counting Butterflies</b>	<b>3</b>
Butterfly Monitoring	5
<b>II From Counts to Flight Curves</b>	<b>7</b>
<b>Data Analyses and Methods</b>	<b>9</b>
Data Generation Processes . . . . .	9
Data Simulation . . . . .	9
Statistical Modelling . . . . .	10
From Counts to Models . . . . .	10
<b>Butterfly Counts</b>	<b>11</b>
Simulation of Butterfly Counts . . . . .	11
Simulation Tool for Butterflies and Other Phenologies . . . . .	11
Simulate Data Sampling Process . . . . .	12
Generalized Additive Models with rbms . . . . .	17
Organising BMS count data . . . . .	18
Fitting a GAM to Butterfly Counts . . . . .	18
Non Gaussian flight curve . . . . .	20
Simple trend case . . . . .	22
<b>III Butterfly Trends</b>	<b>23</b>
<b>Trends Simulation</b>	<b>25</b>
Multisite and trends . . . . .	25
<b>IV Bibliography</b>	<b>29</b>
References	31



# Preface

This project aims to collate documentation of methods used to analyze and work with Butterfly Monitoring Schemes (BMS) data. We focus on methods used to analyze count data associated with Pollard transects that represent the core of the eBMS database.

**NOTE - this book project is a work in progress that only started.**

Reto Schmucki, July 2024



## Part I

# Counting Butterflies





# Butterfly Monitoring

Content under construction



## Part II

# From Counts to Flight Curves



# Data Analyses and Methods

Biodiversity monitoring is essential to assess and understand the status and trends of species and ecosystems. Through active monitoring, we can integrate new information and update the knowledge needed for decision-making. Analyzing monitoring data can range from simple data exploration to the development of complex statistical models. A thorough understanding of both the data and the methods is fundamental to the selection and application of appropriate methods. Such an understanding will also provide important insights into the results enhancing our comprehension and helping us to communicate the results effectively to decision-makers.

## Data Generation Processes

To develop robust monitoring methods and analyze the incoming data appropriately, it is important to understand the processes that generate the data and how the monitoring protocol affects the nature and structure of the data. For example, the life cycle of butterflies affects the number of adult individuals that can be observed in a given location at a given time. The seasonality of the emergence process creates a temporal pattern in the observed and recorded data that must be accounted for in the analysis. Systematic variation may also result from differences in sampling effort (e.g. the area sampled, the time spent recording or the experience of the recorders).

Understanding the influence of both biological and sampling components on the data generation process is crucial for the design and development of methods that can filter out the component of interest while accounting for systematic structures in the data. This also helps in assessing whether the monitoring program accurately captures and reflects the population of interest or whether the sample is biased and unrepresentative. When designing a monitoring program, it is important to identify the population to be monitored and understand how the sampling protocol may affect the representativeness of the data and introduce potential bias.

## Data Simulation

To better understand the influence of species biology and sampling protocols on data generation, we will use data simulation approaches. Data simulation involves generating random data sets based on defined rules and known parameters. This technique is useful not only to illustrate the outcome of ecological and sampling processes but also to test methods and improve our understanding of statistical models and their potential failures.

When carefully designed, data simulation is a powerful tool for testing and validating methods and performing sensitivity analyses to assess their robustness to violations of underlying assumptions. Simulated data sets allow for exploring systematic patterns in the data and evaluating the behavior of models to identify their strengths and limitations.

In the following sections, we will use data simulations to illustrate and explore the different components of butterfly monitoring data generated under different scenarios. Through these simulations, we aim to 1) gain a deeper understanding of the data structure resulting from the ecological and sampling processes involved in

butterfly monitoring schemes, 2) demonstrate and compare different modeling approaches and 3) understand the information that can be derived from these approaches.

## **Statistical Modelling**

### **From Counts to Models**

# Butterfly Counts

## Simulation of Butterfly Counts

We will use simulated data to demonstrate and evaluate methods for calculating butterfly abundance indices, population trends and multi-species indicators such as the European Grassland Butterfly Indicator. To achieve this, we need to generate realistic data sets with known parameters. Data simulation will allow us to apply and test the methods on data generated under different scenarios. This approach will enable rigorous sensitivity analysis and provide crucial insights into the methods and a deeper understanding of their performance and limitations.

```
if(!require("data.table")) install.packages("data.table")
if(!require("ggplot2")) install.packages("ggplot2")
if(!require("devtools")) install.packages("devtools")
if(!require("rbms")) devtools::install_github("RetoSchmucki/rbms")

flc_col <- '#ff8c00'
cnt_col <- '#008b8b'
missing_col <- '#8b0000'
GAM_col <- '#483d8b'
```

## Simulation Tool for Butterflies and Other Phenologies

Because butterflies' life cycle is strongly structured in time, with species-specific phenologies, we must account for this ecological process when simulating individual counts recorded across an entire season. The temporal pattern in the number of adult butterflies (imago) is determined by their emergence rate, the timing of the emergence and the life span of the adult. These parameters will often result in the number of adult individuals increasing over a certain period, up to a peak after when their number starts to decline. If a species can produce more than one generation per season, the number will display additional waves of emergence, each with its respective start, peak and decline periods. For each generation, the hump-shaped temporal pattern in the number of adults can often be described by a function that has a mean (center), variance (width) and a certain level of skewness (asymmetry). When merged, the individual patterns can become hidden under the pattern resulting from the cumulative effect of partly overlapping generations.

To simulate butterfly count data with such phenological patterns, we will use the function `timeseries_sim()` from the R package `butterflyGamSims` developed by Collin Edwards (see Edwards et al. 2023). This package is freely available on [GitHub](#) and will allow us to generate realistic data sets under scenarios with various levels of complexity.

We illustrate how the `timeseries_sim()` function works with a simple case where we simulate butterfly counts for a univoltine species with a Gaussian pattern (i.e. one generation with Normal distribution), where the peak

abundance is observed at day 175 with a standard deviation of 15 days.

```
set.seed(13276)

if(!require("butterflyGamSims")) devtools::install_github("cbedwards/butterflyGamSims")

btfl_data <- timeseries_sim(nsims=1,
  year = c(2023),
  doy.samples = seq(from=1, to=365, by=1),
  abund.type = "exp",
  activity.type = "gauss",
  sample.type = "pois",
  sim.parms = list(growth.rate = 0,
    init.size = 500,
    act.mean = 175,
    act.sd = 15)
)
```

The object produced by the `timeseries_sim()` function contains 1) a `data.frame` `NAME$timeseries` with the time series and 2) a `data.frame` `NAME$parms` with the parameters used for the simulation. In the parameters, you will find the population growth rate (`growth.rate`), the initial population size (`init.size`) measure in number of individuals expected over the season, the peak of the activity curve (`act.mean`) measured in days and the width of the activity curve (`act.sd`) that is measured in standard deviation.

#### **i** Note

Note that not all sampling parameters used for the simulation are included in the `parms` object; the `activity.type` (the distribution function used to define the activity curve), the `sample.type` (the sampling process used to sample random counts along the activity curve) and the `abund.type` (the type of the growth rate, deterministic or with a log-normal process error) are missing.

## Simulate Data Sampling Process

With the `timeseries_sim()` function above, we simulated a regular time series of butterfly counts where the actual number of active butterfly of each day are draw from a Poisson distribution with a given expectation defined by the activity curve along a day-of-year vector  $j$  following the probability density of a normal distribution (Gaussian) with mean equal to the peak day  $\mu$  (`act.mean`) and a standard deviation  $\alpha$  (`act.sd`). Because the integral of the probability density distribution (area under the curve) sum to 1, we can multiply the density by the abundance to retrieve a vector of expected abundance for each day-of-year,  $\lambda_j$ .

$$\lambda_j = abundance * \frac{1}{\alpha\sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{(j-\mu)}{\alpha}\right)^2\right)$$

From the activity curve, we can use the `rpois()` function in R to draw a random value from a Poisson distribution,  $y_i$  representing the count for day  $j$ , where the mean is specified by expected value  $\lambda$  at day  $j$ .

$$y_j = rpois(\lambda_j)$$



From this simulation, we have generated the ecological process for the butterfly counts, for a given abundance distributed over a specific phenology defined by a Gaussian curve with a peak (mean) and a breath (standard deviation).

```
library(knitr)
kable(btfl_data$timeseries[c(1:3,160:163,250:253),])
```

Table 1: Butterfly Simulation Data

	years	doy	count	act	abund.true	onset.true	median.true	end.true	fp.true	sim.id
1	2023	1	0	0.0000000	500	155.8	175	194.2	38.4	1
2	2023	2	0	0.0000000	500	155.8	175	194.2	38.4	1
3	2023	3	0	0.0000000	500	155.8	175	194.2	38.4	1
160	2023	160	9	8.0656908	500	155.8	175	194.2	38.4	1
161	2023	161	11	8.6025942	500	155.8	175	194.2	38.4	1
162	2023	162	10	9.1345489	500	155.8	175	194.2	38.4	1
163	2023	163	5	9.6563851	500	155.8	175	194.2	38.4	1
250	2023	250	0	0.0000496	500	155.8	175	194.2	38.4	1
251	2023	251	0	0.0000354	500	155.8	175	194.2	38.4	1
252	2023	252	0	0.0000252	500	155.8	175	194.2	38.4	1
253	2023	253	0	0.0000179	500	155.8	175	194.2	38.4	1

The additional structure resulting from the monitoring protocol (observation process) can now be added to the simulated time series and replicate a specific protocol. In this first case, we will simulate a protocol with weekly visits and include some missing counts for weeks when the minimal monitoring conditions were not met or the observer was absent. For this, we will write some new R functions. The first function will define the start and end of the monitoring season and sample one monitoring day per week over the season. Then we will write functions to simulate a certain level of missing weekly visits within the season. The likelihood of missing weeks tends to be higher at the beginning and the end of the season and lower in the middle. Let's start with the first function that defines the monitoring season and resamples one day of the simulated time series every week. We will name the function `sim2bms()` as it aligns the simulated time series to the protocol of a specific Butterfly Monitoring Scheme (BMS). The function needs the time series, and additional arguments to define the year that we want to extract `yearKeep`, if the time series must be resampled weekly `weeklySample`, which day should be used for the weekly resampling `weekdayKeep` (this can be a vector of days, e.g. `c(2,3,4,5)`, or a specific day), and finally the monitoring season `monitoringSeason` which correspond to a vector of months, e.g. `c(4,5,6,7,8,9)` represent a season starting in April and ending in September. Note that this function will also add some new variables such as the date, the ISO week number and the day of the week.

```
# data: Time series resulting from the simulation generated for 365 days
# weeklySample: TRUE or FALSE; should the daily count in the time series be resampled weekly?
# weekdayKeep: vector of days c(1,2,..., 7) to be sampled from for the weekly count. If the vector
# contains c(2,3,4), the sampling process will be restricted to Tuesday, Wednesday or Thursday.
# monitoringSeason: vector of months that define the monitoring season (e.g. April to September is c(4,5,6,7,8,9))

sim2bms <- function(data, yearKeep = NULL, weeklySample = FALSE, weekdayKeep = NULL, monitoringSeason = NULL) {
  btfl_ts <- data.table::data.table(data[, site_id := paste0("site_", sim.id)])
  if(!is.null(yearKeep)){
```

```

      btfl_ts <- btfl_ts[years %in% yearKeep, ]
    }
    btfl_ts[, date := as.Date(doy, origin = paste0(years, "-01-01"))-1]
    btfl_ts[, week := isoweek(date)]
    btfl_ts[month(date) != 1 | week < 50, weekday := rowid(week), by = .(site_id, years)]
    if(isTRUE(weeklySample)){
      if(!is.null(weekdayKeep)){
        btfl_ts <- btfl_ts[weekday %in% weekdayKeep, ]
      }
      btfl_ts <- btfl_ts[btfl_ts[,.I[sample(.N, 1)], by = .(week, site_id, years)]]["V"]
    }
    if(!is.null(monitoredSeason)){
      btfl_ts <- btfl_ts[month(date) %in% monitoredSeason, ]
    }
  }
  return(btfl_ts)
}

```

We can apply this function to retrieve a specific year of the simulated time series and add some new variables, leaving all other parameters empty. With the same function, we can also resample weekly counts (e.g. one day from c(2:5)) and restrict the time series to a specific monitoring season (e.g. c(4:9)).

```

set.seed(13276)
y <- c(2023)

btfl_ts <- sim2bms(data = btfl_data$timeseries, yearKeep = y)
btfl_fig1 <- ggplot() +
  geom_point(data=btfl_ts, aes(x=doy, y=count, colour = "count")) +
  geom_line(data = btfl_ts,
    aes(x = doy, y = act, colour = "activity")) +
  xlim(1,365) + ylim(0, max(btfl_ts$count, btfl_ts$act)) +
  scale_colour_manual("",
    breaks = c("count", "activity"),
    values = c(cnt_col, flc_col)) +
  theme_light() +
  theme(legend.position = "inside", legend.position.inside = c(0.9, 0.8)) +
  labs(title = paste0("Simulated butterfly counts (", y, ")"),
    subtitle = "- daily visit",
    x = "Day of Year",
    y = "Count")

btfl_week_smpl <- sim2bms(data = btfl_data$timeseries, yearKeep = y,
  weeklySample = TRUE,
  weekdayKeep = c(2:5),
  monitoredSeason = c(4:9))
btfl_fig2 <- ggplot() +
  geom_point(data=btfl_week_smpl, aes(x=doy, y=count, colour = "count")) +
  geom_line(data = btfl_ts,
    aes(x = doy, y = act, colour = "activity")) +
  xlim(1,365) + ylim(0, max(btfl_ts$count, btfl_ts$act)) +

```

```

    scale_colour_manual("",
      breaks = c("count", "activity"),
      values = c(cnt_col, flc_col)) +
    theme_light() +
    theme(legend.position = "inside", legend.position.inside = c(0.9, 0.8)) +
    labs(title = paste0("Simulated butterfly counts (", y, ")"),
      subtitle = "- weekly visit (random resampled)",
      x = "Day of Year",
      y = "Count")

btfl_fig1
btfl_fig2

```

In the example above, the activity curve represented by the line has a Gaussian shape and counts presented by the points along the curve are independent random samples from a Poisson distribution. Because we sampled a count value for 365 days (day-of-year; doy), the counts are representative of the population of active adult butterflies as if the site was visited every. This implies that a proportion of butterflies are counted more than one day as their lifespan exceeds one day. On Pollard transect, this is how butterfly counts are likely to be counted and reported, but with a different frequency as visits are generally weekly, fortnightly, or even monthly. We can replicate this value by resampling the daily count weekly.

From the weekly visits, counts outside of the monitoring period will not be informed, in many cases these are ‘zeros’ as we expect the monitoring season to align with butterflies’ activity. Some other weeks might be missing from the time series, potentially due to unsuitable weather conditions for monitoring or the recorder’s availability. We can inform and exclude the missing visits by resampling a subset of the weekly counts.

```

missing_prob <- function(data, mu=NULL, alpha = 5, theta = 0.3){
  x_ <- seq_len(nrow(data))
  mu_ <- ifelse(is.null(mu), length(x_) / 2, mu)
  std_ <- sqrt(mu_ / theta)
  y_ <- abs((alpha * exp(-(x_ - mu_)^2) / std_^2)) - alpha) + alpha
  yn_ <- y_ / (sum(y_))
  return(yn_)
}

sample_missing <- function(data, propMissing = 0.25){

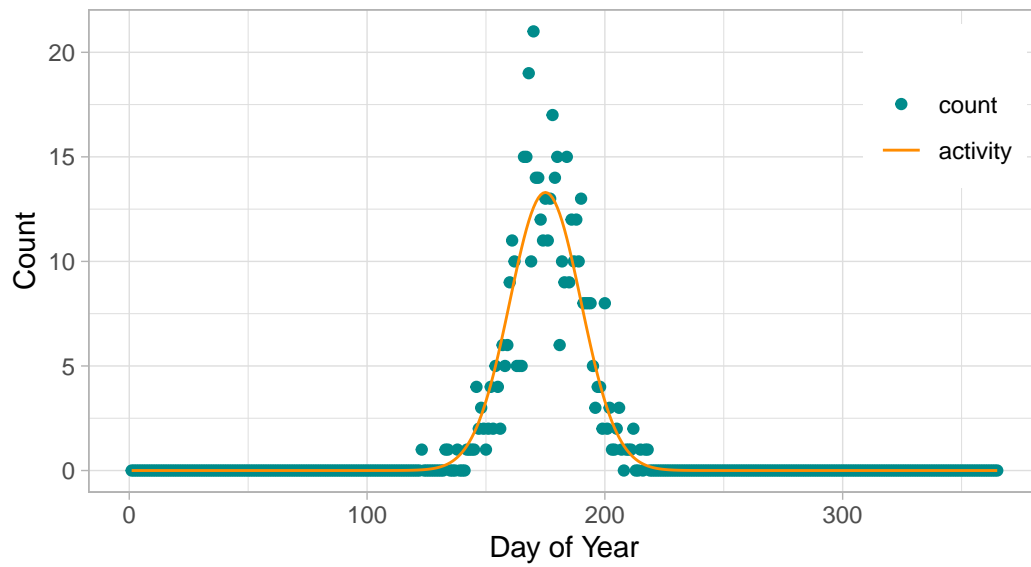
  missing_prob <- data.table::data.table()
  for(i in data[, unique(years)]){
    for(j in data[, unique(site_id)]){
      missing_prob <- rbind(missing_prob, missing_prob(data[years == i & site_id == j, ]))
    }
  }

  missing.week <- data[sample(seq_len(.N), round(propMissing * .N), prob = unlist(missing_prob))
  return(missing.week)
}

```

## Simulated butterfly counts (2023)

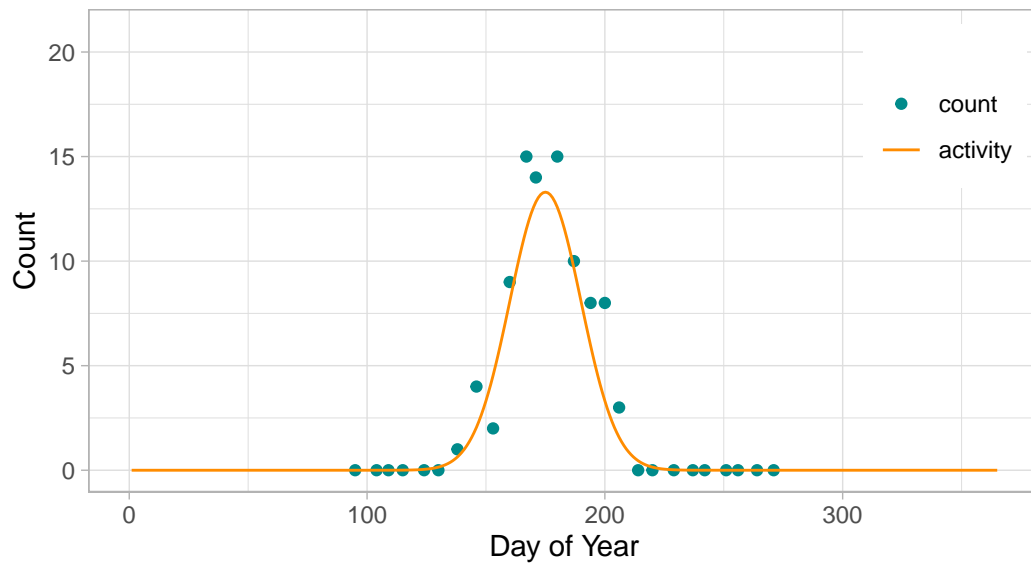
– daily visit



(a) a)

## Simulated butterfly counts (2023)

– weekly visit (random resampled)



(a) b)

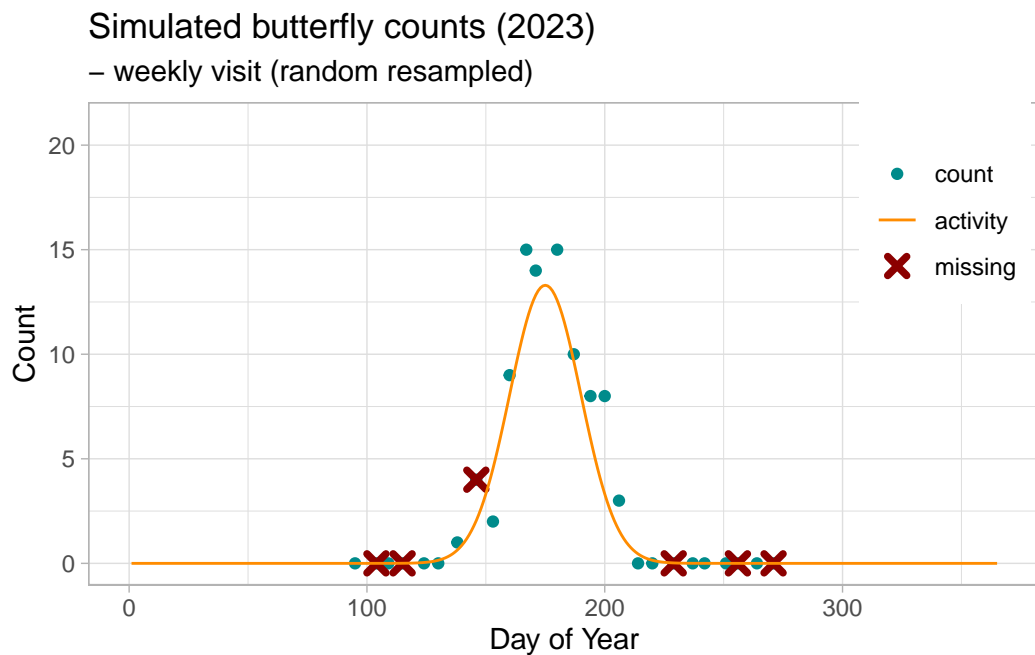
```

btfl_week_missing <- sample_missing(data = btfl_week_smpl, propMissing = 0.25)

btfl_fig3 <- ggplot() +
  geom_point(data=btfl_week_smpl, aes(x=doy, y=count, colour = "count")) +
  geom_point(data=btfl_week_missing, aes(x=doy, y=count, colour = "missing"),
            shape=4, size=2, stroke=2) +
  geom_line(data = btfl_ts,
            aes(x = doy, y = act, colour = "activity")) +
  xlim(1,365) + ylim(0, max(btfl_ts$count, btfl_ts$act)) +
  scale_colour_manual("",
                    breaks = c("count", "activity", "missing"),
                    values = c(cnt_col, flc_col, missing_col)) +
  theme_light() +
  theme(legend.position = "inside", legend.position.inside = c(0.9, 0.8)) +
  labs(title = paste0("Simulated butterfly counts (", y, ")"),
       subtitle = "- weekly visit (random resampled)",
       x = "Day of Year",
       y = "Count")

btfl_fig3

```



## Generalized Additive Models with rbms

We will use the simulation to test the GAM method implemented in the [R package rbms](#) (Schmucki, Harrower A., and Dennis B. 2022). Because recorders only report the number of observed butterflies, zeros are generally not reported but can be derived from the visit dates.

## Organising BMS count data

```
visit_sim <- btfl_week_smpl[!date %in% btfl_week_missing$date, .(site_id, date, count)]
count_sim <- visit_sim[count>=1,][, species := "sp1"]

names(visit_sim) <- toupper(names(visit_sim))
names(count_sim) <- toupper(names(count_sim))

ts_date <- rbms::ts_dwmy_table(InitYear = 2023, LastYear = 2023, WeekDay1 = 'monday')

ts_season <- rbms::ts_monit_season(ts_date,
  StartMonth = 4,
  EndMonth = 9,
  StartDay = 1,
  EndDay = NULL,
  ComplSeason = TRUE,
  Anchor = TRUE,
  AnchorLength = 2,
  AnchorLag = 2,
  TimeUnit = 'd')

ts_season_visit <- rbms::ts_monit_site(ts_season, visit_sim)

ts_season_count <- rbms::ts_monit_count_site(ts_season_visit, count_sim, sp = "sp1")
```

## Fitting a GAM to Butterfly Counts

```
ts_flight_curve <- rbms::flight_curve(ts_season_count,
  NbrSample = 300,
  MinVisit = 5,
  MinOccur = 3,
  MinNbrSite = 1,
  MaxTrial = 4,
  GamFamily = 'nb',
  SpeedGam = FALSE,
  ComplSeason = TRUE,
  SelectYear = NULL,
  TimeUnit = 'd')
```

The flight curve computed by the `rbms::flight_curve()` function is stored in the `...$pheno` object where the days of year are stored under the variable `trimDAYNO` and the standardized flight curve under the variable `NM`. The `NM` variable is scaled to an Area Under the Curve (AUC) that sum to 1. To compare the flight curve derived from the GAM with the activity curve used for the simulation, we must rescale them to the same AUC, in other words, we must rescale the activity curve to have an AUC of 1 or rescale the `NM` to the population size used for the simulation. Here we will rescale the `NM` to match the simulation population size, this will allow us to display the curves and the counts on the same plot with the correct scale.

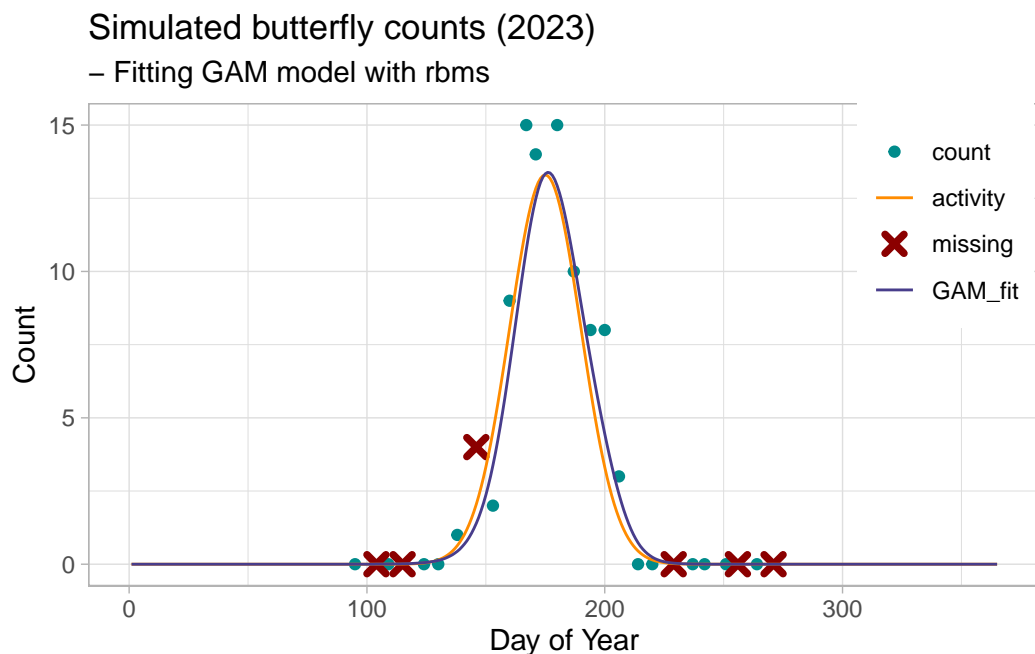
```

pheno <- ts_flight_curve$pheno

btfl_fig4 <- ggplot() +
  geom_point(data=ts_season_count[ANCHOR == 0 & !is.na(COUNT), ], aes(x=DAY_SINCE, y=COUNT)) +
  geom_point(data=btfl_week_missing, aes(x=doy, y=count, colour = "missing"),
    shape=4, size=2, stroke=2) +
  geom_line(data = btfl_ts, aes(x = doy, y = act, colour = "activity")) +
  geom_line(data = pheno,
    aes(x = trimDAYNO, y = btfl_ts[,unique(abund.true)]*NM, colour = "GAM_fit")) +
  xlim(1,365) + ylim(0, max(btfl_ts$act,
    pheno$NM*btfl_ts[,unique(abund.true)],
    btfl_week_missing$count,
    ts_season_count[!is.na(COUNT), COUNT] )) +
  scale_colour_manual("",
    breaks = c("count", "activity", "missing", "GAM_fit"),
    values = c(cnt_col, flc_col, missing_col, GAM_col)) +
  theme_light() +
  theme(legend.position = "inside", legend.position.inside = c(0.9, 0.8)) +
  labs(title = paste0("Simulated butterfly counts (", y, ")"),
    subtitle = "- Fitting GAM model with rbms",
    x = "Day of Year",
    y = "Count")

btfl_fig4

```



To compare the fitted curve with the activity curve, we should use a standard AUC of 1 to enable a fair comparison between models fitted to different population sizes. Using the standardized activity curve and the GAM-generated flight curve (NM), we can calculate the Root Mean Squared Error (RMSE) to estimate the goodness of fit of the flight curve generated with the rbms package.

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \tilde{y}_i)^2}$$

where  $y_i$  is the NM value at time  $i$  and  $\tilde{y}_i$  the value from the standardized activity curve at time  $i$ , from day 1 to  $n$  of the monitoring season.

## Non Gaussian flight curve

The same procedure can be applied to flight curves having more complex shapes. Here we will generate a time series of butterfly counts drawn from a known flight curve (adult activity), using simulation from a Zonneveld model.

```
btfl_data_zn <- timeseries_sim(nsims=1,
  year = c(2023),
  doy.samples = seq(from=1, to=365, by=1),
  abund.type = "exp",
  activity.type = "zon",
  sample.type = "pois",
  sim.parms = list(growth.rate = 0,
    init.size = 500,
    act.mean = 175,
    act.sd = 15,
    #theta = 5,
    zon.theta = 50,
    t0 = 100,
    beta = 5,
    alpha = 0.05)
)

set.seed(13276)
btfl_ts <- sim2bms(data = btfl_data_zn$timeseries, yearKeep = y)

btfl_week_smpl <- sim2bms(data = btfl_data_zn$timeseries, yearKeep = y,
  weeklySample = TRUE,
  weekdayKeep = c(2:5),
  monitoringSeason = c(4:9))

btfl_week_missing <- sample_missing(data = btfl_week_smpl, propMissing = 0.25)

visit_sim <- btfl_week_smpl[!date %in% btfl_week_missing$date, .(site_id, date, count)]
count_sim <- visit_sim[count>=1,][, species := "sp1"]

names(visit_sim) <- toupper(names(visit_sim))
names(count_sim) <- toupper(names(count_sim))

ts_date <- rbms::ts_dwmy_table(InitYear = 2023, LastYear = 2023, WeekDay1 = 'monday')
```



```

ts_season <- rbms::ts_monit_season(ts_date,
  StartMonth = 4,
  EndMonth = 9,
  StartDay = 1,
  EndDay = NULL,
  ComplSeason = TRUE,
  Anchor = TRUE,
  AnchorLength = 2,
  AnchorLag = 2,
  TimeUnit = 'd')

ts_season_visit <- rbms::ts_monit_site(ts_season, visit_sim)

ts_season_count <- rbms::ts_monit_count_site(ts_season_visit, count_sim, sp = "sp1")

# mod_k <- "COUNT ~ s(DAY_SINCE, bs = \"cr\", k = 5) + factor(SITE_ID)"

ts_flight_curve <- rbms::flight_curve(ts_season_count,
  NbrSample = 300,
  MinVisit = 5,
  MinOccur = 3,
  MinNbrSite = 1,
  MaxTrial = 4,
  GamFamily = 'nb',
  SpeedGam = FALSE,
  ComplSeason = TRUE,
  SelectYear = NULL,
  #mod_form = mod_k,
  TimeUnit = 'd')

pheno <- ts_flight_curve$pheno

btfl_fig5 <- ggplot() +
  geom_point(data=ts_season_count[ANCHOR == 0 & !is.na(COUNT), ], aes(x=DAY_SINCE, y=COUNT)) +
  geom_point(data=btfl_week_missing, aes(x=doy, y=count, colour = "missing"),
    shape=4, size=2, stroke=2) +
  geom_line(data = btfl_ts, aes(x = doy, y = act, colour = "activity")) +
  geom_line(data = pheno, aes(x = trimDAYNO, y = btfl_ts[,unique(abund.true)]*NM, colour = "GAM_fit")) +
  xlim(1,365) + ylim(0, max(btfl_ts$act,
    pheno$NM*btfl_ts[,unique(abund.true)],
    btfl_week_missing$count,
    ts_season_count[!is.na(COUNT), COUNT] )) +
  scale_colour_manual("",
    breaks = c("count", "activity", "missing", "GAM_fit"),
    values = c(cnt_col, flc_col, missing_col, GAM_col)) +
  theme_light() +
  theme(legend.position = "inside", legend.position.inside = c(0.9, 0.8)) +
  labs(title = paste0("Simulated butterfly counts - Zonneveld Model (", y, ")"),
    subtitle = "- Fitting GAM model with rbms",

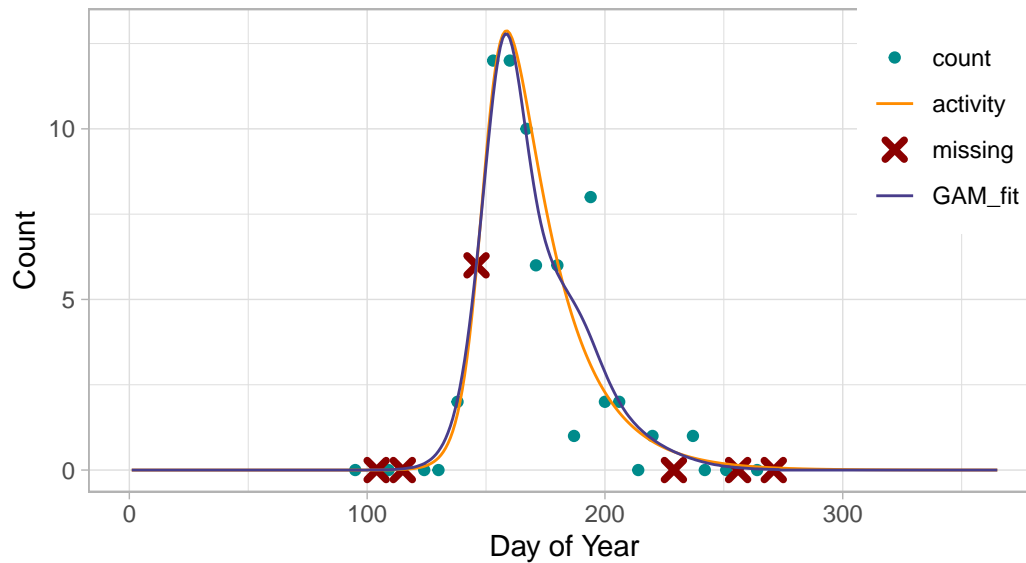
```

```

x = "Day of Year",
y = "Count")
btfl_fig5

```

Simulated butterfly counts – Zonneveld Model (2023)  
– Fitting GAM model with rbms



## Simple trend case

In the first scenario, we will apply the method to a simple case where we have one univoltine species that is monitored over 15 years across 100 sites where the populations follow the same trend with a known growth rate.

## Part III

# Butterfly Trends



# Trends Simulation

**i** Content under construction

## Multisite and trends

Here we will simulate butterfly count data for multiple sites (100), for a univoltine species that has a simple Gaussian flight curve (peak at day 175), where each population has an initial number of 500 individuals (total number of observations accumulated over the season) and a mean population growth rate of 20% over 10 years (2010 to 2020). We introduce some variation (standard deviation of 0.02) in the growth rate observed across the 100 sites.

For this simulation, we will use the `timeseries_sim()` function from the `butterflyGamSims` package and the `sim2bms()` function written above to reformat the simulated data into a format that can be used in the `rbms` package.

```
set.seed(13276)

if(!require("butterflyGamSims")) devtools::install_github("cbedwards/butterflyGamSims")

trend_sim <- data.table()

for(i in 1:100){

  btfl_data <- timeseries_sim(nsims=1,
    year = c(2010:2020),
    doy.samples = seq(from=1, to=365, by=1),
    abund.type = "exp",
    activity.type = "gauss",
    sample.type = "pois",
    sim.parms = list(growth.rate = rnorm(1, mean = 0.2, sd=0.02),
      init.size = 500,
      act.mean = 175,
      act.sd = 15)
  )

  trend_sim <- rbind(trend_sim,
```

```

sim2bms(data = btfl_data$timeseries,
        weeklySample = TRUE,
        weekdayKeep = c(2:5),
        monitoringSeason = c(4:9)),
        growth.rate := btfl_data$parms$growth.rate[,
        site_id := paste0("site_id", i)])
}
knitr::kable(head(trend_sim))

```

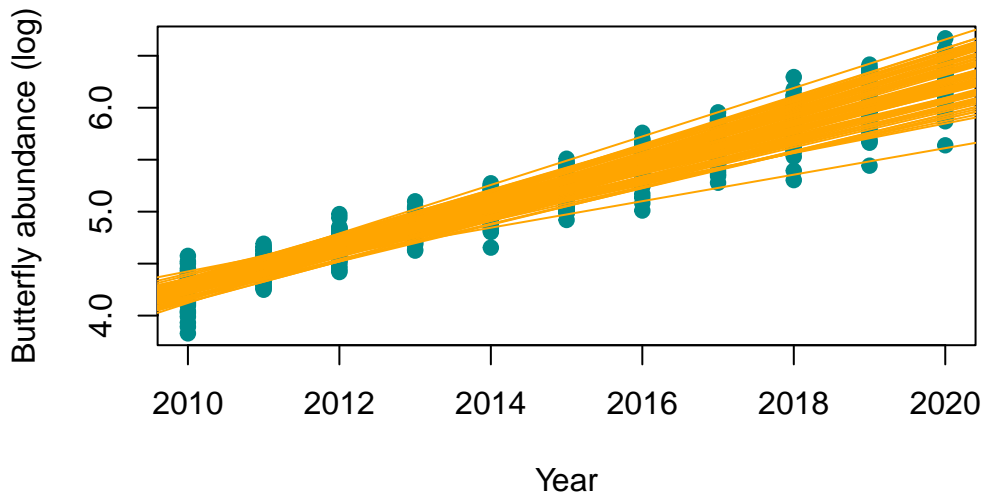
```
## plot and add lm
```

```

trend_sim[, y_count := sum(count), by = .(site_id, years)]
plot(unique(trend_sim[, .(years, log(y_count))]), pch = 19, col = "darkcyan",
     ylab = "Butterfly abundance (log)", xlab = "Year")

for(i in 1:100){
  abline(lm(log(y_count) ~ years, data = trend_sim[site_id == paste0("site_id", i), .(years, y_count)]),
        col = 'orange')
}

```



The overall trend can be calculated by fitting a Generalized Linear Model (GLM), with a Poisson error distribution on the simulated dataset. Note that here we did not include the site effect as all sites had the same abundance (500 individuals) and we excluded the intercept from our model. To illustrate the effect of the variation in the trends across sites on our overall trend estimate, we can bootstrap (resample with replacement) the sites and fit the GLM on the bootstrap sample. We can generate 100 (or more) iterations to illustrate the variation around the trend derived from our 100 sites.

```

plot(unique(trend_sim[, .(years, log(y_count))]), pch = 19, col = "darkcyan",
      ylab = "Butterfly abundance (log)", xlab = "Year")

for(i in 1:100){
  abline(lm(log(y_count) ~ years, data = trend_sim[site_id == paste0("site_id", i), .(years, y_count)],
          col = 'orange')
}

for(i in seq_len(100)){

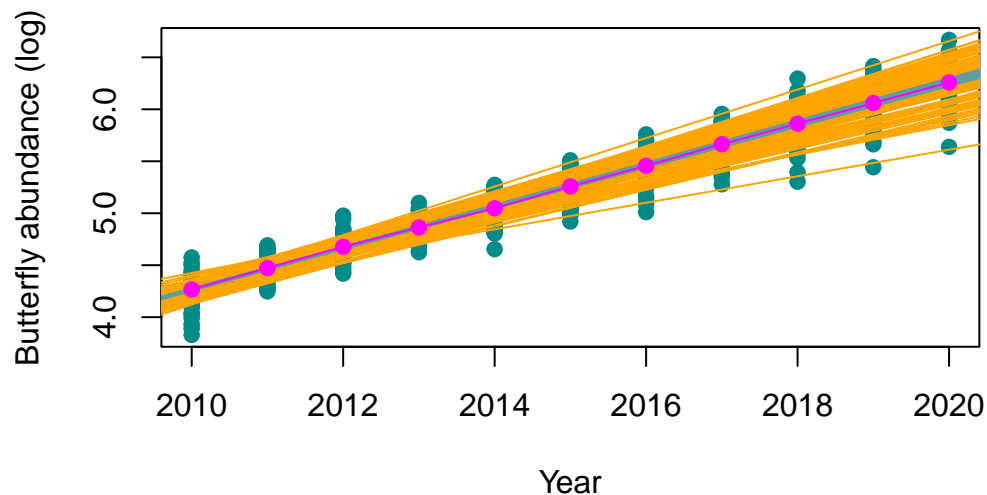
  boot_site <- sample(unique(trend_sim$site_id), trend_sim[, uniqueN(site_id)], replace = TRUE)

  #points(c(2010:2020), as.numeric(coef(glm(y_count~ factor(years) - 1, data = trend_sim[site_id %in% boot_site, ],
  #                                family = "poisson"))), col = 'magenta', pch = 19)

  abline(lm(as.numeric(coef(glm(y_count~ factor(years) - 1, data = trend_sim[site_id %in% boot_site, ],
                                family = "poisson")) ~ c(2010:2020))), col = 'cadetblue', type = 'l')
}

points(c(2010:2020), as.numeric(coef(glm(y_count~ factor(years) - 1, data = trend_sim, family = "poisson"),
  col = 'magenta', pch = 19)
points(c(2010:2020), as.numeric(coef(glm(y_count~ factor(years) - 1, data = trend_sim, family = "poisson"),
  col = 'magenta', type = 'l')

```







Part IV

**Bibliography**



# References

- Edwards, Collin, Cheryl Schultz, David Sinclair, Daniel Marschalek, and Elizabeth Crone. 2023. “Estimating Butterfly Population Trends from Sparse Monitoring Data Using Generalized Additive Models.” December 8, 2023. <https://doi.org/10.1101/2023.12.07.570644>.
- Schmucki, Reto, Colin Harrower A., and Emily Dennis B. 2022. “rbms: Computing generalised abundance indices for butterfly monitoring count data.” <https://github.com/RetoSchmucki/rbms>.

